

אדם, מכונה, מדינה

לקראת אסדרה של בינה מלאכותית

עמיר כהנא | תהילה שוורץ אלטשולר

עיון



המכון הישראלי
לדמוקרטיה



המכון הישראלי
לדמוקרטיה

אדם, מכונה, מדינה לקראת אסדרה של בינה מלאכותית

עמיר כהנא | תהילה שוורץ אלטשולר

Human, Machine, State:
Toward the Regulation of Artificial Intelligence
Amir Cahane | Tehilla Shwartz Altshuler

עיצוב הסדרה והעטיפה: סטודיו Alfabees

ביצוע גרפי: נדב שטכמן פולישוק

עיצוב התרשימים: דנה ברגר

הדפסה: דפוס מאור ולך, ירושלים

התצלום על העטיפה: Shutterstock

מסת"ב 978-965-519-433-3 ISBN

אין לשכפל, להעתיק, לצלם, להקליט, לתרגם, לאחסן במאגר ידע, לשדר או לקלוט בכל דרך או אמצעי אלקטרוני, אופטי או מכני או אחר – כל חלק שהוא מהחומר בספר זה. שימוש מסחרי מכל סוג שהוא בחומר הכלול בספר זה אסור בהחלט אלא ברשות מפורשת בכתב מהמוציא לאור.

© כל הזכויות שמורות למכון הישראלי לדמוקרטיה (ע"ר), 2023

נדפס בישראל, תשפ"ג/2023

המכון הישראלי לדמוקרטיה

רח' פינסקר 4, ת"ד 4702, ירושלים 9104602

טל': 02-5300888

אתר האינטרנט: www.idi.org.il

להזמנת ספרים:

החנות המקוונת: www.idi.org.il/books

דוא"ל: orders@idi.org.il

טל': 02-5300800

כל פרסומי המכון ניתנים להורדה חינם, במלואם או בחלקם, מאתר האינטרנט.

המכון הישראלי לדמוקרטיה

המכון הישראלי לדמוקרטיה הוא מוסד עצמאי א-מפלגתי, מחקרי ויישומי, הפועל בזירה הציבורית הישראלית בתחומי הממשל, הכלכלה והחברה. יעדיו הם חיזוק התשתית הערכית והמוסדית של ישראל כמדינה יהודית ודמוקרטית, שיפור התפקוד של מבני הממשל והמשק, גיבוש דרכים להתמודדות עם אתגרי הביטחון מתוך שמירה על הערכים הדמוקרטיים וטיפוח שותפות ומכנה משותף אזרחי בחברה הישראלית רבת הפנים.

לצורך מימוש יעדים אלו חוקרי המכון שוקדים על מחקרים המניחים תשתית רעיונית ומעשית לדמוקרטיה הישראלית. בעקבותיהם מגובשות המלצות מעשיות לשיפור התפקוד של המשטר במדינת ישראל ולטיפוח חזון ארוך טווח של תרבות דמוקרטית נכונה לחברה הישראלית ולמגוון הזהויות שבה. המכון שם לו למטרה לקדם בישראל שיח ציבורי מבוסס ידע בנושאים שעל סדר היום הלאומי, ליוזם רפורמות מבניות, פוליטיות וכלכליות ולשמש גוף מייעץ למקבלי ההחלטות ולציבור הרחב.

המכון הישראלי לדמוקרטיה הוא זוכה פרס ישראל לשנת תשס"ט על מפעל חיים – תרומה מיוחדת לחברה ולמדינה.

הדברים המתפרסמים בספר זה אינם משקפים בהכרח את עמדת המכון הישראלי לדמוקרטיה.

תוכן העניינים

7	תקציר
17	מבוא
25	פרק ראשון מהי בינה מלאכותית?
57	פרק שני שימושים שונים של בינה מלאכותית
95	פרק שלישי מדיניות אתיקה של בינה מלאכותית
119	פרק רביעי רגולציה של בינה מלאכותית: סקירה השוואתית
161	פרק חמישי בינה מלאכותית וזכויות אדם: חזית חדשה במדיניות טכנולוגיה
179	פרק שישי "אחרי רבים להטות": הטיות אלגוריתמיות
203	פרק שביעי שקיפות אלגוריתמית
219	פרק שמיני פיקוח מוסדי על בינה מלאכותית
231	פרק תשיעי יסודות לרגולציה מכוונת זכויות של בינה מלאכותית בישראל: עקרונות כלליים

פרק עשירי

יסודות לרגולציה מכוונת זכויות של בינה מלאכותית בישראל:

257

ארגז כלים

פרק אחד עשר

המלצות לגבי מוסד מאסדר לבינה מלאכותית בישראל

279

פרק שנים עשר

סיכום

305

iii

Abstract

תקציר

-

בשנים האחרונות מערכות בינה מלאכותית הולכות ונשזרות ברקמת חיי היום-יום. הן ממליצות על נתיבי נסיעה או על השיר הבא שיושמע, תומכות באבחון רפואי, ולאחרונה אף משתתפות באופן פעיל בהכנת שיעורי הבית. גופים ציבוריים ברחבי העולם מטמיעים מערכות אלגוריתמיות שמקבלות החלטות מינהליות הנוגעות להקצאת משאבים, לתכנון, לחיזוי פשיעה או להגנה על המרחב הציבורי, או תומכות בהחלטות כאלה. מעוזרים אישיים דיגיטליים ועד מכוניות אוטונומיות, מרובוטים המבצעים משימות פשוטות ועד מערכות מעקב, זיהוי וחיזוי.

ואולם על אף היתרונות הגלומים במערכות אלגוריתמיות – ללא בקרה על השימוש בהן, על פיתוחן ועל פרישתן הן עלולות לסכן זכויות אדם וחירויות יסוד. סכנות אלו יכולות לבוא לידי ביטוי בכל שרשרת הערך של פיתוחן והשימוש בהן – החל בשלב הגדרת תכליתה של המערכת, עבור בהסתמכות על מידע חלקי, שגוי, מזוהם או מוטה, וכלה באי-החלה של בקרות מאוחרות על תוצרי המערכת. יתר על כן, ככל שהבינה המלאכותית הולכת ומתפתחת כך שמערכותיה נוטות להציג יכולות חדשות שהמפתחים שלהן לא התכוונו להן או לא חזו אותן. מקצת היכולות מעוררות השראה ומקצתן בעלות פוטנציאל נזק, כגון היכולת לבצע פעולות סייבר התקפיות, לתמרן אנשים באמצעות שיחה או להפיץ מידע מלאכותי מוטעה ומטעה. היכולת לזהות יכולות אלו ולהגביל את הסיכונים שהן מביאות איתן נעשתה אפוא אתגר חשוב מאין כמוהו.

מהי בינה מלאכותית, מהם יתרונותיה ואילו חששות היא מעוררת, בעיקר כשהיא משמשת את רשויות השלטון? שאלות אלו הן עניינו של ספר זה.

מקבלי ההחלטות, התעשייה, האקדמיה וארגוני החברה האזרחית ברחבי העולם ובישראל זיהו שהבינה המלאכותית היא טכנולוגיה מתפרצת שיש להיערך אליה באמצעות אסטרטגיה לאומית ומדיניות רגולטורית. בשלהי העשור הקודם פורסמו עשרות מסמכי אתיקה שעסקו בבינה מלאכותית וניסו להניח עקרונות לפיתוח מערכות אלגוריתמיות, לשימוש בהן ולהטמעתן. את הערכים האתיים המוצעים בהם אפשר להעמיד על שבעה עקרונות: שקיפות; הוגנות; מניעת נזקים ובטיחות; אחריות ואחריותיות; פרטיות; קידום הטוב וערך האדם במרכז; חירות ואוטונומיה.

ואולם לא די בעקרונות אתיים. כדי להבטיח שמירה על זכויות אדם וחירויות יסוד יש לעגן עקרונות אלו בחקיקה. ואכן, ברחבי העולם אנו רואים ניצני חקיקה המבקשת לאסדר את השימוש בטכנולוגיות בינה מלאכותית ואת פיתוחן: הן על פי מודל של חקיקה רוחבית החלה על מערכות בינה מלאכותית כאשר הן, כדוגמת תקנות הבינה המלאכותית האירופיות, הן כשמיכת טלאים רגולטורית המאסדרת מגזרים ושימושים קונקרטיים של בינה מלאכותית, כדוגמת הדינים הייעודים שתכליתם להתמודד עם אפליה אלגוריתמית במערכות לגיוס ולקידום עובדים.

ספר זה מציע עקרונות מנחים ליצירת מדיניות בינה מלאכותית מכוונת זכויות וארגז כלים לאסדרת בינה מלאכותית.

המלצות עיקריות

עקרונות מנחים ליצירת מדיניות בינה מלאכותית מכוונת זכויות

האדם במרכז. התכלית המרכזית של הפיתוח של בינה מלאכותית ומערכות לומדות צריכה להיות שירות המין האנושי – הן כקולקטיב הן כפרטים – באופן המיטיב עימו. עיקרון זה של חירות האדם והאוטונומיה שלו מצטרף לעקרון קידום הטוב ועקרון מניעת הנזק ומבסס ביתר שאת את החשיבות של עקרון האדם במרכז, חירות ואוטונומיה. המשמעות המעשית של העמדת האדם במרכז היא שפיתוחן של מערכות נבונות, פרישתן והשימוש בהן ייעשה בפועל על יסוד התפיסה הרואה בהגנה על זכויות יסוד ועל חירות האזרח עיקרון ראשון במעלה, ולא מס שפתיים שנועד להכשיר העדפת עקרונות כגון "קידום חדשנות", אינטרסים כלכליים כגון עידוד מגזר ההייטק או "הפיכת ישראל למובילה טכנולוגית גלובלית", ואף לא ייעול תהליכים במגזר הציבורי.

הדמוקרטיה במרכז. למערכות מבוססות בינה מלאכותית יש פוטנציאל רב לפגוע בדמוקרטיה במובנה הרחב, הן באמצעות השפעה על השיח הציבורי ופיזור של רעיונות, הן באמצעות מכשור לשליטה, למעקב, לזיהוי ולמשטור, הן באמצעות היכולת שלהן לשמש לזריעת ספק ולערעור עצם היכולת לברר את המציאות ולהבחין בין מקור לזיוף ובין אמת לשקר. לכן יש לתת משקל רב לעיקרון "הדמוקרטיה במרכז" גם אם הדבר פוגע לעיתים בהתקדמות טכנולוגיות או בחדשנות.

אוריינות דיגיטלית של מקבלי החלטות. אוריינות דיגיטלית משמעה היכולת לנתח את השוק ולהבין לאן מתפתחת הטכנולוגיה, לפחות בטווח הקצר. למשל: היכן נמצאים כספי המחקר והפיתוח של חברות הענק ומהם הפוטנציאלים שהן רושמות כדי לעגן פיתוחים טכנולוגיים חדשים. היא כוללת גם הבנה של הדרכים המסחריות והרגולטוריות להכוונת פיתוח טכנולוגי, ומכאן גם הבנה של האחריות המוטלת על מעצבי מדיניות להשפיע על התפתחות הטכנולוגיה ולא רק להתבונן בה מן הצד. נדרש רובד ביניים בין הבנת הטכנולוגיה לבין יצירת

מדיניות בעניינה, מסגרת להבנת המשמעויות של מערכות טכנולוגיות, ליכולת לדמיין את האפשרויות החדשות שהן מביאות איתן ולעמוד על השלכותיהן על המוסר החברתי ועל שלד השיטה המשפטית. כיוון שמסגרת זו חסרה לעיתים קרובות, נוצרים פערי הבנה – בייחוד בנושאים בעלי השלכות רחבות כמו בינה מלאכותית.

המשגה חדשה של מערכות ויכולות בתחום הבינה המלאכותית. המושג בינה מלאכותית הוא מטפורה פוליטית-טכנולוגית רבת עוצמה. ההשוואה בין מערכת בינה מלאכותית למוח האנושי יוצרת קרבה ודמיון, והללו מובילים להטמעה חברתית של הרעיון שהמכונות עובדות כמו מוח אנושי, מבצעות פעולות אנושיות בדרך שבה אנשים מבצעים אותן, ולמעשה מתחרות בבני האנוש. אנו מציעים להמשיג את פעולות המכונות ואת התכונות המיוחדות להן בדרך שאינה תלויה בהשוואה זו.

פיתוח זכויות למושאי החלטות של בינה מלאכותית. נדרשת חשיבה מחודשת על זכויות יסוד משתני בחינות. ראשית, יש לצקת משמעות חדשה לתאוריה החוקתית של זכויות האדם; שנית, יש ליצור זכויות דיגיטליות חדשות, שלא היה צורך בהן בעבר, ובעיקר זכויות של פרטים בשעה שהם באים במגע עם מערכות ממוכנות מבוססות אלגוריתמים.

אין להשלים עם מצב שבו ישראל היא חצר אחורית. בשנים האחרונות מקודמות הצעות לאסדרת בינה מלאכותית במדינות מובילות, וגם באיחוד האירופי, ויש להניח שמדובר בתחילתה של מגמה עולמית. גם אם יש שוני ערכי בין האסדרות במקומות שונים – אם בבחירת סוג המערכות המוגדרות מסוכנות ואם במטרות המוצהרות של החקיקה עצמה – יש להן גם הרבה מן המשותף. לכן אין להשלים עם מצב שבו במדינת ישראל לא תהיה חקיקה שתהיה בהרמוניה עם החקיקה המקובלת בעולם.

הרמוניה איננה בהכרח זהות. כאשר הפיתוח מכוון לייצוא ממילא נדרש הסטנדרט הרגולטורי לעמוד לכל הפחות בסטנדרטים הזרים, שהם לרוב גבוהים יותר מהסטנדרטים בישראל. אבל גם כאשר הפיתוח מכוון לשוק המקומי, ומושאי ההחלטות של מערכות הבינה המלאכותית אינם אזרחי חוץ אלא אזרחי המדינה, אין להנמיך את הסטנדרטים, שהרי אזרחי המדינה ראויים אף הם להגנה. אומנם

סטנדרט רגולטורי נמוך מהמקובל בעולם יכול לעורר חדשנות, אך זו לא בהכרח חדשנות רצויה: ישראל עלולה להפוך לחצר אחורית טכנולוגית, כלומר למקום שבו מפתחים מערכות שאסור לפתח אותן, להפיץ אותן או להשתמש בהן לפי הדין הזר.

רגולציה חסינת עתיד. תהליכי חקיקה מפגרים תמיד אחר המציאות. הטכנולוגיה מתקדמת בקצב מהיר, ובמדינות כמו ישראל, שבהן תהליכי החקיקה הם איטיים מאוד, נוצר פער גדול במיוחד. על כן אין לאסדר טכנולוגיה ספציפית, שכן זה מתכון בטוח להתיישנות הרגולציה, אלא להשתדל לקבוע עקרונות מנחים והגדרות כלליות שיקנו גמישות באכיפה עתידית.

מסגרת אסדרה היברידית: עקרונות, זכויות וחקיקה. מסגרת מבוססת עקרונות מציגה מערך עקרונות ליבה אתיים, מסגרת מבוססת זכויות מתמקדת בהגנה על זכויות האדם ועל החירויות של מי שמושפעים מיישומי הטכנולוגיות המבוססות על בינה מלאכותית, ומסגרת מבוססת חקיקה מאפשרת שלא להסתמך אך ורק על אסדרה וולונטרית, המבוססת על הרצון הטוב של השחקנים הכלכליים בשוק. שלוש המסגרות אינן סותרות ויש לשלב ביניהן במסגרת היברידית, המשלבת אסדרה רכה (עקרונות אתיים) עם תפיסת ניהול סיכונים שתתבטא בהוראות החוק ובכללים רגולטוריים נוקשים.

גמישות באשר למועד ההתערבות האסדרתית. במקרים מסוימים יש יתרונות ברורים לרגולציה מוקדמת, כלומר לפני שמוצרים שמבוססים על טכנולוגיה מסוימת חודרים לשוק. אם טכנולוגיה נתפסת מסוכנת במיוחד – פיזית (כמו רכב אוטונומי) או מוסרית (למשל יצירה מלאכותית של תוכן המסית לביצוע מעשי טרור) – הגיוני לאסדר מראש. גם במקרים קיצוניים פחות התערבות מוקדמת יכולה להועיל במה שנוגע לעיצוב כיווני המחקר ולתכנון השקעת המשאבים בפיתוח. יתר על כן, מאחר שבשלב זה הושקעו משאבים מעטים יחסית והעלויות השקועות נמוכות יותר ייתכן שהתערבות רגולטורית תיתקל בהתנגדות מצומצמת יותר מצד בעלי עניין. מנגד, ייתכנו מקרים שבהם עדיף להמתין ולהתמודד עם בעיות כאשר הן מתעוררות במקום לנסות לצפות אותן.

אסדרה ענפית ואסדרה רוחבית. אסדרה רוחבית משיגה מטרות של משילות וודאות מצד השלטון המרכזי, יוצרת הרמוניה רגולטורית ולפיכך יכולה להגביר

את אמון הציבור ואת הוודאות הרגולטורית בעבור התעשייה. מנגד, אסדרה ענפית מאפשרת שימוש ברגולטורים קיימים ובסמכויות האכיפה שלהם, איננה מחייבת להקים מסגרות מוסדיות חדשות, מאפשרת יצירת הסדרים ואמצעי אכיפה המותאמים בצורה מדויקת יותר לענף הרלוונטי, מגבירה את הבהירות והוודאות הרגולטורית ומאפשרת גם לבעלי עניין בתוך כל מגזר להשתתף בחשיבה על ההסדרים האלה. ואולם הקושי באסדרה ענפית הוא בסתירה אפשרית בין ענפים, ביצירת סטנדרטים לא אחידים, בהעמקת הפערים בין הרגולטורים וכן בהותרת מרחבים לא מאוסדרים הנופלים בין הכיסאות. אשר על כן מוצע לנקוט שילוב של אסדרה רוחבית ואסדרה ענפית, דוגמת מודל של מאסדר-על שיש לו סמכויות הנחיה וייעוץ למאסדרים המגזריים.

ארגז כלים לאסדרת בינה מלאכותית

הבנת "מעגל החיים" של מערכות לומדות צריכה לשמש בסיס לאסדרתן. כדי ליצור אסדרה אפקטיבית של מערכות לומדות יש להביא בחשבון את כל רכיבי מעגל החיים שלהן. הואיל ועקרונות האסדרה, כגון הוגנות, פרטיות, שקיפות, אחריותיות וניהול סיכונים, מתבטאים בהקשרים שונים בכל אחד מרכיבי מעגל החיים, היעדר התייחסות לרכיבים אלו עלול ליצור מצב של אסדרת יתר של רכיבים מסוימים והתעלמות מרכיבים אחרים, ולכן לצמצם את האפקטיביות של האסדרה.

לעומת זאת, תפיסה אינטגרטיבית מביאה בחשבון את מכלול רכיבי מעגל החיים ונדרשת גם לקשרי הגומלין ביניהם. למשל: תכלית המערכת ומסגור הבעיה צריכים להשפיע על בחירת המודל (אם לבחור במודל שיש בו עכירות רבה יותר באשר לדרכים שבו הוא מקבל החלטות או שלא לאפשר זאת); תוצאות הערכה בשלב בניית המודל מזינות תהליכים להערכת סיכונים, והם בתורם מחייבים קבלת החלטות הנוגעות לאימון המודל, לפרישתו או לאבטחתו; תכלית המודל משפיעה על הבחירה בממשקי המשתמש (האם מערכת המעניקה ייעוץ רפואי אמורה להציג את האפשרות שהיא טועה? האם נכון לספר למשתמש שהוא מתקשר עם מערכת מלאכותית ולא אנושית?).

חלק חשוב של הבנת מעגל החיים של מערכות לומדות נוגע לצורך לנטר אותן לאחר שיושמו בפועל (post deployment) בעולם האמיתי (למשל הבניה של המערכת בתוך מוצר או בתוך ממשק). זאת משום שמערכת לומדת, שלא כמו מוצרים אחרים (תרופות למשל), יכולה מעצם טיבה להשתנות גם לאחר יישומה בשל המשוב שהיא מקבלת מן המשתמשים.

פיתוח מתודולוגיות לניהול סיכונים. הצורך ליישם מתודולוגיות ניהול סיכונים על מערכות אלגוריתמיות עדיין נמצא בחיתוליו, למרות חיוניותו. אנו מציעים מודל של ניהול סיכונים שלפיו כדי להעריך מסוכנות של מערכת נדרשת התבוננות כפולה – תחילה יש להעריך את פוטנציאל המסוכנות של המערכת כפי שתוכננה; אחר כך יש להעריך את רמת הקשר בין המשימה לתוצאה (alignment), כלומר את האפשרות של המערכת לממש פוטנציאל מסוכנות מחוץ לתפקיד שיייעדו לה מתכנניה.

הערכה כפולה זו תהיה הבסיס לקבלת ההחלטות וכדי שיהיה אפשר להוציאה לפועל יהיה צורך לנסח כללי משילות ובטיחות, כגון כללים לאימון אחראי (אם לאמן מודל חדש שמראה סימנים מוקדמים של סיכון – וכיצד); וכללים ליישום אחראי (אם, מתי וכיצד ליישם מודלים שעלולים להיות מסוכנים); ולקבוע אילו רמות של שקיפות ותיעוד נדרשות במקרה של מודלים שעלול להיות בהם סיכון קיצוני ואילו בקרות ומערכות אבטחת מידע יש ליישם בעניינם.

תיעוד נתונים, משילות נתונים והליכי בקרה (auditing) בדיעבד. המורכבות של מערכות בינה מלאכותית בכלל, ושל מערכות לומדות בפרט, מערימה קשיים מיוחדים על קובעי המדיניות והרגולטורים בבואם לנסח כללי אחריות ולזהות בפועל את השרשרת הסיבתית שהובילה לפגיעה בזכויות, בייחוד בהתחשב בשונות בין מגזר למגזר ובין יישום ליישום.

המשותף לכל אלו הוא שבלי תיעוד ראוי הם בלתי אפשריים. הבסיס לכל בחינה עובדתית של כשל קונקרטי במערכות בינה מלאכותית הוא משילות נתונים ותיעוד קפדני של הליכי עבודה, מקורות מידע, תיוגים, מודלים, תהליכי עיצוב קוד, הערכת סיכונים ובסיסי נתונים, וזיהוי הפערים בכל אלו. עיצוב טוב של התיעוד הוא גם אינטרס של יזמים ומפתחים, שכן הוא מאפשר להם הן לתחקר

בריעבר כשלים ותופעות שהם לא צפו הן לעמוד בחובות תיעוד רגולטוריות ממקורות אחרים.

פיתוח כלים להתמודדות עם הטיות ו"הנדסת הוגנות". אף שאי אפשר למנוע את בעיית ההטיות האלגוריתמיות, בייחוד כשמקור ההטיה הוא במציאות עצמה, המגולמת בנתונים, אפשר לזהות אותה ולצמצם את היקפה. מוצע לנקוט כמה אסטרטגיות לצמצום הטיות אלגוריתמיות, ובהן קיום הליך הוגן מבחינה סטטיסטית, גיוון ההון האנושי בקרב מפתחים של מערכות בינה מלאכותית ויישום הליכי בקרה בריעבר (auditing).

התמודדות עם אתגרי השקיפות האלגוריתמית: איתנות מדעית תהליכית. מוצע לאמץ מודל המבוסס על התפיסה הקלסית של שקיפות, אבל כולל חלופה שמתאמת למגבלות הטכנולוגיות של מערכות אלגוריתמיות שאינן מסוגלות לספק הסבר רציונלי לפלט נקודתי. המודל מיועד למקרים שבהם אין אפשרות לספק פלט, אבל יש נחיצות חברתית לשמור על חובות של שקיפות כדי שלא למנוע פיתוח בטכנולוגיות מסוימות ושימוש בהן.

התמצאות לגבי מוסד מאסדר לבינה מלאכותית בישראל

מוצע להקים מוסד מאסדר לבינה מלאכותית בישראל, שתפקידו יהיו לקדם ולתאם אסדרה של מערכות נבונות בישראל, לרבות תהליך של גיבוש חקיקה רוחבית מתוך מעקב אחר התפתחויות גלובליות מתאימות. המוסד המאסדר יתווה מדיניות בעניין פיתוח מערכות נבונות בישראל, הטמעה שלהן ושימוש בהן; יספק הנחיה מקצועית לרגולטורים ענפיים כדי להבטיח הרמוניזציה של הכללים החלים על פיתוח מערכות אלו, פרישה שלהן ושימוש בהן; וישמש גורם מנחה שיורי, כלומר יהיה ממונה על אסדרת מערכות נבונות פרטיות בתחומים שאין בהם מאסדר מגזרי.

הרגולטור יספק הנחיה מקצועית בתחום הבינה המלאכותית גם לרשויות המדינה. בין השאר הוא ייתן חוות דעת על מסמכי מכרזים ממשלתיים של מערכות בינה מלאכותית כדי להבטיח שהמדינה מצטיידת במערכות בינה מלאכותית העולות בקנה אחד עם הסטנדרטים הישראליים והעולמיים בתחום, וכן כדי לייצר בעבור השוק הפרטי תמריצים "רכים" לעבוד בהתאם לסטנדרטים.

הרגולטור יהיה ממונה גם על ייעוץ בנושאי משפט הבינה המלאכותית, ומוצע שתינתן לו הסמכות להציג עמדה עצמאית בתחומים אלו לכנסת ולבתי המשפט. נוסף על התפקידים הנזכרים, הרגולטור המוצע יהיה מוקד ידע, הדרכה ושיתופי פעולה ויקדם את האוריינות הדיגיטלית בתחומים אלו בקרב גורמי הממשלה. הוא אף ידון עם בעלי עניין מקומיים ובינלאומיים מתחומי התעשייה, הממשל והאקדמיה על ההשפעות האפשריות של טכנולוגיות אלו על החברה.

בעת הזאת יהיה נכון למקום את רגולטור הבינה המלאכותית כיחידה בתוך רשות הרגולציה, שכן תפקידיה הולמים מאוד את העולם המתפתח של אסדרת הבינה המלאכותית – לפחות בשנים הקרובות, עד התייצבותו. ההתמודדות עם אתגרי הגמישות, עם הסביבה הרגולטורית העמוסה ועם הטכנולוגיה מצריכים הקצאת משאבים ראויה, שתאפשר ליצור תקנים הן למומחים טכניים הן למומחים במשפט ובמדיניות. מוצע כי גם בהיעדר חוק בינה מלאכותית שלם תתוקצב הרשות בהחלטת ממשלה כדי שתוכל להקים את היחידה של רגולטור הבינה המלאכותית.

התמצות לתקופת הביניים עד לאסדרה של התחום

חקיקה משלימה וערכוני חקיקה. גם בהיעדר חוק בינה מלאכותית כללי וייעודי יש לחייב כבר בעת הזאת את מקבלי ההחלטות והרגולטורים הייעודיים לעדכן את החקיקה הקיימת ולחוקק חקיקה משלימה. מדובר בעיקר בחוקים כגון חוק התחרות וההגבלים העסקיים, חוק זכויות יוצרים, חוק הגנת הפרטיות, פקודת הראיות וחוק הרכש הממשלתי.

יצירת "אקוסיסטם קדם-אסדרתי". ראוי ליצור אקוסיסטם קדם-אסדרתי שבמסגרתו יפרסמו הרגולטור הייעודי, רגולטורים ענפיים ורגולטורים משלימים (כגון הרשות לפרטיות והרשות לתחרות) הנחיות וגילויי דעת; והם בתורם ישמשו את התעשייה ושחקנים קונקרטיים בה, וגם את בתי המשפט. מסמכי הנחיות אלו יוכלו לשמש את התעשייה כדי לקבוע סטנדרטים של אסדרה עצמית, מתוך הנחה שרגולציה עתידית תהיה דומה להנחיות. הם יוכלו לשמש גם את חברות התקינה כדי ליצור את מסגרות התקינה.

סביב האקוסיסטם הקדם-אסדרתי יתפתחו תהליכים של חידוד דפוסי הבקרה (auditing) בהקשרים שונים, ייערכו תהליכי חשיבה בעניין סטנדרטים של זהירות ראויה וייעשה מאמץ סביר לקדם אחריות בפיתוח וביישום של מערכות לומדות; ייבנו ארגזי חול וייערכו ניסויים רגולטוריים שונים; ויקודם מאגר של מומחים שיוכלו לשמש את התעשייה, את הרגולטורים ואת בתי המשפט בהתמודדות עם סוגיות חדשות ומורכבות. במקביל, התעשייה והשחקנים יוכלו לתת משוב לגופים הרגולטוריים השונים וכך לטייב את ההנחיות. נוסף על כך, בהיעדר אסדרה, בתי משפט יוכלו להשתמש בהנחיות כהשראה לפרשנות בסכסוכים המובאים לפניהם. פסקי הדין יוכלו גם הם להעשיר את גוף הידע ולאפשר לשחקנים וגם לרגולטורים לטייב ולשייף את הנחיותיהם לקראת גיבוש אסדרה.

מבוא

–

מערכות בינה מלאכותית נעשו בשנים האחרונות חלק בלתי נפרד מן החיים. מעוזרים אישיים דיגיטליים ועד מכונות אוטונומיות, מרובוטים המבצעים משימות פשוטות ועד מערכות מעקב, זיהוי וחיזוי. בינה מלאכותית ואלגוריתמים מבוססי למידה נבונה משמשים בין השאר במערכות פיננסיות ממוחשבות,¹ בשירות הסוכנויות לאכיפת החוק והביטחון,² ברפואה,³ וגם בזירה המשפטית. לצד כלים

1 ראו להלן בסעיף 1.9.

2 ראו להלן בסעיף 1.8.2.

3 ראו להלן בסעיף 1.8.2.

חכמים התומכים בעבודתם של משפטנים⁴ יש דיונים תאורטיים על האפשרות שבעתיד יוכלו מערכות אלגוריתמיות "לתפור" חקיקה בהתאמה אישית.⁵

ככל שמערכות הבינה המלאכותית מתפתחות כך הן נוטות להציג יכולות חדשות שהמפתחים שלהן לא התכוונו להן או חזו אותן. יש יכולות שמעוררות השראה ויש יכולות בעלות פוטנציאל להזיק, כגון היכולת לבצע פעולות סייבר התקפיות, לתמרן אנשים באמצעות שיחה או לספק הוראות מעשיות לביצוע פעולות טרור. היכולת לזהות יכולות אלו ולהגביל את הסיכונים הצפונים בהן נעשות אפוא אתגר משמעותי.

ההוצאות של ממשלות ועסקים ברחבי העולם על מוצרים המבוססים על בינה מלאכותית יגיעו ב-2023 ל-500 מיליארד דולר.⁶ עם השימושים האלה מגיעים גם חששות מסוגים שונים: נזקים, אפליה, הוגנות, הטיות, גיוון ופרטיות. על שולחן הדיונים כבר מונחת הצעת חוק של האיחוד האירופי בעניין בינה מלאכותית, הצעה של רשות הסחר ההוגן בארצות הברית, והרבה מסמכי מדיניות והצעות חקיקה במדינות אחרות. גם בישראל מתפרסמים מסמכים שונים בעניין זה.

מהי בינה מלאכותית, מהם יתרונותיה ואילו חששות היא מעוררת, בעיקר כשהיא משמשת את רשויות השלטון? שאלות אלו הן עניינו של ספר זה.

ספר זה נועד לשמש בראש ובראשונה מבוא למקבלי החלטות בכלל ולעוסקים ברגולציה של בינה מלאכותית בפרט. אבל הוא מיועד גם לכל מי שעולם הבינה המלאכותית מעניין אותו ומבקש להבין אותו ואת סביבתו. תחילה נסביר מהי בינה מלאכותית ונסקור יישומים שונים שלה. אחר כך נידרש לאתגרים הייחודיים שהבינה המלאכותית מעמידה לפני הרגולטור – הטיות אלגוריתמיות והסברתיות (explainability), ואף לחשש שמערכות בינה מלאכותית יפגעו בזכויות אדם. ולבסוף נדון בתשתית המשפטית והמוסדית הנדרשת לבינה מלאכותית בטוחה. תשתית כזאת צריכה לתת את הדעת על האדם, מושא ההחלטות האלגוריתמיות,

4 ראו להלן בסעיף 1.8.2.

5 ראו להלן בסעיף 1.8.3.

ועל המשטר הדמוקרטי. מאחר שהן בזירה הטכנולוגית הן בזירה הרגולטורית העולמית תחום אסדרת הבינה המלאכותית הוא דינמי ומתפתח, בחרנו שלא להציע מתווה אסדרה פרטני אלא להמליץ על קווי מתאר, עקרונות וכלים שרצוי לאמץ.

לנוכח רוחב היריעה של הנושא יש להדגיש באילו נושאים ספר זה אינו עוסק.

ספר זה אינו על בינה מלאכותית שבאה לחסל את כולנו, אלא הוא מתבונן במבט מפוכח בצורך לגבש בעשור הקרוב מדיניות טכנולוגיה. התרבות הפופולרית מרבה לתאר תרחישי אימים של מכונות המפתחות תודעה אנושית או על-אנושית ופונות כנגד מפתחיהן האנושיים.⁷ הרעיון של יצורים מלאכותיים, מעשי ידי אדם, בעלי תבונה ואוטונומיה, גירה את הדמיון האנושי לאורך ההיסטוריה. האגדה על הגולם מפראג, שנתנה השראה לשמם של המחשבים הראשונים במכון ויצמן,⁸ מתארת יצור דמוי אדם עשוי חימר או אדמה שעל מצחו כתוב "אמת" והוא ניעור לחיים באמצעות קוד של צירופי אותיות מיסטיים. המפלצת שברא ויקטור פרנקנשטיין, ששבה לרדוף אותו, היא עוד דוגמה לנרטיב הגולם – שבמקרה זה אכן קם על יצורו. ואילו בתרבות הפופולרית בת זמננו התגלגל המיתוס בתרחישי האימים של תבונה ממוחשבת מעשה ידי אדם שמשלתט על הפלנטה ומשעבדת את המין האנושי.⁹

התסריטים האפוקליפטיים האלה נוטים להדהד יתר על המידה באזהרותיהם של מדענים¹⁰ ואנשי טכנולוגיה¹¹ בני זמננו. בסקר שנערך בשנת 2022 בקרב

7 לדוגמאות ראו להלן ה"ש 9.

8 "מחשב אלקטרוני 'גלם' הופעל במכון ויצמן" דבר (30.12.1963).

9 ארתור סי' קלארק "חייג פ' ותקבל את פרנקנשטיין" הרוח הנושבת מן השמש 90 (מתרגם יורם רפפורט, 1978); הרלן אליסון "אין לי פה ואני מוכרח לזעוק" מדע בדיוני: הטוב שבטוב כרך 4, 34 (עורך איזק אסימוב, מחרגמח אילנה דגני בינג, 1981); ראו גם הסרטים אודיסיאה בחלל: 2001 (סטנלי קובריק במאי ומפיק 1968); שליחות קטלנית 2: יום הדין (ג'יימס קמרון במאי ומפיק 1991); מטריקס (במאיכות ליילי ולאנה ואצ'ובסקי 1999).

10 סטיבן הוקינג, למשל, הזהיר מפני פיתוח לא מבוקר של אינטליגנציה מלאכותית. ראו רפאלה גויכמן "סטיבן הוקינג מזהיר בפני הסכנות האמיתיות של המכונות" ynet (15.10.2015).

11 Peter Holley, "Elon Musk's Nightmarish Warning: AI Could Become 'an Immortal Dictator from which We Would Never Escape,'" The Washington Post (6.4.2018)

חוקרי בינה מלאכותית, 36% מהנשאלים השיבו כי יש סבירות שמערכות בינה מלאכותית יוכלו "לגרום במאה זו לקטסטרופה גרועה לפחות כמו מלחמה גרעינית כוללת".¹² ואכן, לאחרונה חתמה שורה של מומחים טכנולוגיים על הצהרה המתריעה על הצורך למתן את הסיכון של הכחדת המין האנושי בשל בינה מלאכותית,¹³ וקדמה לה עצומה של מומחים לבינה מלאכותית ומובילי תעשייה שקראה להשהות לשישה חודשים פיתוח של מערכות בינה מלאכותית מתקדמות מחשש שאלו מסכנות את המין האנושי ועלולות להביא ל"אובדן השליטה על הציוויליזציה שלנו".¹⁴

אכן, קיימים אינטרסים שונים המובילים לאזהרות אלה. החל בניסיון להעביר אחריות לרגולטור על החלטות עסקיות מסוכנות, עובר ברצון לקנות זמן כדי להצליח להתמודד בשוק תחרותי או רצון להשאיר כוח בידי פלטפורמות ולא בידי יחידים, וכלה בחשש כן מפני השלכות חברתיות מסוכנות של שימושים בטכנולוגיה.

ואולם הדימוי הפופולרי של בינה מלאכותית כאיום קיומי על המין האנושי עלול לעוות את הדיון בהשפעת מערכות הבינה המלאכותית בטווח הקרוב והבינוני – בסיכויים הטמונים בה ובסכנות הנשקפות ממנה. יש הטוענים כי הצפת השיח הציבורי בקריאות לרגולציה (לרבות הקמת סוכנות בינלאומית לפיקוח על בינה מלאכותית, לפי הדגם של סבא"א)¹⁵ שתגן עלינו מפני אפוקליפסת בינה מלאכותית נועדה להסיח את הדעת מהבעיות המיידיות שנוגעות לשימוש, פיתוח

Julian Michael et al., *What Do NLP Researchers Believe? Results of the NLP Community Metasurvey* (26.8.2022), available at arXiv

13 אושרי אלקסלטי "סאם אלטמן ובכירים בענף ה-AI מזהירים מפני 'הכחדה' של המין האנושי" גיקטיים (30.5.2023); "Statement on AI Risk: AI Experts and Public Figures Express their Concern about AI risk," CENTER FOR AI SAFETY (30.5.2023)

14 *Pause Giant AI Experiments: An Open Letter*, FUTURE OF LIFE INSTITUTE (22.3.2023)

15 יובל מן "OpenAI מציעה להקים סוכנות בינלאומית לפיקוח על בינה מלאכותית" ynet (23.5.2023).

ופרישה של מערכות אלו,¹⁶ שהרי כבר עתה מערכות נבונות הן חלק ממארג חיי היום-יום של כולנו.

ספר זה אינו על בינה מלאכותית כאישיות משפטית. בשנת 2022 טען בלייק למוין (Lemoine), מהנדס תוכנה בחברת גוגל, שמודל השפה המתקדם של החברה פיתח מודעות עצמית.¹⁷ חברת גוגל הכחישה זאת,¹⁸ אבל השאלה אם מערכות מחשב יוכלו לפתח תודעה עצמאית שתרום אותן מדרגת חפץ דומם לדרגה של יצור תבוני בשר ודם (או מעבר לו) מוסיפה לגרות את הדמיון האנושי.

מבחינה משפטית שאלה זו מעלה סוגיית משנה: האם מערכות מחשב תבוניות ובעלות תודעה ראויות להיחשב לאישיות משפטית?¹⁹ אנלי ניוויץ טוענת שייחוס אישיות משפטית למערכות בינה מלאכותית היא מסך שמגן על מי שייתכן שהם האחראים האמיתיים לפגמים בהתנהלותן של מערכות אלו; על פי עמדתה של ניוויץ, אנשי הפיתוח והמהנדסים הם שלא השכילו לנפות מראש הטיות אלגוריתמיות.²⁰ שאלת האישיות המשפטית מתעוררת לעיתים בהקשרים של זכויות יוצרים – האם ניתן להקנות למערכת בינה מלאכותית זכות בקניין רוחני

16 ראו למשל Matteo Wong, *AI Doomerism is a Decoy*, THE ATLANTIC (2.6.2023).
17 "מהנדס תוכנה של גוגל טען שהבינה המלאכותית של החברה פיתחה תודעה – והושעה" הארץ (13.6.2022).

18 ראו הלן ה"ש 37. מחקרים מאוחרים יותר גורסים שאומנם למודלים מתקדמים של שפה יש מאפיינים של מה שנקרא בספרות הפסיכולוגית "תאוריה של תודעה" (Theory of Mind), כלומר היכולת להבין שלזולת יש מחשבות ואמונות השונות משלך, אך מדובר ביכולת מוגבלת שנשענת על הסתמכות יתר על היריסטיקות ולא על שיקול דעת כללי וחסין. ראו Natalie Shapira et al., *Clever Hans or Neural Theory of Mind? Stress Testing Social Reasoning in Large Language Models* (24.5.2023), available at arXiv

19 ראו לדוגמה הדיון אצל VISA KURKI, *A THEORY OF LEGAL PERSONHOOD* 175–189 (2019);
SIMON CHESTERMAN, *WE, THE ROBOTS? REGULATING ARTIFICIAL INTELLIGENCE AND THE LIMITS OF THE LAW* 114–142 (2021)

Annalee Newitz, *The Curious Case of the AI and the Lawyer*, 255 NEW SCIENTIST 28 (23.7.2022)

שהיא חוללה?²¹ בית המשפט הפדרלי לערעורים בארצות הברית ענה לאחרונה בשלילה על שאלה זו,²² אך הסוגיה עודנה מוסיפה להעסיק ערכאות אחרות.²³

יוער כי יש המתנגדים לשימוש במונח בינה מלאכותית בטענה שהוא מעורר קונוטציות למכונות המפתחות תודעה ורגשות אנושיים ומעודד השוואות שגויות. לבריהם, לא מדובר אלא באוסף של טכניקות אלגוריתמיות סטטיסטיות ולמידת מכונה.²⁴ כך או כך, כפי שנסביר בהמשך החיבור, נדרשת המשגה חדשה שתוכל להרחיק את הדיון מהשוואות שגויות אלו. למשל, החלפת המונח "בינה מלאכותית" ב"למידה חישובית".

ספר זה אינו על בינה מלאכותית במשפט הפרטי, אלא על סוגיות מעולם המשפט הציבורי והמעין ציבורי. מקצת הסוגיות נוגעות לתהליכי הפיתוח של מערכות בינה מלאכותית, הניזונות ממידע רב, ולשאלה מי בעל זכויות היוצרים על המידע,²⁵ או מה מידת הפגיעה בפרטיות הכרוכה בשימוש בו, ומקצתן יכולות להשיק לשאלת האישיות המשפטית ונוגעות לשאלות של אחריות במקום שבו נגרם נזק בשל פעילות של מערכת בינה מלאכותית.²⁶ למשל, מי אשם במקרה של

21 לדין ראו AVIV H. GAON, THE FUTURE OF COPYRIGHT IN THE AGE OF ARTIFICIAL INTELLIGENCE (2021) Ryan Abbott and Elizabeth Shubov, *The Revolution Has Arrived: AI Authorship and Copyright Law*, FLORIDA LAW REVIEW, forthcoming (8.8.2022); Anke Moerland, *AI and Intellectual Property Law*, in THE CAMBRIDGE HANDBOOK OF PRIVATE LAW AND ARTIFICIAL INTELLIGENCE (Philip Morgan and Ernest Lim eds., forthcoming 2023)

22 Thaler v. Perlmutter, No. 1:20-cv-00903-LMB-TCB (E.D. Va. Aug. 5, 2022) 23 Thaler v. Comptroller-General of Patents, Designs and Trade Marks, UKSC [2022] (Decision granting permission to appeal, 12.8.2022)

24 LUC JULIA, THERE IS NO SUCH THING AS ARTIFICIAL INTELLIGENCE (2019); Stefano Quintarelli, *Let's Forget the Term AI. Let's Call Them Systematic Approaches to Learning Algorithms and Machine Inferences (SALAMI)*, QUINTA'S WEBLOG (24.11.2019)

25 לדוגמה ראו משרד המשפטים, שימושים בתכנים מוגנים בזכויות יוצרים לצורך למידת מכונה (18.12.2022).

26 לדוגמה מן העת האחרונה ראו את הצעת דירקטיבת האחריות האזרחית בעניין בינה מלאכותית. הדירקטיבה מציעה כללים של גילוי ראיות בעניין מערכות בינה מלאכותית בסיכון גבוה (ראו להלן בסעיף 4.1.1.2) במסגרת חובענות אזרחיות שאינן על בסיס חוזי, וכללי נטל ההוכחה בעניין חובענות אזרחיות שאינן על בסיס חוזי בשל עליות נזיקיות

תאונה של רכב אוטונומי וכיצד קובעים כי התקיימה רשלנות?²⁷ מי אחראי אם יש נפילה בכורסה בשל שגיאה במערכת האלגוריתמית, הסוחרת במניות במהירות הבזק, או בשל מניפולציות בתמחור מניות, הדרושת יסוד נפשי של כוונה?²⁸

ספר זה אינו על הסכנה הנשקפת לתעסוקה מבינה מלאכותית. לצד היתרונות הכרוכים באוטומציה של תהליכים ואיכות החיים הנלווים לה מתעורר גם חשש שהקדמה תצמצם את הביקוש לידיים עובדות, תייתר מקצועות והתמחויות מסוימים ותגדיל את האבטלה. בסוף הרבעון הראשון של 2023 בחן בית ההשקעות גולדמן סאקס 900 מקצועות ב-22 תעשיות והעריך כי רבע ממשיות העבודה הנוכחיות בתעשיות אלו יכולות להתבצע באופן אוטומטי על ידי בינה מלאכותית.²⁹ הרוח מעריך כי בשל מבנה שוק העבודה הישראלי, בישראל מדובר בכמעט 30% משוק העבודה.

כיום ברור שלא כמו בראשית המהפכה התעשייתית ובתגובת הנגר הלודיטית, המקצועות שנשקפת להם סכנה רבה יותר אינם מקצועות שכרוכים בעבודת כפיים ("צווארון כחול"), אלא מקצועות כגון כתיבה ורעיונאות (קופירייטינג), רפואה ואבחון, ניתוח מידע ונתונים, עיצוב ואומנות, ראיית חשבון ומשפטים, וגם תכנות.³⁰ יש להודות ביושר כי בחלק מן המשימות הנדרשות במקצועות אלו יוכלו מכונות לעלות בביצועיהן על בני אנוש. יתר על כן, ייתכן שהיעלמות מקצועות מסוימים אינה נובעת רק משימוש במערכות אלגוריתמיות, אלא גם

שמקורן במערכות בינה מלאכותית. European Commission, Proposal For a Directive of the European Parliament and of the Council on Adapting Non-Contractual Civil Liability Rules to Artificial Intelligence (AI Liability Directive) COM/2022/496 (28.9.2022) (להלן: הצעה לדיקטיבה אירופית לאחריות בתחום הבינה המלאכותית).

27 Sven Nyholm and Jilles Smids, *The Ethics of Accident-Driving Cars: An Applied Trolley Problem*, 19 *ETHICAL THEORY AND MORAL PRACTICE* 1275 (2016); Karni Chagal-Feferkorn, *The Reasonable Algorithm*, 2018 U. ILL. J.L. TECH. & POL'Y 111 (2018)

28 Gina-Gail S. Fletcher, *Deterring Algorithmic Manipulation*, 74 *VAND. L. REV.* 259 (2021)

29 Jan Hatzius et al., *The Potentially Large Effects of Artificial Intelligence on Economic Growth*, GOLDMAN SACHS ECONOMIC RESEARCH (26.3.2023)

30 Michael Webb, *The Impact of Artificial Intelligence on the Labor Market* (6.11.2019)

ממודלים כלכליים חדשים – דוגמת כלכלת החלטורות (the gig economy) וכלכלת השיתוף (the sharing economy) – שנעזרים בפלטפורמות מקוונות כדי לערער את המערכות המסחריות המסורתיות. כרגע נסתפק בהערכה שגם מקצועות שצופים כי ייעלמו – דוגמאות נפוצות הן הרדיולוגיה, שנעשית כיום באמצעות מערכות אוטומטיות,³¹ או עיצוב, שנעשה באמצעות תוכנות גנרטיביות – גם מקצועות אלו יודקו למומחה אנושי שיחלוק את העבודה עם מכונה כשם שמטוסים בעלי טייס אוטומטי אינם ממריאים בלי קברניט בשר ודם.³² לפיכך פיתוח היכולת לעבוד עם מכונות חכמות ולצידן ילך ויהפוך למיומנות עבודה חיונית ויתבטא במשימות שמכונות היום "לחישת" (שכלול הנוסחה שתזון למכונה כדי לקבל תוצר מיטבי – קוד, טקסט או דימוי חזותי); בהבנת תהליכי היסק והנמקה של מכונות כדי להשלים תהליכי חשיבה אנושיים;³³ ובפיקוח ברמת שונות (קביעת דרגת אוטונומיה של רכב, חוות דעת על אבחון רפואי ועוד).

בינה מלאכותית, לרבות האסדרה שלה, היא תחום דינמי שמשתנה תדיר. ספר זה נועד להיות טקסט מבואי ולהציג למי שאינם בעלי רקע טכנולוגי הן את התחום עצמו הן את אתגרי המדיניות והרגולציה שהוא מעורר. מחקר זה עדכני ליוני 2023 ואינו כולל שינויים או התפתחויות מאוחרים יותר.

אנו מבקשים להודות לפרופ' סוזי נבות, לפרופ' יובל שני, לד"ר רחל ארידור הרשקוביץ ולעו"ד גדי פרל על הערותיהם הטובות לאורך הדרך. לד"ר רועי צזנה ולשירה פטרק על עזרה מחקרית בשלבים הראשונים של המחקר. תודה לאנשי ההוצאה לאור של המכון הישראלי לדמוקרטיה.

31 על רדיולוגיה ראו Ramón Alvarado, *Should We Replace Radiologists with Deep Learning? Pigeons, Error and Trust in Medical AI*, 36 *BIOETHICS* 121 (2022)

32 על הסתגלות של מומחים בתחומי האבחון הרפואי ראו Saurabh Jha and Eric J. Topol, *Adapting to Artificial Intelligence: Radiologists and Pathologists as Information Specialists*, 316 *JAMA* 2353 (2016); לדיון

בנושא טכנולוגיה ופרופסיות ראו DANIEL SUSSKIND AND RICHARD SUSSKIND, *THE FUTURE OF THE PROFESSIONS: HOW TECHNOLOGY WILL TRANSFORM THE WORK OF HUMAN EXPERTS* (2015)

33 להרחבת רעיון זה ראו למשל איתי ברון ותהילה שוורץ אלטשולר "מקום ליד השולחן – האם ה-ChatGPT יכול 'לשבור את הקונספציה'?" המכון לחקר המתודולוגיה של המודיעין (1.5.2023).

פרק ראשון

מהי בינה מלאכותית?

—

הרעיון של ישויות מלאכותיות, אוטומטונים, שאפשר להנדס או להרכיב, נדון ברצינות לאורך השנים.³⁴ בשנות הארבעים של המאה העשרים טבע נורברט וינר את המונח קיברנטיקה כדי לתאר מערכות מבוססות משוב (feedback) והצביע על הדמיון בין פעולתם של מחשבים ובין תהליכים של זיהוי דפוסים במוח האנושי.³⁵ אלן טיורינג, בן תקופתו של וינר, שעסק בתאוריות מתמטיות ובמדעי

34 ראו למשל תומאס הובס לויטן 5 (מתרגם אהרן אמיר 2010).

35 NORBERT WIENER, CYBERNETICS: OR CONTROL AND COMMUNICATION IN THE ANIMAL AND THE MACHINE (1948)

המחשב, מוכר כמי שהציע מבחן היפותטי שנועד לתת מענה אפשרי לשאלה אם מכונה יכולה לחשוב כמו בן אדם.³⁶ תחת שאלה טעונה זו הציע טיורינג משחק שקרא לו משחק החיקוי ולימים נודע בשם מבחן טיורינג. למשחק נדרשים שני בני אדם ומכונה. אדם אחד, המשמש כחוקר, מנהל במקביל דיאלוג בכתב עם האדם השני ועם המכונה, בלי לדעת מי הוא מי. אם החוקר אינו מצליח להבחין בין האדם ובין המכונה אפשר לומר שהמכונה עברה את מבחן טיורינג.

הצירוף בינה מלאכותית נטבע בשנת 1956, בעקבות ההצעה לסדנת המחקר הראשונה בתחום, שנערכה באותה שנה באוניברסיטת דארטמות'. הסדנה התבססה על ההשערה כי "כל היבט של לימוד או כל מאפיין אחר של תבונה יכולים, עקרונית, להיות מתוארים באופן מדויק די הצורך שמכונה תוכל לחקות אותם".³⁷ חרף העניין המחקרי בתחום, לקראת אמצע שנות השבעים דעך העניין בכינה מלאכותית לנוכח האכזבה מההתקדמות בפועל.³⁸

התקופה המכונה "החורף של הבינה המלאכותית" הסתיימה בראשית שנות השמונים, בעקבות פעפוע רעיונות של למידה על ידי חיזוקים (reinforcement learning) לענף מתחום הפסיכולוגיה והתפתחויות בתחום הרובוטיקה.³⁹ לאחר מחזור נוסף של "קיץ" ו"חורף" בשנות השמונים – תנופת השקעות בכינה מלאכותית ואחריה תקופה של אכזבה ומשיכת השקעות מהתחום – פותחו בשנות התשעים מערכות בינה מלאכותית שידעו להתמודד עם בעיות מתמטיות מורכבות (מסוג NP-שלמות). נקודת ציון נוספת היא התבוסה של השחמטאי

Alan Turing, *Computing Machines and Intelligence* 59 MIND 433 (1950) 36

J. McCarthy et al., *A Proposal for the Dartmouth Summer Research Project on Artificial Intelligence* (31.8.1955) 37

38 ראו למשל את סקירתו של המדען הבריטי ג'יימס לייטהיל, שלפיה לא קיים המחקר בתחום את ההבטחות להשפעה רחבת היקף ולכן קטן המימון בבריטניה, ובהמשך באירופה ובארצות הברית. James Lighthill, *Artificial Intelligence: A General Survey*, in ARTIFICIAL INTELLIGENCE: A PAPER SYMPOSIUM (Science Research Council 1973)

39 Sami Haddadin and Dennis Knobbe, *Robotics and Artificial Intelligence*, in ALGORITHMS AND LAW 1, 9-12 (Martin Ebers and Susana Navas eds., 2020)

הרב-אומן גארי קספרוב לכחול עמוק, מערכת בינה מלאכותית שפיתחה איי-בי-אם. ההצלחה של כחול עמוק נשענה על שני רכיבים: הראשון, אלגוריתם חיפוש היריסטי (כלומר התבססות על כללי אצבע כדי לזרז את תוצאות החיפוש); השני, כוח מחשוב רב-עוצמה. רכיב זה הוביל לביקורת שכחול עמוק אינו מבוסס על בינה מלאכותית של ממש, אלא מסתמך על כוח חישובי גס כדי להביס רב-אמן אנושי.⁴⁰

אף שכבר באמצע שנות השישים הציגו ולנטין לאפה ואלכסיי יאבנקו רעיונות בדבר רשתות עצביות עמוקות (deep neural network),⁴¹ נדרשו עוד חמישים שנה עד שהתחום הבשיל. ב-2012 זכה אלגוריתם הלמידה העמוקה של ג'פרי הינטון בפער של ממש בתחרות זיהוי תמונות והגיע לרמת דיוק של 85% (רמת הדיוק האנושית היא 95%).⁴² אפליקציית חלום עמוק של גוגל, שהושקה ב-2015, תרמה גם היא למודעות הציבור לפוטנציאל הטמון בלמידה עמוקה.⁴³

מערכת הלמידה העצמית של מערכת אלפא-גו (AlphaGo) סייעה לה להביס את אלוף העולם במשחק גו (משחק אסטרטגיה מופשט שמקורו בסין). בשנת 2017 הוענקה לה דרגת דאן 9 של כבוד.⁴⁴ המערכת הוננה בעשרות אלפי משחקי גו ושכללה את כישוריה במשחק נגד עצמה. אלפא-גו הציגה אסטרטגיות וסגנונות משחק יצירתיים, ששחקנים אנושיים לא חשבו עליהן בעבר. הגרסה הנוכחית של אלפא-גו, אלפא-גו זירו, לא נדרשה כלל למאגר של משחקים היסטוריים כדי לאמן את עצמה. כללי המשחק הוננו לתוכה ואלפא-גו זירו שיחקה נגד עצמה – בתוך שמונה שעות בלבד הגיעה לרמה שאפשרה לה להביס את אלפא-גו. אלפא-גו זירו התבססה על מנוע ששמו אלפא זירו, שאינו מוגבל רק למשחק הגו

MICHAEL WOOLDRIDGE, *THE ROAD TO CONSCIOUS MACHINES: THE STORY OF AI* (2020) 40

VALENTIN GRIGORYEVICH LAPA AND ALEXEY GRIGORYEVICH IVAKHNENKO, *CYBERNETIC PREDICTING DEVICES* (1965) 41

Alex Krizhevsky, Ilya Sutskever, and Geoffrey E. Hinton, *ImageNet Classification with Deep Convolutional Neural Networks* 25 *ADVANCES IN NEURAL INFORMATION PROCESSING SYSTEMS* 1106 (2012) 42

Christian Zegedy et al., *Going Deeper with Convolutions*, *COMPUTING RESEARCH REPOSITORY* (2014) 43

Google's AlphaGo Gets "Divine" Go Ranking, *THE STRAITS TIMES* (15.3.2016) 44

האסיאתי. אלפא־גו זירו למד שחמט בתוך שעות ספורות וניצח את אחת התוכנות המתקדמות בעולם השחמט באותה עת – כל זאת בלי שהוזן לתוכו ולו משחק אחד של שחמט.⁴⁵ דוגמאות אלו מראות שבינה מלאכותית מסוגלת להגיע לרמת מיומנות על־אנושית בתחומים המצריכים שנים ארוכות של למידה ואימון. יתר על כן, ככל שנוקפות השנים מהירות הפיתוח של הטכנולוגיות המאפשרות זאת הולכת וגדלה.

1.1

הגדרות של בינה מלאכותית

דמיס האסאביס, מייסד DeepMind, החברה אשר פיתחה את אלפא־גו, הגדיר בינה (intelligence) "תהליך שהופך מידע לא

סדר לידע שימושי".⁴⁶ בינה מלאכותית אינה חייבת להיות בעלת נוכחות פיזית במרחב. אין לה צורך ברובוטים או בחומרה ייחודית אחרת (ובלבד שכוח העיבוד העומד לרשותה מאפשר לה להתקיים). אומנם אפשר לדמיין, כפי שעשו סופרי המדע הבדיוני לדורותיהם,⁴⁷ בינה מלאכותית שיש לה נוכחות רובוטית גשמית (כדומה לכינה אנושית). עם זאת, מערכות בינה מלאכותית בנות זמננו, דוגמת מערכת הניווט ויזי או מערכת התרגום של גוגל, הן מערכות אלגוריתמיות שאינן מבוססות רובוטיקה.

ג'ון מקארתי, הנמנה עם האבות המייסדים של דיסציפלינת הבינה המלאכותית, גרס כי המטרה היא לפתח מכונות שמתנהגות כאילו הן בעלות תבונה.⁴⁸ מבחן טיורינג, וניסויי מחשבה אחרים (דוגמת החדר הסיני, שהציע ג'ון סרל כדי למתוח ביקורת על מבחן טיורינג),⁴⁹ מרמזים שבעיית התודעה של הבינה המלאכותית

45 "אינטליגנציה מלאכותית לימדה עצמה שחמט, והפכה לאחת השחקניות הטובות בעולם" הארץ (10.12.2017).

46 Demis Hassabis: *Creativity and AI – The Rothschild Foundation Lecture*, YouTube (10.10. 2018)

47 קארל צ'אפק "R.U.R" דחק ז 424 (תרגום אברהם שלונסקי 2016); איזיק אסימוב *אנוכי הרובוט* (תרגום יאיר שמעוני 1976).

48 McCarthy et al., *לעיל* ה"ש 37.

49 John. R. Searle, *Minds, Brains, and Programs*, 3 BEHAVIORAL AND BRAIN SCIENCES 417 (1980)

היא בעיה פילוסופית בעיקרה ושההגדרה המעשית של בינה מלאכותית צריכה להידרש לדרכי הפעולה שלה. יש הגורסים שהשימוש במונח בינה מלאכותית הוא מתעתע. לדבריהם, תווית זו משמשת לתיאור יכולות טכנולוגיות שאין בהן בינה שדומה לבינה האנושית.⁵⁰

מועצת אירופה הציעה להגדיר בינה מלאכותית "מערכות המדגימות התנהגות תבונית על ידי ניתוח סביבתן ונקיטת פעולות – ברמה מסוימת של אוטונומיה – כדי להשיג מטרות מוגדרות".⁵¹ קבוצת המומחים שכינסה מועצת אירופה (AI HLEG) הדגישה בראש ובראשונה שמערכות אלו הן רציונליות, והוסיפה להגדרה עוד היבטים, בהם היכולת לפרש מידע ולהסתגל לשינויים בסביבה.⁵²

בנוסף המקורי של הצעת תקנות הבינה המלאכותית האירופיות⁵³ הוצע להגדיר בינה מלאכותית "תוכנה שפותחה באמצעות אחת או יותר מהשיטות והגישות המפורטות בתוספת הראשונה [למידת מכונה; גישות מבוססות ידע (מערכות מומחה); או שיטות סטטיסטיות] ויכולה, בהינתן מטרות שהגדירו בני אדם, לייצר פלטים כגון תוכן, תחזיות, המלצות או החלטות המשפיעות על סביבה אמיתית או מדומה". הגדרה דומה מצויה בהצעת החוק הברזילאית לבינה מלאכותית,⁵⁴ ואולם זו מדגישה את היסוד האינטראקטיבי של בינה מלאכותית, שכן היא החריגה במפורש מתחולת החוק מערכות קבלת החלטות המבוססות על פרמטרים

50 ראו לדוגמה Michael I. Jordan, *Artificial Intelligence – The Revolution Hasn't Happened Yet*, 1 HARV. DATA SCI. REV. (2019)

51 ראו European Commission, A EUROPEAN APPROACH TO ARTIFICIAL INTELLIGENCE (להלן: (AI4EU).

52 AI HLEG, A DEFINITION OF AI: MAIN CAPABILITIES AND SCIENTIFIC DISCIPLINES (8.3.2021)

53 PROPOSAL FOR A REGULATION OF THE EUROPEAN PARLIAMENT AND OF THE COUNCIL LAYING DOWN HARMONISED RULES ON ARTIFICIAL INTELLIGENCE (ARTIFICIAL INTELLIGENCE ACT) AND AMENDING CERTAIN UNION LEGISLATIVE ACTS, COM (2021) 206 (21.4.2021) (להלן: הצעת תקנות הבינה המלאכותית האירופיות), בטעיף 3(1) ובחוספת הראשונה.

54 Projeto De Lei N° 21-A de 2020 (29.9.2021). לתרגום אנגלי של הצעת החוק ראו Walter Gaspar, *Non-Official Translation of the Brazilian Artificial Intelligence Bill*, N. 21/2020, CYBERBRICS (25.10.2021) (להלן: הצעת החוק הברזילאית לבינה מלאכותית).

שנקבעו מראש ואין ביכולתן ללמוד, לתפוס, לפרש ולהגיב לסביבה החיצונית על בסיס קלט. הצעת החוק הקנדית לבינה מלאכותית ולמידע⁵⁵ מדגישה את ההיבט האוטונומי של מערכות בינה מלאכותית, המוגדרות מערכות טכנולוגיות שמעבדות מידע שנוגע לפעילות אנושית באופן אוטונומי או אוטונומי למחצה באמצעות אלגוריתם גנטי,⁵⁶ למידת מכונה או טכניקה אחרת במטרה ליצור תוכן, לקבל החלטות, לנסח המלצות או לייצר תחזיות.

ההגדרה המתוקנת בנוסח הפשרה של הצעת תקנות הבינה המלאכותית האירופיות⁵⁷ השמיטה את ההתייחסות לאופן הפיתוח של המערכת, שהייתה בתקנות בהגדרה שבהצעה המקורית. על פי ההגדרה המתוקנת, מערכת בינה מלאכותית היא "מערכת מבוססת מכונה שעוצבה כדי לפעול בדרגות אוטונומיה שונות ויכולה, למטרות מפורשות או לא מפורשות, לייצר פלט כגון תחזיות, המלצות או החלטות המשפיעות על הסביבה הפיזית או הווירטואלית".⁵⁸

מסמך המדיניות (white paper) של משרד הדיגיטל, המדיה, התרבות והספורט של בריטניה מצביע על שני מאפיינים עיקריים של בינה מלאכותית הרלוונטיים

Bill C-27, An Act to Enact the Consumer Privacy Protection Act, 55 the Personal Information and Data Protection Tribunal Act and the Artificial Intelligence and Data Act and to Make Consequential and Related Amendments to Other Acts, 1st Sess, 44th Parl, 2022 (First Reading, 16 June 2022), p. 85-100 (להלן: הצעת חוק הבינה המלאכותית והמידע של קנדה).

56 אלגוריתמים גנטיים הם אלגוריתמים בהשראת רעיון הברירה הטבעית, המשמשים בעיקר להתמודדות עם בעיות של חיפוש ואופטימיזציה. להרחבה ראו, SEYEDALI MIRJALILI, EVOLUTIONARY ALGORITHMS AND NEURAL NETWORKS 43-55 (2019).

57 בחודש מאי 2023 התקבל בוועדת השוק הפנימי ובוועדת חרירות האזרח של הפרלמנט האירופי נוסח מתוקן של הצעת תקנות הבינה המלאכותית האירופיות. Committee on the Internal Market and Consumer Protection and the Committee on Civil Liberties, Justice and Home Affairs, Draft Compromise Amendments on the draft report proposal for a regulation of the European Parliament and of the Council on harmonised rules on Artificial Intelligence (Artificial Intelligence Act) and amending certain Union Legislative Acts Version 1.1 (16.5.2023) (להלן: ההצעה המתוקנת).

58 ס' 3(1) הצעה המתוקנת, שם.

לרגולטור: אדפטיביות ואוטונומיות. אדפטיביות היא היכולת לפעול על בסיס הוראות ודפוסים שלא הוכתבו במפורש על ידי גורם אנושי אלא על ידי "אימון" או למידת מכונה (ראו להלן בסעיף 1.3). מערכות בינה מלאכותית הן גם אוטונומיות – היכולת של בינה מלאכותית לבחור אסטרטגיות פעולה, לנתח נתונים ולהגיב על פיהם היא מה שהופך אותה ל"אינטליגנטית"⁵⁹.

בסיכומו של דבר, המאפיינים המשותפים להגדרות המשפטיות לבינה מלאכותית שהובאו לעיל הם היכולת של מערכות אלו לפעול באופן אוטונומי ולהגיב לשינויים בסביבתן, בהתבסס על טכנולוגיות שונות, שלעיתים מצוינות במפורש בהגדרה. האוטונומיה והאדפטיביות הן גם שמאפשרות למערכות אלגוריתמיות ליצור את הרושם שהתנהגותן היא אנושית.

1.2

דרגות של בינה מלאכותית

מקובל לחלק בינה מלאכותית לשלוש דרגות. הדרגה הנמוכה ביותר נקראת בינה מלאכותית צרה או חלשה. בינה מלאכותית צרה היא מערכת

שיכולתיה משתוות ליכולת האנושית בתחום מוגדר אחד, או בכמה תחומים קשורים, או עולות עליה. כל מערכות הבינה המלאכותית בנות זמננו משתייכות לקטגוריה זו: מערכת בחירת השירים של אפליקציית ספוטיפיי יודעת לקלוע לטעמו של המשתמש כמו די-ג'יי אנושי (יש שחולקים על קביעה זו), אך אינה מסוגלת לנווט במרחב כמו האלוגריתם של אפליקציית ויזו, או לשחק שחמט ברמת המומחיות של אלפא-זירו, ואילו הן אינן יכולות כמוכן לבחור שירים כמו ספוטיפיי.

הדרגה הבאה נקראת בינה מלאכותית כללית או חזקה. סרל נדרש לאפשרות של בינה מלאכותית חזקה כשאמר כי "למחשב מתוכנת כיאות, שיש לו הקלטים והפלטים הנכונים, תהיה תודעה [mind] בדיוק באותו מובן שבו יש לבני אדם

תודעה"⁶⁰ אומנם מרבית המפתחים של בינה מלאכותית לא בהכרח חותרים ליצור מערכת תבונית שעונה על הדרישות התאורטיות הנוקשות של סרל, אבל פיתוח של בינה מלאכותית כללית חותר לעצב מכונות המסוגלות להתמודד עם בעיות רבות, מתחומים שונים, שאינן מוגדרות במדויק.

בינה מלאכותית כללית (AGI) היא ייצוג של יכולות קוגניטיביות אנושיות כלליות בתוכנה. למערכת AGI היכולת למצוא פתרון למשימה לא מוכרת. מערכת כזאת חותרת אפוא לבצע כל משימה שאדם מסוגל לה ברמה של ביצוע אנושי. ההגדרות של AGI משתנות כיוון שמומחים מתחומים שונים מגדירים אינטליגנציה אנושית מנקודות מבט שונות. אבל מדעני מחשב מגדירים את האינטליגנציה האנושית במונחים של יכולת להשיג מטרות, ומקצתם סבורים כי כשתושג בינה מלאכותית כללית יעלו יכולותיה על היכולות האנושיות כיוון שבכוחה לגשת למערכי נתונים עצומים ולעבד אותם במהירות. הדברים אמורים בחשיבה מופשטת, גישה למידע רקע, היגיון פשוט, הבנת סיבה ותוצאה ויכולת ליישם את מה שנלמד במשימה אחת במשימות אחרות וחדשות.

דוגמאות מעשיות ליכולות האלה הן יצירתיות (היכולת לכתוב טקסט, לייצר דימוי חזותי או אורקולי או לכתוב קוד); תפיסה חושית וחזותית (למשל היכולת לדמיין תמונה תלת-ממדית על יסוד חשיפה לתמונה דו-ממדית); יכולות מוטוריות מפותחות (למשל היכולת להוציא צרור מפתחות מתוך כיס, פעולה שייתכן שהיא דורשת גם כושר דמיון); הבנת שפה טבעית (בעיקר את ההיבט ההקשרי של השפה); ויכולות ניווט מפותחות בתוך חללים.

עינינו הרואות שבתקופה האחרונה הולכות ומושגות מקצת היכולות האלה. מגוון מוצרי קצה מבוססים על יכולת למידה עצמית ועל אלגוריתמים של למידה, על בניית מבנים קבועים עבור משימות, על הבנת מערכות סמלים, על שימוש בסוגים שונים של ידע ועל יכולות מטא-קוגניטיביות. לכן מקובל להניח היום כי אין לשרטט הבחנה דיכוטומית בין בינה מלאכותית צרה לבינה מלאכותית כללית אלא יש לרבר על ספקטרום של יכולות. ספקטרום זה מכונה לעיתים Artificial More-general Intelligence (AMGI).

הדרגה האחרונה היא בינה מלאכותית עילאית (superintelligence). ב-1965 חזה המתמטיקאי האנגלי ג'ק גוד את אפשרות קיומן של מכונות שיתעלו מעל הפעילות האינטלקטואלית של החכם באדם. מאחר שמכונות כאלו יוכלו להרכיב מכונות חכמות מהן, יוביל הדבר "בוודאי ל'פיצוץ' של אינטליגנציה"⁶¹. לפי חזון זה, השילוב של יכולות אלו עם הצמיחה המעריכית של יכולות העיבוד של מחשבים יובילו לסינגולריות – עידן עתידי שבו קצב השינויים הטכנולוגיים יהיה מהיר כל כך, והשפעתו עמוקה כל כך, עד שהמין האנושי ישתנה לבלי הכר.⁶²

ואולם ספר זה אינו עוסק בעתידנות אלא במלאכת ההווה. כאמור, מרבית מערכות הבינה המלאכותית הקיימות כיום הן מערכות של בינה מלאכותית צרה ובינה מלאכותית "רחבה" או "כללית". כפי שהולך ומתברר מסוף שנת 2022, חברות ומדינות רבות משקיעות משאבים אדירים בפיתוח יכולות של בינה מלאכותית כללית, אבל עדיין יש מחלוקת באשר ללוח הזמנים הנדרש למימושו. קל וחומר שיש דרך ארוכה, ושנויה במחלוקת אף יותר, עד הגשמת חזון הסינגולריות – בין שהוא רצוי ובין שאין הוא רצוי.

אנו נתרכז אפוא בהשלכות של אותן מערכות בינה מלאכותית שכבר קיימות בשוק על המציאות של היום. מאחר שהמערכות הקיימות משתייכות לקטגוריה של בינה מלאכותית צרה לא נדון בהשלכות החברתיות והמוסריות של בינה מלאכותית "כללית" או "עילאית", כלומר לא נדון למשל בשאלה אם למערכות בינה מלאכותית יש זכויות אדם ובאלו תנאים.⁶³

Irving John Good, *Speculations Concerning the First UltraIntelligent Machine* 6 ADVANCES IN COMPUTERS 31 (1966)

RAY KURZWEIL, *THE SINGULARITY IS NEAR* (2005) 62

63 ראו בעניין זה קביעתו של המשנה לנשיאה מלצר בפס' 31 בפסק דינו בעניין בג"ץ 7846/19 פרקליטות המדינה – יחידת הסייבר נ' לורי שם טוב (12.4.2021), פורסם באר"ש), שלפיה "אין זכויות אדם [...] לרובוטים, מה גם שחלק מאותם רובוטים אפילו לא מופעלים על-ידי בני אדם, אלא על-ידי אינטליגנציה מלאכותית". לדיון מעמיק יותר בשאלה האישיות המשפטית של מערכות נבונות ראו Kurki, לעיל ה"ש 19, בעמ' 175-189.

1.3

למידת מכונה ולמידה חישובית

כבר עתה מחשבון יד פשוט מבצע פעולות אריתמטיות במהירות שכן אנו אינו יכול להשתוות לה. עם זאת, ספק אם אפשר לטעון

שלמחשבון היד יש בינה מלאכותית. המחשבון לא ניחן באוטונומיה, בניתוח רציונלי וביכולת הסתגלות. מכונות בעלות בינה מלאכותית יכולות ללמוד ולהשתכלל.

עד מחצית שנות השמונים הייתה הפרדיגמה הקלסית בינה מלאכותית סימבולית (המוכרת גם בשם Good Old-Fashioned AI, GOFAI). בינה מלאכותית סימבולית מבוססת על קידוד המציאות במכונה באמצעות סמלים. המערכת יכולה לנתח את היחסים שבין הסמלים בהסתמך על מנוע שהחוקים קודדו בו מראש כמארג אלגוריתמי מסועף של הוראות "אם-אז" (if-then). כדי לפתח מערכת שמבוססת על בינה מלאכותית סימבולית נדרש ידע אפריורי בתחום שבו היא עתידה לפעול – המערכת לא תלמד תחום זה באופן עצמאי ותשכלל את האלגוריתם שלה בהתאם, אלא תסתמך על מאגר המידע שהוזן לתוכה. המערכת יכולה לקבל החלטות או המלצות המבוססות על נתונים אלו ברמת מומחיות שמומחה אנושי אינו יכול להגיע אליה ומערכות מומחה כאלו היו הדומיננטיות ביותר בשוק לאורך זמן.

שלא כמו ייצוג סימבולי של מידע מתוך מאגר מידע אפריורי ושל תהליכים קוגניטיביים מסדר גבוה, שעליו מסתמכות מערכות מומחה מטיפוס GOFAI, בגישת הקשרנות (connectionism) לפיתוח בינה מלאכותית⁶⁴ מתחילים מלמטה. גישה זו מבוססת על ההנחה שפיתוח רשתות שמורכבות מיחידות עיבוד פשוטות המכונות נוירונים (אף שאין הם דומים לנוירונים אנושיים), וקישור הרשתות אלו לאלו, יכולים להניב אינטליגנציה – בדומה לאופן שבו המוח האנושי מקושר.

64 Warren S. McCulloch and Walter Pitts, *A Logical Calculus of the Ideas Immanent in Nervous Activity* 5 THE BULLETIN OF MATHEMATICAL BIOPHYSICS 115 (1943); MARVIN MINSKY AND SEYMOUR A. PAPERT, PERCEPTIONS: AN INTRODUCTION TO COMPUTATIONAL GEOMETRY (1969). לקישור מוקדם בין קישוריות לבינה מלאכותית כללית ראו קלארק, לעיל ה"ש 9.

"למידה היא משהו בלתי רצוני שאנחנו, בני האדם, עושים כל הזמן. אנחנו משתמשים במידע שאנחנו קולטים כדי לייצר ידע". כך הסביר בריאיון לעיתון "דה מרקר" פרופ' שי שלו שוורץ, סגן נשיא חטיבת הטכנולוגיה בחברת מובילאיי ומומחה לתחום הלמידה העמוקה מהמחלקה למדעי המחשב באוניברסיטה העברית.⁶⁵ "אם נראה לאדם כלשהו תמונה שמוצגת בה מכונית, הוא יידע שיש בה מכונית, אך לא יידע להסביר איך הוא יודע את זה. זה מה שאני, כחוקר, רוצה לגלות וזה מה שמערכות למידה עמוקה עושות – מנסות להבין כללים. אלה כלים אוטומטיים של הסקה מהפרטים אל הכלל".

מערכות הבינה המלאכותית נועדו להתמודד עם כמויות עצומות של מידע ולהפיק ממנו תובנות מופשטות ברמות תחכום גבוהות. לכן ככל שהן קולטות ומעבדות מידע רב יותר מכל התחומים, כך משתפרת היכולת שלהן ואפשר להטמיע אותן בתחומים רבים יותר. על כן נוסף על ההשפעה המכריעה שיש לאיכות אלגוריתם הלימוד שבבסיס המכונה, יכולת הלמידה שלה תלויה גם באיכותו ובהיקפו של המידע המוזן למכונה.

"למידת מכונה" היא שם של אוסף טכניקות סטטיסטיות, בהן טכניקות הקשורות לרשתות נוירונים, שנועדו לזהות דפוסים או כללים שיכולים להסביר נתונים או לחזות תוצאות עתידיות. זהו תהליך שנעשה באופן אוטונומי, ללא כללים או חוקים ישירים שקודד מפתח המערכת; שלא כמו מערכות מומחה, המבוססות על תובנות אנושיות שהפכו לקוד, למידת מכונה מפיקה תובנות אלו בכוחות עצמה. כשאלגוריתם זיהוי תמונה כגון זה שפיתח ג'פרי הינטון מבחין בין דמות של כלב לדמות של מכונית, אין הוא מתבסס על מערכת של כללים (כגון "למכונית יש ארבעה גלגלים", "לכלב יש ארבע רגליים וזנב"), אלא מייצר מודל מתמטי משלו, בהתאם למאגר הנתונים (במקרה זה, למשל, מאגר הנתונים יכול להיות אוסף של תמונות שכל אחת מהן מוגדרת "כלב", "מכונית" או "לא כלב ולא מכונית"). מאגר הנתונים סופק לצורכי אימון.

65 אלירן רובין | הכירו: Deep Learning – הדבר הכי לוהט בהיי-טק עכשיו" TheMarker (31.3.2016).

באמצעות למידת מכונה תוכנה מסוגלת לייצר בעצמה רצף של כללים על בסיס מאגר נתונים שעומד לרשותה ובהתאם למשוב של מתכנת או גורם אחר. ברמה הבסיסית ביותר למידת מכונה מייצרת עץ החלטות. עץ החלטות הוא בעצם רצף של החלטות "אם-אז" עד לקבלת תוצאה סופית. יתרונה של שיטה זו שהיא עקיבה וברורה להבנה. ואולם שיטת לימוד זו מוגבלת לפעולות שהן לינאריות באופיין וניתנות להמרה לעץ החלטות.

כשמדובר בביצוע משימות מורכבות יותר, יש צורך להפעיל מתודולוגיות למידת מכונה מתוחכמות יותר, כגון רשתות נוירונים. רשתות נוירונים נקראות כך משום שהן בנויות בדומה למבנה של מוח אנושי, שמייצר קשרים בין רצפטורים ונוירונים ומסתמך על הקשרים האלה כדי לבסס פעילות.⁶⁶ בלמידת מכונה מסוג רשת נוירונים המתכנת מזין לתוכנה מאגר נתונים, ונותן לאלגוריתם משוב – מתי הוא הצליח בביצוע הפעולה ומתי לא. הצלחה מתגמלת ומבססת את קשר הנוירונים ולאורך זמן התוכנה מוסיפה לבסס קשרים כאלה עד שהמתכנת סבור שהיא הגיעה לשיעורי הצלחה גבוהים די הצורך. התוצאה הסופית היא רצף של פקודות שמעיד על הצלחה סטטיסטית במשימה. מתודולוגיה זו, הנפוצה בשיטות הלמידה המתקדמות, יעילה מאוד, אך החיסרון שיש בה לענייננו הוא שהיא מתבססת על מתאם ולא על קשר סיבתי.

נוסף על כך, הליך הייצור של אלגוריתמים דורש משאבים ניכרים משני סוגים. ראשית, יש צורך במאגר נתונים גדול, עשיר ומטויב, אשר יאפשר לתוכנה להגיע מהר ככל האפשר להצלחה סטטיסטית במשימתה. ככל שהפעולה הנלמדת מורכבת יותר, כך נדרש מאגר נתונים גדול יותר. לכן אין פלא שהשחקנים המצליחים ביותר בתחום הבינה המלאכותית הם שחקנים שממילא יש בידם מאגרי מידע גדולים.⁶⁷ המשאב השני הוא פעולה אנושית: מתן משוב למכונה כדי שהיא תוכל ללמוד אם הצליחה או לא. פעולה זו אינה רק יקרה, אלא גם מועדת לטעויות; לכן גם בעניין זה יש עדיפות לחברות ענק עתירות משאבים

66 ראו בחלק 1.3 לעיל.

67 ראו בחלק 1.6.4 להלן.

שבאפשרותן לממן את ההשקעה הכבדה, המניבה תוצאות בטווח הארוך. בשל קשיים אלו פותחו דרכים שאינן מחייבות משוב אנושי.

דרך אחת, המתאימה במקרים מסוימים, היא לימוד אדוורסרי. למשל, ליצור מערכת אחת שיש לה מטרה, כגון לפרוץ צופן, ומערכת אחרת שיש לה מטרה הפוכה, כגון לייצר צופן שלא ניתן לפיצוח. בכל סבב לימוד המערכות משתכללות זו כנגד זו ללא צורך במשוב אנושי.

ואולם ברשתות הנוירונים טמון קושי של ממש, שקשור להיותן מעין "קופסה שחורה" – קשה להסביר באופן נקודתי כיצד נוצר הפלט המסוים שאליו הגיעו. המתכנת אינו יודע מה משמעות הקשרים שנוצרו וכיצד הגיעה התוכנה להצלחה, כלומר אין הוא יודע מדוע התוצר הסופי והפלט של התוכנה הם כפי שהם. אומנם מקובל להשוות בין רשתות נוירונים לבין פעולת המוח האנושי, אבל חשוב לציין שהדרך שבה הגיעה התוכנה להצלחה אינה סיבתית, אלא תוצר של פעולות סטטיסטיות, לעיתים נרדומליות. התוכנה אינה מסוגלת להפיק פלט שמסביר כיצד ומדוע התקבלה תוצאה או החלטה מסוימת ולא אחרת. נחזור ונעסוק בעניין זה בהמשך הדברים.

למידה עמוקה היא טכניקה המשמשת ללמד את המכונה לזהות דפוס מסוים על ידי "אימון" – למכונה מוזנת כמות גדולה של נתונים וניתנת לה הוראה לזהות במקבצי הנתונים דפוס או דפוסים מסוימים. לדוגמה, אפשר ללמד את המכונה להבחין בין תמונות שיש בהן כלבים לתמונות שאין בהן כלבים. כדי לאמן את המערכת היא מקבלת מקבצים של נתונים גולמיים – בדוגמה שלנו, את מיקום הפיקסלים בתמונה ואת הגוון שלהם. בדוגמה זו מדובר בלמידה מפוקחת, כך שהנתונים גם מתויגים – כל מקבץ נתונים (תמונה יחידה) מוגדרת (או מתויגת) "כלב" או "לא כלב". אם מספר מקבצי הנתונים יהיה גדול ומגוון די הצורך, המערכת תלמד לזהות טווח מסוים של דפוסים שמאפשרים לה לקשור בין הנתונים הגולמיים לתוצר המופשט הסופי. בדוגמה שלנו היא תלמד לזהות למשל שבתמונות שיש בהן כלבים ייתכנו יחסים מסוימים בין קווים לזוויות או גוונים מסוימים. לאחר שהיא תבחן די תמונות היא תוכל לקבוע על פי מאפיינים אלו שמדובר בתמונה של כלב.

המוח האנושי קולט כמות עצומה של נתונים גולמיים בכל רגע נתון ועל יסוד ניסיון נרכש הוא יודע לעבד אותם ולתת להם משמעות. ניסיון העבר מאפשר לנו לזהות צליל של צופר מכונית, ולהבחין בינו לבין בכי של תינוק, או לזהות אדם מסוים ולא אחר. למידת מכונה פועלת באופן דומה. המכונה, בדומה למערכת העצבים האנושית, יודעת להפיק פשר ממידע גולמי ולא מתויג. מובן שחרף ההשוואה המוח האנושי ומערכות למידה עמוקה אינם זהים ועל אף השימוש במונח נירונים אין הם זהים בפעולתם לנורונים הביולוגיים. עם זאת, בדומה ל"שכבות" הנורונים במוח, המחברות על ידי סינפסות, גם בלמידה עמוקה יש שכבות שונות של עיבוד מידע. שכבות עיבוד אלו מאפשרות רמות שונות של הפשטה בעיבוד המידע – רמת הפיקסל; זיהוי צורות גאומטריות; זיהוי חפצים קונקרטיים (חלון, כביש או לוחית זיהוי); ולבסוף – הקביעה אם מדובר באדם, בעל חיים או רכב.⁶⁸

למידת מכונה יכולה להיות למידה מפוקחת (supervised learning) או למידה לא מפוקחת (unsupervised learning). למידה מפוקחת פירושה שלא לגוריתם מוגדר משתנה

1.4 למידה מפוקחת, למידה לא מפוקחת ולמידה על ידי חיזוקים

מטרה מסוים לחיזוי. למשל, אם המערכת נועדה לחזות רמת סיכון פיננסי של לווה (סיכון גבוה/נמוך),⁶⁹ משתנה המטרה כבר ידוע מראש והמכונה לומדת לחזות אותו בהתבסס על נתוני אימון המוזנים לתוכה. נתונים אלו יכולים לכלול, לדוגמה, את ההיסטוריה הפיננסית של לווים מסוכנים ושל לווים לא מסוכנים ונתונים נוספים כמו מצב משפחתי, צבע עיניים או דת.⁷⁰ אם יוזנו למערכת די נתונים היא תדע לבחור מודל מתמטי שחווה את רמת הסיכון בדיוק יחסי. משתנה המטרה יכול להיות קטגוריאלי (מסוכן/לא מסוכן או רמת סיכון גבוהה/בינונית/נמוכה) והוא יכול להיות רציף (הסיכוי לפשיטת רגל). תיגו מראש ייעשה בהתאם.

68 ש.ס.

69 ראו לדוגמה ס' 16(ב)(4) לחוק נתוני אשראי, התשע"ו-2016, ס"ח 2551 בעמ' 838 (להלן: חוק נתוני אשראי).

70 ייתכן שנחננים אלו יהיו אטורים לשימוש, ראו ס' 51 לחוק נתוני אשראי.

שתיים מן המגבלות של למידה מפוקחת הן שלא־ימון נחוצות כמויות עצומות של נתונים ושלמערכת אין יכולת להתמודד עם סיטואציות חדשות. היכולת להתמודד עם מצבים חדשים היא מיומנות שבני אדם ובעלי חיים רוכשים בשלב מוקדם מאוד בחייהם. למשל, אדם אינו צריך לנסוע במכונית שלו מעבר לצוק כדי להבין שהיא תיפול ותתרוסק; ואין צורך להסביר לאדם שכאשר אובייקט מסתיר אובייקט אחר האובייקט המוסתר ממשיך להתקיים או שכדור שנחבט באלה לכיוון נדנדה עלול לפגוע בילד שמתנדנד. ההתבוננות בעולם והפעולה בו מלמדים בני אדם דברים כמו ממדים, סיבתיות, כוח כבידה וכדומה. אלו הן אבני הבניין שסביבן נצבר ידע מורכב יותר. מערכות הבינה המלאכותית הנוכחיות חסרות את הידע המשותף הזה, ולכן הן רעבות לנתונים ודורשות דוגמאות מתוגות באשר לנתונים שאינם במאגר הידע הבסיסי שלהן.

בלמידה לא מפוקחת לא נעשה אימון כזה. המכונה מקבלת מאגר נתונים לא מתוג ומפיקה ממנו תובנות בכוחות עצמה. כיוון שזהו מאגר לא מתוג, אי־אפשר לאמן את המכונה לזהות משתנה מוגדר (המכונה לומדת לבד להפריד, למשל, בין צורות גרפיות שונות ומחלקת אותן לקטגוריות לפי ה"היגיון" הסטטיסטי שלה; או מבחינה בין רמות סיכון פיננסיות בדרך שונה מזו של מומחה אשראי אנושי).

למידה לא מפוקחת אינה מתאימה לחלוקת הנתונים לקטגוריות שהוגדרו מראש או לזיהוי מתאם בין פרמטרים מוגדרים מראש, אך היא מתאימה למטרות של קיבוץ הנתונים לפי קטגוריות שהאלגוריתם מזהה. למידה לא מפוקחת מתאימה יותר לקיבוץ גורמים יחד על בסיס מאפיינים דומים, לזיהוי חריגות או לזיהוי משתנים שקשורים זה לזה. למידה כזאת מאפשרת לאלגוריתם לזהות מאפיינים בנתונים שמפתחי האלגוריתם לא הגדירו אותם מראש.

האלגוריתם הממליץ למשתמש על השיר הבא שינגן בנגן המוזיקה המקוון שלו, למשל, מבוסס על למידה לא מפוקחת. במערכת מוזן מאגר רשימות השמעה של משתמשים והאלגוריתם מנסה לזהות בהן מאפיינים דומים. כשמשתמש בוחר להאזין לשיר מסוים, האלגוריתם יכול להמליץ לו להאזין לשיר נוסף שמשתמשים אחרים בעלי היסטוריית האזנה דומה בחרו להאזין לו. האלגוריתם יכול להתבסס גם על נתונים נוספים (למשל, תיוגים אקראיים של השירים לפי ז'אנר או שנת הקלטה). בכל מקרה, האלגוריתם מחליט באופן עצמאי מהם הפרמטרים הרלוונטיים לצורך מתן המלצה מתאימה.

גם אלגוריתמים שמזהים שימוש חריג בכרטיס אשראי מבוססים על למידה לא מפוקחת. הם אינם נדרשים לתייג מראש שימושים חריגים או להגדיר פעולות הנחשבות סבירות ופעולות המעוררות חשד אלא רק לזהות שינויים בדפוסי השימוש הרגיל של הלקוח.

טכניקת למידה אחרת היא למידה על ידי חיזוקים (reinforcement learning)⁷¹. בטכניקה זו נבחן הפלט של המכונה וניתן לה משוב. זוהי למידה באמצעות ניסוי וטעייה. האלגוריתם לומד בתהליך חזרור (iterative process) אילו פעולות או פלטים מקבלים משוב חיובי ואילו מקבלים משוב שלילי. באמצעות משוב זה, דמוי מנגנון של שכר ועונש, המכונה לומדת לשכלל ולדייק את עצמה. אם למשל היא לא זיהתה רכישה של תכשיטים יקרים בכרטיס אשראי כשימוש חורג של לקוח, אך בדיעבד התגלתה פעולה זו כשימוש חורג, תלמד המכונה באמצעות המשוב מהטעויות שלה. אפשר לשכלל למידה על ידי חיזוקים בטכניקות אחרות של למידת מכונה (מפוקחת או לא מפוקחת). כך לא זו בלבד שהמערכת לא תחזור על טעויות עבר, אלא היא תלמד לזהות דפוסים חדשים ותימנע משגיאות זיהוי שעדיין לא נתקלה בהן.

מערכת אלפא-גו, למשל, התחילה בלמידה מפוקחת על בסיס נתוני אימון של משחקי גו היסטוריים, אך אחר כך נעזרה גם בתהליך של למידה על ידי חיזוקים. אחרי שהמערכת צברה ניסיון מעיבוד המהלכים במשחקים האנושיים, היא שיחקה משחק רב של משחקים נגד עצמה וקיבלה חיזוק בכל פעם שבחרה אסטרטגיה מנצחת. המערכת זוכרת כמוכן את כל הטעויות שכבר עשתה ונמנעת מלחזור עליהן. אלפא-גו זירו, היורשת של אלפא גו, לא נדרשה ללמידה מפוקחת והתבססה אך ורק על למידה על ידי חיזוקים.

פוטנציאל ההשפעה של הבינה המלאכותית על כל תחום בחיינו ברור. ככל שהטכניקות האלו ישתכללו, ועימן גם סוג הנתונים שהן יכולות לעבד ואיכותם, או שתגדל יכולתן להגיע לתוצרים בלי לעבד כמויות עצומות של מידע, כך יגדל גם הפוטנציאל הטרנספורמטיבי שיביאו איתן.

1.5 בינה מלאכותית יוצרת (גנרטיבית)

מומחים בעלי שם לבינה מלאכותית, לא כל שכן אנשים מן השורה, שפשפו את עיניהם בתדהמה לנוכח מקצת הפיתוחים שיצאו לאוויר העולם בשנת 2022, בעת כתיבת מחקר זה. התעוררה

תחושה שהמכונות חצו סף בלתי נראה ונעשו יצירתיות, תכונה שבעבר הובטח שתישמר לבני אנוש בלבד. DALL-E2 של OpenAI ו-Imagen של גוגל, Midjourney, D-ID, Stable Diffusion – כולם מסוגלים לייצר תמונות וקטעי וידאו על בסיס פקודה טקסטואלית (Prompt), המכונה לעיתים "לחישה"⁷².

אליהם מצטרפים מודלי שפה גנרטיביים שמייצרים טקסטים ברמה גבוהה, עונים על שאלות בצורה סמי-משכנעת ואפילו מבינים מה מצחיק בבדיחה: GPT-4 מבית Open AI, Jurassic-X של AI21 הישראלית ו-PaLM של גוגל. הבנה עמוקה של שפה היא קפיצת מדרגה ביכולות הבינה המלאכותית. כבר בשנת 2017 ידעו מודלי השפה הראשונים, שהתבססו על למידה לא מפוקחת, לומר מהי המילה החסרה או הבאה כשקיבלו סדרה של מילים; ואילו המודלים הנזכרים כאן יודעים לבצע גם מטלות מורכבות יותר, למשל לכתוב טקסט, לנתח סנטימנט בטקסט, לענות על שאלות או אפילו לקבוע אם הטקסט הגיוני או לא. בחודשים הבאים ובשנה הבאה נראה שתפותח גם היכולת לממש מיומנויות מופשטות (למשל לקחת משפט בשפה טבעית ולהכניסה כנוסחה למחשבון; או לקחת שאלות בשפה טבעית ולייצג אותן כנוסחת אקסל או כשאלת SQL).⁷³

מודלים אלו יכולים לסרוק את רשת האינטרנט כדי ליצור תוכן ערוך, החל בחדשות היום, עבור בהמלצות על מוצרים ושירותים וכלה בהמלצות למסלולי טיול, ואפילו לספק טבלאות מסודרות עם קישורים ולוונטיים למקור. הם אף

⁷² סהר מור "בינה מלאכותית יצרה את התמונה הזאת בהסתמך על טקסט בלבד" *ynet* (17.4.2022); יובל מן "הכירו את מחולל התמונות החדש של גוגל" *ynet* (24.5.2022).

⁷³ Nuobei Shi, Qin Zeng, and Raymond Lee, *Language Chatbot: The Design and Implementation of English Language Transfer Learning Agent Apps*, 2020

IEEE 3rd INTERNATIONAL CONFERENCE ON AUTOMATION, ELECTRONICS AND ELECTRICAL ENGINEERING (AUTEEE) (2020); Aditya Ramesh et al., *Hierarchical Text-Conditional Image Generation with CLIP Latents* (13.4.2022), available at [arXiv](https://arxiv.org/abs/2204.06125)

יכולים להדריך כיצד לכתוב תוכן ויראלי לרשתות חברתיות בעזרת ניתוח קוד מן האלגוריתם של רשת חברתית.

שלושת החידושים של המודלים היוצרים לעומת המודלים שקדמו להם:

(1) הם כלליים ולא ספציפיים, כלומר הנתונים שעליהם התאמן המודל אינם ייחודיים למטרה מסוימת (למשל כדי לייצר חיזוי סטטיסטי על מתן הלוואות או לענות על שאלות בנוגע לספרות מדעית בתחום מסוים) אלא נתונים כלליים (כל רשת האינטרנט, כל הערכים בוויקיפדיה וכיו"ב). דבר זה מאפשר מרחב שימושים עצום ואף מספק בסיס לשימושים ספציפיים, שמחייבים הוספת נתונים ייחודיים.

(2) הם יוצרים ולא תיאוריים, כלומר הם יכולים ליצור תוצרים חדשים, החל בכתיבת טקסט או ביצירת דימויים חזותיים ותבניות קול וכלה בכתיבת קוד לביצוע מטלות, ולא רק לשקף תבניות בנתונים הקיימים, כמו למשל זיהוי פנים, קלסיפיקציה של טקסט וחיזוי סטטיסטי מבוסס נתונים.

(3) הם נגישים ולא טכניים – אפשר לתקשר איתם, לשאול אותם שאלות מתוחכמות ולבקש מהם לבצע משימות בשפה פשוטה ("שפה טבעית"), ולא באמצעות כתיבת קוד.

הואיל והשפה הכתובה והמדוברת היא בסיס לכלל הפעילות האנושית, ובכלל זה ליכולת לחשוב חשיבה מופשטת, לפתח רעיונות מורכבים ולמסור אותם מדור לדור וממקום גאוגרפי אחד למשנהו, לחפש, לתקן, לסכם, לסנן, לצטט, לתרגם ולסכם אותם – הואיל וכך, היכולות היוצרות של המודלים של השפה נתפסות טרנספורמטיביות במיוחד לכלל המין האנושי ולכלל המגזרים בו.

גם למעבר למשימות מבוססות שפה טבעית (המכונות גם "נוסחאות" או prompts) יש השלכות. ראשית, העובדה שכל פלט ממודל שפה יכול בתורו לשמש קלט למודל שפה ("טקסט פנימה, טקסט החוצה") מאפשרת להמשיך לאמן ולשכלל את המודלים. שנית, אפשר לבצע פעולות עשירות ומתחכמות באמצעות "שרשור הנחיות". כשמדובר במשימות שדורשות פעולות ביניים או חשיבה מרובת שלבים, הפלט של תת-משימה אחת משמש קלט למשימה הבאה. ולעיתים תת-משימה כוללת גם אחזור מידע ממאגרים חיצוניים (כמו חיפוש מידע במנוע חיפוש או משיכת מידע מכתובת אינטרנטית נתונה).

לכן בטווח הקצר יהיו השימושים של משתמשים אנושיים במודלים הגדולים איטרטיביים ומשותפים. המשתמש האנושי יציג למודל הנחיה ראשונית או שרשרת הנחיות כדי להפיק פלט נתון; יבחן את הפלט ויכוונן את ההנחיה כדי לדייק את הפלט ולשפר את איכותו; יפעיל את המודל פעמים רבות על אותה בקשה כדי לבחור את הגרסאות הרלוונטיות ביותר של הפלט; ולפני השימוש הסופי בתוצר יטייב את הפלט בעצמו. אבל בטווח הארוך יותר, ככל שהמודלים ישתפרו וככל שהמוצרים הבנויים על בסיסם ייעשו נוחים יותר לשימוש ויוטמעו עמוק יותר בשגרות עבודה, יבוצעו המשימות בכל שלב ושלב ללא פיקוח אנושי.

מפתחי מודלים של בינה מלאכותית מתקדמים במהירות ויכולות המודלים לשימוש כללי הולכות וגדלות. המודלים לומדים באמצעות אימון ומציגים יכולות חדשות, לעיתים אף כאלה שהמפתחים שלהם לא צפו. מחקרים מלמדים על "יכולות מתעוררות" שלא נצפו ולא תוכננו מלכתחילה, ככל הנראה בשל היכולת של מודלים גדולים מאוד להשלים משימות, שלא התקיימה כשהמודלים היו קטנים יותר. מדעני מחשב מנסים ליצור רשימה של יכולת חדשות אלו, שלא נדונו בספרות בעבר⁷⁴ ויש המכנים אותן "ניצוצות". נעשים ניסיונות לא רק לזהות יכולות נוספות אלא גם להבין מדוע וכיצד הן צומחות.

ואולם לצד הפוטנציאל הכביר במודלים של שפה שאפשר לממש על בסיסם בינה מלאכותית יוצרת מתגלים גם חששות. ראשית, הואיל והממשקים מבוססים על ניסיון לצפות את רצונו של המשתמש, הם סובלים מתופעה המכונה "הזיות" (הלוצינציות) ומשיבים תשובות שאינן נכונות עובדתית. על אף אזהרותיהן של החברות המעניקות שירותים אלו, הן באמצעות הצהרות הן באמצעות הממשקים עצמם, משתמשים רבים עדיין נופלים בפח ולכן טמונה במודלים אלו סכנה של הפצת מידע מסולף (disinformation) ומידע מוטעה (misinformation) בציבור.

74 ראו Jason Wei et al., *Emergent Abilities of Large Language Models*, TRANSACTIONS ON MACHINE LEARNING RESEARCH (28.3.2023); Sébastien Bubeck et al., *Sparks of Artificial General Intelligence: Early Experiments with GPT-4* (13.4.2023), available at [ARXIV](https://arxiv.org/abs/2303.12712); Deep Ganguli et al., *Predictability and Surprise in Large Generative Models* (15.2.2022)

שנית, מודלים של שפה שעל בסיסם מופעלים ממשקי שיחה עם משתמשים (למשל צ'אטבוטים כמו ChatGPT, בארד ובינג או עוזרים אישיים אוטומטיים כגון AutoGPT) עוקבים אחר הוראות או "הנחיות" של משתמשים ויוצרים תוכן בהסתמך על יכולת חיזוי שרכשו מבסיס נתוני האימון שלהם. השילוב בין היעדר צורך במיומנויות תכנות וכתובת קוד ובין יכולתם של המודלים למלא הוראות הופך אותם למוצרים שימושיים לכלל האוכלוסייה, אך גם עושה אותם פגיעים לשימוש לרעה. דוגמה לשימוש כזה היא "הזרקת הנחיות" המכוונות את מודל השפה להתעלם ממנגנוני ההגבלה והבטיחות שלו עצמו, מה שמכונה לעיתים "פריצה מן הכלא" (jailbreaking). הדבר אפשרי אם מבקשים מהצ'אטבוט ליצור "משחק תפקידים" שבו הוא משחק מודל נוסף שנדרש לעשות מה שהמשתמש מבקש ממנו, ללא מגבלות. כך אפשר להשתמש בממשק כדי לבצע פעולות מזיקות, כגון הפצת תאוריות קונספירציה, או פעולות פליליות, כגון גנבת פרטי אשראי ויצירת הוראות להכנת חומרים מסוכנים. הפלטפורמות העיקריות של הממשקים האלה, כגון OpenAI, מנסות לאסוף דרכי פריצה מן הכלא ולהזיז אותן לנתוני האימון של המערכות כדי שיוכלו להתנגד להן בעתיד, אבל נראה שמדובר בקרב קשה ומסובך.⁷⁵

שלישית, החיבור בין ממשקים המבוססים על המודלים ובין רשת האינטרנט מביא איתו חששות. מודלים גדולים של בינה מלאכותית מתאמנים על כמויות עצומות של נתונים מהאינטרנט. לפי שעה חברות הטכנולוגיה פשוט סומכות על אמיתותם וניקיונם של הנתונים. אבל חוקרים גילו את האפשרות "להרעיל" מאגרי נתונים שמודלים גדולים של שפה מתאמנים עליהם, למשל על ידי רכישת כתובות של אתרי אינטרנט והצפתם בתמונות ובטקסטים שבתורם נכנסים לתוך מאגרי הנתונים האלה. אם לא די בכך, ככל שמשוהו חוזר על עצמו בנתוני האימון של המודל פעמים רבות יותר, כך הוא משפיע יותר על התנהגות המודל ומשפיע על תוצריו.⁷⁶ יתר על כן, ממשקים שחשופים לרשת האינטרנט ושואבים ממנה מידע חשופים גם להתקפה שנקראת "הזרקה מהירה עקיפה", שבה גוף עוין משנה

Will Douglas Heaven, *The Inside Story of How ChatGPT Was Built from the People Who Made It*, MIT Tech. Rev. (3.3.2023)

Nicholas Carlini et al., *Poisoning Web-Scale Training Datasets is Practical* (20.2.2023), available at arXiv

אתר אינטרנט על ידי הוספת הנחיות מוסתרות שגלויות רק לעוזר האוטומטי ונועדו להשפיע על פעולותיו (למשל גורמות לו להעביר אליהם פרטי כרטיסי אשראי, לשלוח בשמו מסרים בדואר אלקטרוני ועוד).

האפשרות ליישם את היכולות של המערכות כדי לגרום נזק, למשל לשכלל מיומנויות הטעיה, לסלף מידע, לבצע עבירות סייבר או לתכנן כלי נשק וטרור, אף שמפתחיהן לא התכוונו לכך, מכונה the alignment problem (בעיית הקשר בין משימה לתוצאה).⁷⁷ היא נתפסת כאחת הבעיות הקשות לפתרון במרחב השימושים במערכות גנרטיביות⁷⁸ ונדון בה בהמשך דברינו.⁷⁹ במקביל, כדי להתגונן מפני היכולת להיעזר בממשק כגון ChatGPT לאיתור חולשות סייבר ולכתובת תוכנה זדונית שמנצלת אותן וגונבת מסמכים, מציעה חברת מיקרוסופט, באמצעות שירות CoPilot, כלים שמטרתם לזהות איומי סייבר ולהגן מפניהם.⁸⁰

יכולות היצירה (היכולות הגנרטיביות) של המכונה ניכרות גם ביכולתה להתבטא בשפות לא טבעיות, כגון שפות קוד שונות. מערכת CoPilot מבית OpenAI ו-GitHub יוצרת כ-40% מן הקוד בחברות המשתמשות בה, ואפשר שבעתיד היא תייתר את הצורך לדעת לכתוב קוד ויהיה די בלימוד "לחישא" ובהזנת הוראות בדבר המוצר המבוקש למערכת.

קרן ההון סיכון סקויה קפיטל הודיעה בסוף 2022 כי היא מאמינה שבינה מלאכותית תיצור "טריליוני דולרים של ערך כלכלי".⁸¹ לדעתה, בינה מלאכותית היא ההתפתחות הטכנולוגית החשובה ביותר בימינו. בתרשים שלהלן מתומצתת התחזית שלה לשנים הבאות.⁸²

Richard Ngo, Lawrence Chan, and Sören Mindermann, *The Alignment Problem from a Deep Learning Perspective* (22.2.2023), available at [ARXIV](#)

Yotam Wolf et al., *Fundamental Limitations of Alignment in Large Language Models* (9.4.2023), available at [ARXIV](#)

79 ראו להלן בפרקים 9, 10.

Microsoft Security Copilot, *Introducing Microsoft Security Copilot*, 80 [YouTube](#) (28.3.2023)

Generative AI: A Creative New World, [SEQUOIA](#) 81

תרשים 1

תחזית יכולות הבינה המלאכותית במרחבי יצירה שונים

	לפני 2020	2020	2022	2023	?2025	?2030
טקסט	איתור ספאם תרגום שאלות ותשובות בסיסיות	יכולות כתיבה בסיסיות "טייטה" ראשונה	יכולות כתיבה ארוכות יותר "טייטה" שנייה	כתיבה משויפת במרחבים ספציפיים (למשל מאמרים מדעיים)	גרסה סופית של טקסט ברמה גבוהה מזו של כותב מקצועי	גרסה סופית של טקסט ברמה גבוהה מזו של כותב מקצועי
קוד	שורה אחת	יכולת יצירה של קוד רבישורות	גרסאות קוד ארוכות דיוק טוב יותר	יותר שפות קוד יותר מרחבים ספציפיים	מעבר ישיר מתיאור טקסטואלי למוצר מבוסס קוד, ברמה סופית טובה מזו של מפתח מיומן	מעבר ישיר מתיאור טקסטואלי למוצר מבוסס קוד, ברמה סופית טובה מזו של מפתח מיומן
תמונות			אומנות צילום	טייטה (מוק-אפ) בתחומי העיצוב, האדריכלות וכד'	מוצר מוגמר (מוצר מעוצב, תוכנית אדריכלית וכד')	מוצר מוגמר ברמה טובה מזו של מעצבים, אומנים, צלמים ואדריכלים מקצועיים
וידאו ותלת-ממד			ניסיונות ראשוניים ליצירת מודלים תלת-ממדיים	"טייטה" ראשונה של תוכני וידאו ותלת-ממד	"טייטה" שנייה של תוכני וידאו ותלת-ממד	משחקי וידאו וסרטים ברמת גימור מקצועית

■ ניסיונות ראשוניים ■ כמעט כאן, בשלבים מתקדמים ■ בשלבים מתקדמים

1.6

הכוחות המניעים את מהפכת הבינה המלאכותית

כפי שראינו לעיל, היסודות התאורטיים של מהפכת הבינה המלאכותית הונחו כבר באמצע המאה הקודמת. תקופתנו אינה אפוא התקופה הראשונה שבה תולים תקוות בהתפתחות תחום זה ומעלים חששות מהשלכותיו – מחזון אופטימי בדבר סינגולריות⁸³ עד אזהרות מפני השתלטות עוינת של בינה מלאכותית עילאית על הפלנטה.⁸⁴ תקופות אלו של התלהבות מבינה מלאכותית – המכונים "קיצים" בעגה המקובלת בתחום – הסתיימו באכזבה, שכן הטכנולוגיה לא עמדה בציפיות שתלו בה. אפשר שהקיץ הנוכחי של הבינה המלאכותית, ועימו השיח על בינה מלאכותית והרצינות שבה נידונות טכנולוגיות אלו בקרב קובעי מדיניות,⁸⁵ יגיעו אל קיצם בחורף נוסף.⁸⁶ עם זאת, כיוון שהדיון הער במהפכת הבינה המלאכותית נובע משילוב של כמה גורמים – מהם טכנולוגיים ומהם כלכליים – אפשר לצפות שאולי הפעם תהיה אחיזתה חזקה וממושכת יותר.

1.6.1. כוח מחשוּב בינה מלאכותית, ובייחוד כזאת הנשענת על טכניקות של למידת מכונה בכלל ועל למידה עמוקה בפרט, נדרשת לכמויות עצומות של מידע כדי לאמן את עצמה, וגם ליכולת הטכנולוגית לעבר את המידע הזה, לאחסן אותו ולאחזר אותו.

עקב מגמות מזעור והתייעלות (הטרנזיסטורים המבוססים על סיליקון קטנו; מהירות העברת הנתונים ברשת גדלה; הייתה עלייה מטאורית ביכולת לאחסן מידע; ועלות האחסון ועיבוד הנתונים פחתה) אנשים פרטיים משתמשים היום במחשבים שלפני עשור או שניים בלבד נחשבו לבעלי יכולת של מחשב-על,

83 Kurzweil, לעיל ה"ש 62.

84 לעיל ה"ש 10-11.

85 ראו להלן בפרק 3.

86 Luciano Floridi, *AI and Its New Winter: From Myths to Realities*, 33 *PHILOSOPHY AND TECHNOLOGY* 1 (2020); Sebastian Schuchmann, *Probability of an Approaching AI Winter*, *TOWARDS DATA SCIENCE* (17.8.2019); Kathleen Walch, *Are We Heading for Another AI Winter Soon?* *FORBES* (20.10.2019)

וכל ארגון מסחרי או מוסד שלטון יכול להרשות לעצמו לאגום ולשמור כמויות עצומות של מידע.

האלגוריתמים אינם פועלים בחלל ריק ונדרש כוח מחשוב כדי להפעיל אותם. חוק מור,⁸⁷ שניסח בשנת 1965 גורדון מור, ממייסדי חברת אינטל, חזה כי צפיפות הטרנזיסטורים במעגלים משולבים במחיר מינימלי תכפיל את עצמה בכל תקופה.⁸⁸ אומנם קצב הגידול המעריכי בכוח המחשוב דעך בשנים האחרונות,⁸⁹ אבל בימינו הגיע כוח זה לבשלות טכנולוגית המאפשרת יישומי בינה מלאכותית מבוססי למידה עמוקה.

ג'נסן הואנג, מנכ"ל חברת השבבים האמריקאית Nvidia, ניסח את חוק הואנג, שקובע כי כוח המחשוב של מעבדים גרפיים צפוי לכל הפחות להכפיל את עצמו מדי שנתיים.⁹⁰ שלא כמו יחידות עיבוד מרכזיות (CPU), שעובדות באופן טורי (יש להן ליבה אחת שמעבדת סדרת מטלות בטור, כלומר זו אחר זו), מעבדים גרפיים (GPU) הם מרובי ליבות ולכן העיבוד נעשה במקביל (כמה ליבות מעבדות בו-זמנית מטלה מסוימת). אומנם מעבדים כאלו יועדו במקור ליישומים גרפיים, אבל נמצא שהם המעבדים המיטביים ליישומים של בינה מלאכותית המבוססת על רשתות נוירונים. כיום מעבדים גרפיים משמשים בפייסבוק, בנטפליקס, בגוגל ובחברות אחרות להרצת מנועים מתקדמים של בינה מלאכותית. לא צפויה האטה בקצב התפתחות המעבדים הגרפיים, ולכן הבינה המלאכותית המבוססת על רשתות עצבים מלאכותיות צפויה להמשיך להתפתח. התפתחות הבינה

87 Gordon E. Moore, *Cramming More Components onto Integrated Circuits*, 87

38 ELECTRONICS 114 (1965)

88 בהתחלה חזה מור שהיא תכפיל את עצמה אחת לשנה; בשנת 1975 הוא עדכן את התחזית וטען שהיא תכפיל את עצמה אחת לכל שנתיים.

89 יש חסמים בסיסיים להמשך הגידול בצפיפות הטרנזיסטורים וביכולת למזער אותם יותר ויותר. Suhas Kumar, *Fundamental Limits to Moore's Law*, available at arXiv (18.11.2015). אינטל עצמה הכריזה שבעתיד לא תוכל לשרוד את מהירות המעבדים, אלא רק את צריכת החשמל שלהם. ראו הראל עילם "החוק הפסיד: אינטל קוברת רשמית את חוק מור" כלכליסט (11.2.2016).

90 Tekla S. Perry, *Move Over, Moore's Law: Make Way for Huang's Law* IEEE 90

SPECTRUM (2.4.2018)

המלאכותית היוצרת בסוף 2022 ותחילת 2023, והעובדה שמעבדים של חברת Nvidia משמשים גם את המודלים הגדולים של השפה, כגון אלה של חברת OpenAI וגוגל, הובילו לעלייה עצומה בשווי החברה עד יותר מטריליון דולר.⁹¹

נתונים שהציגה מעבדת Open AI מראים שמאז שנת 2012 גדל כוח המחשוב ששימש לאמן את מערכות הבינה המלאכותית המתקדמות פי שניים מדי 3.5 חודשים.⁹² היכולות המרשימות של האלגוריתמים שפיתחה חברת DeepMind, למשל, התאפשרו בזכות השקעה כספית ניכרת. ואולם בשנת 2019 הפסידה החברה יותר מחצי מיליארד דולר; ובשנה העוקבת ויתרה חברת האחות שלה, גוגל (שתיהן בבעלות חברת אלפאבית), על סכום של 1.5 מיליארד דולר, שהיא הלוותה ל-DeepMind.⁹³ אפשר שהיסודות שהניחה החברה ישרתו את גוגל בהמשך, אך בגלל עלות הפיתוח ועלות כוח המחשוב הנדרש לאימון אין היא מתאימה לפי שעה לשימוש מסחרי רחב היקף.

כוח המחשוב השתפר מאוד והוא מאפשר לאמן מודלים על יסוד נתונים רבים מאוד, ובכלל זה כל המידע ברשת האינטרנט. כוח החישוב אפשר גם להגדיל את מספר הפרמטרים (הרכיבים במודלי הבינה המלאכותית האחראים על עיבוד המידע). המודלים הגדולים (למשל מודל השפה GPT-4) מכילים מיליארדים של פרמטרים, קפיצה עצומה לעומת המודלים שהיו מבוססים על מיליוני פרמטרים.

אבל כוח החישוב העצום הנדרש בשביל לאמן מודלים של בינה מלאכותית עלול להיות חסם בפני המשך פיתוח ויש הסבורים שמחשוב קוונטי יוכל לסייע לשבור את תקרת הזכוכית בהקשר זה. תכונות קוונטיות קיימות במרחבים מסוימים בטבע, למשל בהולכה חשמלית, בקרינה או בגבישים, ומחשוב קוונטי

91 ליטל סמט "אנבידיה עשתה היסטוריה: חברת השבבים הראשונה בשווי של טריליון דולר" כלכליסט (30.5.2023).

92 Dario Amodei and Danny Hernandez, *AI and Compute*, OPEN AI (16.5.2018). השוו לקצב הגידול (המעריכי אף הוא) של חוק מור - הכפלה מדי שנתיים - כדי לעמוד על הגידול המואץ בכוח המחשוב.

93 Gary Marcus, *DeepMind's Losses and the Future of Artificial Intelligence*, WIRED (14.8.2019); Amy Thomson, *Google Waives \$1.5 Billion DeepMind Loan as AI Costs Mount*, BLOOMBERG (17.12.2020)

פירושו פיתוח יכולות עיבוד המבוססות על מרחב אפשרויות ולא רק על בחירה בינארית בין 0 ל-1. בשנים האחרונות מופנים לתחום זה תקציבי עתק, הן של ממשלות ומוסדות מחקר אקדמיים הן של תאגידים ענקיים כמו גוגל ואייבי-אם. בתחום המחשוב מדובר בטכנולוגיה חדשה שאין לה עדיין ערך של ממש, אבל הפוטנציאל שלה ככל הנראה רב. במאי 2023 חשפה Nvidia מחשב-על ששמו Israel-1, שלדבריה הוא מחשב-העל החזק ביותר בישראל ואחד החזקים בעולם. מחשב-על הוא מחשב המורכב מאלפי מעבדים ומשמש לביצוע משימות מורכבות, בהן פיתוח יישומי בינה מלאכותית יוצרת והרצת סימולציות מדעיות. המטרה המוצהרת של מחשב-העל החדש היא לסייע בשיתופי פעולה עם התעשייה, במחקר ובפיתוח פנימי, וכן לתצוגת תכלית של בניית מחשב-על המבוססים של הפלטפורמה החדשה של Nvidia, Spectrum-X.⁹⁴

1.6.2. קישוריות במרכז המהפכה הגרולה של האינטרנט עומדת הקישוריות – בין שמדובר בקישורים בין פיסות מידע בהיפרטקסט ובין שמדובר בהשלכות התפוצה הנרחבת של הרשת על התקשורת בין בני אדם. ממד נוסף של קישוריות נוגע ליכולת של מכשירים לתקשר אלה עם אלה ועם סביבתם. טכנולוגיה זו, המכונה "האינטרנט של הדברים" (Internet of Things, או IoT), היא ארכיטקטורה המאפשרת חפצים "חכמים", המצוידים בחיישנים שאוספים, מעבדים ומעבירים נתונים כדי ליעל את פעולתם.⁹⁵ לארכיטקטורה זו קשת רחבה של יישומים במגוון תחומים – תכולת המקרר, רשת החשמל (המסגירה את מנהגי השינה של בני הבית), מכשירי הטלוויזיה החכמים (המעבירים בדיוק מרבי מידע על הרגלי צפייה) ואפילו המכוניות האוטונומיות.⁹⁶ החיישנים (מדי תאוצה, חיישנים גירוסקופיים או חיישני תנועה) יכולים למדוד פרמטרים מגוונים – אור וצבע, לחץ ברומטרי, טמפרטורה ולחות, שדות מגנטיים. על אלה אפשר להוסיף גם את מקלט הנווטן

94 רפאל קאהאן "אנבידיה חשפה מחשב על ישראלי – אחד מעשרת החזקים ביותר בעולם" [ynet](https://www.ynet.co.il) (29.5.2023).

95 ראו לדוגמה רועי צזנה, *השולטים בעתיד* 56-64 (2017).

96 על מכוניות אוטונומיות ראו להלן בטעיף 2.1. ראו גם גדי פרל "הבטיחו לנו רכבים אוטונומיים. איפה הם?" *הארץ* (22.2.2022).

(GPS), המצוי בכל מכשיר טלפון חכם, המאפשר לכל מכשיר קצה (וכן ליחידות עיבוד מרכזיות) לקבל מידע בנוגע למקומו במרחב, בכל מקום על פני כדור הארץ.

האינטרנט של הדברים מאפשר איסוף מידע רחב היקף מחיישנים המצויים בכל מקום. ככל שתפוצתם של חיישנים אלו תמשיך לגדול, כך יגדל גם נפח המידע הפוטנציאלי שמערכות בינה מלאכותית יוכלו לעבד אותו ולהתאמן עליו.

1.6.3. בינה מלאכותית איסוף מידע רב ואיכותי, ניתוחו ופיתוח מודלים – מבוטט ענן

ככל שעלות אחסון הנתונים בענן פוחתת, כך עוברות משימות כמו ניתוח מידע באמצעות אלגוריתמים ולמידת מכונה לענן. בשוק הולכות ומתרכבות פלטפורמות של מערכות לומדות המאפשרות לחברות להעביר אליהן מידע ולקבל תוצאות בלי לפתח בעצמן מודלים חישוביים.⁹⁷

כך למשל פיתחה חברת איי-בי-אם עבור מפתחים חיצוניים את מערכת הבינה המלאכותית ווטסון, לשימוש דרך הענן.⁹⁸ משנת 2015 מספקת גוגל למפתחים את Tensorflow, מגוון בינה מלאכותית שנסמך על רשתות נוירונים מבוססות מעבדי GPU,⁹⁹ ובגרסאות מתקדמות על מאיצי בינה מלאכותית – מעגלים משולבים ייחודיים שפיתחה גוגל עבור המערכת.¹⁰⁰ בשנת 2018 החלה גוגל להציע למפתחים להשתמש במאיצים אלו, המכונים (Tensor Processing Units), דרך הענן.¹⁰¹ חברות מסחריות משתמשות ב-Tensorflow ובמנועי

Jeffery Burt, *Learning to Make the Machine Part of AI Invisible and Easy*, THE NEXT PLATFORM (11.12.2020)

Michael del Castillo, *It's 'Entrepreneurs Gone Wild' at IBM's World of Watson Conference*, NEW YORK BUSINESS JOURNAL (5.5.2015)

Cade Metz, *Google Just Open Sourced TensorFlow, Its Artificial Intelligence Engine*, WIRED (9.11.2015)

Norm Jouppi, *Google Supercharges Machine Learning Tasks with TPU Custom Chip*, GOOGLE CLOUD PLATFORM BLOG (18.6.2016)

John Barrus and Zak Stone, *Cloud TPU Machine Learning Accelerators Now Available in Beta*, GOOGLE CLOUD PLATFORM BLOG (12.2.2018)

בינה מלאכותית אחרים בענף לצרכים מגוונים, כגון פענוח סריקות MRI או זיהוי התמונות האפקטיביות ביותר שמעלים משתמשים לפלטפורמות שיתוף כגון Airbnb.¹⁰²

בדרך זו עומדת לרשות כל מפתח אבן בניין טכנולוגית לפיתוח תוכנות ואפליקציות חדשות גם בלי שיידרש למומחיות בארכיטקטורה שביסודה, ובעלות נמוכה (בייחוד כשמביאים בחשבון את עלויות הפיתוח העצומות של בינה מלאכותית).

בשלהי 2020 המליצה ועדה שמינה פורום תל"ם (תשתיות לאומיות למחקר ופיתוח) על הקמת ענף בינה מלאכותית ציבורי בישראל.¹⁰³ בקיץ 2022 החלה בישראל הפעלת פרויקט נימבוס, שתכליתו לספק לזרועות הממשלה שירותי ענף באמצעות חברות גוגל ואמזון.

1.6.4 נתוני ענף
(BIG DATA)
נתוני ענף הוא מונח שנטבע ב-1997 כדי לתאר הדמיית נתונים ואת האתגרים שהיא מעמידה למערכות מחשב.¹⁰⁴ אין הגדרה יחידה לנתוני ענף. יש הרואים בהם מסד נתונים גדול שקשה לנתח או לדמות באמצעים טכנולוגיים מקובלים.¹⁰⁵ אחרים רואים בהם מאגרי מידע בעלי נפח עצום (high volume), המתעדכנים במהירות גבוהה (high velocity), או מאגרים בעלי

Shijing Yao, Qiang Zhu, and Phillippe Siclait, *Categorizing Listing Photos at Airbnb*, THE AIRBNB BLOG (3.5.2018); Jason A. Polzin, *Intelligent Scanning Using Deep Learning For MRI*, TENSORFLOW BLOG (1.3.2019) 102

103 ועדת בינה מלאכותית ומדע נתונים, דוח מסכם (דצמבר 2020), בעמ' 37-38, 83 (להלן: דוח ועדת תל"ם). ראו גם אורי ברקוביץ' "התוכנית הלאומית ל-AI יוצאת לדרך: הקמת ענף במיליארד שקל לחוקרים ולסטארט-אפים" גלובס (22.12.2020).

Michael Cox and David Ellsworth, *Application-Controlled Demand Paging for Out-of-Core Visualization*, PROCEEDINGS OF THE 8TH IEEE CONFERENCE ON VISUALIZATION (1997) 104

FRANK J. OHLHORST, BIG DATA ANALYTICS: TURNING BIG DATA INTO BIG MONEY (2012) 105

מגוון רחב של סוגי נתונים (high variety).¹⁰⁶ גישות אחרות מתעלמות מהנפח של נתוני עתק ומנסות לתת להם הגדרה מהותית – לפי פרמטרים של הקשר, קישוריות וסיבוכיות.¹⁰⁷ חרף ריבוי ההגדרות,¹⁰⁸ אנו נאפיין נתוני עתק לפי נפח, מהירות ומגוון.

כפי שכבר ראינו, בינה מלאכותית – ובפרט כזאת המבוססת על טכניקות של למידת מכונה – זקוקה לנתונים. נתוני עתק הם הדלק המאפשר את פיתוחה ואת שכלולה. כך למשל, עד שנת 2017 אספה חברת טסלה תיעוד מוסרט של יותר מ-1.2 מיליארד קילומטרים של נהיגה.¹⁰⁹ סרטי הווידאו האלה, נוסף על מגוון סוגי נתונים שמקורם בחיישנים אחרים המותקנים ברכב, משמשים אותה לאמן את המערכת לעבור נתיבים ולהימנע מתאונות. נתונים אלו ישמשו את המערכת ללימוד ויהיו הבסיס לרכב האוטונומי של טסלה. חברת גוגל, שאינה משווקת כלי רכב ללקוחות ואין ברשותה נתוני נהיגה שוטפים, הצליחה להשיג באותה עת נתונים רק על כ-2.4 מיליון קילומטרים של נהיגה. היתרון של טסלה ברור ולכן מעריכים שהיא תקדים את גוגל בפיתוח רכב אוטונומי.¹¹⁰

גם ספקי שירותים מקוונים אחרים יכולים להשתמש בנתונים שלהם כדי לאמן את מערכות הבינה המלאכותית שלהם. חשבו על מיליארדי המשתמשים בשירותים של חברות כגון מטא, אפל וגוגל, שעקב היקף הנתונים שהן מייצרות מקנה להן פוטנציאל עצום לאימון מכונות לומדות. לא בכדי קבע מנכ"ל חברת איי-בי-אם, ארבינד קרישנה, כי "הסיבה להתקדמות המואצת של בינה מלאכותית

Mark Beyer and Douglas Laney, *The Importance of "Big Data"* 106
A Definition, GARTNER RESEARCH (21.6.2012)

Vincent Charles and Tatiana Gherman, *Achieving Competitive Advantage Through Big Data: Strategic Implications*, 16 MIDDLE-EAST JOURNAL OF SCIENTIFIC RESEARCH 1069 (2013)

TOBIAS M. SCHOLZ, BIG DATA IN ORGANIZATIONS AND THE ROLE OF HUMAN RESOURCE MANAGEMENT 12-20 (2017)

Michael J. Coren, *Tesla Has 780 Million Miles of Driving Data, and Adds Another Million Every 10 Hours*, QUARTZ (28.5.2016)

Howard Yu, *Google Fumbles While Tesla Sprints Toward a Driverless Future*, FORBES (14.12.2016)

היא שאנחנו מייצרים מדי יום ביומו 2.5 קווינטיליון בייטים של נתונים. זה 2.5 ואחריו 18 אפסים [...] לא די בכסיסי נתונים ובטכניקות אנליטיקה מיושנים. בינה מלאכותית היא הכלי היחיד שיכול לקצור את הנתונים האלו ולרתום אותם להפקת תובנות"¹¹¹.

אפשר לומר בפרפרזה על המשפט הידוע, "איזהו עשיר? בעל הדאטה". מי שמרבה לאסוף נתונים ושולט בהם מחזיק בנכס יקר ערך, שיש שמכנים אותו "הנפט של העולם הדיגיטלי"¹¹². עם זאת, לצד התפתחות מאגרי המידע שנדרשו לצורך למידה מפוקחת, התחזקות יכולות הלמידה העמוקה הלא מפוקחת מאפשרת למודלים להתאמן על כלל המידע ברשת האינטרנט. ככל שהמודלים מתפתחים, כך קטן הצורך שלהם בכמויות מידע עצומות ומתחלף בלמידה מבוססת התנסות (experience-based learning) – התנסות בעולם האמיתי בזמן אינטראקציה עם משתמשים ועם המשימות שהם מתבקשים למלא. באוגוסט 2022 פתחה חברת מטא לסוף שבוע אחד לשימוש חופשי בארצות הברית את הצ'אטבוט החדש שלה, BlenderBot 3. המשתמשים ששוחחו עם המערכת היו אמורים לספק לה משוב לשיפור מודל הבינה המלאכותית שלה.

Stephanie Condon, *IBM Promises a 4,000 Qubit Quantum Computer By 2025: Here's What it Means*, ZDNet (10.5.2022) 111

Michael Palmer, *Data Is the New Oil*, ANA MARKETING MAESTROS (3.11.2006) 112
 SHOSHANA ZUBOFF, *THE AGE OF SURVEILLANCE* וזבוף (2019) 153 CAPITALISM; התרגום שלנו):

דמיינו שיש לכם פטיש. זו למידת מכונה. בעזרת הפטיש טיפסתם על הר גבוה והגעתם לפסגה. זו הדומיננטיות של למידת מכונה בניחות נתונים. על פסגת ההר מצאתם ערמה עצומה של מסמרים, זולה יותר מכל מה שאפשר להעלות על הדעת. זו טכנולוגיית החיישנים החכמים החדשה. כעת מונח לפניכם לוח חלק ועצום [...] ואז אתם לומדים שבכל פעם שאתם נועצים בלוח מסמר בעזרת פטיש למידת המכונה, אתם מפיקים ערך. זו המונטיזציה של נתונים. מה אתם עושים? אתם מתחילים לנעוץ מסמרים ולעולם לא תפסיקו, אלא אם כן מישהו יגרום לכם להפסיק. אבל אין כאן אף אחד שיגרום לכם לעצור. לכן "האינטרנט של כל דבר" הוא בלתי נמנע.

1.6.5 קונסולידציה

של טכנולוגיות

בינה מלאכותית במונח הצר שלה – למידת מכונה מפוקחת, שמתבססת על מאגר – נחלקה באופן מסורתי לכמה סוגים: ראייה חישובית (computational vision), עיבוד שפה טבעית (natural language processing) וזיהוי דיבור (speech recognition). חלוקה זו אפשרה למשל קיום של מערכת שיודעת לכתוב סיפור, אבל אינה יודעת לזהות תמונה של חתול. בשנתיים האחרונות התלכדו סוגי הבינה המלאכותית ונוצרו מכונות שעושות הכול. למשל, מחוללי תמונות שמבינים טקסטים ויודעים להפיק מהם תמונות או מכונות שיודעות להמיר בעיות מילוליות בתרגיל מתמטי שאפשר לפתור במחשבון. מגמה זו של פיתוח מודלים מאוחדים המשלבים יכולות בכמה תחומים (כגון וידאו, תמונות, שפה וקול), הן בארכיטקטורת החומרה הן בתוכנה, מאפשרת את ההתקדמות לעבר בינה מלאכותית כללית יותר: צ'אטבוטים שיכולים להתכתב עם המשתמשים, לדבר איתם בעל־פה ואפילו "לראות" אותם ולהגיב לתנועות הפנים שלהם; רובוטים בעלי יכולת לוגית, כושר ראייה, הבנת שפה, כושר דיבור ומוטוריקה.

זאת ועוד. רוב מערכות הבינה המלאכותית פועלות היום בעולם הווירטואלי. החל במנוע החיפוש של גוגל וכלה באלגוריתמיקה של הרשתות החברתיות, בשירותי הניווט, בעוזרים האישיים הדיגיטליים, במערכות המשמשות לחיזוי ולהמלצה ובמערכות הזיהוי. אבל יש לשער שבשנים הבאות יעברו מערכות הבינה המלאכותית מן העולם הווירטואלי לעולם המעשי. דוגמאות מוכרות הן רכב אוטונומי ומכשור רפואי. אבל יש דוגמה מסקרנת ומסעירה יותר: הבינה המלאכותית עתידה לחלחל מהעולם הווירטואלי לעולם האמיתי באמצעות רובוטים הומנואידים (humanoids), שיהיו להם יכולות לוגיות ויכולות ראייה, הבנת שפה ודיבור; שילמדו בצורה לא מפוקחת באמצעות חיזוקים ולא יודקו לכמויות גדולות של נתונים לצורך אימון; ושתוכנות כגון מנועי משחק יוכלו לבחון אותם באמצעות סימולציות שידמו את העולם האמיתי ואחרי כן ישתמשו בטכניקות כגון sim2real כדי לבחון בעולם האמיתי רובוט שנבדק בסימולציה.

1.6.6. התפתחות של בינה מלאכותית באמצעות מערכות קוד יוצרות (גנרטיביות) בשימוש הנלל

הכוחות שיוסיפו להניע את מהפכת הבינה המלאכותית הם הפונקציונליות של מערכות בינה מלאכותית שאינן תלויות ברמת המיומנות הטכנית של המשתמש וזמינותן של מערכות המסוגלות ליצור קוד ואלגוריתמים, כלומר חיזוק. היכולת של משתמשים שאינם מוצאים יישום שהולם את צורכיהם ליצור מוצר משלהם. תחילה יוצרו מוצרים ששירתו את מהנדסי התוכנה, כגון Tabnine ו-DeepCode, שבדומה לנעשה בשפה טבעית חזו מה תהיה שורת הקוד הבאה או ידעו לאתר באגים בקוד שנכתב. היום הולך וגדל מספר הפלטפורמות שאינן מחייבות ידע טכני בתכנות; כמעט כל אחד יכול ליצור, לבדוק וליישם פתרונות המבוססים על בינה מלאכותית באמצעות ממשקים פשוטים של גרירה ושחרור או ממשקי תכנות יישומים (application programming interface [API]).¹¹³ אפשר להביא לדוגמה את SwayAI, המשמשת לפיתוח יישומי בינה מלאכותית ארגוניים, ואת Akkio, שיכולה ליצור כלי חיזוי וקבלת החלטות. מערכת CoPilot מבית OpenAI ו-GitHub יוצרת כ-40% מן הקוד בחברות שמשתמשות בה ואפשר שבעתיד היא תיתר את הצורך לדעת לכתוב קוד ודי יהיה ללמוד כיצד לנסח "לחישה" מתאימה (prompt) ולומר למערכת מהו המוצר המבוקש. אפשרויות אלו יעניקו לעסקים ולארגונים את היכולת להתגבר על המחסור במדעני נתונים ומהנדסי תוכנה, כך שהפקת התוצרים של מוצרים מבוססי בינה מלאכותית תהיה פשוטה, זולה ובהישג ידם של רבים.

1.6.7. סיכום

ככל שעלות הרכישה, האחסון והעיבוד של נתונים תלך ותפחת, כך תלך ותשתכלל כלכלת האלגוריתמים ותציע עוד ועוד אבני בניין של שירותי בינה מלאכותית בענף. בעקבות המחשוב המתעצם והיכולת ליצור מוצרים מבוססי בינה מלאכותית באמצעות "לחישה" למערכת הבינה היוצרת, מתחזקות התחזיות שלפיהן בעשורים הבאים ישתלבו רובוטים ובינה מלאכותית כמעט בכל חלק בחיינו. משימושים שונים של מערכות בינה מלאכותית שנסקור בפרק הבא עולה כי לא מדובר עוד בעתידנות, שכן תחזיות אלו קרובות מאוד להתממש.

113 ממשק תכנות יישומים (API), הוא ערכה של פונקציות ופרוצדורות קוד מן המוכן המאפשרת יצירה פשוטה של יישומים שמתממשים ליישום קונקרטי.

פרק שני

שימושים שונים
של בינה מלאכותית

—

בינה מלאכותית אינה מוגבלת לתחומים צרים כגון משחקי שחמט, גו או שש-כש. היא אינה מוגבלת גם ליישומים תאורטיים במעבדות הפיתוח כגון הבחנה בסיסית בין כלבים לחתולים. לעיתים נדמה שלא חולף שבוע בלי שכותרות העיתונים מבשרות על יכולת כזאת או אחרת של מערכות בינה מלאכותית – החל בכתיבת ספרים ומאמרים¹¹⁴ ובהפקת דימויים גרפיים או יצירות קומיקס

114 יובל פלוטקין "קראתם פעם ספר שנכתב בעזרת בינה מלאכותית?" ynet
John K. Waters, *AI21 Labs Emerges from Stealth with the First* ;(7.12.2020)
"AI-based Writing Companion," Pure AI (28.10.2020)

מתיאורים מילוליים¹¹⁵ (לעיתים ברמה גבוהה מזו של אמנים בשר ודם);¹¹⁶ עבור בזיהוי רגשות של כלבים¹¹⁷ (או של בני אדם, חרף הספקות באשר לבשלות המדעית של טכנולוגיה זו),¹¹⁸ בניחוש עמדותיה של השופטת העליונה המנוחה רות ביידר גינזבורג בכל נושא שבעולם¹¹⁹ ובהמלצות על העברת שחקנים לקבוצות כדורגל;¹²⁰ וכלה בפיתוחים יישומיים יותר, כגון ניהול מרכזים לוגיסטיים,¹²¹ זיהוי צורת חלבונים,¹²² קבלת החלטות בתהליכי הפריה מלאכותית¹²³ או זיהוי גלקסיות והבחנה ביניהן.¹²⁴ תקצר היריעה מלסקור את כל השימושים האפשריים בבינה מלאכותית, ועל כן נתמקד בפרק זה בכמה תחומים מובילים המדגימים את הפוטנציאל היישומי הטמון במערכות נבונות.

115 אושרי אלקסלטי, "מפחיד ומרשים בו זמנית: מודל ה-AI של OpenAI יוצר עכשיו תמונות מהמילים שלכם" *Geektime* (6.1.2021); מור, לעיל ה"ש 72; Rich Johnston, *Abolition Of Man, First Comic Book Entirely Drawn By A.I. Algorithm*, BLEEDING COOL (14.6.2022)

116 Kevin Roose, *An A.I.-Generated Picture Won an Art Prize. Artists Aren't Happy*, THE NEW YORK TIMES (2.9.2022)

117 "זו העם שלך, או שאתה שמח לראות אותי? קולר מיוחד מזהה רגשות אצל כלבים", *הארץ* (14.1.2021).

118 "Immature Biometric Technologies Could Be Discriminating Against People" Says ICO in Warning to Organisations, ICO (26.10.2022); Ifeoma Ajuwa, *Automated Video Interviewing as the New Phrenology*, 36 BERKELEY TECHNOLOGY LAW JOURNAL 1173 (2021)

119 *Israeli Company AI21 Labs Creates AI Model of Ruth Bader Ginsburg*, THE JERUSALEM POST (18.6.2022)

120 "טכנולוגיה ישראלית מחאימה שחקנים לקבוצות כדורגל באמצעות בינה מלאכותית" *New-Tech Magazine* (16.8.2022).

121 נטע-לי בינשטוק "אמזון, מאחוריך: צה"ל מחליף את האפסנאים בבינה מלאכותית" *כלכליסט* (14.4.2021).

122 קייד מץ "צעד שיאיץ פיתוח תרופות: חוקרים פתרו בעיה בת 50 שנה עם בינה מלאכותית" *הארץ* (1.12.2020).

123 אושרי אלקסלטי "במקום ילדי מבחנה, ילדי AI: בזכות האלגוריתם של Embryonics הישראלית, הושגו כבר 6 הריונות מוצלחים" *Geektime* (21.1.2021).

124 זאב בליזובסקי "בינה מלאכותית זיהתה 80,000 גלקסיות ספירליות" *הידען* (27.8.2020).

2.1 בינה מלאכותית במערכות תחבורה

כלי רכב אוטונומיים הם מיישומי הבינה המלאכותית המוכרים ביותר לציבור. המונח הלטיני למכונית, אוטומוביל, מבטא את יכולת ההנעה העצמית של הרכב, להבדיל מכרכה המונעת מכוח סוסים. למכוניות ולאוטוסטרדות נקשרה בתרבות הילה של חירות ושל שחרור האדם באמצעים טכנולוגיים, אך הסכנה לתאונות לא נעלמה.¹²⁵ כלי הרכב האוטונומיים – הניחנים לכאורה ביכולת לנוע במרחב באופן עצמוני, ללא התערבות אנושית – אמורים בין השאר לצמצם עד מאוד את הסיכונים הכרוכים בנהיגה אנושית.¹²⁶

את ההבחנה המקובלת בין דרגות שונות של אוטונומיה של כלי רכב פיתח ארגון מהנדסי הרכב (SAE).¹²⁷ על פי הבחנה זו (ראו לוח 1), יש שש דרגות של אוטונומיה. בשלהי 2015 העריך אלון מאסק, מנכ"ל חברת טסלה, שבתוך שנתיים יפותחו כלי רכב ברמת אוטומציה מלאה.¹²⁸ חמש שנים לאחר מכן המשיך מאסק להשמיע תחזיות בדבר הגעתן של מכוניות אוטונומיות לחלוטין לכביש בעגלא ובזמן קריב.¹²⁹ עם זאת, רמת הבשלות של טכנולוגיות הנהיגה הגיעה עד 2020 לרמת אוטונומיה 3 (נהיגה עצמית מותנית). חרף בשלותן של מערכות אלו, כמה מיצרניות הרכב מהססות לשווק אותן.¹³⁰

125 HENRI LEFEBVRE, *EVERYDAY LIFE IN THE MODERN WORLD* 98–109 (1971)

126 לסקירה עדכנית על סיכונים אלו ראו מתנאל בן אבי, *מגמות בבטיחות בדרכים בישראל 2013–2019* (חטיבת מידע ומחקר, הרשות הלאומית לבטיחות בדרכים 2020).

127 *Automated Driving: Levels of Driving Automation are Defined in New SAE International Standard J3016*, SOCIETY OF AUTOMOTIVE ENGINEERS INTERNATIONAL (18.6.2018). יוער כי עד 2016 השתמש משרד התחבורה האמריקאי בטקסונומיה משלו, אך דומה לזו של SAE. לטקסונומיה זו ראו *NHTSA Lays Out Ground Rules for Autonomous Vehicles*, THE CAR CONNECTION (7.6.2013)

128 Fred Lambert, *Tesla CEO Elon Musk Drops His Prediction of Full Autonomous Driving from 3 Years to Just 2*, ELECTREK (21.12.2015)

129 "מאסק מבטיח: מכונית אוטונומית ב־25 אלף דולר תוך 3 שנים" כלכליסט (23.9.2020); "אלון מאסק: בעוד שנה טסלה תוציא מיליון רכבים אוטונומיים" גלובס (23.4.2019).

130 תומר הדר "דילמת רמה 3: מדוע חברות הפאר הגרמניות לא ממהרות לפתח רכב אוטונומי?" כלכליסט (13.11.2020).

לוח 1

רמות אוטונומיה של כלי רכב אוטונומיים

רמה	שם	הגדרה
רמה 0	אין אוטומציה	נהג מיומן מבצע את כל משימות הנהיגה, אך מסתייע במערכות התרעה.
רמה 1	סיוע לנהג	מערכות אוטומטיות לסיוע לנהג מתקנות פעולות האצה והאטה או היגוי במצב נהיגה נתון ¹³¹ בהסתמך על מידע על סביבת הנהיגה ובהנחה שנהג מיומן מבצע את כל שאר משימות הנהיגה.
רמה 2	נהיגה עצמית חלקית	מערכות אוטומטיות לסיוע לנהג מבצעות פעולות האצה, האטה והיגוי במצב נהיגה נתון בהסתמך על מידע על סביבת הנהיגה ובהנחה שנהג מיומן מבצע את כל שאר משימות הנהיגה.
רמה 3	נהיגה עצמית מותנית	מערכות אוטומטיות מבצעות את כל משימות הנהיגה במצב נהיגה נתון בהנחה שנהג מיומן יגיב אם יתעורר צורך להתערב בנהיגה (כלומר יוכל לחזור לשליטה מלאה ברכב בתוך פרק זמן סביר מקבלת ההתרעה).
רמה 4	נהיגה עצמית ברמה גבוהה	מערכות אוטומטיות מבצעות את כל משימות הנהיגה במצב נהיגה נתון גם אם נהג אנושי אינו יכול להגיב תגובה סבירה לבקשת המערכת שישלוט שליטה מלאה ברכב.
רמה 5	אוטומציה מלאה	מערכות אוטומטיות מבצעות את כל משימות הנהיגה בכל תנאי הסביבה שנהג אנושי יכול להתמודד איתם.

מה שמבחינן בין השיח בן זמננו על כלי רכב אוטונומיים ובין פנטזיות שהוצגו בסרטי מדע בדיוני¹³² הוא יכולת המימוש. המימוש מתאפשר בזכות ארכיטקטורה

131 מצב נהיגה נתון הוא תרחיש שיש לו משימות נהיגה אופייניות – נהיגה בכביש הפתוח, השתרכות בפקק וכו'.

132 החל במכוניות אוטונומיות שמופיעות כבדרך אגב בסרטי הז'אנר, כגון זיכרון גורלי (פול ורהובן במאי 1990) או דוח מיוחד (סטיבן שפילברג במאי 2002), ועד לכאלו המשמשות דמויות מרכזיות, כגון בסדרה אביר על גלגלים (ג'ון לרסון יוצר 1982-1986).

המשלבת כמה טכנולוגיות – הנעה מכנית קלסית, חיישנים, קישוריות ובינה מלאכותית.¹³³ מצלמות, טכנולוגיית לידאר (LiDAR, light detection and ranging – זיהוי אור וטווח), מכ"מים, מערכות אינפרה-אדום, נווטנים (מערכות GPS) וכלים אחרים – כל אלה מאפשרים למערכת האוטונומית לאסוף באופן שוטף מידע רב על סביבתה המיידית ולהעריך אילו אובייקטים נמצאים בקרבתה ואילו סכנות עלולות להתממש.

ואולם לא מדובר בהכרח ביחידה אוטונומית אחת ויחידה (ונוסף על כך, רכב אוטונומי אינו בהכרח "מכונית").¹³⁴ ארכיטקטורות שונות של קישוריות – רכב לרכב (V2V); רכב לתשתיות (V2I), דוגמת סנסורים קרקעיים בצדי הדרך; רכב לענן (V2C); ואף רכב להולכי רגל (V2P) – כל אלו מניחים את התשתית הרעיונות ל"אינטרנט של הרכבים" (IoV), שבו תאפשר הקישוריות זרימת מידע בין כל החיישנים וכלי הרכב במרחב. בדרך זו יהיה אפשר למזער סיכונים ולטייב את איכות התחבורה אף יותר.¹³⁵

בינה מלאכותית היא רכיב מרכזי של כל טכנולוגיה של נהיגה אוטונומית. בין שמדובר במימוש של אוטונומיה באמצעות מערכות מומחה שמסתמכות על כללי נהיגה קבועים מראש ובין שמדובר בטכניקות של למידה עמוקה¹³⁶ – היכולת המעשית של רכב להתמצא במרחב בהתבסס על קלט שמתקבל מחיישנים מותנית בפיתוח של מערכות נבונות.¹³⁷

במצב שבו חלק ניכר מהתנועה תעבור לשליטה של מערכות נבונות, כלומר כששיעור הרכבים האוטונומיים מכלל הרכבים בשוק יהיה ניכר, יהיה אפשר

133 ראו למשל Hannah Yee Fen Lim, *Autonomous Vehicles and the Law* 7–15 (2018)

134 ראו יובל מן "זהירות, ספינה אוטונומית חוצה את האטלנטי" *ynet* (17.6.2021).

135 Zaigham Mahmood, *Connected Vehicles in the IoV: Concepts, Technologies and Architectures*, in *Connected Vehicles in the Internet of Things: Concepts, Technologies and Frameworks for the IoV* 3–18 (Zaigham Mahmood ed., 2020)

136 ראו Coren, *לעיל* ה"ש 109.

137 ראו Alex Kendall et al., *Learning to Drive in a Day*, in *International Conference on Robotics and Automation (ICRA)* 8248 (2019)

לצפות לעלייה ברמת הבטיחות בדרכים, לעלייה בקיבולת הכבישים¹³⁸ ולהפחתה בפליטת מזהמים.¹³⁹ מקצת התרחישים צופים מעבר לשימוש ברכבים אוטונומיים במערכות היסעים ציבוריות ויותר על רכבים פרטיים.¹⁴⁰

עם זאת, רכבים אוטונומיים אינם היישום היחיד של בינה מלאכותית במערכות תחבורה. מערכות נבונות יכולות לייעל את התנועה באמצעות הסדרת פעולת הרמזורים¹⁴¹ או ביזור מערכות הניווט. מערכות הניווט עצמן מכילות רכיבים של בינה מלאכותית ויכולות להציע לכל נהג נתיב זהה או לנסות לווסת את כלל התנועה של המשתמשים בהן.¹⁴² מערכות נבונות יכולות גם לנהל יישומים תחבורתיים של כלכלה שיתופית, כגון מיזמי שיתוף נסיעות (car pooling/ride sharing)¹⁴³ או שירותי היסעים שיתופיים (כגון השירותים שמציעות חברות Uber ו-Lyft או Bubble ו-GetTaxi).

RESHAPING URBAN MOBILITY WITH AUTONOMOUS VEHICLES: LESSONS FROM THE CITY OF BOSTON 138
(World Economic Forum in collaboration with the Boston Consulting Group, 2018)

Jeffery B. Greenblatt and Samveg Saxena, *Autonomous Taxis Could Greatly Reduce Greenhouse-Gas Emissions of US Light-Duty Vehicles*, 5 NATURE CLIMATE CHANGE, 860 (2015)

Alex Davies, *China's Self-Driving Bus Shows Autonomous Tech's Real Potential*, WIRED (7.10.2015)

Cathy Wu et al., *Framework for Control and Deep Reinforcement Learning in Traffic*, IEEE 20TH INTERNATIONAL CONFERENCE ON INTELLIGENT TRANSPORTATION SYSTEMS (ITSC) 1-8 (2017); Abu Salman Shaikat et al., *An Image Processing and Artificial Intelligence Based Traffic Signal Control System of Dhaka*, ASIA PACIFIC CONFERENCE ON RESEARCH IN INDUSTRIAL AND SYSTEMS ENGINEERING (APCORISE) 1-6 (2019)

Shoshana Vasserman, Michal Feldman, and Avinatan ראו לדוגמה Hassidim, *Implementing the Wisdom of Waze*, PROCEEDINGS OF THE TWENTY-FOURTH INTERNATIONAL JOINT CONFERENCE ON ARTIFICIAL INTELLIGENCE (IJCAI 15) (2015)

אודי עצינן "מוביט רווייזו ישתפו פעולה כדי לעקוף את הפקקים" כלכליסט 143 Xusen Cheng et al., *The Good, the Bad, and the Ugly: Impact of Analytics and Artificial Intelligence-Enabled Personal Information Collection on Privacy and Participation in Ridesharing* 31 EUROPEAN JOURNAL OF INFORMATION SYSTEMS 339 (2021)

2.2 זיהוי פנים

בפרק הקודם תיארונו מערכת שמשמשת בטכניקות של למידה עמוקה כדי להבחין, למשל, בין תמונות שיש בהן כלבים לתמונות שאין בהן כלבים. באותו אופן פותחו גם אלגוריתמים לזיהוי פנים. זיהוי פנים הוא התליך אוטומטי שמשווה בין שתי תמונות שונות של פנים כדי לבחון אם מדובר באותו אדם.

לאלגוריתמים של זיהוי פנים, המבוססים על בינה מלאכותית, יש שלל יישומים.¹⁴⁴ מערכות לזיהוי פנים מאפשרות למיין ולמפתח את שלל התמונות המצולמות דרך קבע במכשירי הטלפון החכמים של כולנו, לתייג חברים ברשתות החברתיות, לאמת זהות של אדם באמצעות פרמטרים ביומטריים ואף להצפין או לנעול מכשירי קצה ולשלוט בגישה לחשבונות מקוונים.¹⁴⁵ מובן שיש הברל בין

United States Government Accountability Office, FACIAL RECOGNITION 144 TECHNOLOGY: PRIVACY AND ACCURACY ISSUES RELATED TO COMMERCIAL USES, 11-13 GAO-20-522 (13.7.2020); Sharon Nakar and Dov Greenbaum, *Now You See Me. Now You Still Do: Facial Recognition Technology And The Growing Lack Of Privacy* 23 J. SCI. & TECH. 889, 97-100 (2017); IAN BERLE, FACE RECOGNITION TECHNOLOGY: COMPULSORY VISIBILITY AND ITS IMPACT ON PRIVACY AND THE CONFIDENTIALITY OF PERSONAL IDENTIFIABLE IMAGES 17-23 (2020); איחן לשם "משמר הלילה" של סין: מערכת לזיהוי פנים תמנע מקטינים לשחק כשאסור" הארץ (11.7.2021); לדוגמאות של שימוש בטכנולוגיית זיהוי פנים בישראל ראו מיכל רז-חיימוביץ "למנוע שימוש של קטינים: אפליקציות הקורקינטים של ווינד תפעיל טכנולוגיה לזיהוי פנים" גלובס (14.10.2020); ניצן שפיר "שימוש בטכנולוגיית זיהוי פנים לבעלי חו ירוק בהבימה ובאצטדיון בלומפילד פוגע בפרטיות ואינו חוקי" גלובס (20.5.2021). יוער כי כותביו של דוח של מרכז המחקר והמידע של הכנסת מציינים כי למיטב ידיעתם אין גורם בגופי המטה הממשלתיים שהם פנו אליו שציינו כי נציגיו משתמשים היום בטכנולוגיות זיהוי פנים. רועי גולדשמידט השימוש בטכנולוגיות זיהוי וניטור במרחב הציבורי 8 (מרכז המחקר והמידע של הכנסת 2020). עם זאת, הן הצבא הן המשטרה סירבו לבקשה של האגודה לזכויות אדם לפי חוק חופש המידע לקבל מידע על השימוש שלהם בטכנולוגיות זיהוי פנים. ראו "שימוש בטכנולוגיות לזיהוי פנים על ידי המשטרה והצבא", האגודה לזכויות האזרח בישראל (2.5.2021).

145 כבר בשנת 2008 הציגה חברת לנובו מחשבים ניידים שאפשרו להחליף את השימוש בסיסמה בזיהוי פנים של משתמשים מורשים. FACIAL RECOGNITION TECHNOLOGY AND THE CULTURE OF SURVEILLANCE 125 (2011) KELLY A. GATES, OUR BIOMETRIC FUTURE.

שימוש בטכנולוגיות זיהוי פנים למטרות פרטיות, כמו אשרור גישה למכשיר קצה פרטי, ובין פרישה של מערך מצלמות לזיהוי פנים במרחב הציבורי.¹⁴⁶

השימוש בטכנולוגיות לזיהוי פנים לביקורת גבולות¹⁴⁷ ולאמצעי תשלום מקוונים הולך וגובר. גם סוכנויות ביטחון ואכיפת חוק מביעות עניין גובר בטכנולוגיות לזיהוי פנים,¹⁴⁸ שכן הן מאפשרות להן לזהות במרחב מרושת במצלמות איומים ועבריינים,¹⁴⁹ ואף מרתיעות בני אדם (כמעין פנאופטיקון ממשטר) מביצוע פשעים.¹⁵⁰ לכן יש לתת את הדעת על האופן שבו שימוש בטכנולוגיות אלו עלול לפגוע בזכויות אדם,¹⁵¹ ואין פלא שבתקנות הבינה המלאכותית האירופיות ניתנה תשומת לב מיוחדת לטכנולוגיות אלו (בהקשר הרחב של מערכות ביומטריות).¹⁵²

Evan Selinger and Brenda Leong, *The Ethics of Facial Recognition Technology*, in *THE OXFORD HANDBOOK OF DIGITAL ETHICS* 101-122 (Carissa Véliz ed., 2021)

Berle, לעיל ה"ש 144, בעמ' 17-19.

148 שם, בעמ' 19-20; Clare Garvie, Alvaro M. Bedoya, and Jonathan Frankle, *The Perpetual Line-Up - Unregulated Police Face-Recognition in America*, CENTER ON PRIVACY AND TECHNOLOGY AT GEORGETOWN LAW (2016); Katelyn Ringrose, *Law Enforcement's Pairing of Facial Recognition Technology with Body-Worn Cameras Escalates Privacy Concerns* 105 VA. L. REV. ONLINE 57 (2019)

149 Kyriakos N. Kotsoglou and Marion Oswald, *The Long Arm of the Algorithm? Automated Facial Recognition as Evidence and Trigger for Police Intervention* 2 FORENSIC SCIENCE INTERNATIONAL: SYNERGY 86-89 (2020)

150 ראו לדוגמה את הנחיית רשם מאגרי המידע מס' 4/2012, "שימוש במצלמות אבטחה ומעקב ובמאגרי התמונות הנקלטות בהן", פס' 2.2-2.5.

151 ראו לדוגמה את הסקירה הזאת: IN FOCUS: FACIAL RECOGNITION TECH STORIES AND RIGHTS HARMS FROM AROUND THE WORLD (International Network of Civil Liberties Organizations (INCLCLO) 2021). כמו כן ראו את הדיווחים מן הזמן האחרון על כוונת המשטר האיראני לאכוף את החוק המחייב נשים לעטות חג'אב באמצעות טכנולוגיה לזיהוי פנים: "בעקבות 'התרופפות' בצניעות של נשים, איראן פונה לטכנולוגיות זיהוי פנים" (7.9.2022).

152 ס' 3(33) להצעת תקנות הבינה המלאכותית האירופיות (ראו לעיל ה"ש 53) מגדיר מידע ביומטרי מידע שמקורו בהליכי עיבוד טכניים הנוגעים בין השאר לזהותו היחודית של אדם, לדוגמה מראהו (בהצעה המחוקקת, ס' 3(33) מפנה ישירות לס' 4(14) ל-GDPR, שממנו נלקחה הגדרה זו). ס' 5(1)(d) להצעה המחוקקת אוסר על שימוש במערכות זיהוי

אלגוריתמים המשמשים לזיהוי מבוססים על מערך כללים או על תהליך לחישוב או ניתוח מאפיינים של פנים אנושיות. איכותם של אלגוריתמים כאלו תלויה באופיו של בסיס הנתונים שהם מתאמנים עליו; למשל, מאגר התמונות ששימש לפיתוחם יקבע מה יצליחו לזהות טוב יותר – פנים של בני אדם ממוצא אסיאתי או פנים מערביות.¹⁵³ לאחרונה ציין הממונה על היישומים הביומטריים כי יש קשיים בהרכשה (acquisition) של תמונות פנים של תושבים בעלי גון עור כהה; בהתאם לכך, היעדר כיוול ראוי של מערכות ביומטריות עלול להשפיע על הדיוק ועל מרווח הטעות בעת שימוש בהן.¹⁵⁴ לנוכח מגוון היישומים הפוטנציאליים של טכנולוגיה זו יש להטות אלו, גם אם נעשו בתום לב, השלכות הרות גורל והן עלולות להנציח פערים חברתיים ואפליה של מיעוטים.¹⁵⁵ למשל, אלגוריתם זיהוי פנים שהתאמן על מאגר תמונות של אנשים ממוצא מערבי ומשמש לזיהוי חשודים בפשע עלול להגיע לתוצאות מדויקות פחות בזיהוי חשודים מרקע אתני אחר.

ואכן, יישומים של מערכות זיהוי פנים למטרת אכיפת חוק או ביטחון לאומי שנויים במחלוקת – לא רק בשל הפגיעה בפרטיות עקב פרישת מצלמות מעקב

ביומטריות בזמן אמת במרחבים ציבוריים. ס' 52(2) מטיל על מערכות אלו חובות שקיפות ייחודיות. כמו כן, מערכות זיהוי ביומטריות מוגדרות כמערכות בסיכון גבוה (ס' 1 בחוספת III לטייטה). ההצעה המתוקנת מגדירה כמערכות בסיכון גבוה גם מערכות מבוססות ביומטריקה (המסיקות בהתבסס על נתונים ביומטריים מאפיינים אישיים של אדם, ובכלל זה מערכות לזיהוי רגשות).

153 כך למשל, מחקר מצא שרמת הדיוק של אלגוריתמים לזיהוי פנים שפותחו במדינות מזרח אסיה לזיהוי תווי פנים אסיאתיים גבוהה יותר מרמת הדיוק של אלגוריתמים שפותחו במערב, ולהפך. P. Jonathon Phillips et al., *An Other-Race Effect For Face Recognition Algorithms*, ACM TRANSACTIONS ON APPLIED PERCEPTION (2009)

154 דוח הממונה על היישומים הביומטריים מס' 11 (21.6.2021), בעמ' 6. חששות דומים הועלו בעקבות אישור חוק המצלמות המיוחדות בוועדת השרים לענייני חקיקה (טייטת חוק לחיקון פקודת המשטרה (מערכות צילום מיוחדות), התשפ"ב-2022. לנוסח מיום 10.4.2022 שאושר בוועדת השרים בחודש מאי 2022). ראו לדוגמה רוני (פנטנש) מלכאי "מצלמות לזיהוי פנים יסכנו את יוצאי אתיופיה" הארץ (11.5.2022).

Fabio Bacchini and Ludovica Lorusso, *Race, Again: How Face Recognition Technology Reinforces Racial Discrimination* 17 JOURNAL OF INFORMATION, COMMUNICATION & ETHICS IN SOCIETY 321 (2019), Selinger and Leong

לעיל ה"ש 146.

במרחב הציבורי, אלא גם בשל העיבוד האוטומטי של מידע ביומטרי רגיש,¹⁵⁶ שמעורר קשיים בדיוק ובזיהוי,¹⁵⁷ והגברת עוצמתן של תופעות המקשות על זיהוי אנושי רגיל של חשודים במסגרת מסדרי זיהוי.¹⁵⁸

בערים רבות בעולם המערבי הוגבל במידה ניכרת השימוש במערכות כאלה,¹⁵⁹ והפרלמנט האירופי קרא לחרם על מערכות זיהוי פנים למטרות שיטור.¹⁶⁰ חברות הטכנולוגיה הגדולות, ובראשן אמזון, אייבי-אם ומייקרוסופט, הודיעו שיפסיקו לייצר מערכות זיהוי פנים לתכליות שיטור ולהשקיע בחברות המפתחות אותן, או שיימנעו מלספקן לסוכנויות אכיפת החוק.¹⁶¹ בשלהי 2021 הודיעה מטא (פייסבוק)

156 R (Bridges) v. Chief Constable of South Wales Police, ראו לדוגמה, [2020] EWCA (Civ) 1058, at para. 88–89

157 לפי דיווח ברשת סקיי הבריטית, שיעור הטעויות בזיהוי של המערכת לזיהוי פנים שהייתה בשימוש משטרת המטרופולין של לונדון הגיע ל-81%. Rowland Manthorpe and Alexander J. Martin, 81% of "Suspects" Flagged By Met's Police Facial Recognition Technology Innocent, *Independent Report Says*, SKY NEWS (4.7.2019); מחקר אחר הצביע על שיעור טעויות נמוך יותר, אך עדיין גבוה – 38% ו-63% בשתי בדיקות שונות: Pete Fussey and Daragh Murray, INDEPENDENT REPORT ON THE LONDON METROPOLITAN POLICE SERVICE'S TRIAL OF LIVE FACIAL RECOGNITION TECHNOLOGY 69–72 (Human Rights Centre, University of Essex 2019)

158 Laura Moy, *Facing Injustice: How Face Recognition Technology May Increase the Incidence of Misidentifications and Wrongful Convictions*, 30 WM. & MARY BILL RTS. J. 337, 350 (2021)

159 ראו הסקירה של תהילה שוורץ אלטשולר ועמיר כהנא חוות דעת בעניין תזכיר חוק לחיקוק פקודת המשטרה (נוסח חדש) (מערכות צילום מיוחדות), החשפ"א-2021-12-15 (חוות דעת, המכון הישראלי לדמוקרטיה (29.7.2021).

160 European Parliament Resolution of 6 October 2021 on Artificial Intelligence in Criminal Law and its Use by the Police and Judicial Authorities in Criminal Matters, para. 27 (2020/2016(INI))

161 Alex Hern, *IBM Quits Facial-Recognition Market over Police Racial-Profiling Concerns*, THE GUARDIAN (9.6.2020); Jay Greene, *Microsoft Won't Sell Police its Facial-Recognition Technology, Following Similar Moves by Amazon and IBM*, THE WASHINGTON POST (11.6.2020); Jeffrey Dastin, *Amazon Extends Moratorium on Police Use of Facial Recognition Software*, REUTERS (18.5.2021)

על הפסקת פעולתן של מערכות זיהוי הפנים שלה.¹⁶² מנגד, בישראל קודם לאחרונה תזכיר חוק לאסדרת מערכת עין הנץ, מערכת משטרתית לניטור לוחיות רישוי, והונחה התשתית המשפטית להחילו בעתיד גם על מערכות לזיהוי פנים.¹⁶³

פרישה רחבת היקף של מערכות זיהוי פנים תרחיב את השימוש בקטגוריות חדשות של מידע. ייתכן שמערכות נבונות המנתחות נתונים חזותיים יכולות לייצר באופן אוטומטי מידע על רגשות של פרטים הנקלטים בעדשות המצלמות¹⁶⁴ או על מאפייני התנהגותם במרחב. מערכות כאלו יכולות לנתח את טיב האינטראקציה החברתית של יחידים וקבוצות ואת אופי הקשרים ביניהם. יכולות אלו לא היו קיימות בעבר ושילובן בפרקטיקות של פיקוח ומשטור במרחב הציבורי עלול להעמיק את אופי פרופיל האישיות שניתן לבנות באופן אוטומטי וחסר הבחנה לאזרחים החולפים בעל כורחם על פני עדשת המצלמה, וממילא להגדיל את מידת הפגיעה בפרטיות.

יש לתת את הדעת שמערכות זיהוי פנים ניחנות ביכולות שחורגות מזיהוי פנים גרידא. מערכות צילום נבונות מסוגלות להפיק מן הנתונים הוויזואליים מידע על מאפייני תנועה, על מגדר, על שיוך אתני¹⁶⁵ או על רגשות (חרף הספקות בעניין

Kashmir Hill and Ryan Mac, *Facebook, Citing Societal Concerns, Plans to Shut Down Facial Recognition System*, THE NEW YORK TIMES (2.11.2021); Jerome Pesenti, *An Update on Our Use of Face Recognition*, FACEBOOK (2.11.2021)

163 בעת כתיבת שורות אלו נדונה שאלת ההסמכה של המשטרה להשתמש במערכת "עין הנץ" בבג"ץ 641/21 האגודה לזכויות אזרח נ' משטרת ישראל. לעיון בעתירת האגודה לזכויות אזרח ראו "להפסיק להשתמש במערכת שעוקבת אחרי נהגים" האגודה לזכויות אזרח (2.11.2022). לתזכיר החוק ראו תזכיר חוק לתיקון פקודת המשטרה (נוסח חדש) (תיקון מס') (מערכות צילום מיוחדות), התשפ"א-2021 (8.7.2021). כן ראו שורר אלטשולר וכהנא, לעיל ה"ש 159.

164 לאחרונה, למשל, דווח כי חברת המצלמות קאנון החקינה במשרדה שבבייג'ינג מערכת זיהוי פנים המחירה כניסה רק לעובדים "שמחים". James Vincent, *Canon Put AI Cameras in its Chinese Offices that only Let Smiling Workers Inside*, THE VERGE (17.6.2021)

165 אושרית גן-אל "ניו יורק טיימס: רשויות בסין משתמשות בזיהוי פנים לניטור מיעוטים" גלובס (15.4.2019); עודד ירון "לא יכול להסביר את זה; בכיר בוואווי התפטר בעקבות פיתוח 'אזעקת אויגורים' " הארץ (16.12.2020).

בשלותה של טכנולוגיה זו).¹⁶⁶ למשל, אמן בלגי אימן אלגוריתם לזהות בזמן אמת מתי דעתם של פוליטיקאים בלגים מוסחת מדיונים פרלמנטריים והם עסוקים במכשיר הטלפון הנייד שלהם. אותה מערכת אף פרסמה את הסרטון ברשתות החברתיות ותייגה אוטומטית את הפוליטיקאים הסוררים.¹⁶⁷ מערכות זיהוי פנים אינן מוגבלות כמוכן לכיוש נבחרי ציבור וניתן לרתום יכולות דומות לניתוח מידע בהקשר רחב יותר – למשל לזיהוי מאפייני תנועה של התקהלויות חשודות או של אירועים אלימים.

2.3

בינה מלאכותית

במערכות אכיפת החוק והצדק

2.3.1 לעורך דין מְכָנִי
שלשלתי אסימון:
בינה מלאכותית
בליגל-טק

עריכת דין היא פרופסיה שמרנית¹⁶⁸ הנרתעת מאימוץ חידושים טכנולוגיים.¹⁶⁹ על כן חדירתן של טכנולוגיות בינה מלאכותית למגזר המשפטי – ובייחוד הפרטי – נעשית עקב בצד אגודל. ואולם פוטנציאל ההשפעה של הטכנולוגיות האלה רב. לצד השמרנות המקצועית, חסם כניסה נוסף שמקשה על חדירת הבינה המלאכותית לעולם המשפט הוא מרכזיות השפה בו.

166 לעיל ה"ש 118.

167 *The Flemish Scrollers, 2021–2022*, DRIES DEPOORTER (2021)

168 להתייחסות כללית לפרופסיה המשפטית כאל פרופסיה שמרנית ראו יניב רוזנאי וקארין פאר פרידמן "עריכת דין מהפכנית" המשפט כ 303, 310–316 (2015); גל סידס, שמוליק בכר, אופיר נוה "החינוך המשפטי: על מה שבפנים ועל מה שבחוץ" עיוני משפט כה 243 (2001); Edgar Bodenheimer, *The Inherent Conservatism of the Legal Profession*, 23 *IND. L. J.* 221 (1948)

169 Chay Brooks, Cristian Gherhes, and Tim Vorley, *Artificial Intelligence in the Legal Sector: Pressures and Challenges of Transformation*, 13 *CAMBRIDGE JOURNAL OF REGIONS, ECONOMY AND SOCIETY* 135, 146 (2020)

אף שהמילה קוד פירושה בין השאר אסופת חוקים, קוד משפטי אינו קוד מכונה.¹⁷⁰ השפה המשפטית המרכיבה את החוק – הקוד המשפטי – יכולה להיות עמומה ורב-משמעית.¹⁷¹ גם בהיעדר עמימות, עיבוד מסמכים משפטיים דורש פיתוח טכניקות לעיבוד שפה – היכולת להבחין בין מונח משפטי לשפה טבעית או היכולת לזהות תבניות משפטיות והקשרים.¹⁷²

מערכות בינה מלאכותית מסוגלות כיום למיין ולסווג מאות מסמכים לצורך הליכי גילוי (discovery) במהלך משפט, או הליכי גילוי נאות (due diligence) בעסקאות ענק בין חברות – משימות שבעבר היו מנת חלקם של צוותים גדולים של עורכי דין.¹⁷³ החברה הישראלית LawGeex, למשל, פיתחה מערכת שמסוגלת לזהות בהסכם סודיות היבטים שונים שיש לתת עליהם את הדעת מבחינה משפטית. בחודש פברואר 2018 ערכה החברה ניסוי שבחן את ביצועי המערכת בהשוואה לכיצועים של 20 עורכי דין אנושיים. המשתתפים התבקשו לזהות היבטים משפטיים בחמישה הסמכי סודיות. האלגוריתם הגיע לממוצע של 94% הצלחה בתוך 26 שניות; עורכי הדין, לעומת זאת, הגיעו לממוצע של 85% הצלחה, אבל היו זקוקים לשעה וחצי בממוצע כדי לסקור את כל חמשת החוזים.¹⁷⁴

170 על היחס בין קוד לחוק ראו למשל Primavera De Filippi and Samer Hassan, *Blockchain Technology as a Regulatory Technology: From Code Is Law to Law Is Code* (2018), available at arXiv

171 ראו לדוגמה Diana Raffman, *Vagueness in Law: Placing the Blame where It's Due*, in *VAGUENESS IN THE LAW – PHILOSOPHICAL AND LEGAL PERSPECTIVES* 49 (Geert Keil and Ralf Poscher eds., 2016); H. L. A. HART, *THE CONCEPT OF LAW* 124–136 (2nd edition, 1994); BRIAN BIX, *LAW, LANGUAGE, AND LEGAL DETERMINACY* (1993)

172 ראו להלן בסעיף 2.5.

173 John Markoff, *Armies of Expensive Lawyers, Replaced by Cheaper Software*, *THE NEW YORK TIMES* (4.3.2011); Nicholas Barry, *Man Versus Machine Review: The Showdown between Hordes of Discovery Lawyers and a Computer-Utilizing Predictive-Coding Technology*, 15 *VAND. J. ENT. & TECH. L.* 343 (2013); Samuel Sahagian, *iLawyer*, 3 *UCLA J. L. & TECH. DIGEST* 8 (2021)

174 *Comparing the Performance of Artificial Intelligence to Human Lawyers in the Review of Standard Business Contracts*, LAWGEEX (2018)

לכאורה אפשר לומר שתוצאות אלו מבשרות את "קץ הפרופסיה"¹⁷⁵, אך אפשר לומר גם שיש בהן כדי להראות שלבינה מלאכותית יש פוטנציאל להיות טכנולוגיה משבשת ולשנות סדרי עולם בתחום המשפט. האוטומציה של הסכמים פשוטים (כמו הסכמי סודיות) מקילה על עורכי הדין את העבודה ומאפשרת להם להתמחות בסוגיות מורכבות יותר ובשירותים אישיים לפי צורכי הלקוח. אפשר שהשימוש במערכות ממוחשבות לעיבוד חוזים אף יוביל לניסוחם בשפה פשוטה יותר ולצמצום הנטייה המקצועית לנקוט לשון מעורפלת ומסובכת כדי לחפות על בעיות מהותיות בחוזה. לעת עתה, גם בשל הכללים החלים על שירותים משפטיים,¹⁷⁶ נדרש אישורו של עורך דין אנושי לתוצאה שהפיק האלגוריתם.

לצד עריכת חוזים, היבט אחר של העבודה המשפטית הוא פרשנות הדין, ובכלל זה הפסיקה. לעיתים מזומנות, חוות הדעת של עורכי הדין על חוקיות מצב דברים עובדתי מסוים אינן אלא תחזית בדבר האופן שבו יפסוק בית המשפט בהינתן אותן נסיבות. תחזית זו נסמכת על פרשנות החקיקה ועל פרשנות ההלכה המשפטית הרלוונטית, ולעיתים אף מסתייעת בהערכה של עמדותיהם האפשריות של שופטים מסוימים. את חוות הדעת המשפטית הלא-ממוכנת, המבוססת על ניסיונו המקצועי של עורך הדין המנסח אותה, ניתן להחליף או למצער להשלים בתחזית ממוכנת. אחד היישומים של מערכות נבונות הוא חיזוי, ומערכת בינה מלאכותית שתדע לעבד מסמכים משפטיים תוכל לנתח מאות או אלפי תיקים דומים, לפנות לחקיקה המתאימה ולנבא את התשובה המשפטית הצפויה בשאלה מסוימת.

מערכות כאלו אינן פרי הרמיון. חברת ההזנק הקנדית Blue J Legal, למשל, משתמשת בלמידת מכונה לחיזוי החלטות מיסוי.¹⁷⁷ מחקרים אחרים מציעים

175 ראו למשל RICHARD SUSSKIND, THE END OF LAWYERS? RETHINKING THE NATURE OF LEGAL SERVICES (2010).

176 ס' (3)20-(4) לחוק לשכת עורכי הדין, חשכ"א-1961, ס"ח 347 בעמ' 178.

177 Benjamin Alarie, Anthony Niblett, and Albert H. Yoon, *Using Machine Learning to Predict Outcomes in Tax Law*, 58 CANADIAN BUSINESS L. J. 231 (2016); idem, *How Artificial Intelligence Will Affect the Practice of Law*, 68 UNIVERSITY OF TORONTO L. J. 106 (2018). ראו למשל כיצד אלגוריתם מבוטס למידת מכונה מעריך כיצד יחליט בית המשפט למיטתם לטובת הוצאות משפטיות על קניין רוחני: Benjamin Alarie and Kimberly Condon, *Deducting Legal Expenses: Unpacking the IRS's Appeal in Mylan*, TAX NOTES (26.8.2022).

מודלים לחיזוי החלטות של טריבונלים בעניין סכסוכים בין משכירים לשוכרים,¹⁷⁸ ואף החלטות של בית המשפט העליון של ארצות הברית¹⁷⁹ או של בית הדין האירופי לזכויות אדם (ECHR).¹⁸⁰

ואולם כשמערכות בינה מלאכותית מתמודדות עם חיזוי של החלטות שיפוטיות אין הן נדרשות להסתפק בנתונים משפטיים גרידא (נסיבות המקרה, תקדימים ולשון החוק), אלא יכולות להשתמש בכל מידע לבר-משפטי שאפשר להוסיף למאגר נתוני האימון. מחקרים אקדמיים מראים כי תוצאות משפטיות יכולות להיות מושפעות מגורמים שאינם רלוונטיים לתיק הנדון. שעת הדין וסמיכותה לארוחת הצוהריים למשל – עלולה לכאורה להשפיע על החלטת השופט.¹⁸¹ גם מוצאם של המתדיינים או השופט עלול להשפיע על ההחלטה בתיק.¹⁸²

Hannes Westermann et al., *Using Factors to Predict and Analyze Landlord-Tenant Decisions to Increase Access to Justice*, in PROCEEDINGS OF THE SEVENTEENTH INTERNATIONAL CONFERENCE ON ARTIFICIAL INTELLIGENCE AND LAW (ICAIL '19) 133 (2019)

Daniel Martin Katz, Michael J. Bommarito II, and Josh Blackman, *A General Approach for Predicting the Behavior of the Supreme Court of the United States*, PLoS ONE (12.4.2017). יוער שהניסיונות למדל סטטיסטי את החלטות בית המשפט העליון של ארצות הברית אינם חדשים, אך לא כולם התבססו על למידת מכונה.

Ilias Chalkidis, Ion Androutsopoulos, and Nikolaos Aletras, *Neural Legal Judgment Prediction in English*, PROCEEDINGS OF THE 57TH ANNUAL MEETING OF THE ASSOCIATION FOR COMPUTATIONAL LINGUISTICS 4317 (2019)

181 במחקר של שי דנציגר, יונתן לבב וליאורה אבנעים-פסו נמצא מחאם בין מועד הדיון בוועדת השחרורים להחלטה – ככל שהזמן ממועד ארוחת הצוהריים לדיון היה ארוך יותר, כך פחת הסיכוי לשחרור מוקדם. Shai Danziger, Jonathan Levav, and Liora Avnaim-Pesso, *Extraneous Factors in Judicial Decisions*, 108 PROCEEDINGS OF THE NATIONAL ACADEMY OF SCIENCES 6889 (2011). ואולם בחינה חוזרת של המחקר תלתה את המחאם הסטטיסטי בהיבטים פרוצדורליים אחרים: אסירים שיוצגו על ידי עורך דין (ועל כן ההיתכנות לשחרורם גבוהה יותר) נדונו בוועדה לפני האסירים הלא מיוצגים. Kerem Weinshall-Margel and John Shapard, *Overlooked Factors in the Analysis of Parole Decisions* 108 PROCEEDINGS OF THE NATIONAL ACADEMY OF SCIENCES E833 (2011). ראו גם סמדר רייספלד "מחוץ לחוק: מה משפיע על החלטות של שופטים?" הארץ (2.12.2015). לביקורת נוספת על השימוש במחקר "השופטים הרעבים" כהצדקה להחלפת שופטים בשר ודם במערכות נבונות ראו Konstantin Chatziathanasiou, *Beware the Lure of Narratives: "Hungry Judges" Should not Motivate the Use of Artificial Intelligence*, 23 GERMAN L. J. 452 (2022)

182 ראו לדוגמה אורן גזל-אייל ואח' "ערבים ויהודים בהליכי הארכת מעצר ראשוני" משפטים לח 627 (2009).

ביולי 2019 נחקק בצרפת חוק האוסר להשתמש בנתוני הזיהוי של שופטים לצורכי ניתוח משפטי.¹⁸⁵ הנימוק: שימוש כזה, המתיימר לחזות מראש הכרעה משפטית של כל שופט ושופט, עלול להביא להפעלת לחץ על שופטים או לנקיטת טקטיקות של ברירת סמכות מקומית (forum shopping).¹⁸⁴ מבקרי האיסור גורסים כי מדובר במהלך שנועד להגן על השופטים מאחריות (accountability) מוגברת.¹⁸⁵ אפשר גם שהמהלך הזה נועד לשמר מראית עין של שפיטה עקבית, שאינה מושפעת מ"גודל כף רגלו של הצ'נסלור".

מבלי לקבוע מסמרות בסוגיה נראה אפוא כי רתימתה של הבינה המלאכותית למטרות של חיזוי תוצאות משפטיות או לבחינה אקדמית אמפירית היא אכן טכנולוגיה משבשת שמצריכה טיפול עדין – הן כדי למנוע הנצחה של הטיות קיימות המצויות בנתונים (וכן של פרשנות משפטית, שעלולה לקפוא על שמריה),¹⁸⁶ הן כדי שהיא עצמה תתקבל בגופים מוסדיים קיימים, המנסים לשמר את הסטטוס קוו.

יש כמה מופעים אפשריים של הממשק בין בינה מלאכותית לעולמות התוכן המשפטיים. כפי שראינו לעיל, כבר היום אלגוריתמים שמבוססים על בינה מלאכותית יכולים לחזות החלטות שיפוטיות במגוון תחומים ובדיוק גבוה יחסית, ועורכי דין יכולים להיעזר בהם כדי לעצב את האסטרטגיה המתאימה ביותר ללקוחות שהם מייצגים. מערכות חיזוי של החלטות שיפוטיות יכולות לתרום גם לעבודת השיפוט עצמה, בהציען בסיס

2.3.2. שופט אוטומטי שלל לי רישיון:

בינה מלאכותית
ככלי עזר בשיפוט
ובאכיפת החוק

LOI n° 2019-222 du 23 mars 2019 de programmation 2018-2022 et de réforme pour la justice (1) – Article 33

France, Conseil Constitutionnel, Loi de programmation 2018-2022 et de réforme pour la justice, Décision n° 2019-778 DC

Jena McGill and Amy Salzyzn, *Judging by Numbers: How Will Judicial Analytics Impact the Justice System and Its Stakeholders?* 44 DAL. L. J. 249 (2021)

Daniel Maggen, *Predict and Suspect: The Emergence of Artificial Legal Meaning*, 23 N. C. J. L. & Tech 67 (2021)

אמפירי לרפלקסיה על תהליכי קבלת ההחלטות של השופטים. מערכות בינה מלאכותית יכולות לשמש מערכות תומכות החלטה לשופטים ולסייע להם להגיע באופן יעיל יותר לתובנות מדויקות, עקביות וטובות יותר.

המערכת השיפוטית אינה חפה מטעויות אנוש, והשאיפה לנטרל אותה ככל האפשר מהגורם האנושי אינה חדשה. מערכות חיזוי תומכות החלטה במערכת המשפט תורמות לאובייקטיביות של ההחלטה,¹⁸⁷ אך יש הטוענים כי לנוכח קיומן של הטיות אלגוריתמיות¹⁸⁸ אין כאן אלא מראית עין של אובייקטיביות מדעית.¹⁸⁹ תומס ג'פרסון, מאבות האומה האמריקאית, קבע שיש "להניח לשופט להיות מכונה גרידא".¹⁹⁰ ייתכן שבשנים האחרונות מתחיל חזונו של ג'פרסון לקרום עור וגידים, שכן מדינות שונות משלבות מערכות בינה מלאכותית בעבודת רשויות השיפוט שלהן. אסטוניה הכריזה על כוונתה להשתמש בבינה מלאכותית כשופט הדין בתביעות קטנות (שסכומן אינו עולה על 7,000 אירו), בכפוף לזכות לערער על ההחלטה הממוכנת לפני שופט אנושי.¹⁹¹ גם סרביה שוקלת להשתמש בטכנולוגיות דומות.¹⁹² בבריטניה העלו כמה איגורים מקצועיים והאקדמיה הצעות להפחתת העומס השיפוטי באמצעות מערכות בינה מלאכותית.¹⁹³

187 ראו למשל Jennifer L. Skeem and Christopher Lowenkamp, *Risk, Race, & Recidivism: Predictive Bias and Disparate Impact*, 54 *CRIMINOLOGY* 680 (2015); Anthony W. Flores, Kristin Bechtel, and Christopher T. Lowenkamp, *False Positives, False Negatives, and False Analyses: A Rejoinder to "Machine Bias: There's Software Used Across the Country to Predict Future Criminals. and it's Biased against Blacks,"* 80 *FED. PROBATION J.* 38 (2016).

188 ראו להלן בפרק 6.

189 ראו למשל John Logan Koepke and David G. Robinson, *Danger Ahead: Risk Assessment and the Future of Bail Reform*, 93 *WASH. L. REV.* 1725 (2018).

190 Thomas Jefferson, *Thomas Jefferson to Edmund Pendleton* (26.8.1776).

191 Eric Niler, *Can AI Be a Fair Judge in Court? Estonia Thinks So*, *WIRED* 191 (25.3.2019).

192 Milica Stojanovic, *Serbia Eyes Artificial Intelligence in Courts, but Experts See Dangers*, *BALKANINSIGHT* (25.1.2021).

193 YADONG CUI, *ARTIFICIAL INTELLIGENCE AND JUDICIAL MODERNIZATION* 23 (2020).

בכל הנוגע להטמעת טכנולוגיות בינה מלאכותית במערכת המשפט, סין נמצאת בחזית. מערכת לזיהוי תיקים דומים היא רכיב מרכזי בפרויקט "בית המשפט החכם" של סין. מערכת זו מתבססת על אלגוריתמים המוצאים דפוסים בתיקי בית המשפט (שלא כמו מערכות המנסות לחלץ אוטומטית את ההלכה שנפסקה).¹⁹⁴

מערכת המשפט הסינית משתמשת בעוד מערכת תומכת, "מערכת 206". מערכת זו מסייעת לשופטים ולבאי כוחם של הצדדים בהצגת הראיות הרלוונטיות במשפט. היא מסוגלת לעבד ולהציג באופן ברור ומקיף את מארג הראיות הנוגעות לסוגיה מסוימת ולאתר פערים ואי-התאמות בין ראיות שונות. היא מסוגלת גם לתמלל שפה מדוברת לטקסט, ובכך מייתרת את הצורך בקלדנים אנושיים.¹⁹⁵

בשנת 2019 החל פיילוט של יישום רדיקלי יותר של בינה מלאכותית במערכת המשפט. ההכרעה המשפטית נעשית באמצעות אלגוריתם, אך בהשגחת שופט אנושי.¹⁹⁶ עדיין מוקדם להעריך אם השיפוט האלגוריתמי הזה יוסיף לשמש בסין כלי תומך לשופטים, או שבעתיד יוסר הפיקוח האנושי ויוקנו לו סמכויות עצמאיות. ראוי לציין שבכל פעם שהשופט רוחה החלטה של המערכת הוא נדרש לנמק בכתב מדוע,¹⁹⁷ תנאי שעלול להשפיע על העצמאות השיפוטית ולמצער

Li Xiaohui, *Research on the Building of China's Smart Court in the Internet Era*, 8 CHINA LEGAL SCI. 30, 44 (2020); George G. Zheng, *China's Grand Design of People's Smart Courts*, 7 ASIAN J. L. & Soc. 561 (2020); Jinting Deng, *Should the Common Law System Welcome Artificial Intelligence? A Case Study of China's Same-Type Case Reference System* 3 GEO. L. TECH. REV. 223 (2019)

195 שם. ראו גם Cui, לעיל ה"ש 193, בפרק 7-8.

Nyu Wang and Michael Yuan Tian, "Intelligent Justice": AI Implementations in China's Legal AI Transformation, in ARTIFICIAL INTELLIGENCE AND ITS DISCONTENTS: CRITIQUES FROM THE SOCIAL SCIENCES AND HUMANITIES 197, 212 (Ariane Hanemaayer, ed., 2022); Bin Wei et al., *A Full-Process Intelligent Trial System for Smart Court*, 23 FRONTIERS OF INFORMATION TECHNOLOGY & ELECTRONIC ENGINEERING 186 (2022)

Stephen Chen, *China's Court AI Reaches Every Corner of Justice System, Advising Judges and Streamlining Punishment*, SOUTH CHINA MORNING POST (13.7.2022)

לשנות את מרכז הכובד של ההחלטות המשפטיות.¹⁹⁸ בשנת 2012 פרסמו החוקרים ברדרן ומקנטייר מודל לחיזוי מסוכנות של עצירים לפני הליכים. החוקרים העריכו שיישום המודל היה מייתר רבע מהמעצרים לפני הליכים ושיעור הפשיעה (ובפרט הפשיעה האלימה) לפני משפט היה פוחת.¹⁹⁹ מחקר מאוחר יותר בחן החלטות של שופטים בעיר ניו יורק בעניין שחרור בערבות ופיתח על סמך ממצאיו כלל שחרור אלגוריתמי מקביל. נמצא כי אילו השתמשו באלגוריתם היה אפשר להפחית את רמת הפשיעה בעיר או את מספר המעצרים בשיעור ניכר.²⁰⁰

ואולם גם בארצות הברית החלטות שיפוטיות הנתמכות באלגוריתמים אינן נחלתן של תאוריות אקדמיות. מערכת הצדק האמריקאית משתמשת במערכות אוטומטיות להערכת מסוכנות בכל אחד משלביה – שיטור, מתן גזר דין, כליאה ושחרור ממאסר.²⁰¹

סרטו של שפילברג משנת 2002, "דוח מיוחד",²⁰² המבוסס על סיפור קצר של פיליפ ק' דיק,²⁰³ מתאר מחלקה במשטרה שתפקידה למנוע פשעים עתידיים בטרם יתרחשו בהסתמך על דיווח של מוטנטים הניחנים ביכולת חיזוי על-טבעית. רשויות אכיפת החוק ברחבי העולם ישמחו לזכות ביכולת דומה לחזות פשיעה

198 מנגד, הגישה הסינית רואה באוטומציה כלי להתמודדות עם שחיתות ולהבניית שיקול דעת אחיד והרמוני. על פי גישה זו כל החלטה שיפוטית שסוטה מהחלטה האלגוריתמית חייבת בהנמקה מטעמים של צדק פרוצדורלי והוגנות.

199 Shima Baradaran and Frank L. McIntyre, *Predicting Violence*, 90 Tex. L. Rev. 497, 500 (2012)

200 מאחר שיש מתאם בין רמת הפשיעה בעיר לשיעור העצורים שלא שוחררו בערבות חישובו החוקרים בנפרד תוצאות של שני תרחישים שונים. כשהוקפאה רמת הפשיעה בעיר (ביחס לנתוני האימון) הפחית האלגוריתם את מספר העצורים ב-42%; כשהוקפא שיעור העצורים (ביחס לנתוני האימון) הביא האלגוריתם לירידה של כ-25% ברמת הפשיעה. ראו Jon Kleinberg et al., *Human Decisions and Machine Prediction*, 133 The Quarterly Journal of Economics 237 (2017)

201 RUHA BENJAMIN, RACE AFTER TECHNOLOGY 63 (2019)

202 לעיל ה"ש 132.

203 PHILIP K. DICK, THE MINORITY REPORT (1956)

באמצעות בינה מלאכותית. לצד הצורך להבחין בין נבואה לחיזוי,²⁰⁴ או השאלה העולה בסרט בדבר הנורמטיביות של ההנחות הדטרמניסטיות הגלומות בקביעה האקטוארית שפלוני עתיד לעבור עבירה,²⁰⁵ מערכות הבינה המלאכותית של ימינו מעלות שאלות נוספות באשר לרמת הדיוק של קביעות כאלו ולאפשרות שהן מכילות הטיות אלגוריתמיות.

בכואם לגזור את דינו של נאשם נעזרים מקצת בתי המשפט האמריקאיים בדוחות מ"מערכת לניהול פרופילים לסנקציות חלופיות עבור קציני מבחן" (COMPAS, Correctional Offender Management Profiling for Alternative Sanction). מערכת זו מבצעת, בין היתר, הערכת מסוכנות של הנאשם, כלומר בוחנת את סיכוייו לחזור לסורו (רצידיביזם לפני המשפט, רצידיביזם כללי ורצידיביזם אלים).²⁰⁶ בשנת 2016 פורסמה באתר ProPublica סדרת כתבות שהעלו חשש שמערכת COMPAS סובלת מהטיות גזעניות²⁰⁷ לאחר שנמצאו הבדלים בין התפלגות הערכת הסיכונים שהפיקה המערכת עבור נאשמים שחורים להתפלגות שהפיקה עבור נאשמים לבנים. באותה שנה, בעניין Loomis, דחה בית המשפט העליון של ויסקונסין את הטענה שעצם השימוש בהערכת מסוכנות אלגוריתמית מפרה כללים של הליך נאות, וקבע שההחלטה השיפוטית יכולה

Dr. Lance B. Eliot, *Quandary of Precrime Detection Getting Nearer* 204
Via *Advances in AI* (28.12.2021)

PETER MARKS, *IMAGINING SURVEILLANCE* 101-103 (2015); Jackson ראו למשל
Polansky and Henry Fradella, *Does "Precrime" Mesh with the Ideals of U.S. Justice? Implications For the Future of Predictive Policing* 15
(2016-2017) CARDOZO PUB. L. POL'Y & ETHICS J. 253. לעיסוק בשאלת האוטונומיה של
הפרט בדין הישראלי "לכתוב את סיפור חייו" ראו ע"א 10064/02 מגדל חברה לביטוח
נ' אבו־חנא, פ"ד (ס) 13 (2005). ראו גם אליעזר ריבלין וגיא שני "תפיסה עשירה של
עקרון השבח המצב לקדמותו בתורת הפיצויים הנזיקיים - המקרה של שימוש בסטטיסטיקה
בפסיקה פיצוי לנפגע קטין" *משפט ועסקים* י 499 (2009).

State v. Loomis, 881 N.W.2d 749, 754 (Wis. 2016) 206

Julia Angwin et al., *Machine Bias: There's Software Used across the Country to Predict Future Criminals. and It's Biased against Blacks*,
PROPUBLICA (23.5.2016) 207

להשתמש במערכת COMPAS כמקור מידע נוסף – ובתנאי ששאר הגורמים שהובאו בחשבון יירשמו.²⁰⁸

אימון בינה מלאכותית בנתונים היסטוריים שמגולמת בהם (במשתנים חברתיים-כלכליים למשל) הטיה מערכתית לחובתה של קבוצת מיעוט עלולה אפוא לייצר מערכת חיזוי בעל אפקט גרסיבי המנציח פערים.²⁰⁹ יש הטוענים שההפך הוא הנכון, ושאפשר להסתייע בבינה מלאכותית כדי לזהות מופעים של אפליה.²¹⁰

מבחינות רבות אין שוני רב בין מערכות לחיזוי החלטות שיפוטיות ובין מערכות תומכות החלטה שיפוטית: אלו ואלו מנסות להגיע לתוצאה המשפטית הצפויה. עם זאת, מערכות לחיזוי החלטות שיפוטיות הן חיצוניות למערכת המשפט והניתוח שלהן שואף להגיע לתוצאה האמפירית הרצויה (מה יחליט בית המשפט), ואילו מערכות תומכות החלטה שיפוטית מספקות תוצאה נורמטיבית (מה צריך בית המשפט להחליט). משכך, מערכות אלו נדרשות לעמוד בסטנדרטים נוקשים יותר, הכפופים למגבלות הדין הפרוצדורלי ולכללי הצדק הטבעי.

הדין הצרפתי אוסר על מערכות תומכות החלטה במערכת המשפט. לפי חוק הגנת הפרטיות של צרפת, "החלטה משפטית הכרוכה בהערכת התנהגותו של אדם לא תתבסס על עיבוד אוטומטי של נתונים שנועד להעריך היבטים מסוימים באישיותו של אותו אדם".²¹¹ ואולם אפשר לפרש איסור זה ולומר שאין הוא אוסר

208 State v. Loomis, לעיל ה"ש 206. בית המשפט העליון של ארצות הברית דחה בקשת רשות ערעור של הנאשם Loomis v. Wisconsin, 137 S. Ct. 2290 (2017) (denying cert.) להרחבה בעניין לומיס ראו John Villasenor and Virginia Foggo, *Artificial Intelligence, Due Process and Criminal Sentencing*, 2020 Mich. St. L. Rev. 295 (2020); Katherine Freeman, *Algorithmic Injustice: How the Wisconsin Supreme Court Failed to Protect Due Process Rights in State v. Loomis*, 18 N. C. J. L. & Tech. 75 (2016)

209 בהקשר זה ראו למשל Sandra G. Mayson, *Bias In, Bias Out*, 128 YALE L. J. 2122 (2019)

210 ראו למשל Jon Kleinberg et al., *Algorithms as Discrimination* (28.7.2020) *Detectors*; בהקשר של הפרטום ב-ProPublica על ההטיות במערכת COMPAS (לעיל ה"ש 208) ראו Flores, Bechtel, and Lowenkamp, לעיל ה"ש 187.

211 Loi 78-17 du 6 janvier 1978 relative à l'informatique, aux fichiers et aux libertés [Law 78-17 of January 6 1978 on Information Technology,

שימוש במערכות תומכות החלטה אלא אוסר להסתמך עליהן בלבד בהחלטות משפטיות, בדומה לעניין Loomis.²¹²

מערכות נבונות נמצאות בשירות מנגנוני אכיפת החוק והביטחון גם בשלבים הקודמים להליך המשפטי – מניעת פשע וחיזוי טרור. האתגר של חיזוי הטרור – יותר מהאתגר של מניעת הפשע – מצריך יכולות איסוף מידע רב-ממדיות,²¹³ כלי ניתוח ועיבוד חזקים, ובכלל זה בינה מלאכותית, ומודלים תאורטיים החוזים רדיקליזציה ומסוכנות בהסתמך על פרופילים ברשתות חברתיות.²¹⁴ שירותי הביטחון כבר משתמשים בכלים אלו כדי לחזות פעילות טרור.²¹⁵

גם גופי אכיפת החוק ברחבי העולם מאמצים טכנולוגיות דומות. את הדימוי של המפה העירונית המעטרת את קירות תחנות המשטרה ועליה מסומנים

Data Files and Civil Liberties] JOURNAL OFFICIEL DE LA RÉPUBLIQUE FRANÇAISE [J.O] [OFFICIAL GAZETTE OF FRANCE], Jan. 6, 1978, p. 7, Art. 47

Florence G'sell, *AI Judges*, in THE CAMBRIDGE HANDBOOK OF ARTIFICIAL INTELLIGENCE, GLOBAL PERSPECTIVES ON LAW AND ETHICS (Larry A. Dimatteo, Cristina Poncibo, and Michal Cannarsa eds., 2022)

213 בהקשר של איסוף, ריכוז ועיבוד מידע פנים ממשלתי ראו 115 (2006) BverGE 320; וכן עמיר כהנא ויובל שני רגולציה של מעקב מקוון בדין הישראלי ובדיו המשווה 193-191, 261-267 (מחקר מדיניות 123, המכון הישראלי לדמוקרטיה 2019).

214 ראו לדוגמה - Michael Wolfowicz et al., *Faces of Radicalism - Differentiating between Violent and Non-Violent Radicals by Their Social Media Profiles* 116 COMPUTERS IN HUMAN BEHAVIOR (2021); Abhishek Sachan and Devshri Roy, *TGPM: Terrorist Group Prediction Model For Counter Terrorism*, 44 INTERNATIONAL JOURNAL OF COMPUTER APPLICATIONS 49 (2012); Robert Pelzer, *Policing of Terrorism Using Data from Social Media*, 3 EUR. J. SECUR. RES. 163 (2018)

215 ראו למשל ג'ון בראון " ידו"ח מיוחד": על סמך מה עוצרת ישראל מאות פלסטינים על פשעים שלא ביצעו?" הארץ (19.04.2017); גילי כהן "ראש השב"כ: איחרנו יותר מ-2,000 מפגעים בודדים פוטנציאליים מתחילת 2016" הארץ (27.06.2017); אור הירשאווגה והגר שיזף "סיכול ממוקד: השיטה החדשה להתמודדות עם הטרור נחשפת" הארץ (26.05.2017); עמוס הראל "בינה מלאכותית תתפוס את מקומם של אנשי המודיעין: העתיד לפי מפקד Badi Hasisi, Simon Perry, and Michael (1.10.2021). ראו גם Wolfowicz, *Counter-Terrorism Effectiveness and Human Rights in Israel*, in INTERNATIONAL HUMAN RIGHTS AND COUNTERTERRORISM 409, 419-422 (Eran Shor and Stephen Hoadley eds., 2019)

אזורי הפשיעה בחוטים וסיכות החליפו זה מכבר מערכות מידע גאוגרפי (GIS) המעבדות נתוני פשיעה ממוכנים ומשרטטות על מסך המחשב "נקודות חמות" (hot spots) או "אזורים עתירי פשיעה" (high-crime areas).²¹⁶ זוהי גישת חיזוי משטרתי מבוססת מיפוי (predictive mapping).²¹⁷

מערכות חיזוי פשיעה בשירות ארגוני שיטור יכולות להצביע על אזורים שבהם יש סבירות גבוהה לעבירות רכוש. זאת, בין השאר, בהסתמך על תאוריות קרימינולוגיות שמראות שבתנאים מסוימים צפויה עלייה במספר עבירות הרכוש – פריצה לבית בשכונה מסוימת שעולה יפה יכולה למשל לעודד עבריינים לנסות לשחזר את ההצלחה באותו אזור.²¹⁸ עיבוד היסטורי של אירועי עבריינות בתא שטח גאוגרפי מסוים יכול לזהות אזורים שבהם העבריינים צפויים לפעול שוב. הפניית כוחות שיטור לפטרול באותן "נקודות חמות" עשויה להרתיע עבריינים מלבצע שם עבירות. ואכן, כוחות משטרה ברחבי ארצות הברית דיווחו על ירידה ניכרת בהיקף עבריינות הרכוש בעקבות השימוש במערכות חיזוי פשיעה אלו.²¹⁹ עם זאת, מחקרים נוספים הצביעו על עלייה ברמת העבריינות לאחר מכן,²²⁰ ואחרים לא מצאו כל שיפור מובהק סטטיסטית ברמת הפשיעה במחוזות שבהם השתמשו בטכניקות החיזוי האלו.²²¹

216 על טכנולוגיות מיפוי קרימינולוגיות ראו Andrew Guthrie Ferguson, *Crime Mapping and the Fourth Amendment – Redrawing "High-Crime Areas,"* 63 *HASTINGS L. J.* 179 (2011). ראו גם Andrew G. Ferguson, *Policing Predictive Policing*, 94 *WASH. U. L. REV.* 1109 (2017).

217 Rosamunde van Brakel, *Pre-Emptive Big Data Surveillance and its (Dis)Empowering Consequences: The Case of Predictive Policing*, in *EXPLORING THE BOUNDARIES OF BIG DATA* 117, 120-122 (Bart van der Sloot, Dennis Broeders, and Erik Schrijvers eds., 2016).

218 ראו למשל Daniel B. Neill and Wilpen L. Gorr, *Detecting and Preventing Emerging Epidemics of Crime*, 4 *ADVANCES IN DISEASE SURVEILLANCE* (2007); Jerry H. Ratcliffe and George F. Rengret, *Near-Repeat Patterns in Philadelphia Shootings* 21 *SECURITY J.* 58 (2008).

219 ראו Ferguson, *Policing Predictive Policing* בעמ' 1130.

220 ראו למשל *Richmond Police Chief Says Department Plans to Discontinue "Predictive Policing" Software*, *THE RICHMOND STANDARD* (24.6.2015).

221 PRISCILLIA HUNT, JESSICA SAUNDERS, AND JOHN S. HOLLYWOOD, *EVALUATION OF THE SHREVEPORT PREDICTIVE POLICING EXPERIMENT* (RAND CORP., 2014).

מערכות חיזוי אלגוריתמיות נמצאות בשימוש המשטרה גם כדי להתמודד עם פשיעה אלימה, מתוך הנחה שבאזורים גאוגרפיים מסוימים תנאי השטח נוחים יותר לביצועה – סמטאות ספציפיות שתנאי התאורה בהן לקויים או שנתבי המילוט מהן נוחים, או שכונות שנמצאות בשליטת כנופיה זו או אחרת. חברת Hunchlab פיתחה מערכת המבוססת על גישה זו. היא חוזה פשיעה באמצעות למידת מכונה, בהסתמך על נתוני פשיעה היסטוריים ועל פרמטרים גאוגרפיים, תחבורתיים ומטאורולוגיים.²²² נראה שבחודשי הקליטה שלה בשיקגו הצליחה המערכת להפחית את שיעורי הפשיעה האלימה.²²³

חיזוי פשיעה נעשה גם באמצעות מערכת דירוג חברתי (social scoring).²²⁴ אין הכוונה בהכרח למערכת דרקונית דוגמת זו המשמשת בסין, שיש לה מטרות אחרות,²²⁵ אך העיקרון שביסוד המערכות להערכת מסוכנות וחיזוי רצידיביזם – שתיהן טכניקות חיזוי משטריות שתכליתן זיהוי (predictive identification) – הוא עיקרון זהה: הרשת החברתית של אדם, כלומר מכריו, חבריו ומשפחתו, היא אינדיקציה לרמת העבריינות הפוטנציאלית שלו. כאן מודל השיטור אינו פטרולים באזורים מועדי פשיעה לצורכי הרתעה, אלא זיהוי ישיר של חשודים באמצעים סטטיסטיים. גישה זו עורנה בחיתוליה ויישומה מצומצם יחסית, אך ברחבי ארצות הברית אפשר לראות ניצנים של פרויקטים בתחום.²²⁶

Aaron Shapiro, *Reform Predictive Policing* 541 NATURE 458 (2017) 222

Timothy McLaughlin, *As Shootings Soar, Chicago Police Use Technology to Predict Crime*, REUTERS (5.8.2017) 223

ראו Ferguson, *Policing Predictive Policing* לעיל ה"ש 216, בעמ' 1137. 224

225 להגדרת דירוג חברתי בחקנות הבינה המלאכותית האירופיות, ראו ס' 3(45c) להצעה המחוקקת, לעיל ה"ש 57. כן ראו כהנא ושני, לעיל ה"ש 213, בעמ' 22; הרשות להגנת הפרטיות, דירוג חברתי בראי הזכות לפרטיות: סקירת רקע בעניין שימוש במערכות לדירוג חברתי (2020.4.21). בס' 5(1)(I) להצעה חקנות הבינה המלאכותית האירופיות, לעיל ה"ש 53, מוצע לאסור שימוש במערכות דירוג חברתי בנסיבות מסוימות. על מדינת המעקב הסינית ראו GEOFFREY CAIN, *THE PERFECT POLICE STATE: AN UNDERCOVER ODYSSEY INTO CHINA'S TERRIFYING SURVEILLANCE DYSTOPIA OF THE FUTURE* (2021). על מערכת הדירוג החברתי בסין ראו גם עמדתו הזהירה יותר של Werbach, לעיל ה"ש 239.

226 לסקירה ראו Ferguson, *Policing Predictive Policing*, לעיל ה"ש 216, בעמ' 1143-1138.

2.3.3. חוק בתפירה אישית: רגולציה תבוססת בינה מלאכותית

בינה מלאכותית, כפי שראינו לעיל, יכולה להתבסס על קוד ויכולה להתבסס על למידת מכונה.²²⁷ מערכות מהסוג הראשון נקראות מערכות מומחה. רגולציה המסתייעת במערכות מומחה נשענת על כללי החלטה דטרמיניסטיים ("אם x אז y"). כללים אלו מאפשרים לשרטט עצי החלטה מורכבים מאוד, שיכולים להתאים כלל משפטי למקרה פרטני. מערכות מומחה יכולות לסייע בקבלת החלטות מינהליות בעניין אכיפת דיני מיסים מורכבים²²⁸ או לפשט ולבאר חקיקה מורכבת ולהתאים את התוצאה המשפטית הנכונה לנסיבות ספציפיות (למשל, בבחינת זכויות סוציאליות). ניסוח חקיקה שמתאימה למערכות מומחה דורשת בהירות והיגיון (ולא, הקוד לא יעבוד).²²⁹ הפעלת מערכות כאלו מחייבת קידום של שקיפות ואחריותיות.²³⁰

גישה אחרת היא חקיקה מונחית נתונים (data-driven law). בשנת 2016 הושג מס על בעלי דירה שלישית ויותר.²³¹ הנוסחה שלפיה חושב המס, שתוארה באריכות בחקיקה,²³² התבססה על שני פרמטרים – שטח הדירה ומקומה הגאוגרפי (שיש לו השפעה על כמה מדדים שהלשכה המרכזית לסטטיסטיקה מפרסמת). נוסחה זו הצריכה חישוב שכולל פעולות מתמטיות מורכבות, שמעידות שהנוסחה היא ככל הנראה תוצר של גרסיה לוגיסטית של מדדים אזוריים של הלשכה המרכזית לסטטיסטיקה עם נתוני שווי דירות לפי שטח. הנוסחה הסופית ששימשה בחוק

Mirelle Hildenbrandt, *Algorithmic Regulation and the Rule of Law* 227
376 PHIL. TRANS. R. Soc A (2018)

228 ראו לדוגמה עת"מ 2663-05-19 הר שמש מושב שיחופי בע"מ נ' הממונה על חופש המידע ברשות המיסים (5.9.2019).

229 ראו לדוגמה 70 Tax L. Rev. 377 Sarah B. Lawsky, *Formalizing the Code* (2017)

230 ראו לדוגמה את עניין הר שמש, לעיל ה"ש 228. ראו גם את עמדת ידיד בית המשפט של התנועה לחופש מידע בערעור עע"מ 6782/19 הר שמש מושב שיחופי בע"מ אגש"ח נ' הממונה על חוק חופש המידע ברשות המיסים.

231 ס' 120-148 והחוספת לפרק יב לחוק ההתייעלות הכלכלית (תיקוני חקיקה ליישום המדיניות הכלכלית לשנת התקציב 2017 ו-2018), התשע"ז-2016, ס"ח 2592 בעמ' 17.

232 שם, בחוספת לפרק יב.

הייתה תוצאה של כמה גלגולים, שכן חוקרים מצאו כי גרסאות קודמות שלה אינן מנבאות את שווי הדירה ברמת דיוק סבירה.²³³ בסופו של דבר ביטל בית המשפט העליון בפרשת קוונטינסקי את הוראות החוק המטילות מס זה,²³⁴ בנימוק שהן לא נגעו לתוכנו של החוק אלא לפגמים בשורש הליכי החקיקה שלו בכנסת. עם זאת, מס ריבוי דירות משמש דוגמה לחקיקה מונחית נתונים, שאף עברה טיוב בעקבות הערות של מומחים. ואולם מורכבות הנוסחה ושינויים בערכי המקדמים שלה, שנעשו בשלבי החקיקה האחרונים, מרמזים שעדיין יש מקום לשיפור.

כשחקיקה מונחית נתונים אינה מנוסחת על יסוד מודלים סטטיסטיים שפיתחו אנליסטים אנושיים, וניסוחה מופק בטכניקות של למידת מכונה, קשה מאוד להתחקות בדיעבד על הרציונל שהנחה אותה. בשני המקרים הנוסחה תלויה בטיב הנתונים המשמשים לפיתוח המודל ובפונקציית המטרה שלה.²³⁵

קייסי וניבלט גורסים כי בגלל פערי המידע בין המחוקק לאזרח נדרשת חקיקה המנוסחת ככללים וסטנדרטים. אין המחוקק יכול לחזות מראש כל מערך נסיבות אפשרי שהאזרח עשוי להיקלע אליו, וגם אילו היה יכול הייתה חקיקה כזאת כנראה מסורבלת ויקרה מאוד. לשיטתם של קייסי וניבלט, הבינה המלאכותית תטייב את החוק באמצעות מילוי פערי המידע.²³⁶ הם סבורים שמערכות בינה מלאכותית יכולות להוביל לתום עידן החקיקה המנוסחת ככללים וסטנדרטים ולהחלפתה ב"מיקרו־דירקטיבה".²³⁷

במודל שמתארים קייסי וניבלט, החוק יכיל קטלוג של דירקטיבות תפורות לפי מידה, המתמחות בדיוק במה מותר בכל סיטואציה פרטנית. האזרח יעקוב אחר דירקטיבות שעברו אופטימיזציה עבורו ויתעדכן באופן אוטומטי כשחוקים ישתנו.

233 רוני גולן ודני בן שחר "נוסחת המס על דירה שלישיית מסובכת ומעוותת" *TheMarker* (23.9.2016).

234 בג"ץ 10042/16 צחי קוונטינסקי נ' כנסת ישראל (6.8.2017).

235 Hildenbrandt, לעיל ה"ש 227, בעמ' 3.

236 Anthony J. Casey and Anthony Niblett, *Self-Driving Laws*, 66 U. TORONTO L.J. 429 (2016)

237 *Idem*, *Self-Driving Contracts*, 43 J. CORP. L. 1 (2017)

לדברי קייסי וניבלט, מוקד קבלת ההחלטות יעבור מקביעת כללים וסטנדרטים לניסוח המטרה שהמחוקק מבקש להשיג.

עמרי בן שחר ואריאל פורת מתארים בספרם מודל דומה לזה של קייסי וניבלט.²³⁸ הם משרטטים חזון – שלאחרים הוא ייראה אוטופי ולאחרים מעורר פלצות²³⁹ – של מערכת דינים מותאמת אישית, הנעזרת במערך מורכב של חיישנים, נתוני עתק ואלגוריתמיקה משוכללת כדי להכווין התנהגות אנושית באופן אופטימלי ומדויק. כך למשל, במקום לקבוע הגבלה קשיחה של המהירות המרבית המותרת בכביש יוכל המחוקק לקבוע פונקציית מטרה של "מזעור תאונות דרכים", והמיקרו־דירקטיבות יתאימו באופן אלגוריתמי את המהירות המרבית המותרת לכל נהג ונהג בהתאם לניסיונו, לתנאי הכביש או לפרמטרים אחרים.

שחקנים רבים בתחום הפיננסיים והכלכלה מאמצים מערכות בינה מלאכותית.²⁴⁰ בנקים, למשל, משתמשים במערכות לניתוח חוזי הלוואה מסחריים, החוסכות להם שעות עבודה

2.4 יישומים פיננסיים של בינה מלאכותית

OMRI BEN-SHAHAR AND ARIEL PORAT, PERSONALIZED LAW: DIFFERENT RULES FOR DIFFERENT PEOPLE (2021). ראו גם סימפוזיון מקוון בנושא באתר כתב העת THE UNIVERSITY OF CHICAGO LAW REVIEW

239 ראו בהקשר זה את מאמרו של קווין וורבך על מערכת הדירוג החברתי הסינית, שנחפסה במערב כפנאופטיקון דיסטופי מבעית. וורבך מציע להסתכל על מערכת זו, בניכוי ההקשרים הדכאניים של המשטר האוטוריטרי הסיני, כארכיטקטורה של אסדרה אלגוריתמית המאפשרת האחדה ושיוויון במענה הרגולטורי לאזרחים. Kevin Werbach, *Orwell that Ends Well? Social Credit as Regulation for the Algorithmic Age*, 2022 U. Ill. L. Rev. 1417 (2022)

240 ראו גם הסקירה ב־"ARTIFICIAL INTELLIGENCE AS A DISRUPTIVE TECHNOLOGY" 36-39. כמו כן, ראו הסקירה בארי אחיעז ואח' בינה מלאכותית במגזר הפיננסי: שימושים נפוצים, אתגרים וסקירה השוואתית של התמודדות רגולטורית (מוגש למחלקה הכלכלית בייעוץ וחקיקה [משרד המשפטים], הפקולטה למשפטים, אוניברסיטת תל אביב, 18.7.2022).

אנושיות רבות מספור;²⁴¹ עוד הם משתמשים במערכות לחיזוי התנהגות לקוחות של בנקאות קמעונאית,²⁴² במערכות ממוכנות לייעוץ בנקאי,²⁴³ במערכות להתראה על ניסיונות להלבנת הון²⁴⁴ ובמערכות לדירוג אשראי.²⁴⁵

הן הבנקים הן גורמים חוץ-בנקאיים משתמשים במערכות לדירוג אשראי. מערכות אלו מתבססות על מודלים סטטיסטיים המשקללים את ההיסטוריה הפיננסית של הלקוח (בעיקר עמידה בהחזר הלוואות) ועל פרמטרים אחרים שנמצאים באופן מסורתי במאגרי המידע הבנקאיים²⁴⁶ כדי לתת לו ציון דירוג אשראי המשקף את "הסיכוי שלקוח יעמוד בפרעון התשלומים שבהם הוא מתחייב".²⁴⁷

לצד טכניקות מסורתיות אלו של דירוג אשראי, בשנים האחרונות הולך וגובר השימוש במודלים שפותחו בטכניקות של למידת מכונה, על בסיס מידע שלא משמש במודלים הקלאסיים. קשרים ברשתות חברתיות של הלקוח, סוג הרכב שהוא מחזיק בו, רמת השכלה ואפילו סגנון הכתיבה של הודעות SMS – כולם

Hugh Son, *JPMorgan Software Does in Seconds What Took Lawyers 360,000 Hours*, BLOOMBERG (27.2.2017)

242 חיזוי בנקאי (predictive banking) מבוסס על חיזוי ההתנהלות הפיננסית של לקוחות בנקאיים – ניבוי תשלומים חודשיים חוזרים, אתרעות על הוצאות חריגות, המלצות על מכשירים פיננסיים מתאימים וחוכניות חיטכון. ראו E. Y. Barakina and I.S. Ismailov, *Legal Regulation of Using the Artificial Intelligence Technology in the Banking*, in *Economic Systems in the New Era: Stable Systems in an Unstable World* 40, 43-44 (Svetlana Igorevna Ashmarina et al. eds., 2021)

243 ראו Lee Reiners, *Regulation of Robo-Advisory Services*, in *FINTECH: LAW AND REGULATION* 353 (Jelena Madir ed., 2019)

244 Filip Koprivec et al., *Screening Tool for Anti-Money Laundering Supervision*, in *Big Data and Artificial Intelligence in Digital Finance: Increasing Personalization and Trust in Digital Finance Using Big Data and AI* 233 (John Soldatos and Dimosthenis Kyriazi eds., 2022)

245 Nikita Aggarwal, *The Norms of Algorithmic Credit Scoring*, 80 *CAMBRIDGE LAW JOURNAL* 42 (2021)

246 השוו לס' 16(ב)4 לחוק נתוני אשראי.

247 ס' 2 לחוק נתוני אשראי.

יכולים להיות רלוונטיים בדירוג אשראי מבוסס למידת מכונה.²⁴⁸ מידע מעין זה מתחיל לזרום גם מהרשתות החברתיות עצמן: פייסבוק רשמה לאחרונה פטנט על טכנולוגיה המאפשרת לה להגדיר דירוג אשראי של כל משתמש בהסתמך על דירוגי האשראי של כל חבריו, מתוך הנחה שאדם שחבריו הם בעלי דירוג אשראי נמוך, או שאין הם נוטים להחזיר הלוואות, סיכוי נמוך יותר להחזיר הלוואות במועד.²⁴⁹

ספק אם הבנקים הגדולים יימנעו מלשלב בינה מלאכותית שמבוססת על למידה חכמה ומשתמשת באופן יצירתי במידע חלופי²⁵⁰ במנגנונים המסייעים לקבלת החלטות בנוגע לכראיות מתן הלוואות (להבדיל ממערכות מומחה סטטיות שמבוססות על מודל סטטיסטי שפיתחו אנליסטים אנושיים); על אחת כמה וכמה כשחברות אינטרנטיות מתחילות לנשוף בעורפם של הבנקים ולספק הלוואות בטוחות ורווחיות יותר לכל דורש. אפשר לשער אפוא שבשלב מסוים יאמצו גם הבנקים את הבינה המלאכותית – אף שאולי שימושם בה יהיה מוגבל יותר ויתחשב בפרמטרים מעטים מאלו שמשמשים את חברות ההזנק.

לבינה מלאכותית יש יישומים גם בתחום ההשקעות. האגדה מספרת שהמקור להונה העצום של משפחת רוטשילד הוא פערי מידע. לנתן רוטשילד נודע על ניצחון האנגלים בקרב ווטרו ווהוא מיהר למכור את שטרי החוב האנגליים שברשותו ויצר את הרושם שאנגליה הפסידה. סוחרים אחרים מיהרו להיפטר משטרי החוב שלהם במחירי הפסד, ואנשיו של רוטשילד רכשו אותם בזיל הזול.²⁵¹ פערי המידע נסגרו רק לאחר כמה ימים. בתחילת המאה ה-19 אפשר

Matthew Adam Bruckner, *The Promise and Perils of Algorithmic Lenders' Use of Big Data*, 93 Chi. Kent L. Rev. 3, 15 (2018)

Mark Sullivan, *Facebook Patents Technology to Help Lenders Discriminate against Borrowers Based on Social Connections*, VENTURE BEAT (4.8.2015)

250 בכפוף למגבלות הרגולטוריות הקיימות. ראו לעיל, ה"ש 246.

251 לדיון באגדה זו ובמקורותיה ראו Niall Ferguson, *The House of Rothschild, Volume 1: Money's Prophets, 1798–1848* (1999)

פיתוחה של רשת הטלגרף לסוחרים בעלי גישה לטכנולוגיה זו יתרון תחרותי דומה: היא אפשרה להעביר ידיעות מעיר לעיר בתוך חצי שעה.²⁵²

מאז תחילת האלף הנוכחי מאפשר המסחר הרבובטי (אלגוריתמי) קבלת החלטות בתוך שבריר שנייה וניצול פערים זעירים במידע: למכור ולקנות ניירות ערך בפרקי זמן שסוחרים אנושיים אינם יכולים להתחרות בהם.²⁵³ חלק מהאלגוריתמים מסתמכים על מידע שמגיע גם מחוץ לבורסה, כגון מסרים ברשתות חברתיות כמו טוויטר ופייסבוק.²⁵⁴ סביבת המסחר רוויה באלגוריתמים מתחרים – כל אלגוריתם מנסה למצוא דרכים מיטביות להשאת רווחים. יש אסטרטגיות מסחר אלגוריתמיות שבהן אלגוריתם אחד מנסה לשטות באלגוריתמים אחרים באמצעות הודעה על רצון לרכוש מניות מסוימות, המבוטלת לאחר שברירי שנייה.

²⁵⁵ המסחר האלגוריתמי בבורסות לניירות ערך תופס נתח ניכר ממחזור העסקאות.²⁵⁵ חרף הרווח האינפניטיסימלי מכל עסקה, היכולת לבצע יותר מאלף עסקאות מדי שנייה הופכת רווחים שוליים אלו לסכומים לא מבוטלים. עם זאת, יש ראיות

Bob Pisani, *Plundered by Harpies: An Early History of High-Speed Trading*, 20 FINANCIAL HISTORY (2014)

Neil Johnson et al., *Abrupt Rise of New Machine Ecology Beyond Human Response Time*, 3 SCIENTIFIC REPORTS, art. no. 2627 (11.9.2013); Brendan Conway, *Wall Street's Need for Trading Speed: The Nanosecond Age*, THE WALL STREET JOURNAL (14.6.2011)

Steve Y. Yang and Sheung Yin Kevin Mo, *Social Media and News Sentiment Analysis for Advanced Investment Strategies*, in SENTIMENT ANALYSIS AND ONTOLOGY ENGINEERING: AN ENVIRONMENT OF COMPUTATIONAL INTELLIGENCE 237 (Witold Pedrycz and Shyi-Ming Chen eds., 2016); Nuno Oliveira, Paulo Cortez, and Nelson Areal, *Some Experiments on Modeling Stock Market Behavior Using Investor Sentiment Analysis and Posting Volume from Twitter*, WIMS '13: PROCEEDINGS OF THE 3RD INTERNATIONAL CONFERENCE ON WEB INTELLIGENCE, MINING AND SEMANTICS (2013); Raúl Gómez Martínez, Miguel Prado Román, and Paola Plaza Casado, *Big Data Algorithmic Trading Systems Based on Investors' Mood*, 20 J. BEHAVIORAL FINANCE 22 (2019)

Kevin O'Connell, *Has Regulation Affected the High Frequency Trading Market?*, 27 CATH. U. J. L. & TECH 145, 150-151 (2019)

שרווחי המסחר האלגוריתמי פוחתים מעט בשנים האחרונות, ככל הנראה כיוון שהתחרות אינה עם בני אדם איטיים בעולם הפיזי אלא עם אלגוריתמים מהירים באותה מידה.²⁵⁶ שאלת מעמדן המשפטי של חברות המספקות שירותי מסחר אלגוריתמיים נדונה גם בישראל. בית המשפט המחוזי נדרש לקבוע אם חברות אלו חייבות ברישיון לניהול תיקי השקעות.²⁵⁷

בינה מלאכותית יכולה להיות לעזר גם בהשקעות חוץ-בורסאיות (private equity). משקיעים פרטיים וממשלתיים ערים היטב לקושי לזהות אילו חברות הזנק יעלו לגדולה. כך למשל, מתוך יותר מ-500 חברות הזנק שהתקבלו לחממה היוקרתית Y Accelerator, 37 בלבד הצליחו להימכר או זכו להערכת שווי גבוהה דייה (יותר מ-40 מיליון דולרים). לשון אחר, כ-93% מחברות הזנק שעברו את מסלול הסינון והתקבלו לחממה לא נחלו הצלחה של ממש. בהתחשב בכך ששיעור הקבלה לחממה זו הוא זעום, נראה כי גם משקיעים מנוסים עלולים להתקשות בזיהוי חברות הזנק בעלות פוטנציאל.²⁵⁸ לכן בינה מלאכותית שביכולתה לזהות חברות ומיזמים בעלי פוטנציאל טוב להשקעה היא בעלת ערך רב. ואכן, גופי השקעה שונים פיתחו מערכות אלגוריתמיות לזיהוי השקעות מעין אלו.²⁵⁹

M. Philips, *How the Robots Lost: High-Frequency Trading's Rise and Fall*, BLOOMBERG BUSINESS (6.6.2013). ראו גם O'Connell, שם, בעמ' 152-154.

257 ראו ח"צ (כלכלית) 47119-12-15 אפרימוב נ' יו אס ג'י קפיטל ישראל בע"מ (פורסם בנבו, 24.6.2019); רשות ניירות ערך, הוראה לבעלי רישיון בקשר למחן שירותים תוך שימוש באמצעים טכנולוגיים (2016).

258 Henry Blodget, *DEAR ENTREPRENEURS: Here's How Bad Your Odds Of Success Are*, INSIDER (28.5.2013)

259 Vlad Savov, *Telefónica Will Let an Algorithm Decide which Startups to Invest In*, THE VERGE (12.8.2015); Kirk Kardashian, *Could Algorithms Help Create a Better Venture Capitalist?* FORTUNE (5.8.2015). ראו למשל באחר CORRELATION VENTURES.

2.5

עיבוד שפה

בפרק הקודם ראינו כי כבר בעת ניסוח מבחן טיורינג השתרש הרעיון שהיכולת לעבד שיג ושיח מילולי עם בני אדם ולהשתתף בו היא

אמת מידה לסיווג טכנולוגיה כבינה מלאכותית.²⁶⁰ העיסוק בטכניקות של עיבוד שפה טבעית (NLP) החל עוד בשנות החמישים של המאה העשרים. עיבוד שפה טבעית הוא תחום שקשור לבלשנות ולאינטליגנציה מלאכותית, ומלבד יומרתיו לפענח את הקשר שבין השפה לתודעה האנושית²⁶¹ יש לו שלל יישומים מעשיים, המאפשרים עיבוד אוטומטי ומשוכלל של טקסט.²⁶²

כדי לפענח את הרובד הסמנטי של משפט, כלומר להבין את המשמעות הבסיסית שלו, עלינו לעבד תחילה את הרובד התחבירי שלו. רובד שלישי, מתקדם יותר, הוא הרובד הפרגמטי – ניתוח משמעות המשפט בהקשר של השיח שבו הוא נאמר. טכניקות עיבוד שפה טבעית מכירות במורכבות השפה הטבעית – הן מבחינת רמת ההפשטה הנדרשת כדי לעבד אותה (תחביר, סמנטיקה, הקשר) הן מבחינת הרזולוציה של הניתוח (משפט יחיד לעומת שיח). העמימות של השפה הטבעית – אפילו של משפטים פשוטים – הם מהאתגרים המיוחדים העומדים לפני העוסקים בתחום.

טכניקות של עיבוד שפה טבעית מפרקות את התהליך המורכב למשימות: זיהוי משפטים ומילים ברצף של טקסט, אפיון תחבירי של משפטים, זיהוי שמות של אנשים ומקומות, המרה של טקסט לקול וקול לטקסט ואפילו ניתוח רגשות (sentiment analysis), סיכום מסמכים או תרגום מכונה.

260 לעיל ה"ש 36.

Noam Chomsky, *Language and Nature*, 104 MIND 1, 2 (1995); Patricia 261 Smith Churchland, *Can Neurobiology Teach Us Anything about Consciousness?* 67 PROCEEDINGS AND ADDRESSES OF THE AMERICAN PHILOSOPHICAL ASSOCIATION 23 (1994)

262 לדוגמאות של עיבוד שפה טבעית בתחום המחקר המשפטי האקדמי ראו: LAW AS DATA: COMPUTATION, TEXT, AND THE FUTURE OF LEGAL ANALYSIS (Michael A. Livermore and Daniel N. Rockmore eds., 2019)

אלגוריתמים לעיבוד שפה טבעית נדרשים להתאמן על בסיס נתונים גדול ככל האפשר. בסיס נתונים זה אינו מאגר טקסטואלי ותו לא, אלא מצבור טקסטים שזוהו ותויגו בהם פונמות, מורפמות, מבנים תחביריים, רגשות ועוד – הכול לפי מטרת המודל.²⁶³

לנוכח המורכבות של שפות טבעיות, תרגום מכונה הוא אחד היישומים המאתגרים ביותר של עיבוד שפה טבעית. למילים רבות יש יותר ממשמעות אחת, בין שפות אין התאמה חד-חד-ערכית ולא לכל מילה או ביטוי יש מילה או ביטוי מקביל בשפה אחרת. גם היכולת למיפוי תחבירי של היגדים בשפת התרגום לא תמיד קיימת.²⁶⁴ בתחילת הדרך הייתה הנחת העבודה שמתודולוגיות של פענוח צפנים (קריפטוגרפיה) הן המפתח לתרגום מכונה יעיל,²⁶⁵ אך ניסיונות אלו לא צלחו. עד שנת 2000 התבססה הגישה הפרדיגמטית לתרגום מכונה על מנוע חוקים, שהסתמך לצורך התרגום על קידוד של כללים תחביריים ומילון. מתחילת המאה הגישה המובילה לתרגום מכונה היא שימוש בטכניקות סטטיסטיות של ניתוח קורפוס שפה ובלמידת מכונה.

טכניקות של עיבוד שפה אף מאפשרות לצ'אט בוטים לשוחח ברשת עם בני אנוש. מאז הוצגה בשנת 1966 תוכנת השיחה הראשונה, הפסיכיאטרית האלקטרונית אלייזה,²⁶⁶ התפתחה רמת השיחה בין מערכות ממוחשבות לבני אדם התפתחות ניכרת. אומנם עדיין אפשר להבחין בין צ'אטבוטים לבני אדם, אך כיום הם פרוסים באתרים מקוונים רבים כדי לספק שירות מקוון בהתכתבות ולחסוך המתנה למפעיל אנושי. גם עוזרים דיגיטליים דוגמת סירי ואלכסה משתמשים

263 ראו לדוגמה "קורפוס השפה העברית – תיג מורפולוגי", Data Gov.

264 LYNE BOWKER AND JAIRO BUITRAGO CIRO, MACHINE TRANSLATION AND GLOBAL RESEARCH 37-53 (2019)

265 John Hutchins, *Warren Weaver and the Launching of MT, in EARLY YEARS* 265 IN MACHINE TRANSLATION 17-20 (John Hutchins ed., 2000)

266 Eleni Adamopoulou and Lefteris Moussiades, *Chatbots: History, Technology, and Applications*, 2 MACHINE LEARNING WITH APPLICATIONS (2020) בשנות השבעים פנה המשורר דוד אבדן למעבדות איי-בי-אם בבקשה שיאפשרו לו לשוחח עם אלייזה. ראו דוד אבדן כל השירים כרך ג 9-104 (מוסד ביאליק 2010); ראו גם גל מור "הבינה של אבדן ואלייזה המלאכותית" ynet (2.7.2001).

בפונקציות מורכבות של עיבוד שפה – מהמרת קול לטקסט עד עיבוד בקשות שונות של משתמשים.

לניטור תוכן מקוון יש תכליות מגוונות – זיהוי תכנים פורנוגרפיים שמפירים את כללי הקהילה, זיהוי הפרות של זכויות יוצרים והתמודדות עם ביטויי שנאה ועם סילוף מידע. מרבית הפלטפורמות המקוונות נשענות על שילוב של אלגוריתמים ובקרה אנושית לניטור תוכן, ולפי שעה נראה שאין הן יכולות להסתמך על אלגוריתמים בלבד.²⁶⁷ אומנם משימות של ניטור תוכן אינן מוגבלות לניטור טקסטים, אבל פיתוח יכולות עיבוד של שפה טבעית הן קריטיות להתמודדות עם האתגרים העולים מביטויי שנאה ומסילוף מידע. אתגרים אלו קשים במיוחד לנוכח היעדר הגדרות אחידות לביטויי שנאה ולסילוף מידע,²⁶⁸ לחשש מפני ההשפעות השליליות של ניטור יתר על זכויות אדם²⁶⁹ ולקשיים הטכניים הכרוכים בזיהוי ההקשר שבו נכתבו הדברים.²⁷⁰

יישומי בינה מלאכותית צפויים לשנות את מערכת הרפואה בכל שלביה, משלב המניעה עד שלב הטיפול, ולחולל תמורות במגוון רחב של תחומים – מרפואה פרטית עד תכנון מדיניות בריאות הציבור.

2.6 בינה מלאכותית ברפואה

267 ראו לדוגמה את ההתנצלות שפרסמה טוויטר לאחר שנאלצה להורות לחלק ניכר מאנשי בקרת התוכן שלה שלא להגיע למקום עבודתם בעקבות התפרצות מגפת הקורונה באפריל 2020. Vijaya Gadde and Matt Derella, *An Update on our Continuity Strategy During COVID-19*, TWITTER BLOG (16.3.2020)

268 ראו לדוגמה ROTEM MEDZINI and TEHILLA SHWARTZ ALTSHULER, DEALING WITH HATE SPEECH ON SOCIAL MEDIA (2019); Alexander Brown, *What Is Hate Speech? Part 1: The Myth of Hate*, 36 LAW AND PHILOSOPHY 419 (2017); idem, *What Is Hate Speech? Part 2: Family Resemblances*, 36 LAW AND PHILOSOPHY 561 (2017); ALEXANDER BROWN AND ADRIANA SINCLAIR, THE POLITICS OF HATE SPEECH LAWS 11–19 (2020); Ronan Ó. Fathaigh, Natali Helberger, and Naomi Appelman, *The Perils of Legally Defining Disinformation*, 10 INTERNET POL'Y. REV. (2021)

269 ראו לדוגמה להלן, בסעיף 5.5. United Nations General Assembly, REPORT OF THE SPECIAL RAPPORTEUR ON THE PROMOTION AND PROTECTION OF THE RIGHT TO FREEDOM OF OPINION AND EXPRESSION (להלן: HRC 2018) (A/HRC/38/35, 2018), para 29.

טעויות רפואיות²⁷¹ הן הגורם השלישי למוות בקרב האוכלוסייה בארצות הברית.²⁷² לא מן הנמנע שמצב הדברים דומה גם בישראל. מערכות בינה מלאכותית יכולות לצמצם את היקף הטעויות הרפואיות באמצעות דיאגנוזה שמבוססת על אנליטיקה של ההיסטוריה הרפואית של רכבות ומיליוני חולים אחרים, שרמת הדיוק שלה גבוהה יותר.²⁷³ באמצעות זיהוי אנומליות בסריקות MRI או צילומי רנטגן, ואפילו סיוע בהליכים כירורגיים;²⁷⁴ באמצעות זיהוי הפרעות פסיכיאטריות על פי שינויים בדפוסי הדיבור,²⁷⁵ באמצעות אפליקציות לזיהוי נגעים סרטניים בעור;²⁷⁶ באמצעות חיזוי כוונות התאבדות על סמך היסטוריה פסיכיאטרית או פוסטים ברשתות החברתיות (הללו רלוונטיים למקרי התאבדות שאינם ניתנים לחיזוי על רקע היסטוריה); ועוד.²⁷⁷

כבר עתה יש מערכות נבונות שיכולות האבחון שלהן עולות על אלו של רופא אנושי במגוון תחומים, ובהם דרמטולוגיה, קרדיולוגיה ואונקולוגיה.²⁷⁸ גם

271 בהקשר זה, טעות רפואית פירושה פעולה לא מכוונת (מעשה או מחדל), פעולה שאינה משיגה את התוצאה שלשמה היא בוצעה, פעולה שבוצעה בצורה לא נכונה, פעולה שגויה שאינה מתאימה לנסיבות העניין או חריגה מההליך הטיפולי המקובל – בין שנגרם נזק לחולה ובין שלא.

Martin. A. Makary and Michael Daniel, *Medical Error – the Third* 272
Leading Cause of Death in the U.S., BMJ (2016)

Simon Parkin, *The Artificially Intelligent Doctor Will Hear You* 273
Now, MIT TECHNOLOGY REV. (9.3.2016)

Adam Bohr and Kaveh Memarzadeh, *The Rise of Artificial* 274
Intelligence in Healthcare Applications, in ARTIFICIAL INTELLIGENCE IN
HEALTHCARE 25, 34–35 (Adam Bohr and Kaveh Memarzadeh eds., 2020)

Joseph Frankel, *How Artificial Intelligence Could Help Diagnose* 275
Mental Disorders, THE ATLANTIC (23.8.2016)

Madhumita Murgia, *Google Launches AI Health Tool for Skin* 276
Conditions, FINANCIAL TIMES (18.5.2021)

Mason Marks, *Artificial Intelligence-Based Suicide Prediction*, 21 277
YALE J. L. & TECH 98, 106–111 (2019)

Dayong Wang et al., *Deep Learning for Identifying Metastatic* 278
Breast Cancer (2016), available at ARXIV; Gongning Luo et al., *A Novel Left*
Ventricular Volumes Prediction Method Based on Deep Learning Network
in Cardiac MRI, 43 COMPUTING IN CARDIOLOGY CONFERENCE (CINIC) (2016); Andre Esteve

בתחומים שבהם יש למומחיות אנושית יתרון – קבוע או זמני – על פני מערכות נבונות, הסתייעות של המומחה האנושי בהן יכולה להפחית את הסיכוי לטעויות אנוש.²⁷⁹ מערכות נבונות יכולות לסייע בזיהוי חולים שנדבקו בנגיף קורונה בהסתמך על ניתוח צילומי רנטגן של החזה²⁸⁰ או שכלול שיטת האבחון הקלינית של הנגיף.²⁸¹ מערכות בינה מלאכותית יכולות לתרום גם בשלב הפרוגנוזה או להציע תחזיות באשר להתפתחות מחלות או לסיכוי שתחול הידרדרות במצבם של מטופלים.²⁸²

גם ברמת השירות למטופלים, בינה מלאכותית יכולה להיות לעזר. קופת חולים כללית הכניסה לאחרונה לשימוש מערכת בינה מלאכותית המסייעת לרופא לקבל תמונה מגוונת על מצבו הבריאותי של המטופל, לרבות המלצות להמשך טיפול; המתמללת את השיח בין הרופא למטופל; המנתחת דפוסים של אי-הגעה לתורים שנקבעו מראש; ואפילו מנתחת את מצבו הנפשי של המטופל בזמן אמת.²⁸³ היתרונות של מערכות בינה מלאכותית ברפואה אינן מתמצות בתפקידן ככלי עזר דיאגנוסטי. מחקרים מראים שבעיני המטופלים ייעוץ רפואי באמצעות אפליקציית צ'אטבוטים רגיש וקשוב יותר מייעוץ של רופא בשר ודם.²⁸⁴

et al., *Dermatologist-Level Classification of Skin Cancer with Deep Neural Networks*, 542 NATURE 115 (2017)

279 ראו Wang et al., ש.ס.

Tulin Ozturket et al., *Automated Detection of COVID-19 Cases Using Deep Neural Networks with X-ray Images* 121 COMPUTERS IN BIOLOGY AND MEDICINE (2020)

Liping Sun et al., *Combination of Four Clinical Indicators Predicts the Severe/Critical Symptom of Patients Infected COVID-19*, 128 JOURNAL OF CLINICAL VIROLOGY 104431 (2020); Yazeed Zoabi, Shira Deri-Rozov, and Noam Shomron, *Machine Learning-Based Prediction of COVID-19 Diagnosis Based on Symptoms*, 4 NPJ DIGITAL MEDICINE 3 (2021)

S. SCOTT GRAHAM, *THE DOCTOR AND THE ALGORITHM: PROMISE, PERIL AND THE FUTURE OF HEALTH AI* 46-48 (2022)

283 ישראל וולמן "הטכנולוגיה שתשדרג את מערכת הבריאות" *ynet* (29.10.2022).

John W. Ayers et al., *Comparing Physician and Artificial Intelligence Chatbot Responses to Patient Questions Posted to a Public Social Media Forum* 183 JAMA INTERN. MED. 589 (2023)

– ²⁸⁵(precision medicine) מדויקת רפואה מאפשרת רפואה מותאמת אישית לכל חולה בהתבסס על כמות גדולה של נתונים משלל מקורות, ובהם תיקים רפואיים אלקטרוניים, מידע שהחולים עצמם משתפים ברשתות חברתיות, נתונים המגיעים ממחשוב לכיש, נתונים סביבתיים ועוד. המטרה: לנבא לכל אדם, חולה או בריא, אילו סיכונים רפואיים צפויים לו בעתיד ומה הדרכים היעילות ביותר למנוע את התממשותם. כלומר למצוא מה הטיפול המתאים ביותר עבורו בהינתן ההיסטוריה הרפואית שלו, הסביבה שהוא חי בה, המאפיינים הגנטיים שלו וכדומה.²⁸⁶ כל זה מתוך הנחה ששיטת הטיפול השכיחה או המצליחה ביותר עבור כלל האוכלוסייה לא בהכרח תתאים לכל חולה.

יתר על כן, בכוחן של מערכות בינה מלאכותית להשפיע לחיוב על בריאות הציבור – למשל באמצעות סיוע בניטור התפרצות מגפות.²⁸⁷ כבר בשנת 2002 היו למערכת המודיעין הרפואית הגלובלית GIPHIN אינדיקציות ראשונות להתפרצות מגפת הסארס בסין, בהסתמך על ידיעה עיתונאית על עלייה במכירות של תרופות אנטי-ויראליות.²⁸⁸ מערכות בינה מלאכותית יכולות לסייע

Bertalan Mesko, *The Role of Artificial Intelligence in Precision Medicine*, 2 EXPERT REVIEW OF PRECISION MEDICINE AND DRUG DEVELOPMENT 239 (2017);

J. Larry Jameson and Dan L. Longo, *Precision Medicine – Personalized, Problematic, and Promising*, 372 N. ENGL. J. MED 2229 (2015)

Po-Yen Wu et al., *Omic and Electronic Health Record Big Data Analytics for Precision Medicine*, 64 IEEE TRANSACTIONS ON BIOMEDICAL ENGINEERING 263 (2017)

287 לסקירת שימושים של יישומי בינה מלאכותית לתכליות אפידמיולוגיות ראו Said Agrebia and Anis Larbi, *Use of Artificial Intelligence in Infectious Diseases*, in ARTIFICIAL INTELLIGENCE IN PRECISION HEALTH 415, 422-423 (Debmalya Muhammad Barh ed., 2020); לבינה מלאכותית בהתפרצות הקורונה ראו לדוגמה Awaiz et al., *Overview of IoT and Machine Learning For E-Healthcare in Pandemics and Health Crises*, in DATA SCIENCE ADVANCEMENTS IN PANDEMIC AND OUTBREAK MANAGEMENT 16, 30-33 (Eleana Asimakopoulou and Nik Bessis eds., 2021)

M. Dion, P. AbdelMalik, and A. Mawudeku, *Big Data and the Global Public Health Intelligence Network (GPHIN)*, 41 CANADA COMMUNICABLE DISEASE REPORT 209 (2015)

בחיזוי התפרצויות מקומיות וכך להביא להקצאה יעילה של משאבים רפואיים ולסייע למקבלי ההחלטות.²⁸⁹ בשלבים מוקדמים של התפרצות נגיף קורונה בישראל הוצע להשתמש במערכות נבונות כדי להקצות לכל תושב פרופיל רפואי המשקף את הסיכוי שיידבק בנגיף.²⁹⁰ עם זאת, יש מחקרים שמצאו שכלי בינה מלאכותית רבים כשלו בסיוע בהתמודדות עם הקורונה עקב הסתמכות על בסיסי נתונים "מלוכלכים" וחלקיים או עקב טעויות באימון המודל.²⁹¹

289 Agrebi and Larbi, **לעיל** ה"ש 287. בהקשר של זיהוי התפרצויות מקומיות של נגיף הקורונה ראו Matheus Henrique Dal Molin Ribeiro et al., *Short-Term Forecasting COVID-19 Cumulative Confirmed Cases: Perspectives for Brazil*, 135 CHAOS, SOLITONS & FRACTALS (2020)

290 רפאלה גויכמן "משרד הבריאות חבר ל-NSO כדי לדרג את הסיכוי שחידבקו בקורונה" **TheMarker** (29.3.2020). ראו גם הרשות להגנת הפרטיות, **לעיל** ה"ש 225.

291 לסקירה ראו Will Douglas Heaven, *Hundreds of AI Tools Have Been Built to Catch Covid. None of them Helped*, MIT TECHNOLOGY REVIEW (30.7.2021); Sean Mann et al., *Artificial Intelligence Applications Used in the Clinical Response to COVID-19: A Scoping Review*, 1 PLoS DIGITAL HEALTH (2022)

פרק שלישי

מדיניות אתיקה
של בינה מלאכותית

—

העניין הגובר בטכנולוגיות בינה מלאכותית, לצד הבנת הסיכויים והסיכונים הטמונים בהן, הביאו לפרסום עשרות מסמכים המבקשים להציע מתווי מדיניות, אסטרטגיות וקווים מנחים אתיים הנוגעים לתחום זה,²⁹² ולאחרונה גם להגשת כמה הצעות חקיקה באירופה ובארצות הברית.

292 לסקירות חלקיות ראו, Anna Jobin, Marcello Ienca, and Effy Vayena, *The Global Landscape of AI Ethics Guidelines*, 1 NATURE MACHINE INTELLIGENCE 389 (2019); Yi Zeng, Enmeng Lu, and Cunqing Huangfu, *Linking Artificial Intelligence Principles* (2018), available at ARXIV; Thilo Hagendorff, *The Ethics of AI Ethics: An Evaluation of Guidelines*, 30 MINDS AND MACHINES 99

מקורן של כמה יוזמות כאלו בגופים על-לאומיים, כגון הארגון לשיתוף פעולה ולפיתוח כלכלי (OECD), שלצד סקירה מחקרית של בינה מלאכותית²⁹³ פרסם מסמך עקרונות בנושא.²⁹⁴ את המסמך אימץ גם פורום עשרים שרי הכלכלה ונגידי הבנקים המרכזיים (G-20).²⁹⁵ בשנת 2018 פרסם הדָּח המיוחד של האומות המאוחדות (UN Special Rapporteur) לענייני חופש ביטוי דוח העוסק בהשפעות הבינה המלאכותית על חופש הביטוי וזכויות אחרות.²⁹⁶ כדי להתוות את אסטרטגיית הבינה המלאכותית של אירופה²⁹⁷ הקימה מועצת אירופה קבוצת מומחים (AI HLEG) שניסחו קווים מנחים לבינה מלאכותית אמינה²⁹⁸ והגדרה מקובלת לבינה מלאכותית.²⁹⁹ במהלך שנת 2020 פרסמה המועצה מסמך מדיניות על התייעצות ציבורית בנושא בינה מלאכותית "לקראת מערכת אקולוגית של מצוינות ואמון".³⁰⁰

מדינות רבות אחרות שקדו אף הן על ניסוח מסמכי מדיניות הנוגעים לבינה מלאכותית, אך לרוב במסגרת אסטרטגיה לאומית בעניין בינה מלאכותית ולא מצע לדיון נפרד באתיקה של בינה מלאכותית. בשנת 2016 למשל פרסם ממשל

(2020): Jessica Fjeld et al., *Principled Artificial Intelligence: Mapping Consensus in Ethical and Rights-Based Approaches to Principles For AI*, BERKMAN KLEIN CENTER (2020)

ARTIFICIAL INTELLIGENCE IN SOCIETY (OECD publishing 2019) 293

OECD COUNCIL RECOMMENDATION ON ARTIFICIAL INTELLIGENCE (22.5.2019) (להלן: OECD (2019)). 294

G20 MINISTERIAL STATEMENT ON TRADE AND DIGITAL ECONOMY (2019) 295

ראו HRC 2018, לעיל ה"ש 270. 296

ראו AI4EU, לעיל ה"ש 50. 297

ETHICS GUIDELINES FOR TRUSTWORTHY AI (AI HLEG, 18.12.2018) (להלן: AI HLEG) 298

ETHICS GUIDELINES FOR TRUSTWORTHY AI (AI HLEG 8.4.2019) (להלן: AI HELG); (2018 (2019)). 299

A DEFINITION OF AI: MAIN CAPABILITIES AND SCIENTIFIC DISCIPLINES (AI HLEG, 8.4.2019) 300

European Commission, ON ARTIFICIAL INTELLIGENCE – A EUROPEAN APPROACH TO EXCELLENCE AND TRUST (European Commission 2020) (להלן: European Commission); WHITE PAPER ON ARTIFICIAL INTELLIGENCE: PUBLIC CONSULTATION TOWARDS A EUROPEAN APPROACH FOR EXCELLENCE AND TRUST (17.7.2020) (להלן: EU White Paper 2020). 300

אובמה דוח בנושא בינה מלאכותית בעקבות סדרת סדנאות מומחים שקיים,³⁰¹ ובשנת 2019 פרסם ממשל טראמפ את הצו הנשיאותי לשמירה על מעמדה המוביל של ארצות הברית בבינה מלאכותית, המתווה את האסטרטגיה האמריקאית בעניינה.³⁰² הצו הנשיאותי אינו מכיל את המונח אתיקה, אלא מתרכז בצורך לפתח בינה מלאכותית ולהיות חלוצים בתחום זה.³⁰³ בהתאם להוראות הצו פרסם משרד הניהול והתקציב בבית הלבן (OMB) בשנת 2020 חוזר הכולל עקרונות של מדיניות אסדרה של בינה מלאכותית. מטרת החוזר הייתה בעיקר להנחות סוכנויות פדרליות בשיקולים הנוגעים לאסדרה של התחום באופן שיקדם חדשנות, הסרת חסמים והגנה על ערכים אמריקאיים כמו פרטיות, חירות וזכויות אדם.³⁰⁴ באוקטובר 2022 פרסמה המחלקה המייעצת לנשיא ארצות הברית בנושאי מדע וטכנולוגיה (OSTP, the Office of Science and Technology Policy) את "המתווה למגילת זכויות בבינה מלאכותית – מערכות אוטומטיות בשירות העם האמריקאי".³⁰⁵ על אף הכותרת, מדובר במסמך מדיניות בלתי מחייב, המציע קווי מתאר לפיתוח מדיניות ולעיצוב פרקטיקות שיגנו על הזכויות האזרחיות, יקדמו ערכים דמוקרטיים בכנייה ובפיתוח של מערכות אוטומטיות ויודאו

National Science and Technology Council, Executive Office of the President, PREPARING FOR THE FUTURE OF ARTIFICIAL INTELLIGENCE (2016) 301

US President, EXECUTIVE ORDER ON MAINTAINING AMERICAN LEADERSHIP IN ARTIFICIAL INTELLIGENCE (2019) 302

על פי תקציב אמריקה לשנת 2019, הבינה המלאכותית היא אחד מתחומי העניין העיקריים בתחום המחקר והפיתוח. Office of Management and Budget, Executive Office of the President, BUDGET OF THE UNITED STATES GOVERNMENT, FISCAL YEAR 2019, בעמ' 36. 303

Office of Management and Budget, Executive Office of the President, MEMORANDUM FOR THE HEADS OF EXECUTIVE DEPARTMENTS AND AGENCIES, GUIDANCE FOR REGULATION OF ARTIFICIAL INTELLIGENCE APPLICATIONS, M-21-06 (17.11.2020) 304

The White House, Office of Science and Technology Policy, BLUEPRINT FOR AN AI BILL OF RIGHTS: MAKING AUTOMATED SYSTEMS WORK FOR THE AMERICAN PEOPLE (2022) (להלן: AI BOR). ראו גם עמיר כהנא ודפנה דרור-שפוליאנסקי "זכויות אדם בבינה מלאכותית? על המתווה האמריקאי למגילת זכויות בבינה מלאכותית ומשמעותו" אחר משפט ועסקים (18.5.2023). 305

שלא תהיה פגיעה בזכויות הציבור, בהזדמנויות שלו או בנגישות שלו לצרכים חיוניים.³⁰⁶

רוח הוועדה של בית הלורדים לבינה מלאכותית נדרש להיבטים עקרוניים של אתיקה בבינה מלאכותית.³⁰⁷ גרמניה, שאימצה אסטרטגיה לאומית בנושא בינה מלאכותית,³⁰⁸ הקימה ועדת מומחים לאתיקה של מידע. תפקיד הוועדה לייעץ לממשל הפרדלי בעניין שימושים אתיים בבינה מלאכותית, והדוח הראשון שלה משרטט קווים מנחים בתחום זה.³⁰⁹ במאי 2019 פרסמה האקדמיה של בייג'ין לבינה מלאכותית (ארגון שמימונו בא ממשרד המדע והטכנולוגיה של סין ומהממשלה המוניציפלית של בייג'ין) את מסמך עקרונות הבינה המלאכותית של בייג'ין,³¹⁰ ויש הסבורים שפרסומו מאותת על נכונותה של סין להצטרף לשיח הבינלאומי על שימוש אתי בטכנולוגיה.³¹¹ המשרד לענייני פנים ותקשורת של יפן פרסם ב-2017 טיוטת הנחיות למפתחי מערכות נבונות.³¹²

בישראל, ועדת המשנה של המיזם הלאומי למערכות נבונות בנושא אתיקה ורגולציה של בינה מלאכותית (להלן: ועדת המשנה) הגישה בשלהי שנת 2019 דוח מסכם. ברוח נכללו העקרונות האתיים שנמצאו בעיניה בעלי חשיבות

306 AI BOR, לעיל ה"ש 305, בעמ' 3.

307 House of Lords Select Committee on Artificial Intelligence, AI IN THE UK: READY, WILLING AND ABLE? 125 (2018)

308 Bundesministerium für Bildung und Forschung; Bundesministerium für Wirtschaft und Energie; Bundesministerium für Arbeit und Soziales, STRATEGIE KÜNSTLICHE INTELLIGENZ DER BUNDESREGIERUNG (15.11.2018)

309 Datenethikkommission der Bundesregierung, OPINION OF THE DATA ETHICS COMMISSION 43 (22.1.2020)

310 Beijing Academy of Artificial Intelligence, BEIJING AI PRINCIPLES (2019)

311 Will Knight, *Why Does Beijing Suddenly Care about AI Ethics?* MIT TECHNOLOGY REVIEW (31.5.2019)

312 The Conference toward AI Network Society, Ministry of Internal Affairs and Communications, DRAFT AI R&D GUIDELINES FOR INTERNATIONAL DISCUSSIONS (28.7.2017)

עבור מקבלי החלטות.³¹³ בשנת 2021 גיבש הצוות הבין-משרדי לרגולציה ואתיקה בתחום הבינה המלאכותית בישראל, במסגרת התוכנית הלאומית לבינה מלאכותית, רשימת עקרונות אתיים שמן הראוי שהתוכנית תפעל לאורם. בגיבוש הרשימה הושם דגש על עקרונות ה-OECD ועל המלצות ועדת המשנה. המסמך "עקרונות מדיניות, רגולציה ואתיקה בתחום הבינה המלאכותית", שפרסם משרד החדשנות, המדע והטכנולוגיה בשלהי שנת 2022 כולל גם המלצות לעקרונות אתיים.³¹⁴

לא רק גופים על-לאומיים ומדינות עוסקים בסוגיה. עוד עוסקים בה תאגידים בינלאומיים, ובהם אינטל,³¹⁵ מייקרוסופט,³¹⁶ גוגל,³¹⁷ איי-בי-אם³¹⁸ וסוני;³¹⁹ קבוצות חשיבה מקרב התעשייה, במסגרת התאגדויות אד-הוק דוגמת "שותפות לבינה מלאכותית" (Partnership on AI), שחברות בה גוגל, פייסבוק, אמזון ואפל ועוד;³²⁰ או גופים ותיקים כמו הארגון הבינלאומי של מהנדסי האלקטרוניקה והחשמל (IEEE).³²¹ מובן שייתכן שמסמכי מדיניות אלו, מאת חברות פרטיות,

313 המיזם הלאומי למערכות נבונות, "ועדת המשנה בנושא אחיקה ורגולציה של בינה מלאכותית, דוח מסכם" (2019) (להלן: המיזם הלאומי למערכות נבונות (2019)).

314 משרד החדשנות, המדע והטכנולוגיה, מדיניות רגולציה ואתיקה בתחום הבינה המלאכותית בישראל (2022) (להלן: מסמך מדיניות רגולציה של משרד החדשנות). ראו גם תהילה שוורץ אלטשולר, עמיר כהנא, גדי פרל ואורי פריימן אסדרת הבינה המלאכותית בישראל מחייבת קווים אדומים שימנעו פגיעה בזכויות יסוד (חוות דעת, המכון הישראלי לדמוקרטיה 30.12.2022).

315 *Artificial Intelligence – The Public Policy Opportunity*, INTEL (2017)

316 *Responsible AI*, MICROSOFT

317 *Our Principles*, GOOGLE AI

318 Adam Cutler and Milena Pribić, *EVERYDAY ETHICS FOR ARTIFICIAL INTELLIGENCE* (IBM Design Program Office, 2018)

319 *Sony Group AI Ethics Guidelines*, SONY (2018)

320 *Partnership on AI*, AI (2016)

321 IEEE, *ETHICALLY ALIGNED DESIGN: A VISION FOR PRIORITIZING HUMAN WELLBEING WITH AUTONOMOUS AND INTELLIGENT SYSTEMS* (2019)

אינם אלא ניסיון לבלום רגולציה עתידית באמצעות הצהרת כוונות עמומה שעיקרה איתות סגולה (virtue signaling) לא מחייב.³²²

גם האקדמיה וארגוני החברה האזרחית פרסמו ניירות עמדה בנושא אתיקה ורגולציה של בינה מלאכותית. עם אלו אפשר למנות, בין השאר, את הצהרת מונטריאול של אוניברסיטת מונטריאול;³²³ את הצהרת טורנטו, שמשותפת לארגונים אמנסטי ו-Access Now;³²⁴ את הדוח "פרטיות וחופש ביטוי בעידן הבינה המלאכותית", שפרסמו הארגונים Privacy International ו-19 Article;³²⁵ ופרסומים נוספים מאת מרכזי מחקר, איגודים מקצועיים ומכונים אקדמיים.³²⁶

ריאן קיילו העיר כי מסמכי מדיניות אתיקה מאפשרים פרשנות גמישה ומתאפיינים ברקמה פתוחה יחסית, כיוון שרמת ההפשטה בהם גבוהה. עוד ציין קיילו שגם

322 ראו למשל עמדתו של חומס מצינגר, חבר קבוצת המומחים האירופית, הרואה אף במסמך המומחים האירופי (AI HELG 2019), לעיל ה"ש 298) דוגמה לטיוח אחי – דיונים אחרים כאמצעי לשיהוי אסדרה אפקטיבית. Thomas Metzinger, *Ethics Washing*. Inga Ulmancic, *Artificially Made in Europe*, DER TAGESSPIEGEL (8.4.2019). *Intelligence in the European Union – Policy, Ethics and Regulation*, in THE ROUTLEDGE HANDBOOK OF EUROPEAN INTEGRATIONS, 254, 263 (Thomas Hoerber, Gabriel Weber, and Ignazio Cabras eds., 2022); Paul Nemitz, *Constitutional Democracy and Technology in the Age of Artificial Intelligence*, 376 PHILOSOPHICAL TRANSACTIONS OF THE ROYAL SOCIETY A (2018).

323 THE MONTREAL DECLARATION FOR A RESPONSIBLE DEVELOPMENT OF ARTIFICIAL INTELLIGENCE (להלן: הצהרת מונטריאול). (2017).

324 THE TORONTO DECLARATION: PROTECTING THE RIGHT TO EQUALITY IN MACHINE LEARNING (להלן: הצהרת טורונטו). (2018).

325 PRIVACY AND FREEDOM OF EXPRESSION IN THE AGE OF ARTIFICIAL INTELLIGENCE (Privacy International and Article 19, April 2018) (להלן: PI&A19).

326 ראו למשל *Principles for Accountable Algorithms and a Social Impact Statement for Algorithms*, FAT\ML (2016); *Asilomar AI Principles*, FUTURE FOR LIFE INSTITUTE (2017); Peter Cihon, STANDARDS FOR AI GOVERNANCE – INTERNATIONAL STANDARDS TO ENABLE GLOBAL COORDINATION IN AI RESEARCH & DEVELOPMENT (Future of Humanity Institute, Oxford University, 2019); US Public Policy Council, Association for Computing Machinery, STATEMENT ON ALGORITHMIC TRANSPARENCY AND ACCOUNTABILITY (2017) (להלן: ACM 2017).

כשיש קונצנזוס בנוגע למוסריות של טכנולוגיות בינה מלאכותית, אין בנמצא מנגנונים לאכיפת הפרות. לכן חברות מסחריות העוסקות בפיתוח טכנולוגיות בינה מלאכותית, ובפרט החברות המובילות בשוק, מעדיפות סטנדרטים אתיים שאפשר להתעלם מהם לפי הצורך על פני כללים מחייבים מבחינה משפטית.³²⁷

לא נתאר בפרוטרוט את כל המסמכים האלו – שנוסחו בזמנים שונים למטרות שונות – אלא נתרכז בעקרונות שחוזרים ומופיעים בהם. צוות חוקרים במרכז ברקמן קליין באוניברסיטת הרווארד סקר שלושים ושישה מסמכים מסוג זה ואיתר שמונה תמות שחוזרות ברובם.³²⁸ ז'ובן, ינקה וואיינה דגמו שמונים מסמכים כאלו וזיהו באמצעות טכניקות כמותיות כ-11 אשכולות של עקרונות שחוזרים ומופיעים בהם.³²⁹ זנג, לו והונגפאו דגמו 27 מסמכים כאלו, השתמשו בניתוח סמנטי כדי לאפיין נושאים חוזרים וקישרו בין נטיית המסמכים לציין נושאים אלו ובין טיב הגורם שחיבר אותם (ממשלתי, פרטי או מגזר שלישי).³³⁰ אומנם ייתכן שאפשר להעלות השגות מתודולוגיות על הממצאים הכמותיים של מחקרים אלו ומחקרים דומים, אבל ההיריסטיקה האיכותנית של העקרונות השכיחים באתיקה של הבינה המלאכותית העולה ממחקרים אלו מדויקת די הצורך לצרכים שלנו.

תמת השקיפות³³¹ שכיחה במסמכי אתיקת בינה מלאכותית.³³² פיילד ועמיתיה כוללים בה, נוסף על שקיפות, גם הסברתיות, קוד פתוח ומידע פתוח, וכן אינטראקציה מיועדת עם בינה מלאכותית.³³³ לעיתים עקרון השקיפות

3.1 שקיפות

- Ryan Calo, *Artificial Intelligence Policy: A Primer and Roadmap*, 51 327
U.C. Davis L. Rev. 399, 408 (2017)
- ראו Fjeld et al., לעיל ה"ש 292. 328
- ראו Jobin, Ienca and Vayena, לעיל ה"ש 292. 329
- ראו Zeng, Lu and Huangfu, לעיל ה"ש 292. 330
- ראו בהרחבה גם בפרק 7 להלן; וכן ראו פרל ושוורץ אלטשולר, להלן ה"ש 881. 331
- 332 94% מן המסמכים שדגמו Fjeld et al., וכ-87% מן המסמכים שדגמו Jobin, Ienca and Vayena (לעיל ה"ש 292) מתייחסים לשקיפות.
- 333 Fjeld et al., לעיל ה"ש 292, בעמ' 41.

מתואר כתנאי ולעיתים הוא נכלל בעקרונות אחרים, כגון עקרון ההוגנות³³⁴ או עקרון האחריותות.³³⁵ יש המדגישים שהשקיפות חשובה כיוון שהיא מאפשרת פיקוח של בעלי עניין מגוונים (ובהם ארגוני החברה האזרחית) על מערכות בינה מלאכותית, והלכה למעשה היא תנאי לאחריותות.³³⁶

דרך אחרת להמשיג את עקרון השקיפות בהקשר של בינה מלאכותית הוא לפנות לעקרון המפורשות (explicability). המודל האתי שהציעו לוצ'אנו פלורידה ועמיתיו לבינה מלאכותית נסמך על ארבעה נדבכים ביו-אתיים כלליים, שעליהם נוסף נדבך המפורשות.³³⁷ לשיטתם, עקרון המפורשות הוא מזיגה בין ההיבט האפיסטמולוגי – ההבנה של אופן פעולת מערכת הבינה המלאכותית,³³⁸ להיבט האתי – שאלת האחריות לפעולתה.³³⁹ גם כשמערכת בינה מלאכותית פועלת לכאורה בשקיפות מלאה ורבבות שורות הקוד שלה פתוחות לעיון הציבור, אין ערובה שהדיוטות יצליחו להבין כיצד היא פועלת ומי אחראי לתוצריה.³⁴⁰

ההבחנה בין שקיפות גרידא ובין מפורשות מופיעה במסמך המומחים של האיחוד האירופי (להלן: מסמך המומחים). על פי המסמך, עקרון המפורשות נמנה עם עקרונות-העל של בינה מלאכותית אתית. כדי שתתקיים מפורשות נדרשת שקיפות של תהליכי הבינה המלאכותית, היכולות והמטרות של מערכות בינה

334 ראו למשל National Science and Technology Council, לעיל ה"ש 301, בעמ' 41; European Commission 2020, לעיל ה"ש 300, בעמ' 11.

335 ראו המיזם הלאומי למערכות נבונות (2019), לעיל ה"ש 313, בעמ' 16; Margot E. Kaminski, *Understanding Transparency in Algorithmic Accountability*, in THE CAMBRIDGE HANDBOOK OF THE LAW OF ALGORITHMS 121 (Woodrow Barfield ed., 2021).

336 ראו למשל PI&A19, לעיל ה"ש 325, בעמ' 28; OECD 2019, לעיל ה"ש 294, בעמ' 8.

337 Luciano Floridi et al., *Ai4people – An Ethical Framework for a Good AI Society: Opportunities, Risks, Principles, and Recommendations*, 28 MINDS AND MACHINES 689 (2018).

338 עקרונות הבינה המלאכותית של חברת מיקרוסופט, למשל, מגדירות שקיפות כיכולת להבין מערכת בינה מלאכותית. ראו *Responsible AI*, MICROSOFT, לעיל ה"ש 316.

339 השוו לפירוק מושג השקיפות אצל PAULA BODDINGTON, TOWARDS A CODE OF ETHICS FOR ARTIFICIAL INTELLIGENCE 29 (2017).

340 CHRISTOPH BARTNECK, CHRISTOPH LÜTGE, ALAN WAGNER, AND SEAN WELSH, AN INTRODUCTION TO ETHICS IN ROBOTICS AND AI 36 (2021).

מלאכותית צריכים להיות גלויים וההחלטות שהן מקבלות צריכות להיות ניתנות להסבר. כשמדובר בקופסאות שחורות שאי-אפשר להסביר את פעולתן נדרשים אמצעים משלימים, בהם היכולת לבקר את המערכת (auditability)³⁴¹ ולעקוב אחר התהליכים בה (נעקבות, traceability)³⁴².

היבט ההסבריות של עקרון השקיפות (או למצער היכולת לבקר מערכות בינה מלאכותית ולעקוב אחר התהליכים בהן) נועד, בין השאר, לזהות את הנזקים שהן עלולות להסב³⁴³ ולאפשר בקרה אנושית על החלטות ממוכנות בעלות חשיבות.³⁴⁴ היבט נוסף של שקיפות הוא אינטראקציה מיועדת עם בינה מלאכותית.³⁴⁵ מערכות בינה מלאכותית נדרשות להציג את עצמן למשתמשים בתור כאלה ולא לגרום להם לחשוב בטעות שמדובר באנשים בשר ודם.³⁴⁶ במסמך מדיניות הרגולציה של משרד החדשנות נכללים ההסבריות והאינטראקציה המיועדת בעקרון השקיפות.³⁴⁷ ב"מתווה האמריקאי למגילת זכויות בכינה מלאכותית" הם נכללים בעיקרון "יידוע והסבר".³⁴⁸

יש מסמכים שנדרשים להיבט הציבורי של שקיפות מערכות המידע. ההסתמכות על קוד פתוח מאפשרת פיתוח אלגוריתמים משותפים, פיקוח משותף רחב עליהם

341 ראו גם ACM 2017, לעיל ה"ש 326, בעמ' 2 (עיקרון 6).

342 AI HELG 2019, לעיל ה"ש 298, בעמ' 13. על נעקבות ראו גם: AI PRINCIPLES: RECOMMENDATIONS ON THE ETHICAL USE OF ARTIFICIAL INTELLIGENCE BY THE DEPARTMENT OF DEFENSE (Defense Innovation Board 2019), עיקרון 8, עיקרון 3.

343 ראו *Asilomar AI Principles*, לעיל בה"ש 326, עיקרון 7. לביקורת ראו BODDINGTON, לעיל ה"ש 339, בעמ' 107.

344 ראו *Asilomar AI Principles*, לעיל ה"ש 326, עיקרון 8; Commission Regulation 2016/679 of 27 Apr. 2016 on the Protection of Natural Persons with Regard to the Processing of Personal Data and on the Free Movement of Such Data, and Repealing Directive 95/46/EC (General Data Protection Regulation), 2016 O.J. (L 119) 1 (EU), Article 22

345 ראו לדוגמה AI HELG 2019, לעיל ה"ש 298, בעמ' 18; ס' 32 להצהרת טורונטו, לעיל ה"ש 324; *Guidelines for Artificial Intelligence*, DEUTSCHE TELEKOM (2018).

346 ראו לדוגמה ס' 52(1) להצעת תקנות הבינה המלאכותית האירופיות, לעיל ה"ש 53; ס' 5(V)(a) להצעת החוק הברזילאית לבינה מלאכותית, לעיל ה"ש 54.

347 מסמך מדיניות הרגולציה של משרד החדשנות, לעיל ה"ש 314, בעמ' 106.

348 AI BOR, לעיל ה"ש 305, בעמ' 40-45.

ואיגום משאבים. כך גם אפשר למנוע היווצרות של מונופולים בתחום – גופים החולשים לברם על שירותי בינה מלאכותית מסוימים או על מידע.³⁴⁹

עקרון ההוגנות נידון אף הוא במסמכים רבים.³⁵⁰ הטיפולוגיה של ז'ובן, ינקה וואיינה קושרת את ההוגנות גם לצדק,³⁵¹ ובפרט למניעה, לצמצום ולניטור של הטיות בלתי רצויות ואפליה.³⁵² במסמכים אלו ההיבטים הדומיננטיים של הוגנות הם אי-אפליה ומניעת הטיות,³⁵³ אך עקרון ההוגנות מוזכר לעיתים גם בדיונים על שוויון,³⁵⁴ צדק חלוקתי,³⁵⁵ הכלה (inclusiveness)³⁵⁶ וגיוון (diversity)³⁵⁷ בתכנון מערכות הבינה המלאכותית, וכן בהסתמכות על נתונים אמנים ומייצגים ללמידת מכונה. כאשר מדובר במערכות נבונות המקבלות החלטות הרות גורל בעניין אזרחים, יש להבטיח כי אלו תפעלנה בהוגנות ובשוויון.³⁵⁸

3.2 הוגנות

- 349 ראו לדוגמה את עיקרון 1.7 אצל: Beijing Academy of Artificial Intelligence, לעיל ה"ש 310; OECD 2019, לעיל ה"ש 294, בעמ' 8.
- 350 100% מהמסמכים שדגמו Fjeld et al., וכ-76% מהמסמכים שדגמו Jobin, Ienca and Vayena (לעיל ה"ש 292) מתייחסים לשקיפות.
- 351 ראו למשל Floridi et al., לעיל ה"ש 337, בעמ' 698-699; הצהרת מונטריאל, לעיל ה"ש 323, עיקרון 6.
- 352 Jobin, Ienca and Vayena, לעיל ה"ש 292.
- 353 ראו Google, לעיל ה"ש 317, עיקרון 2; Defense Innovation Board, לעיל ה"ש 342, עיקרון 2; AI BOR, לעיל ה"ש 305, בעמ' 23-29.
- 354 למשל, בהצהרת טורונטו מצוין בפירוש "הצהרה זו מתמקדת בזכות לשוויון ולהיעדר אפליה" (ס' 12 להצהרת טורונטו, לעיל ה"ש 324), ובהמשכה מתוארים קווים מנחים שלאורם מדינות צריכות לפעול כדי לצמצם ולמתן את הסיכון לאפליה שמקורה בטכנולוגיות למידת מכונה הן במגזר הציבורי (ס' 30-41) הן במגזר הפרטי (ס' 42-51).
- 355 ראו Datenethikkommission, לעיל ה"ש 309, בעמ' 46-47.
- 356 ראו למשל Responsible AI, Microsoft, לעיל ה"ש 316.
- 357 Cutler and Pribić, לעיל ה"ש 318, בעמ' 32; Sony Group AI Ethics Guidelines, לעיל ה"ש 319, עיקרון 5. לדיון השוואתי בניסוח עקרון ההוגנות התוצאתית במסמכים שונים ראו Fjeld et al., לעיל ה"ש 292, בעמ' 51.
- 358 ראו למשל National Science and Technology Council, לעיל ה"ש 301, בעמ' 34.

מסמך המומחים מבחין בין הוגנות מהותית להוגנות פרוצדורלית.³⁵⁹ הוגנות מהותית משמעה מחויבות לחלוקה הוגנת ושווה של כל היתרונות הטמונים במערכות בינה מלאכותית ושל כל העלויות שלה, ולהיעדר הטיות לא-הוגנות, אפליה וסטיגמות. מסמך המומחים מדגיש בהקשר זה גם את עקרון המידתיות ואת חשיבות האיזון בין אינטרסים מתחרים. הוגנות פרוצדורלית, לעומת זאת, משמעה האפשרות לערער על החלטות של מערכות בינה מלאכותית ולקבל סעדים.

אימון או הפעלה של מערכות בינה מלאכותית עלולים להתבסס על בסיסי נתונים לא שלמים או כאלה שמשקפים הטיות היסטוריות או ממשל לא תקין. משום כך הם עלולים להביא לידי אפליה עקיפה של פרטים או קבוצות מסוימים.³⁶⁰ מסמך המומחים מונה אפוא דרכים שונות למימוש בינה מלאכותית ראויה לאמון, ובהן גיוון, אי-אפליה והוגנות.³⁶¹ דרכים אחרות הן שיתוף מגוון בעלי עניין בתהליכי התכנון, ההטמעה וההפעלה השוטפת של מערכות בינה מלאכותית, וכן הבטחת נגישות של מגוון משתמשים למערכות אלו. מובן שפרקטיקות לא הוגנות, כגון אפליה אלגוריתמית בתמחור צרכנים,³⁶² יכולות להיות גם מכוונות, אך הן אינן מענייננו כאן.

רוח ועדת המשנה נדרש להיבטים של שוויון תוצאתי וממליץ להימנע ממצבים שבהם תוצרי הבינה המלאכותית מעצימים רק קבוצות מסוימות באוכלוסייה. דוח זה, כמו מסמך המומחים, נדרש גם לצורך לייצג מגוון אוכלוסיות – הן בתוך מאגרי המידע שהמערכות מסתמכות עליהן הן בשלבי הפיתוח שלהן – ורואה בו היבט חשוב של עקרון ההוגנות.³⁶³ על פי מסמך מדיניות

359 AI HELG 2019, לעיל ה"ש 298, בעמ' 12.

360 ראו למשל HRC 2018, לעיל ה"ש 270, בפס' 38: *BigData: Discrimination*, EUROPEAN UNION AGENCY FOR FUNDAMENTAL RIGHTS *in Data-Supported Decision Making*, (30.5.2018)

361 AI HELG 2019, לעיל ה"ש 298, בעמ' 18.

362 Michal S. Gal and Niva Elkin-Korren, *Algorithmic Consumers*, 30 HARV. J. L. & TECH. 309 (2017)

363 ראו המיזם הלאומי למערכות נבונות (2019), לעיל ה"ש 313, בעמ' 15.

הרגולציה של משרד החדשנות, "שוויון ומניעת אפליה פסולה" הוא עיקרון אתי עצמאי.³⁶⁴

3.3

מניעת נזקים ובטיחות

בעיקרון זה נשמעים במוכן מסוים הדיו של חוק הרובוטיקה הראשון של אסימוב (וחוק האפס),³⁶⁵ וכן של עקרון היסוד של האתיקה

הרפואית, *primum non nocere* (ראשית, אל תזיק).³⁶⁶ אחת מהמלצות דוח הוועדה של בית הלורדים על בינה מלאכותית היא שלא להעניק למערכות בינה מלאכותית סמכות אוטונומית לפגוע, להשמיד או להטעות בני אדם.³⁶⁷ עקרונות אסילומר נוקטים לשון של היעדר חתרנות (non-subversion): סמכויות השליטה במערכות בינה מלאכותית אל להן לחתור תחת תהליכים חברתיים ואזרחיים החיוניים לבריאותה של החברה.³⁶⁸ מסמך ההמלצות של משרד ההגנה של ארצות הברית מדבר על משילות (governable AI): עיצובן של מערכות נבונות צריך לכלול את היכולת לזהות נזקים לא צפויים ולמנוע אותם.³⁶⁹ המתווה האמריקאי למגילת זכויות כבינה מלאכותית כולל בעקרונותיו "מערכות בטוחות ויעילות".³⁷⁰

364 מסמך מדיניות הרגולציה של משרד החדשנות, לעיל ה"ש 314, בעמ' 105.

365 אסימוב, לעיל ה"ש 47, בעמ' 6. ראו גם Robin Murphy and David D. Woods, *Beyond Asimov: The Three Laws of Responsible Robotics*, 24 IEEE INTELLIGENT SYSTEMS 14 (2009)

366 TOM L. BEAUCHAMP AND JAMES F. CHILDRESS, *PRINCIPLES OF BIOMEDICAL ETHICS*, 150-201 (7th ed., 2013); גיל סיגל "ביו־אתיקה בעולם המערבי: ישראל שובן אירופה לארה"ב" הרפואה 142, 143 (2004). ראו גם התייחסות מפורשת למקור הביו־אתי של הכלל אצל: Floridi et al., לעיל ה"ש 337, בעמ' 697.

367 House of Lords Select Committee on Artificial Intelligence, לעיל ה"ש 307, בפס' 417(5).

368 *Asilomar AI Principles*, לעיל ה"ש 326, עיקרון 17.

369 Defense Innovation Board, לעיל ה"ש 342, עיקרון 5.

370 AI BOR, לעיל ה"ש 305, בעמ' 23-29.

ז'ובן, ינקה וואיינה מציינים שעקרון מניעת הנזק שכיח יותר מעקרון עשיית הטוב וקושרים בינו ובין עקרון הבטיחות.³⁷¹ המיפוי של פיילרד ועמיתיה מזכיר את עקרון מניעת הנזק רק במשתמע ומתמקד בבטיחות ובביטחון.³⁷² ההתמקדות בהיבטים אלו מסייעת לקונקרטיזציה של עקרון מניעת הנזק וליישומו.³⁷³ על פי מסמכים אחדים, ניהול סיכונים³⁷⁴ ויכולת חיזוי (predictability)³⁷⁵ הם פנים נוספים של בטיחות במערכות נבונות.

מסמך המומחים נדרש לקונקרטיזציה של עקרון מניעת הנזק. המסמך מסביר שלמנוע נזק פירושו לוודא שמערכות נבונות אינן גורמות נזקים חדשים ואינן מחריפות נזקים קיימים – בין שמדובר בנזקים פרטיים ובין שמדובר בנזקים קבוצתיים, לרבות נזקים לא מוחשיים לסביבות פוליטיות, חברתיות ותרבותיות;³⁷⁶ ואף אינן משפיעות לשלילה בכל דרך אחרת על בני אדם. הדבר מחייב שמירה על שלומם הנפשי והפיזי של בני אדם ועל כבודם.

מסמך המומחים נדרש גם להיבטים של בטיחות וביטחון ומציין כי מערכות נבונות צריכות להיות חסינות מבחינה טכנית – עליהן לפעול באופן בטוח, מאובטח ואמין, ואל להן לגרום לנזק בלתי מכוון.³⁷⁷ יש לתת את הדעת גם על

371 ראו Jobin, Ienca and Vayena, לעיל ה"ש 292, בעמ' 394. כ-71% מהמסמכים שהם דגמו מחייחסים לעיקרון של מניעת נזקים, לעומת 41% שמחייחסים לעיקרון של עשיית הטוב.

372 Fjeld et al., לעיל ה"ש 292, בעמ' 37. במונח ביטחון נדרשת ה-Datenschutzkommission (לעיל ה"ש 309, בעמ' 45) להיבטים של אבטח סייבר, הגנת מידע ובטיחות פיזית ונפשית, וכן להיחכנות של נזקים קינטיים-ממשיים בשל פעילות מערכות בינה מלאכותית.

373 כפי שמעריה Boddington, לעיל ה"ש 339, בעמ' 107. טיוטת כללי האחיקה של ממשלת אוסטרליה לבינה מלאכותית כללה את העיקרון של מניעת נזק: ARTIFICIAL INTELLIGENCE: AUSTRALIA'S ETHICS FRAMEWORK (a discussion paper, Australian Government (2019); אך הוא נעדר מכללי האחיקה הסופיים שפורסמו לאחר מכן. הכללים שפורסמו, לעומת זאת, כוללים עקרונות של בטיחות וביטחון: Australia's AI Ethics Principles, DEPARTMENT OF INDUSTRY, SCIENCE AND RESOURCES (2020).

374 ראו Beijing Academy of Artificial Intelligence, לעיל ה"ש 310.

375 AI HELG 2019, לעיל בה"ש 298, בעמ' 22.

376 שם, בעמ' 12.

377 ראו גם Google, לעיל ה"ש 317, עיקרון 3; Asilomar AI Principles, לעיל ה"ש 326, עיקרון 6; OECD 2019, לעיל ה"ש 294, בעמ' 8.

אבטחת סייבר – חוסנן של המערכות נמדד בין השאר בעמידותן בפני פצחנים המבקשים להשתמש בהן לרעה.³⁷⁸ דוח ועדת המשנה נדרש להיבטים של בטיחות וביטחון (הגנת סייבר ואבטחת מידע) בנפרד.³⁷⁹ כך גם מסמך ההמלצות היפני.³⁸⁰ מסמך מדיניות הרגולציה של משרד החדשנות כולל המלצה על עיקרון אתי דומה: "אמינות, עמידות, אבטחה ובטיחות".³⁸¹

3.4

אחריות ואחריותיות

שאלת האחריותיות של מערכות נבונות משיקה לשאלת האחריות של אלגוריתמים.³⁸² טכנולוגיות של למידת מכונה עלולות להגביר את רמת חוסר היכולת לצפות את פעולתן של מערכות אלו בהשוואה למערכות שמסתמכות על פיתוח אנושי,³⁸³ אך אלו ואלו מעלות שאלות זהות בדבר ייחוס אחריות להחלטותיהן.

מסמכי מדיניות של בינה מלאכותית מזכירים תדיר עקרונות של אחריות ואחריותיות.³⁸⁴ ואולם במקרים רבים אין הם מוגדרים די הצורך,³⁸⁵ ומסמכים

378 ראו למשל National Science and Technology Council, לעיל ה"ש 301, בעמ' 37-36; Datenethikkommission, לעיל ה"ש 309, בעמ' 45; US President, לעיל ה"ש 302, בס' (d)2.

379 ראו המיזם הלאומי למערכות נבונות (2019), לעיל ה"ש 313, בעמ' 21.

380 The Conference toward AI Network Society, לעיל ה"ש 312, בעמ' 7 (עקרונות מס' 4-5).

381 מסמך מדיניות הרגולציה של משרד החדשנות, לעיל ה"ש 314, בעמ' 107-108.

382 לשאלת האחריותיות של מערכות מחשב, ראו Helen Nissenbaum, *Accountability in a Computerized Society*, 2 SCIENCE & ENGINEERING ETHICS 25 (1996); Joshua A. Krolle et al., *Accountable Algorithms*, 165 U. PA. L. REV. 633 (2016-2017).

383 Paul B. de Laat, *Algorithmic Decision-Making Based on Machine Learning from Big Data: Can Transparency Restore Accountability?*, 31 PHILOSOPHY & TECHNOLOGY 525 (2018).

384 97% מהמסמכים שדגמו Fjeld et al., וכי-71% מהמסמכים שדגמו Jobin, Ienca and Vayena, (לעיל ה"ש 292), מתייחסים לאחריותיות או אחריות.

385 ראו הערתם של Jobin, Ienca and Vayena, לעיל ה"ש 292, בעמ' 394. FAT/ML, לעיל ה"ש 326, מציינים במפורש שהם נקטו לשון כללית ולא מוגדרת כדי לאפשר יישום רחב של העקרונות שהם מציעים.

שונים ממקמים אותם ברמות שונות בהיררכיה. כך למשל, דוח ועדת המשנה נדרש לרכיבים שונים של אחריות – שקיפות, הסברתיות ואחריות – בלי להגדיר אחריות מהי.³⁸⁶ לעומת זאת, פלורדי ועמיתיו רואים באחריותות רכיב של הסברתיות.³⁸⁷ יש גם הקושרים בין אחריותות לעקרון ההוגנות³⁸⁸ או לעקרון השקיפות.³⁸⁹

לפי עמדת הארגון (Fairness, Accountability, and Transparency in Machine Learning), נקודת המוצא של אתיקה של מערכות נבונות היא שאי-אפשר לתלות את האשמה לפעולות של אלגוריתמים באלגוריתמים, אלא רק בבני אדם.³⁹⁰ ערכים אתיים כגון אחריות, הסברתיות או הוגנות נועדו לסייע בזיהוי הגורמים שיש לייחס להם אחריותות לתוצאות ההחלטות של מערכות נבונות. אומנם על פי מסמכים אחרים אחריותות אינה העיקרון המארגן של אתיקה של מערכות נבונות, אך גם הם רואים בו את אחד העקרונות שלה ולרוב מצביעים על הצורך באחריותות אנושית לתוצאות של פעילות מערכות לכינה מלאכותית.

מסמך האיגוד למכונות חישוביות גורס כי על מוסדות להיות אחראים להחלטות המתקבלות על ידי האלגוריתמים שהם משתמשים בהם גם כשאיי-אפשר לספק הסבר לפעולתם.³⁹¹ הצהרת מונטריאול מזהירה מהפחתת האחריות האנושית עקב שימוש במערכות נבונות.³⁹² מסמך המלצות של משרד ההגנה של ארצות הברית קובע שעל בני אדם לקבל אחריות על הפיתוח של מערכות נבונות ועל תוצאות פעולתן.³⁹³ טיוטת ההנחיות למפתחי מערכות נבונות שפרסם המשרד

386 ראו המיזם הלאומי למערכות נבונות (2019), לעיל ה"ש 313, בעמ' 16–17.

387 Floridi et al., לעיל ה"ש 337, בעמ' 698–700.

388 AI HELG 2019, לעיל ה"ש 298, בעמ' 19.

389 ס' 50 להצהרת טורונטו, לעיל ה"ש 324.

390 FAT/ML, לעיל ה"ש 326.

391 ראו ACM 2017, לעיל ה"ש 326, עיקרון 3.

392 הצהרת מונטריאול, לעיל ה"ש 323, עיקרון 9: "אחריות: אל לו לפיתוח מערכת בינה מלאכותית להפחית מאחריותם של בני אדם במקום שבו יש לקבל החלטות".

393 Defense Innovation Board, לעיל ה"ש 342, עיקרון 1.

לענייני פנים ותקשורת של יפן מציינת כי המפתחים של מערכות נבונות נושאים באחריות (accountable) כלפי בעלי עניין שונים, לרבות המשתמשים בהן.³⁹⁴ הקווים המנחים של חברת דויטשה טלקום מכריזים בפשטות "אנחנו אחראים".³⁹⁵ ה-OECD, לעומת זאת, נמנע מלהצביע על הגורם שלו מיוחסת אחריות וקובעים כי "שחקני בינה מלאכותית", כלומר "מי שממלאים תפקיד פעיל במחזור החיים של מערכות בינה מלאכותית, לרבות פרישתן או תפעולן",³⁹⁶ צריכים לשאת באחריות לתפקודן התקין של מערכות בינה מלאכותית (בהתחשב בתפקידיהם, בהקשר ובחזית הידע).³⁹⁷

מסמך המומחים אינו מסביר מהם אחריות ואחריותיות, אלא מונה כמה תנאים להתקיימותן.³⁹⁸ יכולת לבקר את המערכת (auditability),³⁹⁹ מזעור השפעות שליליות ודיווח עליהן, גישה רציונלית לאילוצי שקלול תמורות (trade-offs) ביישום עקרונות אתיים, ואפשרות לסעד כאשר פעולת מערכות הבינה המלאכותית אינה תקינה.⁴⁰⁰

מסמך מדיניות הרגולציה של משרד החדשנות נמנע אף הוא מלבאר את המושג ולפיו "מפתחי בינה מלאכותית, מפעיליה או המשתמשים בה יגלו אחריות לתפקודה התקין, ולקיום העקרונות האחרים בפעילותם".⁴⁰¹ לפי מסמך זה, עקרון האחריות "נועד לבסס את קיומו של 'גורם אחראי' לבינה מלאכותית [...]".

- 394 The Conference toward AI Network Society, לעיל ה"ש 312, עיקרון 9.
- 395 Deutsche Telekom, לעיל ה"ש 345, בס' 1.
- 396 OECD 2019, לעיל ה"ש 294, בעמ' 7.
- 397 שם, בעמ' 8. ראו גישה דומה בהצהרה הכנס הבינלאומי השנתי של נציבי הגנת הפרטיות בעולם משנת 2018: DECLARATION ON ETHICS AND DATA PROTECTION IN ARTIFICIAL INTELLIGENCE, 4 International Conference of Data Protection and Privacy Commissioners (ICDPP) (2018).
- 398 AI HELG 2019, לעיל ה"ש 298, בעמ' 19-20.
- 399 שם. ראו גם ACM 2017, לעיל ה"ש 326, עיקרון 6.
- 400 לעוד מופעים של עקרון הסעד ראו לעיל ה"ש 317, עיקרון 2; FAT/ML, לעיל ה"ש 326; HRC 2018, לעיל ה"ש 270, בפס' 39-41.
- 401 מסמך מדיניות הרגולציה של משרד החדשנות, לעיל ה"ש 314, בעמ' 108.

מערכות בינה מלאכותית נוצרות בידי אדם, ולכן כברירת מחדל האדם הוא שאחראי להן.⁴⁰²

3.5 פרטיות

מופע ראשוני של רגולציית מערכות נבונות אפשר למצוא בדיני הגנת המידע האירופיים, הכוללים את הזכות שלא להיות מושא להחלטות אוטומטיות.⁴⁰³ הקשר בין פרטיות למערכות בינה מלאכותית, ובפרט למערכות למידה נבונה, המסתמכות על מידע רב, אינו מפתיע.⁴⁰⁴ ואכן, במסמכים רבים נידון נושא הפרטיות והצורך בשמירה עליה במסגרת הדיון בעקרונות האתיים הבסיסיים של הבינה המלאכותית.⁴⁰⁵ מסמך המומחים מונה שלושה היבטים של ההגנה על הפרטיות במערכות נבונות: שמירה על פרטיות והגנת מידע לאורך כל מחזור חיי מערכת הבינה המלאכותית; איכותו ושלמותו של המידע; והגבלת הרשאות הגישה למידע במערכות נבונות.⁴⁰⁶ לפי המתווה האמריקאי למגילת זכויות בבינה מלאכותית, יש לעצב ככל האפשר מערכות אוטומטיות לפי עקרונות של הנדסת פרטיות ולהתבסס על הסכמה של מושאי המידע; לפי עיקרון זה, אין לבצע מעקב רציף (continuous surveillance) במקרים שבהם טכנולוגיות מעקב עלולות להגביל זכויות, הזדמנויות או גישה.⁴⁰⁷

402 שם, שם.

403 ראו ס' 22 של התקנוח הכלליות בדבר הגנת מידע (GDPR), שאת שורשיו אפשר לאתר בס' 15 לדיקטיב הגנת המידע, Directive 95/46/EC of 24 October 1995 on the Protection of Individuals with regard to the Protection of Personal Data and on the Free Movement on such Data [1995] OJ L'281/31 הרשקוביץ ותהילה שוורץ אלטשולר הצעת חוק הגנת הפרטיות התשע"ט 2019-69-71 (המכון הישראלי לדמוקרטיה 2019); כהנא ושני, לעיל ה"ש 213, בעמ' 263-265. ברמת החקיקה הלאומית, כבר ב-1978 נדרשו דיני הגנת הפרטיות של צרפת לעיבוד אוטומטי - ראו ס' 2-3 אצל, Loi no. 78-17 du 6. janvier 1978 relative à l'informatique, à la protection des fichiers et aux libertés

404 ראו למשל HRC 2018, לעיל ה"ש 270, בפס' 34.

405 97% מהמסמכים שדגמו Fjeld et al., וכ-56% מהמסמכים שדגמו Jobin, Ienca and Vayena (לעיל ה"ש 292), מתייחסים לפרטיות.

406 AI HELG 2019, לעיל ה"ש 298, בעמ' 17.

407 AI BOR, לעיל ה"ש 305, בעמ' 30-39.

לעיתים פרטיות מוזכרת כערך שיש לכבד או לקדם⁴⁰⁸ ולעיתים כזכות שיש להגן עליה.⁴⁰⁹ יש מסמכים שקושרים בין פרטיות להגנת מידע;⁴¹⁰ ויש אחרים שקושרים בין פרטיות לשליטה ומציינים את זכותו של המשתמש לשלוט במידע שלו.⁴¹¹ מסמכים רבים מדגישים את החשיבות של הנדסת פרטיות (או "עיצוב לפרטיות", *privacy by design*), ושל פרטיות כברירת מחדל (*privacy by default*),⁴¹² כלומר הטמעה יזומה של עקרונות הגנת הפרטיות כבר בשלב התכנון של מערכות הבינה המלאכותית כפעולת מנע, ולא רק לאחר מעשה.⁴¹³ במסמך מדיניות הרגולציה של משרד החדשנות מוטמע עקרון הפרטיות בעיקרון שכותרתו "האדם במרכז". לפי עיקרון זה, "פיתוח בינה מלאכותית, או שימוש בה, ייעשו תוך כיבוד שלטון החוק, זכויות יסוד ואינטרסים ציבוריים, ובפרט תוך שמירה על כבוד האדם ופרטיותו".⁴¹⁴

כדי לטפח אמון בכינה מלאכותית, המסתמכת על נתוני עתק, יש להבטיח שמירה על פרטיות.⁴¹⁵ יש מתח בין פרטיות ובין שקיפות, שהוא עיקרון מרכזי

408 ראו Partnership on AI, לעיל ה"ש 320; *Responsible AI*, Microsoft, לעיל ה"ש 316.

409 ראו *Sony Group AI Ethics Guidelines*, לעיל ה"ש 319, עיקרון 4: *Asilomar AI*, לעיל ה"ש 326, עיקרון 13; *The Conference toward AI Network Society*, לעיל ה"ש 312, עיקרון 6; HRC 2018, לעיל ה"ש 270, בפס' 33-35.

410 ראו Cutler and Pribić, לעיל ה"ש 318, בעמ' 40: *Sony Group AI Ethics Guidelines*, לעיל ה"ש 319, עיקרון 4; ראו Google, לעיל ה"ש 317, עיקרון 5; Intel, לעיל ה"ש 315.

411 ראו Datenethikkommission, לעיל ה"ש 309, בעמ' 45; Cutler and Pribić, לעיל ה"ש 318, בעמ' 40: *Asilomar AI Principles*, לעיל ה"ש 326, עיקרון 12.

412 בשנת 2010 אומצו בכנס נציבי הגנת הפרטיות בעולם עקרונות הנדסת הפרטיות והפרטיות כברירת מחדל. גם הדין האירופי הכיר בעקרון הנדסת הפרטיות (ס' 25 לתקנות הכלליות בדבר הגנת מידע, GDPR). להרחבה ראו ארידור הרשקוביץ ושוורץ אלטשולר, לעיל ה"ש 403, בעמ' 89-90; ראו גם מיכאל בירנהק "פרטיות במשבר: הנדסה חוקתית והנדסת פרטיות" משפט וממשל כד 149 (2022), בחלק ה.

413 ראו Google, לעיל ה"ש 317, עיקרון 5; Intel, לעיל ה"ש 315, בעמ' 7-8; Datenethikkommission, לעיל ה"ש 309, בעמ' 123.

414 מסמך מדיניות הרגולציה של משרד החדשנות, לעיל ה"ש 314, בעמ' 104.

415 ראו למשל National Science and Technology Council, לעיל ה"ש 301, בעמ' 1(d); US President, לעיל ה"ש 302, בס' 1(d).

אחר של אתיקת הבינה המלאכותית. שקיפות אלגוריתמית עלולה לעיתים להתנגש בעקרונות של פרטיות, שכן בחינת תפקודן של מערכות בינה מלאכותית מחייבת פעמים רבות גישה לנתונים שאימנו את האלגוריתם.⁴¹⁶ כך למשל, מודל השפה של ChatGPT פותח בין השאר בהסתמך על בסיסי נתונים הכוללים מידע על אנשים פרטיים. גם אם מידע זה נלקח ממקורות גלויים ברשת הוא עלול, בשל הצטברותו, להסגיר מידע רגיש אגב שיחה עם המודל.⁴¹⁷ יתר על כן, הטמעת נתונים אלו במודל עלולה להתנגש בזכות להישכח (the right to be forgotten).⁴¹⁸

3.6

קידום הטוב

מקצת המקורות מונים את קידום הטוב (beneficence) עם התכליות האתיות של בינה מלאכותית.⁴¹⁹ קידום הטוב מופיע בראש

חמשת העקרונות שניסחו פלורידיו ועמיתיו,⁴²⁰ ומקורו בין השאר בכללי האתיקה הרפואית.⁴²¹ בטיוטות המוקדמות של מסמך המומחים נכלל עיקרון זה במפורש,⁴²² אך אחר כך הוא הושמט.⁴²³ עם זאת, מסמך המומחים בגרסתו הנוכחית מונה "רווחה חברתית וסביבתית" עם הדרכים למימוש בינה מלאכותית אמינה.⁴²⁴ גם

416 ראו לדוגמה ACM 2017, לעיל ה"ש 326, עיקרון 5; המיזם הלאומי למערכות נבונות (2019), לעיל ה"ש 313, בעמ' 19.

417 Melissa Heikkilä, *What Does GPT-3 "Know" About Me?* MIT TECHNOLOGY REVIEW (31.8.2022)

418 Eric Wallace et al., *Does GPT-2 Know Your Phone Number?* BAIR (20.12.2020)

419 כ-48% מהמסמכים שדגמו Jobin, Ienca and Vayena (לעיל ה"ש 292) מתייחסים לעקרון קידום הטוב (beneficence).

420 Floridi et al., לעיל ה"ש 337, בעמ' 697-698.

421 ראו לדוגמה BEAUCHANP AND CHILDRESS, לעיל ה"ש 366, בעמ' 202-248.

422 AI HELG 2018, לעיל ה"ש 298, בעמ' 8-9.

423 Andrea Renda, *Europe: Toward a Policy Framework for Trustworthy AI*, in THE OXFORD HANDBOOK OF ETHICS IN AI 651, 656 (2020)

424 AI HELG 2019, לעיל ה"ש 298, בעמ' 21.

הצהרת מונטריאול משתמשת במונח רווחה,⁴²⁵ ואילו עקרונות אסילומר ורוח הוועדה של בית הלורדים לנושא בינה מלאכותית נוקטים את המונח "טוב משותף".⁴²⁶ יש מסמכים שמדברים על קידום הטוב עבור "כל יצור תבוני",⁴²⁷ עבור "בני אדם רבים ככל האפשר"⁴²⁸ או לצורך "בניית חברה טובה יותר";⁴²⁹ ואילו מסמכים שמקורם בחברות פרטיות מדגישים לעיתים את קידום הטוב עבור לקוחותיהם.⁴³⁰

צמיחה בת קיימא היא העיקרון הראשון המוזכר בעקרונות האתיקה של מסמך מדיניות הרגולציה של משרד החדשנות. המסמך מעמיד בראש עקרונות האתיקה שהוא מציע את העיקרון של "בינה מלאכותית לקידום צמיחה, פיתוח בר קיימא ומובילות ישראלית בחדשנות". כאן למסגור הכלכלי של עקרון קידום הטוב יש גם היבט של חוסן לאומי ("יש לתת את הדעת לחשיבות של שימור עוצמה טכנולוגית ויכולת טכנולוגית").⁴³¹ עם זאת, ספק אם "מובילות ישראלית בחדשנות" היא אכן עיקרון אתי.⁴³²

- 425 הצהרת מונטריאול, לעיל ה"ש 323, עיקרון 1; OECD 2019, לעיל ה"ש 294, בעמ' 7.
- 426 House of Lords Select Committee on Artificial Intelligence, לעיל ה"ש 307, בפס' 417(1).
- עקרונות אסילומר גורסים כי יש לפתח חבונת-על (superintelligence) רק לאור אידאלים אתיים מקובלים (widely shared). ראו Asilomar AI Principles, לעיל ה"ש 326, עיקרון 23. בודינגטון מעירה, ובצדק, כי נוסח זה עמום וחלוי באופן שבו אנו מזהים "אידיאלים אתיים מקובלים". BODDINGTON, לעיל ה"ש 339, בעמ' 110.
- 427 הצהרת מונטריאול, לעיל ה"ש 323, עיקרון 1.
- 428 ראו Asilomar AI Principles, לעיל ה"ש 326, עיקרון 14 (לביקורת ראו BODDINGTON, לעיל ה"ש 339, בעמ' 109); Partnership on AI, לעיל ה"ש 320.
- 429 ראו Sony Group AI Ethics Guidelines, לעיל ה"ש 319, עיקרון 1; ראו גם GOOGLE, לעיל ה"ש 317, עיקרון 1.
- 430 ראו למשל DEUTSCHE TELEKOM, לעיל ה"ש 345, בס' 3.
- 431 מדיניות הרגולציה של משרד החדשנות, לעיל ה"ש 314, בעמ' 103-104.
- 432 נוסח העיקרון מזכיר את ס' 2(ב) לתזכיר חוק הגנת הסייבר, להלן ה"ש 461, שמונה עם מטרות מערך הסייבר את "קידום ישראל כמובילה עולמית בתחום הסייבר". הגדרת מובילות גלובלית בתחום טכנולוגי מסוים כמטרה מוסדית ועיקרון אתי גם יחד מעוררת תמיהה מסוימת, שכן עקרונות אתיים רכים או מטרות סטטוטוריות של גופים מוסדיים הם חלק מהמסגרת הפרשנית התוחמת את מרחב הסמכויות של הרגולטור.

3.7

**האדם במרכז,
חירות ואוטונומיה**

עקרונות של חירות ואוטונומיה, ובהם גם עקרונות המעמידים את האדם במרכז, נכללים בכמה ממסמכי המדיניות והאתיקה של בינה מלאכותית.⁴³³ לפי מסמכים אלו, בינה מלאכותית

צריכה לקדם את חירויות הפרט (כגון חופש הביטוי,⁴³⁴ הגדרה עצמית מידעית או אחרת⁴³⁵ ואת החופש מפני מניפולציה⁴³⁶), לשמור על חירויות אלו ולראוג לאוטונומיה שלהן.⁴³⁷ תמות של חירות ואוטונומיה מבטאות בין השאר גישה אנתרופוצנטרית, המעמידה את האדם במרכז.⁴³⁸ גישה זו באה לידי ביטוי בבקרה אנושית על פעולות אוטומטיות,⁴³⁹ במתן האפשרות שלא להשתמש במוצרים המבוססים עליה⁴⁴⁰ או בקריאה כללית לשליטה אנושית בטכנולוגיות.⁴⁴¹

אוטונומיה היא עיקרון יסוד באתיקה רפואית,⁴⁴² ואחד העקרונות המארגנים של אתיקת הבינה המלאכותית שניסחו פלורידו ועמיתיו.⁴⁴³ אוטונומיה היא זכותו

- 433 69% מהמסמכים שדגמו Fjeld et al. כוללים התייחסויות להיבטים של שליטה אנושית בטכנולוגיה. כ-40% מהמסמכים שדגמו Jobin, Ienca and Vayena מתייחסים לחירות ואוטונומיה (לעיל ה"ש 292).
- 434 PI&A19, לעיל ה"ש 325; HRC 2018, לעיל ה"ש 270.
- 435 AI HELG 2019, לעיל ה"ש 298, בעמ' 12; Datenethikkommission, לעיל ה"ש 309, בעמ' 43-44. להגדרה עצמית מידעית (informational self-determination) ראו (1983) BverfGE 65, 1, וכן כהנא ושני, לעיל ה"ש 213, בעמ' 187.
- 436 AI HELG 2019, לעיל ה"ש 298, בעמ' 10.
- 437 שם, בעמ' 12.
- 438 OECD 2019, לעיל ה"ש 294, בעמ' 7; ראו Asilomar AI Principles, לעיל ה"ש 326, עיקרון 11 (לביקורת ראו BODDINGTON, לעיל ה"ש 339, בעמ' 108); HRC 2018, לעיל ה"ש 270, בפס' 47-52.
- 439 ראו ס' 22 של התקנות הכלליות בדבר הגנת מידע (GDPR); AI HELG 2019, לעיל ה"ש 298, בעמ' 12; AI BOR, לעיל ה"ש 305, בעמ' 46-52.
- 440 House of Lords Select Committee on Artificial Intelligence, לעיל ה"ש 307, בעמ' 27.
- 441 Asilomar AI Principles, לעיל ה"ש 326, עיקרון 16 (לביקורת ראו BODDINGTON, לעיל ה"ש 339, בעמ' 109).
- 442 ראו לדוגמה BEAUCHAMP AND CHILDRESS, לעיל ה"ש 366, בעמ' 101-149. בדין הישראלי, ראו ע"א 2781/93 דעקה נ' בית החולים "כרמל", חיפה, פ"ד נג(4) 526, 573-571 (1999).
- 443 Floridi et al., לעיל ה"ש 337, בעמ' 697-698.

של הפרט להחליט עבור עצמו, ורוח ועדת המשנה אינו נדרש רק ליכולת הפיזית של אדם לבחור בין אפשרויות, אלא גם להיבטים אפיסטמולוגיים המאפשרים בחירה מושכלת.⁴⁴⁴ הצהרת מונטריאול קוראת לפתח מערכות בינה מלאכותית מתוך כיבוד האוטונומיה האנושית במטרה להגביר את יכולתם של פרטים לשלוט בחייהם.⁴⁴⁵ על פי המלצת רוח ועדת בית הלורדים, אין להעניק למערכות בינה מלאכותית את הזכות האוטונומית לבחור "לפגוע, להשמיד או להטעות בני אדם".⁴⁴⁶ מסמך המומחים של האיחוד האירופי רואה בכיבוד האוטונומיה של האדם וחירויותיו את אחת התכליות של זכויות היסוד שעליהן נכון האיחוד. משכך, על מערכות בינה מלאכותית לאפשר למשתמשים בהן לממש באופן רציף את זכותם להגדרה עצמית ולהשתתפות בתהליכים הדמוקרטיים.

עוד עולה ממסמכים שונים שבינה מלאכותית אמורה לקדם סולידריות אנושית⁴⁴⁷ או ערכים של קיימות ושל שמירה על הסביבה.⁴⁴⁸ הטמעתם של עקרונות אתיים בפיתוח, בפרישה, בהטמעה ובתפעול של מערכות נבונות נועדה לטפח אמון ציבורי במערכות אלו ולקדם את הטוב הכללי (תהיה הגדרתו אשר תהיה).

מסמך מדיניות הרגולציה של משרד החדשנות מונה עם המלצותיו לעקרונות אתיים את העיקרון "האדם במרכז – כיבוד זכויות יסוד ואינטרסים ציבוריים".⁴⁴⁹ לפי עיקרון זה, "בעת פיתוח ושימוש בטכנולוגיות חדשניות, האדם ויחסי הגומלין שלו עם החברה בה הוא פועל, מהווים נקודת מכוון נורמטיבית". בעיקרון זה כלולים, לפי מסמך זה, אינטרסים ציבוריים, עקרון שלטון החוק וזכויות יסוד (לרבות כבוד האדם ופרטיות). את כל אלו יש לכבד אגב פיתוח בינה מלאכותית ושימוש בה.

444 ראו המיזם הלאומי למערכות נבונות (2019), לעיל ה"ש 313, בעמ' 20.

445 הצהרת מונטריאול, לעיל ה"ש 323, עיקרון 2.

446 House of Lords Select Committee on Artificial Intelligence, לעיל ה"ש 307, בפס' 417(5). ראו גם להלן בסעיף 5.7.

447 Miguel Luengo-Oroz, *Solidarity Should Be a Core Ethical Principle of AI*, 1 NATURE MACHINE INTELLIGENCE 494 (2019)

448 AI HELG 2019, לעיל ה"ש 298, בעמ' 12; Floridi et al., לעיל ה"ש 337, בעמ' 697-696. ראו גם ס' 3(6) להצעת חוק הבינה המלאכותית הברזילאית, לעיל ה"ש 54.

449 מדיניות הרגולציה של משרד החדשנות, לעיל ה"ש 314, בעמ' 104-105.

הספרות הענפה על אתיקה של בינה מלאכותית, שנסקרה בפרק זה, צמחה מן התפיסה שטכנולוגיה חדשה זו איננה מאוסדרת, ואף אין הסכמה אוניברסלית בעניין דרכי אסדרתה, ולכן נכון ליצור עבורה מסגרת של כללי אתיקה לא פורמליים. לצד זאת עלתה גם הטענה שהיצמדות לכללי אתיקה היא בבחינת "טיוח אתי" (ethics washing), כלומר ניסיון ליצור מראית עין של אחריות אגב התחמקות ממחויבות מעשית לאחריות זו.⁴⁵⁰ מהלך זה כולל הימנעות מכוונת מהגדרות בהירות, מוחשיות ומפורטות של העקרונות האתיים; אי-קביעת מנגנוני אכיפה ברורים, שענישה בצידם, במקרה של אי-ציות לכללי האתיקה; היעדר התייחסות למגוון החוליות בשרשרת הערך של יצירת המערכות והמוצרים; והתמקדות בעיקר בחוליה של המהנדסים ושל כותבי הקוד.

לטיוח אתי כמה תכליות. תכלית אחת היא לחזק את אמון המשתמשים במוצרים מבוססי בינה מלאכותית (ואחת היא אם מדובר באמון בעניין בטיחות המוצרים, בהוגנות הליך הייצור שלהם או בהיבט אחר). תכלית אחרת היא לשחרר את התעשייה ממגבלות ולאפשר התחמקות מרגולציה ממשלתית או בינלאומית מחייבת. אחת התוצאות של הפרקטיקה הזאת היא האפשרות לבחור לאילו כללי אתיקה לציית ומאילו להתעלם.

אף על פי כן, ולזול יתר בכללי האתיקה (מה שמכונה לעיתים ethics bashing), מתוך תפיסה שטקסטים אלו אינם שווים את הנייר שעליו הם כתובים ואין ביכולתם להתמודד עם המציאות שבה מתפתחת הבינה מלאכותית, גם הוא אינו מוצדק. ראשית, כללי אתיקה יכולים לשמש גשר במצב הביניים שבין התפרצות של טכנולוגיה לבין יצירת רגולציה מדינתית או בינלאומית מותאמת. שנית, מחויבות מוסרית של עובדים בתוך תאגידים ושל תאגידים ללקוחות ולמשתמשים יכולה לשמש לעיתים שוט חזק לא פחות מרגולציה ממשלתית מחייבת במודל של ציווי ושליטה (command and control).⁴⁵¹

Meredith Whittaker et al., AI Now REPORT 2018 29–30 (AI Now, 2018) 450

Elettra Bietti, *From Ethics Washing to Ethics Bashing: A View on Tech Ethics from Within Moral Philosophy*, PROCEEDINGS OF THE 2020 CONFERENCE ON FAIRNESS, ACCOUNTABILITY, AND TRANSPARENCY (January 2020) 451

פרק רביעי

רגולציה של
בינה מלאכותית:
סקירה השוואתית

—

לצד מסמכי המדיניות שתוארו לעיל הגיעה בשנת 2021 גם סנונית ראשונה של רגולציית בינה מלאכותית – הצעת תקנות הבינה המלאכותית האירופיות.⁴⁵² אורסולה פון דר ליין, נשיאת הנציבות האירופית, הכריזה עם בחירתה ביולי 2019 שבתוך מאה ימים היא תפרסם הצעת חוק בתחום הבינה המלאכותית,⁴⁵³

452 ראו לעיל ה"ש 53.

453 *Von der Leyen's Real 100-day Challenge*, POLITICO (28.11.2019)

וכעת הגיעה השעה. בהמשך השנה פורסמו הצעות חקיקה לאסדרת תחום הבינה המלאכותית גם בסין⁴⁵⁴ ובברזיל⁴⁵⁵.

גם ארצות הברית אותתה על עניין רגולטורי בבינה מלאכותית בפרסום ההמלצות של ועדת הסחר הפדרלית (FTC) לשימוש הוגן בבינה מלאכותית.⁴⁵⁶ המלצות אלו מצביעות על כמה מקורות סטטוטוריים שמהם הרשות הפדרלית יכולה לינוק את סמכותה לאסור הטיות בבינה מלאכותית ואף פרקטיקות נאותות (שקיפות, הוגנות, הסברתיות, חסינות ואחריותיות).⁴⁵⁷

4.1

תקנות הבינה המלאכותית האירופיות

נוסח ההצעה של תקנות הבינה המלאכותית האירופיות, שעתיד להיכנס לתוקף בשנת 2024 וגם לאחר פרסומו הראשוני יהיה כפוף לשינויים,⁴⁵⁸ נועד להתמודד עם אחד האתגרים

המשמעותיים העומדים לפני מחוקק שניגש לאסדר טכנולוגיה – אי-הוודאות בנוגע להתפתחות העתידית שלה. לנוכח קצב התפתחות הטכנולוגיה ואיטיות המחוקק – שנדרש ללמוד את הטכנולוגיות החדשות, להעריך בזמן אמת מה תהיה השפעתן על הסדר הקיים ולגזור מהערכה זו דפוס אסדרה אופטימלי – אסטרטגיה אפשרית להתמודדות עם אתגרים אלו היא אסדרה "חסינת עתיד"

454 ראו להלן בסעיף 4.2.

455 ראו להלן בסעיף 4.3.

456 Elisa Jillson, *Aiming for Truth, Fairness, and Equity in your Company's Use of AI*, FEDERAL TRADE COMMISSION BUSINESS BLOG (19.4.2021)

457 ראו גם Andrew Smith, *Using Artificial Intelligence and Algorithms*, FEDERAL TRADE COMMISSION BUSINESS BLOG (8.4.2021)

458 בשלבי הבאח הספר לדפוס, במאי 2023, אכן פורסם נוסח מחוקן של הצעת תקנות הבינה המלאכותית האירופיות (ראו ההצעה המחוקנת, לעיל בה"ש 57). ההפניות להצעת תקנות הבינה המלאכותית האירופיות הן להצעה המקורית, אך במקומות שבהם היא עודכנה אנחנו מפנים להצעה המחוקנת. ביוני 2023 אושר הנוסח המחוקן בפרלמנט האירופי. ראו רפאל קאהאן "הפרלמנט האירופי אישר טיוטת חוק לרגולציה על בינה מלאכותית" Ynet (14.6.2023).

(future proof).⁴⁵⁹ אסדרה חסינת עתיד מאמצת דגמי חקיקה שאינם קשורים אך ורק לטכנולוגיה מסוימת אלא מבוססים על עקרונות-על גמישים שאפשר להתאימם להתפתחויות עתידיות. כאלה הם, למשל, חוק סמכויות החקירה של בריטניה⁴⁶⁰ או הדגם שהוצע בתזכיר חוק הגנת הסייבר בישראל.⁴⁶¹

נוסח התקנות האירופיות התבסס על מסמכים קודמים⁴⁶² ומטרתו לעצב מסגרת משפטית רגולטורית לבינה מלאכותית אמינה (trustworthy). הוא מגלם בתוכו זכויות יסוד אירופיות ונועד לאפשר לאנשים להשתמש בביטחון במערכות בינה מלאכותית, ובה בעת לעודד עסקים לפתח מערכות כאלה. בדברי ההסבר לנוסח התקנות האירופיות נכתב שמערכות בינה מלאכותית נועדו להיות כלי לקידום הטוב החברתי ולהביא לשגשוגם ורווחתם של בני אדם; משכך, הכללים החלים על מערכות נבונות צריכים להתרכז באדם.

רצוי להזכיר כי נוסף על נוסח התקנות האירופיות, שיידון להלן בהרחבה, גם התקנות הכלליות בדבר הגנת מידע (General Data Protection Regulation, GDPR), שנכנסו לתוקף בשנת 2018, כוללות הוראות הנוגעות למערכות אוטומטיות של קבלת החלטות.⁴⁶³ לפי מסמך זה, למושאי מידע מוקנית הזכות שלא להיות מושאן של החלטות המבוססות על עיבוד אוטומטי גרידא,

Sofia Ranchordas and Mattis Van 't Schip, *Future-Proofing* 459
Legislation for the Digital Age, in TIME, LAW AND CHANGE: AN INTERDISCIPLINARY
STUDY 347 (Yaniv Roznai and Sofia Ranchordas eds., 2020). לביקורת על גישה
חסינת עתיד לרגולציה של טכנולוגיה, בייחוד בהקשר של מעקב מקוון, ראו MARIA HELEN
MURPHY, *SURVEILLANCE AND THE LAW* 79-81 (2019); Paul Ohm, *The Argument against*
Technology-Neutral Surveillance Laws, 88 *TEX. L. REV.* 1685 (2010)

460 Investigatory Powers Act 2016, c.25 (Eng.). על חוק סמכויות החקירה של
בריטניה ראו כהנא ושני, *לעיל* ה"ש 213, בעמ' 138-186. על עיצובו בדגם חסין עתיד ראו
Murphy, *שם*, וכן Graham Smith, *Future-Proofing the Investigatory Powers*
Bill, *CYBERLEAGUE* (15.4.2016)

461 תזכיר חוק הגנת הסייבר ומערך הסייבר הלאומי, התשע"ח-2018 (להלן: תזכיר חוק
הגנת הסייבר); ראו גם רחל ארידור הרשקוביץ ותהילה שוורץ אלטשולר מהו סייבר? *חלק*
ב: *אחגרי האסדרה של הגנת הסייבר* 160-192 (2023).

462 ובהם AI HELG 2019, לעיל ה"ש 298; AI HELG 2018, לעיל ה"ש 298; AI4EU,
לעיל ה"ש 50.

463 ס' 22 של התקנות הכלליות בדבר הגנת מידע (GDPR).

שיש להן השלכות משפטיות משמעותיות. בדברי המבוא לתקנות הכלליות בדבר הגנת מידע נאמר כי "בכל מקרה יהיה העיבוד הנזכר כפוף לבקורות מתאימות, ובכללן [העברת] מידע ספציפי למושא המידע והזכות להתערבות אנושית בהחלטה; [למושא המידע] תינתן הזכות להביע את עמדתו, לקבל הסבר של ההחלטה שהתקבלה לאחר הערכה כזאת ולערער עליה".⁴⁶⁴

4.1.1. ניהול סיכונים כיוון שתקנות הבינה המלאכותית של אירופה

עוצבו כך שיהיו חסינות עתיד,⁴⁶⁵ הן נוקטות גישה של ניהול סיכונים ושומרות על ניטרליות טכנולוגית. למשל, מערכות בינה מלאכותית אינן מוגדרות בנוסח המקורי של התקנות המוצעות כבעלות זיקה לטכנולוגיה מסוימת, אלא כמערכות שפותחו באמצעות אחת או יותר מהגישות המפורטות בתוספת הראשונה לתקנות (תוספת שאפשר לתקן ולהרחיב בהתאם להתפתחויות טכנולוגיות עתידיות). לגישות הכלולות ברשימה יש יכולת לייצר פלט (תחזיות, המלצות, תוכן או החלטות) בהינתן מטרות שהגדירו בני אדם.⁴⁶⁶ פלורדי בירך על הגדרה זו, שמשקפת בעיניו גישה מפוכחת הרואה בבינה מלאכותית טכנולוגיה לפתרון בעיות (ולא למשל ישות תבונית של ממש, שעלולה לקום על יוצריה או להיחשב אישיות משפטית נפרדת).⁴⁶⁷ פטריק גלאונר, לעומת זאת, טען כי ההגדרה רחבה מדי, שכן לכאורה לפי הגדרה זו כל תוכנה שמשמשת בממוצע (פרוצדורה סטטיסטית) נחשבת לבינה מלאכותית.⁴⁶⁸

464 ס' 71 לדברי המבוא לתקנות הכלליות בדבר הגנת מידע (GDPR). ראו גם Gianclaudio Malgieri and Giovanni Comandé, *Why a Right to Legibility of Automated Decision-Making Exists in the General Data Protection Regulation*, 7 INT'L DATA PRIVACY L. 243, 248 (2017)

465 ראו דברי ההסבר, לעיל ה"ש 51, בעמ' 3, 12 ו-15. ראו גם המבוא לתקנות, בפס' 6, 71.

466 ס' 3(1) להצעת תקנות הבינה המלאכותית האירופיות, לעיל ה"ש 53. Luciano Floridi, *The European Legislation on AI: A Brief Analysis of its Philosophical Approach*, 34 PHILOSOPHY & TECH. 215 (2021)

468 Patrick Glauner, *An Assessment of the AI Regulation Proposed By the European Commission*, in THE FUTURE CIRCLE OF HEALTHCARE: AI, 3D PRINTING, LONGEVITY, ETHICS, AND UNCERTAINTY MITIGATION (Sepehr Ehsani, Patrick Glauner, Philipp Plugmann, and Florian M. Thieringer eds., 2022)

ההגדרה שבהצעה המתוקנת משככת את החששות שהביע גלאונר, שכן היא משמיטה את ההתייחסות לאופן הפיתוח של המערכת ושמה דגש על אוטונומיה. בינה מלאכותית היא "מערכת מבוססת מכונה שעוצבה כדי לפעול ברמות שונות של אוטונומיה והיא יכולה לייצר, למטרות מפורשות או לא מפורשות, פלטים כגון תחזיות, המלצות או החלטות המשפיעות על הסביבה הפיזית או הווירטואלית".⁴⁶⁹ הגדרה זו מתיישבת עם הגדרת ה-OECD למערכות בינה מלאכותית,⁴⁷⁰ אבל יש לתת את הדעת על העמימות שבנוסח "רמות שונות של אוטונומיה". פרשנות מרחיבה מדי של ההגדרה הזאת עלולה לכלול גם מערכות אוטומטיות פשוטות, כגון דלתות אוטומטיות מבוססות גלאי נפח.

גישת ניהול הסיכונים שבתקנות מחלקת את המערכות הנבונות לארבע רמות סיכון: מערכות אסורות, מערכות בסיכון גבוה, מערכות שחלה עליהן חובת שקיפות מיוחדת ומערכות בסיכון נמוך. יוער כי התקנות המוצעות לא יחולו על מערכות בשימוש צבאי בלעדי (דוגמת מערכות נשק אוטונומיות), גם אם רמת הסיכון בהן גבוהה.⁴⁷¹

4.1.1.1. מערכות אסורות

על פי נוסח התקנות האירופיות יש לאסור שימוש במערכות בינה מלאכותית שיש בהן סכנה קיצונית והן אינן עולות בקנה אחד עם ערכי היסוד האירופיים.⁴⁷² התקנות המוצעות אוסרות על מערכות שמשמשות בטכניקות תת-הכרתיות או מניפולטיביות כדי לעוות באופן מהותי את ההתנהגות האנושית באופן שמקשה במידה ניכרת על יכולתו של אדם לקבל החלטה מיודעת, וגורמים לו לקבל

469 ס' 3(1) להצעה המתוקנת, לעיל ה"ש 57.

470 על פי עקרונות ה-OECD מערכת בינה מלאכותית היא "מערכת מבוססת מכונה שיכולה, בהינתן קבוצת מטרות שהגדירו בן אנוש, לערוך תחזיות, לתת המלצות או לקבל החלטות המשפיעות על סביבות ממשיות או וירטואליות ברמות משתנות של אוטונומיה". ראו OECD 2019, לעיל ה"ש 294, בעמ' 3.

471 ס' 2(3) להצעת תקנות הבינה המלאכותית האירופיות, לעיל ה"ש 53.

472 שם, ס' 5(1). ראו גם Rostam Josef Neuwirth, *Prohibited Artificial Intelligence Practices in the Proposed EU Artificial Intelligence Act* (29.10.2022)

החלטה שאחרת לא היה מקבל, באופן שסביר שיסב לאדם או לקבוצת אנשים נזק פיזי או פסיכולוגי משמעותי, וזאת למעט מערכות בינה מלאכותית לתכליות תרפויטיות.⁴⁷³

עוד אוסרות התקנות המוצעות מערכות שמנצלות חולשות של קבוצות שונות באוכלוסייה (למשל חולשות שקשורות לגיל, ליכולות פיזיות או ליכולות קוגניטיביות) באופן שסביר שיסב נזקים פיזיים או פסיכולוגיים מהותיים;⁴⁷⁴ מערכות סיווג ביומטריות המסווגות אנשים בהסתמך על פרמטרים רגישים או באמצעות הסקה של פרמטרים רגישים (למעט מערכות שנועדו לתכליות תרפויטיות, בכפוף להסכמה מדעת של מי שנחשף אליהן);⁴⁷⁵ וכן מערכות נבונות שנועדו לבצע דירוג חברתי⁴⁷⁶ או לסווג או להעריך את מידת האמון שיש לתת באדם על בסיס התנהלותו החברתית או תכונות אופי שלו באופן שעלול להביא להחלטות לרעתם של יחידים או קבוצות בהקשרים חברתיים שאין להם זיקה להקשרים שבהם המידע נאסף או נוצר במקור;⁴⁷⁷ מערכות שיטור מונחה עתיד (predictive policing), כלומר מערכות המעריכות את הסיכון שאדם יבצע עברה פלילית או מינהלית (לרבות הערכת סיכויי רצידיביזם);⁴⁷⁸ מערכות המייצרות או מרחיבות מאגרי נתונים לזיהוי פנים בהסתמך על איסוף גורף מרשת האינטרנט או על חומר מצולם ממערכות טלוויזיה במעגל סגור;⁴⁷⁹ מערכות בינה מלאכותית שתכליתן לזהות רגשות לתכליות של אכיפת חוק; או ניהול מעברי גבול, או לשמש במוסדות חינוכיים ובהקשרים תעסוקתיים.⁴⁸⁰

473 ס' 5(1)(a) להצעה המתוקנת, לעיל ה"ש 57.

474 ס' 5(1)(b), שם.

475 ס' 5(1)(ba), שם.

476 ס' 3(45c), שם, מגדיר דירוג חברתי כהערכה או סיווג של אנשים בהסתמך על התנהגותם החברתית, על מעמדם החברתי-כלכלי או על תכונות האופי שלהם (הידועות או החזויות).

477 ס' 5(1)(I), שם.

478 ס' 5(1)(da), שם. ראו גם סעיף 2.3.2 לעיל.

479 ס' 5(1)(db), שם. נראה שסעיף זה נועד לאסור על מאגרי זיהוי פנים כגון אלו שמפתח חברת Clearview Ai. ראו גם היל, להלן ה"ש 1059.

480 ס' 5(1)(dc) להצעה המתוקנת, לעיל ה"ש 57. ראו גם סעיף 2.3.2 לעיל.

התקנות אוסרות גם שימוש במערכות זיהוי ביומטריות בזמן אמת במרחבים ציבוריים.⁴⁸¹

4.1.1.2. מערכות בסיכון גבוה

נוסח התקנות האירופיות מבחין בין שני סוגים עיקריים של מערכות בינה מלאכותית בסיכון גבוה: מערכות שמשמשות רכיב בטיחות כמוצר⁴⁸² ומערכות שלפעולתן יש השלכות על זכויות יסוד מתוך רשימה מפורשת,⁴⁸³ אם אלו מסכנות באופן משמעותי את בריאותם, בטיחותם או זכויותיהם הבסיסיות של אנשים.⁴⁸⁴

רשימה זו מפורטת בתוספת השלישית לתקנות וכוללת כ-20 קטגוריות של מערכות בינה מלאכותית: מערכות זיהוי ביומטריות; מערכות שנועדו לבקרה ובטיחות בתשתיות קריטיות; מערכות שמשמשות לקבלת מועמדים ללימודים או למתן ציונים; מערכות שמשמשות לגיוס עובדים ולהערכתם; מערכות דירוג אשראי; מערכות לקביעת זכאות לקבלת שירותים ציבוריים או לקבלת זכויות סוציאליות; מערכות לקביעת סדרי עדיפות בטיפול של שירותי חירום והצלה ציבוריים

481 ס' (d)(1)5, שם, אוסר לחלוטין פרישת מערכות אלו בלי להידרש לתכליתן, ואילו הנוסח המקורי של ס' (d)(1)5 להצעת תקנות הבינה המלאכותית האירופיות, לעיל ה"ש 53, קובע סייגים בעניין זה. מערכת זיהוי ביומטריות "בזמן אמת" אסורות, אלא אם כן הן נחוצות לחיפוש ממוקד אחר קורבנות פוטנציאליים מסוימים של פשע (כגון ילדים חטופים); למניעת סכנה ספציפית מוגדרת, משמעותית ומיידית לחייהם של אנשים או לשלומם; לסיכול מתקפת טרור; או לאיתורם וזיהוים של עבריינים בעבירות חמורות. ס' (2)5 לתקנות הוסיף סייגים לשימוש במערכות זיהוי ביומטריות גם במקרים חריגים אלו. לפי אותם סייגים יש להביא בחשבון את עוצמת הסכנה הפוטנציאלית ואת השלכות הרוחב החברתיות הטמונות בהפעלתן. כן נדרש להבטיח שכל הסתייעות פרטנית במערכת זיהוי ביומטרית תיעשה בכפוף לצו של גורם שיפוטי או גוף מינהלי עצמאי. על פיקוח אקט אנטה לצורך מעקב מקוון ראו עמיר כהנא ויובל שני פיקוח על מעקב מקוון בישראל 44-46 (מחקר מדיניות 149, המכון הישראלי לדמוקרטיה 2020).

482 לפי ס' (14)3 להצעת תקנות הבינה המלאכותית האירופיות, לעיל ה"ש 53, רכיב בטיחות במוצר או במערכת הוא רכיב של מוצר או של מערכת שיש לו פונקציה של בטיחות או שתקלה בו עלולה להיות סיכון בריאותי או בטיחותי לאנשים או לרכוש.

483 ס' 6 להצעת תקנות הבינה המלאכותית האירופיות, לעיל ה"ש 53.

484 ס' (2)6 להצעת המחוקקת, לעיל ה"ש 57. הנציבות הפרסם קווים מנחים המגדירים את הנסיבות שבהן פלט של מערכות בינה מלאכותית עלול להיות איום מסוג זה בתוך שישה חודשים מיום כניסת התקנות לתוקף.

(טריאז); מערכות שונות בשימוש רשויות אכיפת החוק (בהן מערכות לזיהוי דיפ-פייק,⁴⁸⁵ להערכת מצב נפשי, להערכת סיכון, לחיזוי רצידיביזם על בסיס פרופילים, לחיזוי פשיעה על בסיס פרופילים,⁴⁸⁶ לאנליטיקה קרימינולוגית); מערכות בשירות רשויות ההגירה והגבולות;⁴⁸⁷ וכן מערכות בשירות מערכת המשפט, שנועדו לסייע במחקר, בפרשנות עובדות וחוקים וביישום החוק במערך נסיבות מוגדר. אפשר להוסיף על כל אלו עוד מערכות בינה מלאכותית בסיכון גבוה, ובתנאי שרמת הסיכון הנשקפת מהן שקולה לזו של המערכות ברשימה או גבוהה מהן.⁴⁸⁸

נוסח התקנות האירופיות הנוגעות למערכות בסיכון גבוה מושתת על הדגם האירופי לבטיחות מוצרים.⁴⁸⁹ התקנות המוצעות מטילות על המערכות הללו חובות בעניין נתונים, משילות נתונים (data governance), כלומר היכולת

485 דיפ-פייק הוא כינוי לטכנולוגיה מבוטסת בינה מלאכותית המאפשרת יצירת סרטוני וידאו, תמונות או קטעי קול מזויפים כך שיראו אמיתיים. האסדרה המשפטית של טכנולוגיה זו עודנה בחיתוליה. בשנת 2019 התקבלה במדינת קליפורניה חקיקה האוסרת על שימוש בטכנולוגיות אלו כדי להונות מצביעים או לפגוע בפוליטיקאים במסגרת תעמולה פוליטית, וכן על שימוש בטכנולוגיות אלו כדי להשתיל בסרטונים פורנוגרפיים דמויות של אנשים שלא הביעו את הסכמתם לכך. ראו גם "קליפורניה יוצאת נגד ה-Deep Fakes, חוקקה שני חוקים נגד החופעה" כלכליסט (7.10.2019); תהילה שורץ אלטשולר ואיתי ברון פייק ניוז: הדור הבא (מאמר דעה, המכון הישראלי לדמוקרטיה 14.7.2019); הרשות להגנת הפרטיות "פרטיות ואבטחת מידע בשימוש בטכנולוגיות Deepfake (זיוף עמוק)" (28.8.2022); תב"כ 9-24 יש עתיד - בראשות יאיר לפיד נ' עמותת "כן לשלום" Anna Pesetski, *Deepfakes: A New Content Category for a* (18.1.2021) *Digital Age*, 29 Wm. & Mary Bill Rts. J. 503 (2020); Mamta Sareen, *Threats and Challenges by DeepFake Technology, in DEEPFAKES: CREATION, DETECTION, AND IMPACT* 100 (Loveleen Gaur ed., 2023)

486 ראו גם ס' 71 לדברי המבוא לתקנות הכלליות בדבר הגנת מידע (GDPR).

487 למערכות אלגוריתמיות בשירות ביקורת הגבולות באיחוד האירופי ראו לדוגמה Amanda Musco Eklund, *Frontex and "Algorithmic Discretion"* (Part I), *VERFASSUNGSBLOG* (10.9.2022); idem, *Frontex and "Algorithmic Discretion"* (Part II), *VERFASSUNGSBLOG* (10.9.2022)

488 ס' 7 להצעה הממוקנת, לעיל ה"ש 57.

489 Michael Veale and Frederik Zuiderveen Borgesius, *Demystifying the Draft EU Artificial Intelligence Act: Analysing the Good, the Bad, and the Unclear Elements of the Proposed Approach*, 4 *COMPUTER L. REV. INT'L* 97 (2021)

לדעת איזה מידע נמצא ברשותך ולשלוט בניהולו), תיעוד, יידוע, בקרה אנושית, חסינות, דיוק ובטיחות – הכול בהתבסס על עבודת קבוצת המומחים האירופית.⁴⁹⁰ עמידה בחובות אלו – מהן מהותיות ומהן פרוצדורליות – נבחנת במסגרת בחינת תאימות (conformity assessment),⁴⁹¹ שמשמשת תנאי לשחרורן של מערכות אלו לשוק האירופי. אחרי שמערכות אלו משוחררות לשוק הן נתונות לפיקוחן של רשויות פיקוח לאומיות (market surveillance authorities), שתפקידן לבחון אם מערכות אלו עדיין פועלות בהלימה לתקנות.⁴⁹²

לאורך כל מחזור החיים של מערכות בינה מלאכותית בסיכון גבוה יש להפעיל מערכת לניהול סיכונים. מערכת זו נועדה לזהות ולנתח סיכונים צפויים ולהעריך אותם, וכן לבחון מתווים צפויים של שימוש לרעה, להעריך סיכונים על בסיס נתונים ממערכות ניטור אחרות (post market),⁴⁹³ ולאמץ אמצעים נאותים לניהול סיכונים בהתאם להוראות התקנות.⁴⁹⁴ בדיקות של מערכות בינה מלאכותית בסיכון גבוה ייערכו לפני שחרורן לשוק.⁴⁹⁵

בהצעה האירופית מושם דגש מיוחד על תיעוד, מידע ונתונים מתוך הבנה שמשילות נתונים אפקטיבית הכרחית לשמירה על איכות גבוהה של מידע, שהיא עצמה תנאי הכרחי לתפקודן של מערכות בינה מלאכותית.⁴⁹⁶ יש להבטיח שנתוני אימון, תיקוף (validation) ובדיקה יהיו כפופים לפרקטיקות מתאימות של משילות נתונים, המביאות בחשבון היבטים מסוימים של עיצוב המערכת, איסוף מידע, פעולות הכנה של הנתונים, פערים בנתונים ופוטנציאל להטיות.⁴⁹⁷ בסיסי נתונים בשימוש מערכות בינה מלאכותית יביאו בחשבון היבטים מיוחדים

490 ראו AI HELG 2019, לעיל ה"ש 298.

491 ס' 16e, 19, 49 להצעת תקנות הבינה המלאכותית האירופיות, לעיל ה"ש 53.

492 ראו להלן, וכן ס' 62-65 להצעת תקנות הבינה המלאכותית האירופיות, לעיל ה"ש 53.

493 שם, ס' 61.

494 שם, ס' (2)9.

495 שם, ס' (7)9.

496 שם, פס' 44 למבוא.

497 שם, ס' (2)10.

הנוגעים לאופן פעולתן של מערכות בינה מלאכותית בסיכון גבוה.⁴⁹⁸ לצורך ניטור הטיות שליליות⁴⁹⁹ רשאים הספקים של מערכות אלו לעבד קטגוריות מיוחדות של מידע אישי, כגון מידע שיש בו כדי לזהות מוצא אתני או עמדות פוליטיות.⁵⁰⁰

יש להקפיד על תיעוד טכני של מערכות בינה מלאכותית בסיכון גבוה. על מלאכת התיעוד להתחיל בטרם שחרורן לשוק והתיעוד צריך להישאר מעודכן גם לאחר מכן.⁵⁰¹ התיעוד נועד לאפשר להדגים כי המערכת פועלת בהתאם להוראות התקנות.⁵⁰² נוסף על התיעוד הטכני נדרשות המערכות לתעד את פעולתן השוטפת (log) לצורכי נעקבות.⁵⁰³

על מערכות בינה מלאכותית בסיכון גבוה חלה גם חובת שקיפות. יש לעצב ולפתח מערכות שיהיו שקופות די הצורך לאפשר למשתמשים בהן לפרש את הפלט שלהן ולהשתמש בו כראוי.⁵⁰⁴ מנגנוני השקיפות יכללו בין השאר הנחיות רלוונטיות למשתמשים,⁵⁰⁵ שמקצתן מוגדרות בתקנות.⁵⁰⁶ פינק ופינק מעירות שחובה זו חלה על מפתחיהן של מערכות בינה מלאכותית בסיכון

498 שם, ס' 10(4).

499 על הטיות ראו להלן בפרק 6.

500 ס' 10(5) להצעת תקנות הבינה המלאכותית האירופיות, לעיל ה"ש 53. ס' 9 של התקנות הכלליות בדבר הגנת מידע (GDPR) אוסר על עיבוד סוגים שונים של מידע אישי רגיש שלא בהסכמת מושא הנחונים (או אם חריגים אחרים אינם מתקיימים). ס"ק (g)-(5) 10 להצעה המתוקנת, לעיל ה"ש 57, מונים תנאים מיוחדים שבכפוף להם יותר להשתמש במידע אישי רגיש למטרות זיהוי הטיות שליליות: אי־אפשר לזהות או לתקן את ההטיה באמצעות מידע סינחטי או מותמם; המידע עבר פסידו־אנונימיזציה; נוהלי האבטחה נאותים; המידע הרגיש לא יועבר לצד שלישי; המידע הרגיש נמחק לאחר תיקון ההטיה או עם תום תקופת השימור שלו.

501 ס' 11(1) להצעת תקנות הבינה המלאכותית האירופיות, שם. בהצעה המתוקנת, לעיל ה"ש 57, יש כמה הקלות בדרישות אלו לעסקים בינוניים וקטנים (SME), וכן לחברות הזנק.

502 שם, ס' 11(2).

503 שם, ס' 12.

504 שם, ס' 13(1).

505 שם, ס' 13(2).

506 שם, ס' 13(3).

גבוה ואין בה כדי להקנות זכויות או סעדים למשתמשי הקצה או למושאייהן של החלטות אוטומטיות. הן מבהירות כי אין בחובה זו כדי לייצר זכות להסבר – לא כלפי משתמשי הקצה ולא כלפי מושאייהן של החלטות מבוססות בינה מלאכותית; החובה נוגעת לפיתוח המערכות באופן שמאפשר הסבר, מבלי להורות עליו.⁵⁰⁷

על פי התקנות המוצעות, מערכות בסיכון גבוה יפותחו באופן שיאפשר פיקוח אנושי אפקטיבי לכל אורך מחזור החיים שלהן כדי למזער או למנוע סיכונים לבריאות, לבטיחות או לזכויות יסוד עקב שימוש סביר במערכות אלו (לרבות שימוש צפוי לרעה).⁵⁰⁸ בדרך זו תהיה למפקח האנושי האפשרות לעמוד על היכולות והמגבלות של המערכות המפוקחות, להיות מודע לנטייה להסתמך או להסתמך יתר על המידה על התוצרים שלהן, לפרש נכונה את הפלט שלהן, להחליט בנסיבות מסוימות שלא להשתמש בו (או להתעלם ממנו או להפוך את התוצאות שלהן), וכן להתערב במהלך פעולתן או לעצור אותה.⁵⁰⁹

עיצוב ופיתוח של מערכות בינה מלאכותית בסיכון גבוה נדרש לפי התקנות המוצעות כדי להבטיח שמערכות אלו יעמדו ברמה נאותה של דיוק, חוסן (robustness) ואבטחה.⁵¹⁰ בין השאר מוצע שמערכות בינה מלאכותית בסיכון גבוה שלמירת המכונה שלהן נמשכת גם לאחר שחרורן לשוק יפותחו באופן שיבטיח מיתון של משוברים המגבירים הטיות קיימות במערכת.⁵¹¹

מהחובות החלות על מערכות בינה מלאכותית בסיכון גבוה נגזרות החובות החלות על הספקים שלהן.⁵¹² ספקי מערכות אלו נדרשים בין השאר להקים

Melanie Fink and Michèle Finck, *Reasoned A(I)administration: 507 Explanation Requirements in EU Law and the Automation of Public Administration*, 47 Eu. L. Rev 376 (2022)

508 ס' (1) 14 להצעת תקנות הבינה המלאכותית האירופיות, לעיל בה"ש 53.

509 שם, ס' (4) 14.

510 שם, ס' (1) 15.

511 שם, ס' (3) 15.

512 שם, ס' 16-25.

מערכת בקרת איכות שתבטיח ציות להוראות התקנות,⁵¹³ לדאוג לתיעוד הטכני של המערכות,⁵¹⁴ לשמור על התיעוד המתמשך (log) של פעולתן,⁵¹⁵ להבטיח כי טרם שחרורן לשוק יעברו המערכות בחינת תאימות,⁵¹⁶ וכן לרשום מערכות אלו במרשם האירופי למערכות בינה מלאכותית בסיכון גבוה.⁵¹⁷

חובה מרכזית בתקנות המוצעות היא בחינת התאימות (conformity). בחינת התאימות יכולה להיעשות בהסתמך על בקרות פנימיות של הספק⁵¹⁸ או בהסתמך על הערכה של גורם ביקורת חיצוני מוסמך (notified body).⁵¹⁹ בחינת התאימות נועדה לברוק את הלימתן של מערכות בינה מלאכותית בסיכון גבוה להוראות התקנות וכוללת בחינה של מערכת בקרת האיכות,⁵²⁰ של התיעוד הטכני,⁵²¹ של התהליכי העיצוב והפיתוח ושל אופן ניטורן של המערכות לאחר ששוחררו לשוק. גורמי הביקורת החיצוניים ינפיקו תעודה שתוקפה לא יעלה על חמש שנים, המעידה על הלימת המערכת לתקנות.⁵²² יואנידיס וגוטסופולו מראים שהמונח תאימות, שמתייחס למערך מוגדר של הוראות, צר יותר

513 שם, ס' 16b, 17.

514 שם, ס' 16(c). ס' 50 מורה כי על תיעוד זה להיות נגיש לרשויות המוסמכות הלאומיות לתקופה של עשר שנים מיום שחרורה של מערכת בינה מלאכותית בסיכון גבוה לשוק.

515 שם, ס' 16(d), 20.

516 שם, ס' 16(e), 49.

517 שם, ס' 16(f), 51, 60. נוסף על הרישום על הספק להמציא לרשות הלאומית המוסמכת הצהרת תאימות אירופית (EU declaration of conformity, ראו ס' 48, וכן חוספה V לתקנות). על הספק לשמור מסמכים שונים, ובהם ההצהרה והתיעוד הטכני של המערכת ושל מערכת בקרת האיכות שלה, ולהעבירם לרשות הלאומית המוסמכת לתקופה של עשר שנים (ס' 50).

518 שם, ס' 43(1)(A). ראו גם ההוראות שבתוספת VI לתקנות.

519 שם, ס' 43(1)(B). ראו גם ההוראות שבתוספת VII לתקנות.

520 שם, ס' 16(b), 17. וכן ס' 2 בחוספה VI לתקנות; ס' 3 בחוספה VII לתקנות.

521 שם, ס' 16(c), 18. וכן ס' 3 בחוספה VI לתקנות; ס' 4 בחוספה VII לתקנות.

522 שם, ס' 44.

מהמונח החלופי ציות (compliance), שמתייחס לעמידה כללית בהוראות הרגולטוריות.⁵²³

לאחר שחרורן של מערכות בינה מלאכותית בסיכון גבוה לשוק נדרשים הספקים שלהן להקים ולתעד מערכת ניטור אחרת שתאסוף ותנתח באופן שיטתי נתונים רלוונטיים על ביצועיהן כדי לאפשר לספק לאמוד את הלימתן הרציפה להוראות התקנות.⁵²⁴ מערכת הניטור נדרשת לדווח לרשויות הפיקוח הלאומיות על השוק על כל תקרית חריגה⁵²⁵ במערכות אלו שעולה כדי הפרה של דינים אירופיים שנועדו להגן על זכויות אדם.⁵²⁶

בתביעות נזיקין אזרחיות נגד ספקי בינה מלאכותית, אי-עמידה בחובות החלות על ספקים על פי נוסח התקנות האירופיות עלולה להשפיע על הנטל הראייתי להראות שיש קשר סיבתי בין הפרת התקנות לבין הנזק שגרמה מערכת בסיכון גבוה.⁵²⁷

נוסח התקנות האירופיות מחיל חובות לא רק על ספקי מערכות בינה מלאכותית בסיכון גבוה, אלא גם על שחקנים אחרים באקוסיסטם, ובכלל זה על יבואנים, על מפיצים⁵²⁸ ואף על המשתמשים במערכות אלו. יבואנים של מערכות בינה מלאכותית בסיכון גבוה נדרשים להבטיח שמערכות אלו מתועדות כראוי ועברו

Nikolaos Ioannidis and Olga Gkotsopoulou, *The Palimpsest of Conformity Assessment in the Proposed Artificial Intelligence Act: A Critical Exploration of Related Terminology*, EUROPEAN LAW BLOG (2.7.2021)

524 ס' 61 להצעת תקנות הבינה המלאכותית האירופיות, לעיל ה"ש 53.

525 תקריות חריגות מוגדרות תקריות שהביאו או היו עלולות להביא, במישרין או בעקיפין, למוחו של אדם או לפגיעה חמורה בבריאותו, לשיבוש משמעותי של הניהול והתפעול של תשתיות קריטיות, לפגיעה בזכויות יסוד המוגנות בדין האירופי או לפגיעה משמעותית בסביבה או ברכוש. ראו ס' 3(44) להצעה המתוקנת, לעיל ה"ש 57.

526 ס' 62(1) להצעת תקנות הבינה המלאכותית האירופיות, לעיל ה"ש 53.

527 ס' 4(2) להצעה לדירקטיבה אירופית לאחריות בתחום הבינה המלאכותית, לעיל ה"ש 26.

528 ס' 27 להצעת תקנות הבינה המלאכותית האירופיות, לעיל ה"ש 53.

בחינת תאימות;⁵²⁹ ואם אין הלימה ביניהן ובין הוראות התקנות עליהם להימנע מהפצתן.⁵³⁰ כשספק, מפיץ או צד שלישי אחר משחררים לשוק מערכת בינה מלאכותית בסיכון גבוה תחת שמם או סימן מסחרי שלהם, או כשהם משנים את תכליתה המקורית של מערכת בינה מלאכותית שכבר קיימת בשוק, או כשהם משנים באופן משמעותי את המערכת עצמה, חלות עליהם החובות החלות על ספקי בינה מלאכותית בסיכון גבוה.⁵³¹

משתמשים של מערכות בינה מלאכותית בסיכון גבוה – למשל משתמשים מוסדיים המעבדים מידע אישי של רכבות לקוחות ויותר – נדרשים בין השאר להשתמש במערכות אלו בהתאם להוראות ותיעודן, להבטיח שהמידע שמוזן למערכות אלו רלוונטי לתכלית שהן מיועדות לה, לתעד ולנטר את פעילותן ולדווח לספק או למפיץ כשמתעורר חשש לסכנה של ממש.⁵³² כתביעות נזיקין אזרחיות, אי-עמידה בחובות החלות על משתמשים על פי נוסח התקנות האירופיות עלולה להשפיע על הנטל הראייתי לקיומו של קשר סיבתי בין הפרתן לבין הנזק שגורמת מערכת בסיכון גבוה.⁵³³

4.1.1.3. מערכות שחלה עליהן חובת שקיפות מיוחדת

נוסח התקנות האירופיות מגדיר כמה סוגים של מערכות בינה מלאכותית שחלה עליהן חובת שקיפות מיוחדת. כך למשל, מערכות שנמצאות במגע עם בני אדם, ולפי נוסח ההצעה המתוקנת גם מערכות בינה מלאכותית גנרטיביות,⁵³⁴ נדרשות להיות מעוצבות כך שבני אדם ידעו שהם באינטראקציה עם מכונה;⁵³⁵ יש ליידע

529 שם, ס' 26(1).

530 שם, ס' 26(2).

531 שם, ס' 28.

532 שם, ס' 29.

533 ס' 4(3) להצעה לדירקטיבה אירופית לאחריות בתחום הבינה המלאכותית, לעיל ה"ש 26.

534 ס' 28b(4)(a) להצעה המחוקנת, לעיל ה"ש 57.

535 ס' 52(1) להצעת תקנות הבינה המלאכותית האירופיות, לעיל ה"ש 53. השוואה AI HELG 2019, לעיל ה"ש 298, בעמ' 18; ס' 32 להצהרת טורונטו, לעיל ה"ש 324.

משתמשים במערכות לזיהוי רגשות אנושיים⁵³⁶ או לסיווג ביומטרי⁵³⁷ בנוגע לפעולת המערכת.⁵³⁸

נוסף על כך, מערכות שעושות מניפולציות בתצלומים, בקטעי קול או בסרטוני וידאו עד כדי יצירת תוכן מזויף (deep fake) נדרשות לגלות שהתוכן הוא תוצאה של שינוי או שהוא נוצר באופן מלאכותי. חובה זו לא תחול על מערכות לזיהוי עבריינות ומניעתה ואף לא במקרים שבהם הדבר נחוץ למימוש חופש הביטוי, האומנות והמדע כפי שהוא מוגדר באמנה האירופית לזכויות אדם וחירויות יסוד (Convention for the Protection of Human Rights (and Fundamental Freedoms, CFR)⁵³⁹.

אין בהחלת חובת השקיפות המיוחדת החלה על מערכות מסוימות כדי לגרוע מחובות אחרות שיכולות לחול על אותן מערכות אם אלו מסווגות כמערכות בסיכון גבוה.⁵⁴⁰

4.1.1.4. מערכות בסיכון נמוך

מערכות בינה מלאכותית בסיכון נמוך הן קטגוריה שירית, כלומר כל מערכת שאינה מערכת בסיכון גבוה. בנוסח התקנות האירופיות מוצע לעודד מפתחים

536 מערכות לזיהוי רגשות מוגדרות בס' 3(34) להצעת תקנות הבינה המלאכותית האירופיות, לעיל ה"ש 53, כמערכות בינה מלאכותית שתכליתן היא זיהוי או הסקה של רגשות, מחשבות, מצבי תודעה או כוונות של בני אדם או קבוצות על בסיס מידע ביומטרי. טכנולוגיה כזאת חוכל לשמש, למשל, לגילוי שקרים על בסיס ניתוח הבעות פנים, אבל נכון לכתיבת שורות אלו אין היא מניבה תוצאות מדויקות די הצורך. בשנת 2021 רמזה חברת הביטוח הדיגיטלית הישראלית למונייד בציוץ שפרסמה שהיא משתמשת ביכולות כאלה, אך דבריה עוררו סערה והחברה נסוגה מהציוץ. ראו שגיא כהן "חברת הביטוח למונייד אמרה שהיא משתמשת ב־AI כדי לזהות שקרנים – וחוללה סערה" *TheMarker* (27.5.2021).

537 מערכות לסיווג ביומטרי מוגדרות בס' 3(35) (שם) כמערכות בינה מלאכותית שפועלות על בסיס מידע ביומטרי ותכליתן סיווג בני אדם לקטגוריות ספציפיות, כגון מגדר, גיל, צבע שיער, צבע עיניים, קעקועים, מוצא אתני או נטייה מינית.

538 שם, ס' 52(2).

539 שם, ס' 52(3).

540 שם, ס' 52(4).

של מערכות בסיכון נמוך לאמץ מרצונם את החובות החלות על מערכות בסיכון גבוה. כן מוצע שהנציבות האירופית, מדינות חברות והוועד האירופי לבינה מלאכותית (EU AI Board) יעודדו פיתוח של תקן רצוני למערכות בינה מלאכותית בתעשייה, בשיתוף משתמשים ובעלי עניין אחרים.⁵⁴¹

4.1.2 "ארגז חול" רגולטורי
 לצד החשיבות של התקנות האירופיות המוצעות להגנה על זכויות יסוד יש גם חשש שרגולציית יתר תצנן את יצר חדשנות ותפגע אפוא בפיתוחים הטכנולוגיים.⁵⁴² התקנות מתמודדות עם אתגר זה באמצעות הסדרים המאפשרים הקמה של "ארגזי חול" רגולטוריים – סביבה מבוקרת לפיתוח, בדיקה ותיקוף של מערכות חדשניות טרם שחרורן לשוק.⁵⁴³

4.1.3 סנקציות
 בהצעת התקנות האירופיות מוצע להשית קנסות גדולים בגין הפרתן, ואלו האמירו עם פרסום נוסח ההצעה המתוקנת. לפי הצעה זו, אם תפר חברה את האיסור על פרישת מערכות בינה מלאכותית מסוכנות ושימוש בהן, יוטל עליה קנס מינהלי של עד 40 מיליון אירו או 7% ממחזור הפעילות העסקית שלה, הגבוה מהם.⁵⁴⁴ בגין הפרת חובות הנוגעות למשילות נתונים יוטל קנס מינהלי של עד 20 מיליון אירו או 4% ממחזור הפעילות העסקית של החברה.⁵⁴⁵ בגין הפרות אחרות יוטלו קנסות מינהליים של עד 10 מיליון אירו או 2% ממחזור הפעילות העסקית השנתי של החברה, הגבוה מהם.⁵⁴⁶

541 שם, ס' 69.

542 ראו לדוגמה, Gönenç Gürkaynak, İlay Yılmaz, and Güneş Haksever, *Stifling Artificial Intelligence: Human Perils*, 32 COMPUTER L. AND SEC. REV. 749 (2016).

543 ס' 53-54 להצעת תקנות הבינה המלאכותית האירופיות, לעיל ה"ש 53.

544 ס' (3) 71 להצעה המחוקנת, לעיל ה"ש 57.

545 שם, ס' (4) 71.

546 שם, ס' (5)-(4) 71.

4.1.4. אסדרת מודלי יסוד אומנם בין פרסום הצעת תקנות הבינה המלאכותית האירופיות באפריל 2021 לפרסום (Foundation models)

נוסח ההצעה המתוקנת במאי 2023 חלפו כשנתיים, אך נראה כי במונחים טכנולוגיים חלפו שנות דור. לא היה אפשר להתעלם מהפריצה של המערכות הגנרטיביות ומהתפוצה הרחבה שלהן, וההצעה המתוקנת כללה אפוא הוראות פרטניות בעניין מודלי יסוד.

מודל יסוד הוא "מודל בינה מלאכותית שאומן על בסיס נתונים רחב, תוכנן להפיק פלט כללי ואפשר להתאימו למגוון רב של משימות נפרדות".⁵⁴⁷ כדי להבהיר את מצבם המשפטי של ספקי מודלי יסוד, ושל ספקי מערכות בינה מלאכותית הנסמכות עליהם, נדרש הסדר פרטני.⁵⁴⁸ ההסדר בהצעה המתוקנת כולל הוראות החלות על מודלי יסוד בבינה מלאכותית גנרטיבית, המוגדרים "מודלי יסוד שמשמשים במערכות בינה מלאכותית כדי לחולל תוכן ברמות משתנות של אוטונומיה, ובכלל זה מלל מורכב, תמונות, צליל או וידאו".⁵⁴⁹

ספק של מודל יסוד נדרש להבטיח בטרם שחרורו לשוק שהמודל מציית להוראות החוק, בין שהוא מופץ כמודל עצמאי ובין שהוא מוטמע במערכת בינה מלאכותית או במוצר אחר. ההוראות חלות על כל ערוצי ההפצה – גם אם הגישה למודל חינמית (קוד פתוח).⁵⁵⁰ הספק נדרש להראות שננקטו אמצעים (עיצוב, בדיקה וניתוח) למזעור סיכונים לבריאות, לבטיחות, לזכויות יסוד, לסביבה, לרמוקרטיה ולשלטון החוק; שבסיסי הנתונים עובדו בהתאם לממשל נתונים המתאים למודלי יסוד ונבחנה מידת התאמתם של מקורות מידע, אפשרות קיומן של הטיות ודרכי מיתון אפשריות לפגיעה בזכויות;⁵⁵¹ שמודל היסוד עוצב ופותח באופן שמבטיח רמות מתאימות של תפקוד, צפיות, בטיחות, אבטחת מידע ויכולת תיקון לכל אורך מחזור החיים שלו, מתוך שימוש בתקנים סביבתיים מתאימים הממזערים צריכה אנרגטית;⁵⁵²

547 שם, ס' (dc)(1c)3.

548 שם, פס' 60g לדברי המבוא.

549 שם, ס' (4)28b.

550 שם, ס' (1)28b.

551 שם, ס' (b)(2)28b.

552 שם, ס' (E)-(I)28b(2).

שיש תיעוד טכני מקיף והנחיות ברורות לשימוש כדי לאפשר לספקי מערכות הבינה המלאכותית הנסמכות על מודל היסוד לציית לחובותיהן לפי תקנות הבינה המלאכותית לתקופה של עשור לאחר שחרור מודל היסוד לשוק;⁵⁵³ שהוקמה מערכת בקרת איכות שתבטיח ציות להוראות התקנות ותתעד אותן; וכן שמודל היסוד נרשם במרשם האירופי המתאים.⁵⁵⁴

לפי ההצעה המתוקנת, על ספקים של מודלי יסוד שמשמשים בבינה מלאכותית גנרטיבית יחולו חובות שקיפות מיוחדות.⁵⁵⁵ עליהם לאמן, לעצב ולפתח את מודל היסוד באופן שיבטיח הגנה נאותה מיצירת תוכן המפר את הדין האירופי,⁵⁵⁶ וכן לתעד את השימוש בנתוני אימון המוגנים בזכויות יוצרים.⁵⁵⁷

4.1.5. התקנות האירופיות: אסדרתן של טכנולוגיות מתעוררות אינה פשוטה, הערכה ראשונית

שכן יכולותיהן ודפוסי השימוש בהן יכולים להשתנות תוך כדי תהליך האסדרה. פלוריד ציין לחיוב שמנסחי ההצעה לא ראו לנגד עיניהם חזון דיסטופי של מערכות בינה מלאכותית המשמידות את האנושות, אלא נקטו גישה פרגמטית שלפיה טכנולוגיות של בינה מלאכותית אינן אלא מערכות סטטיסטיות לפתרון בעיות.⁵⁵⁸

ההגדרה הרחבה של מערכות בינה מלאכותית בנוסח התקנות האירופיות, גם בנוסח שבהצעה המתוקנת, מקשה עלינו להבחין בין בינה מלאכותית לאלגוריתמים, אך היות שבקרב המומחים אין תמימות דעים באשר להגדרה המדויקת של בינה מלאכותית,⁵⁵⁹ אפשר שאכן מוטב לאמץ הגדרה רחבה. חרף החשש מהחלת התקנות האירופיות על מערכות שהגדרתן כבינה מלאכותית גבולית או שנויה במחלוקת, יש לשים לב שהגישה שנוקטות התקנות האירופיות בעניין ניהול סיכונים אינה מחילה חובות כלשהן על מערכות בסיכון נמוך.

553 שם, ס' 28b(2)(c) וכן 28b(3).

554 שם, ס' 28b(2)(g).

555 שם, ס' 28b(4)(a), וכן ראו סעיף 4.1.1.3 לעיל.

556 שם, ס' 28b(4)(b).

557 שם, ס' 28b(4)(c).

558 Floridi, לעיל ה"ש 467, בעמ' 219.

559 ראו לעיל בסעיף 1.1.

הרגש שניתן בתקנות להיבטים של משילות נתונים וניהול מידע ממחיש את ההבחנה בין אלגוריתמים "פשוטים" לבינה מלאכותית, שסם החיים שלה הוא נתונים. עוד בטרם פרסום התקנות המוצעות הכירה החקיקה האירופית בקשר שבין נתונים לבינה מלאכותית. סעיף 22 של ה-GDPR, המקנה למושא המידע את הזכות שלא יתקבלו בעניינו החלטות על בסיס עיבוד אוטומטי של המידע האישי שלו, קושר בין אוטומציה לנתונים.⁵⁶⁰ ההקלות המוצעות בארגז החול הרגולטורי בסעיפים 53-55 של התקנות האירופיות (שעיקרן פטור מעקרון צמידות המטרה כדי לפתח מערכות בינה מלאכותית על בסיס נתונים שנאספו למטרות אחרות), גם הן מעידות על ההכרה שהמידע הכרחי לקידומן ולפיתוחן של מערכות נבונות.

אחת ממטרות התקנות האירופיות היא להציע פרקטיקות מומלצות לניהול מידע ולפיתוח מערכות בינה מלאכותית. על פי סעיף 10 של תקנות אלו, נתוני אימון, תיקוף ובריאות של מערכות בינה מלאכותית יהיו כפופים לפרקטיקות מתאימות של משילות נתונים.⁵⁶¹ שיעור הקנס המקסימלי המוטל בגין הפרות של סעיף 10 זה לשיעור המוטל בגין פרישת מערכות בינה מלאכותית אסורות ושימוש בהן, עדות נוספת לחשיבות שמנסחי התקנות האירופיות מייחסים למדיניות נתונים ולמשילות נתונים נאותות להגנה על זכויות יסוד.

הרצון להבטיח שהתקנות יהיו רלוונטיות גם לאחר שינויים עתידיים במערכות הבינה המלאכותית ייצר לעיתים ניסוחים עמומים שיהיה צורך להבהיר בעתיד. האם למשל אפשר להגדיר רשתות חברתיות (או למצער, את הבינה המלאכותית החולשת על ה"פיד") "מערכות שמפעילות טכניקות תת-סיפיות כדי לשנות התנהגות אנושית", שנמנות עם מערכות הבינה המלאכותית האסורות לשימוש?⁵⁶²

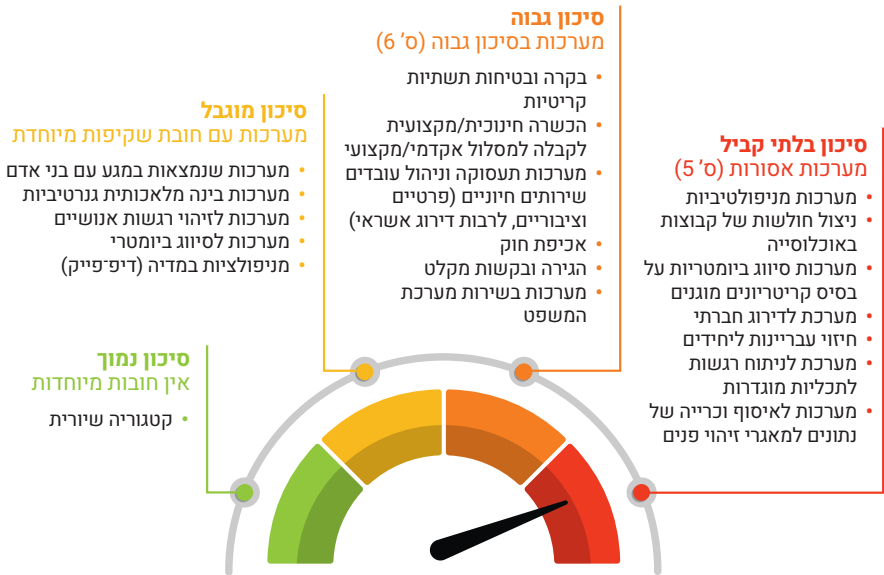
560 ראו לעיל ה"ש 403.

561 ראו לעיל בסעיף 4.1.1.2.

562 ס' (א) 5 להצעת תקנות הבינה המלאכותית האירופיות, לעיל ה"ש 53, Floridi, לעיל ה"ש 467, בעמ' 219, מביע חשש שמא עמימות זו תגדיר שלא לצורך מערכות בינה מלאכותית לחליליות תמימות כמערכות אסורות. גם תומס ברי ופרדריק פון-בוטמר סוברים כך. ראו Thomas Burri and Fredrik von Bothmer, *The New EU Legislation on Artificial Intelligence: A Primer* (21.4.2021)

הצורך לייצר ודאות משפטית בתחומים מסוימים שיש בעניינם שיח ציבורי ער (דירוג חברתי, דיפ־פייק וזיהוי פנים)⁵⁶³ תרם לנוסח המקורי של התקנות בעיקר באמצעות יצירת איים של החרגות – בפרט במה שנוגע לשימוש בטכנולוגיות אלו למטרות ביטחוניות ולאכיפת חוק.⁵⁶⁴ בנוסח ההצעה המתקנת המסתמן נראה שהחרגות אלו נעלמו, וכי לפנינו מדרג סדור יותר של רמות סיכון, כמתואר בתרשים 2.

תרשים 2 סיווג רמות הסיכון של מערכות בינה מלאכותית בתקנות הבינה המלאכותית האירופיות



563 ראו לדוגמה את השיח הער על מערכות לזיהוי פנים וחוקיותן באירופה ומחוץ לה. Joe Devansan, *EU Privacy Debate Rages Regarding Facial Recognition Firm Clearview AI*, T_HQ (1.6.2021); *Legality of Collecting Faces Online Challenged*, BBC News (27.5.2021); Amir Ali, *RCMP Violated Canadian Privacy Act with Facial Recognition Technology*, DHNews (10.6.2021)

564 ראו לדוגמה ההוראות המיוחדות לשימוש במערכות סיווג ביומטריות בידי גורמי אכיפת חוק שבס' (1)(d)6 להצעת תקנות הבינה המלאכותית האירופיות, לעיל ה"ש 53.

יש הסבורים למשל שהסיווג המוצע של מערכות בינה מלאכותית שחל עליהן איסור קטגורי, החשוב כשהוא לעצמו, אינו כולל את כל מערכות הבינה המלאכותית שראוי שיהיו אסורות, ובייחוד מערכות העוסקות במעקב המונים.⁵⁶⁵ עם זאת, כשבוע לפני פרסום התקנות האירופיות דלפה טיוטה שלהן, שכללה איסור מפורש על מעקב חסר הבחנה למטרות שאינן ביטחון לאומי.⁵⁶⁶ ההצעה המתוקנת הרחיבה את סוגי המערכות האסורות. כך למשל, הביקורת על סיווגן של מערכות לזיהוי רגשות כמערכות שחלה עליהן אך ורק חובת שקיפות מיוחדת, ולא כמערכות בסיכון גבוה,⁵⁶⁷ זכתה למענה מסוים בהצעה המתוקנת.⁵⁶⁸ כך גם הקולות שקראו שלא להסתפק באיסורים הפרטניים על שימוש במערכות זיהוי פנים במרחב הציבורי לצורך אכיפת חוק, ותחת זאת להחיל איסור גורף על כל שימוש במערכות זיהוי פנים במרחב הציבורי,⁵⁶⁹ לרבות על שחקנים פרטיים.⁵⁷⁰ ההצעה המתוקנת הרחיבה את האיסורים על שימוש במערכות ביומטריות בזמן אמת,⁵⁷¹ והוסיפה הגנות מפני סיווג ביומטרי על בסיס קריטריונים מוגנים.⁵⁷² ואולם מכל הדוגמאות האלה עולה שחסרה הגדרה כללית לסיכון קטגורי חלף הרשימה הסגורה או נוסף עליה.⁵⁷³

- 565 ראו לדוגמה Sarah Chander, *EU's AI Law Needs Major Changes to Prevent Discrimination and Mass Surveillance*, EDRI (28.4.2021)
- 566 James Vincent, *The EU Is Considering a Ban on AI For Mass Surveillance and Social Credit Scores*, THE VERGE (14.4.2021)
- 567 Gianclaudio Malgieri and Marcello Ienca, *The EU Regulates AI but Forgets to Protect our Mind*, EUROPEAN LAW BLOG (7.7.2021)
- 568 ס' 5(1)(dc) להצעה המחוקנת, לעיל ה"ש 57.
- 569 *EDPB & EDPS Call for Ban on Use of AI For Automated Recognition of Human Features in Publicly Accessible Spaces, and Some Other Uses of AI that Can Lead to Unfair Discrimination*, EUROPEAN DATA PROTECTION SUPERVISOR (21.6.2021)
- 570 *People, Risk and the Unique Requirements of AI*, ADA LOVELACE INSTITUTE (31.3.2022)
- 571 השוו את ס' 5(1)(d) להצעה המחוקנת, לעיל ה"ש 57, לס' (2)-(1)(d) להצעת תקנות הבינה המלאכותית האירופיות, לעיל ה"ש 53.
- 572 ס' 5(1)(ba) להצעה המחוקנת, לעיל ה"ש 53.
- 573 ראו למשל Lilian Edwards, *REGULATING AI IN EUROPE: FOUR PROBLEMS AND FOUR SOLUTIONS 3* (Expert opinion, Ada Lovelace Institute, 31.3.2022)

ייתכן גם שההתייחסות הפרטנית להיבטים אקטואליים של בינה מלאכותית עלולה להביא לידי רגולציה מבוססת סיכון, "רפלקס רגולטורי עקב סכנה" (risk regulation reflex). במקרים של נזק נדיר אך רחב היקף יש נטייה ציבורית לחוץ על הרגולטור להגיב בחומרה יתרה כדי למנוע את הישנותו.⁵⁷⁴

ביקורת נוספת על התקנות האירופיות נוגעת לבחירה המתעתעת במונח "משתמשים" לתיאור מי שפורשים או מפעילים מערכות בינה מלאכותית, אף שמונח זה מורה בדרך כלל על משתמשי הקצה של מערכות אלו. יש המעירים כי כך קל יותר לתת לאותם משתמשים פטור מהחובה לבצע הערכת תאימות.⁵⁷⁵

יתר על כן, יש בהבחנה בין "משתמשים" של מערכות בינה מלאכותית ל"ספקים" שלהן כדי לעודד את המגזר הציבורי להימנע מפיתוח פנימי של מערכות בינה מלאכותית ליישומי ממשל. כך למשל, אם רשויות הממשלה יפתחו מערכת לאיתור הונאות בתחום הרווחה דוגמת SyRI ההולנדית,⁵⁷⁶ יהיה להן מעמד של ספקי בינה מלאכותית.⁵⁷⁷ ובתור ספקים, יחולו עליהן החובות שבתקנות למיתון הסיכונים הפוטנציאליים של המערכת.⁵⁷⁸ לעומת זאת, אם יבחרו הרשויות בחלופה של מיקור חוץ, וירכשו מספק חיצוני מערכת מן המוכן, לא יחולו עליהן אלא החובות המוטלות על "משתמשים",⁵⁷⁹ כלומר הן יהיו רשאיות להפעיל שיקול דעת נרחב בהאצלת סמכויות הפיקוח לספקי המערכת.⁵⁸⁰

Nicolas Petit and Jerome De Cooman, *Models of Law and Regulation for AI*, in THE HANDBOOK OF AI 199, 205 (Anthony Elliott ed., 2022) 574

Martin Ebers et al., *The European Commission's Proposal for an Artificial Intelligence Act: A Critical Assessment by Members of the Robotics and AI Law Society (RAILS)*, 4 J 589 (2021) 575

ראו להלן בטעיף 5.2. 576

סי' 3(2), 28(1) להצעת תקנות הבינה המלאכותית האירופיות, לעיל ה"ש 53. 577

שם, סי' 16. 578

שם, סי' 29. 579

ראו בהקשר זה *How the EU's Flawed Artificial Intelligence Regulation Endangers the Social Safety Net: Questions and Answers*, HUMAN RIGHTS WATCH (10.11.2021) 580

יש הגורסים כי המגרעת העיקרית של התקנות טמונה בהתמקדותן בחובות החלות על תאגידים ולא בזכויות של בני אדם.⁵⁸¹

כדי להתגבר על הבלבול המושגי במונח "משתמשים" לא די להבחין בין משתמשי קצה (למשל רופא המשתמש במערכת דיאגנוסטיקה) למפעילים ריכוזיים של מערכות (למשל קופת החולים המנהלת את מערכת הדיאגנוסטיקה); יש להגדיר קטגוריה נוספת של "אנשים מושפעים" (affected persons) – בני אדם (או ישויות משפטיות) שמושפעים מפעילות מערכות הבינה המלאכותית.⁵⁸²

ההצעה של התקנות לאסדרת הבינה המלאכותית באירופה היא מהלך שאפתני, המזהה ריק רגולטורי כלל-עולמי ומנסה לייצר הובלה אירופית בתחום.⁵⁸³ לתקנות יש תחולה אקסטרטריטוריאלית,⁵⁸⁴ כלומר הן חלות גם על ספקים ממדינות מחוץ לאיחוד האירופי שמספקים מערכות בינה מלאכותית בתחומיו, וכן על ספקים ומשתמשים של מערכות בינה מלאכותית ממדינות שאינן שייכות לאיחוד האירופי אם הם משתמשים בפלט של מערכות אלו בתוך האיחוד. אקס-טריטוריאליות זו, המזכירה את "האימפריאליזם הרגולטורי" של ה-GDPR,⁵⁸⁵

581 ראו למשל Hannah Van Kolfschooten, *EU Regulation of Artificial Intelligence: Challenges for Patients' Rights*, 59 COMMON MARKET LAW REVIEW 81, 106 (2022).

582 *People, Risk and the Unique Requirements of AI*, לעיל ה"ש 570.

583 ואגליס פפונסטנטינו ופול דה הרט מעירים שהתקנות המוצעות הן דוגמה נוספת ל"ברוטליות" של דיני האיחוד, המזהים ואקום רגולטורי בתחומים טכנולוגיים וכופים את דין האיחוד על המדינות החברות בו בלי שהללו ניסו להתאים את הדין הלאומי שלהן למציאות הטכנולוגית החדשה. ראו Vagelis Papakonstantinou and Paul De Hert, *EU Lawmaking in the Artificial Intelligent Age: Act-ification, GDPR Mimesis, and Regulatory Brutality*, EUROPEAN LAW BLOG (8.7.2021).

584 ס' 2 להצעת תקנות הבינה המלאכותית האירופיות, לעיל ה"ש 53.

585 ראו Graham Greenleaf, *The "Brussels Effect" of the EU's "AI Act" on Data Privacy Outside Europe*, 171 PRIVACY LAWS & BUSINESS INTERNATIONAL REPORT He Li, Lu Yu, and Wu He, *The Impact of GDPR on Global Technology Development*, 22 J. OF GLOBAL INFORMATION TECH. MANAG'NT 1 (2019); Cedric Ryngaert and Mistale Taylor, *The GDPR As Global Data Protection Regulation?*, 114 AMERICAN J. INT'L. L. 5 (2020); Graham Greenleaf, *Global Data Privacy Laws 2021: Despite COVID Delays, 145 Laws Show GDPR Dominance*, 169 PRIVACY L. & BUS. INT'L REPORT 1 (2021); Huw

שתחולתו האקס-טריטוריאלית משפיעה על דיני הגנת הפרטיות במדינות רבות המעוניינות לסחור עם האיחוד האירופי,⁵⁸⁶ עשויה לתמרץ שחקנים רבים בתעשייה מחוץ לאיחוד האירופי לפעול בהתאם לדין האירופי ואף לעודד מדינות אחרות לפתח את דיני הבינה המלאכותית שלהם לפי המודל האירופי.⁵⁸⁷ מומחים מעריכים – עוד בטרם נחקקו התקנות – כי לכל הפחות לבחינת התאימות תהיה השפעה מחוץ לגבולות האיחוד.⁵⁸⁸

בשלהי חודש אוגוסט 2021 פרסמה מינהלת הסייבר של סין את טיוטת התקנות לניהול המלצות אלגוריתמיות במערכות מידע אינטרנטיות,⁵⁸⁹ ובמרץ 2022 נכנסו התקנות לתוקף.⁵⁹⁰ חקיקה זו היא חלק ממהלך אסדרה מקיף של ענף ההייטק שנקט

4.2 אסדרת בינה מלאכותית בסין

-
- Roberts et al., *Achieving a "Good AI Society": Comparing the Aims and Progress of the EU and the US*, 27(68) SCIENCE AND ENGINEERING ETHICS 7 (2021)
- 586 לרבות ישראל. ראו למשל עמיר כהנא "המטאטא יאה לנו – הכלי של השב"כ והחלטת הנאותות האירופית" **אתר משפט ועסקים** (29.3.2021); ארידור הרשקוביץ ושוורץ אלטשולר, **לעיל ה"ש** 403, בעמ' 8.
- 587 בהקשר זה ראו גם ANU BRADFORD, *THE BRUSSELS EFFECT: HOW THE EUROPEAN UNION RULES THE WORLD* 78–81, 131–169 (2020)
- 588 Charlotte Siegmann and Markus Anderljung, *The Brussels Effect and Artificial Intelligence: How EU Regulation Will Impact the Global AI Market*, CENTRE FOR THE GOVERNANCE OF AI (16.8.2022). לגישה הטבורה ש"אפיק בריסל" של הצעת תקנות הבינה המלאכותית האירופיות יהיה מתון ראו, Alex Engler, *The EU AI Act Will Have Global Impact, but a Limited Brussels Effect*, BROOKINGS (8.6.2022)
- 589 *Translation: Internet Information Service Algorithmic Recommendation Management Provisions (Opinion-Seeking Draft)*, DIGICHINA (14.10.2021)
- 590 *Translation: Internet Information Service Algorithmic Recommendation Management Provisions: Effective March 1, 2022*, DIGICHINA (10.1.2022) (להלן: תקנות מערכות ההמלצה האלגוריתמיות הסיניות). לרקע עליהן ולהרחבה בעניינן ראו Gilad Abiri and Xinyu Huang, *The People's (Republic) Algorithms*, 12 NOTRE DAME J. INT'L & COMP. L. 16 (2022)

הממשל הסיני כדי לבסס את כוחו מול הענף המתחזק.⁵⁹¹ התקנות הסיניות חלות על מערכות המלצה אלגוריתמיות מבוססות רשת המספקות תכנים למשתמשים – מערכות המלצה אישיות על תכנים, מערכות קבלת החלטות ומערכות חיפוש וסינון.⁵⁹²

התקנות מורות בין השאר לספקי שירותי מערכות המלצה אלגוריתמיות "לשמור על מוסר הציבור ועל אתיקה מסחרית ומקצועית", וגם לכבד עקרונות של הוגנות, צדק, פתיחות ושקיפות.⁵⁹³ הן מנחות לפעול לשימוש במערכות אלגוריתמיות לקידום הטוב ואוסרות להשתמש בהן לפעילויות שיש בהן משום פגיעה בביטחון הלאומי, שיבוש הסדר הכלכלי והחברתי או הפרה של כל דין.⁵⁹⁴ הן אף אוסרות על שימוש לרעה במידע ומניפולציה שלו במטרה להשפיע על דעת הקהל⁵⁹⁵ ועל הפצת ידיעות כוזב (fake news) באמצעות ספקי מערכות המלצה אלגוריתמיות.⁵⁹⁶ מערכות אלגוריתמיות המפירות את הנוהג והסדר הציבוריים, ובפרט כאלו המעודדות התמכרות של משתמשים לשימוש בהן או מקדמות צריכת יתר, אסורות במפורש.⁵⁹⁷ אסורות במיוחד מערכות המעודדות התמכרות כזאת בקרב קטינים.⁵⁹⁸ התקנות אף מחייבות ספקי שירותי המלצות אלגוריתמיות לוודא שמערכות שנועדו לשימוש קשישים (elderly) אינן מאפשרות הונאות והן נוחות לקהל היעד שלהם.⁵⁹⁹

591 ראו בהקשר זה Angela Huyue Zhang, *Agility over Stability: China's Great Reversal in Regulating the Platform Economy*, 63 HARV. INT'L. L. J. (2022); Chang Che, *China's "Big Tech Crackdown": A Guide*, THE CHINA PROJECT (2.8.2021).

592 ס' 2 לתקנות מערכות ההמלצה האלגוריתמיות הסיניות, לעיל ה"ש 590.

593 שם, ס' 4.

594 שם, ס' 6.

595 שם, ס' 14.

596 שם, ס' 13.

597 שם, ס' 8.

598 שם, ס' 18.

599 שם, ס' 19.

התקנות מפרטות שורה של הגנות על משתמשים, ובהן הזכות לעיין במידע עליהם, לערוך אותו ולקבל הסברים כשהם סבורים שפעולת המערכת השפיעה באופן מהותי על הזכויות והאינטרסים שלהם.⁶⁰⁰ על ספקי מערכות ההמלצה האלגוריתמיות חלה חובת שקיפות כללית ועליהם לתת פומבי לעקרונות הפעולה של המערכות.⁶⁰¹

התקנות מנחות את הספקים של מערכות ההמלצה האלגוריתמיות לנקוט אמצעים טכניים כדי למנוע מן האלגוריתמים שלהם להשפיע השפעה שלילית על המשתמשים או לעורר מחלוקות ושערוריות ציבוריות.⁶⁰² המערכות נדרשות לאפשר התערבות אנושית ומאפשרים למשתמשים בחירה אוטונומית, אך עליהן להציג בערוצי התוכן המרכזיים שלהם "תוכן התואם ערכים מקובלים".⁶⁰³

בתקנות מוצע להקים מרשם מרכזי של מערכות המלצה אלגוריתמיות במשרד הגנת הסייבר והמידע של סין.⁶⁰⁴ משרד זה ישמש רגולטור ויפקח על המערכות לצד משרדים רלוונטיים נוספים.⁶⁰⁵ אם יימצא שספק מערכות המלצה אלגוריתמיות הפר את הוראות התקנות ינחה המשרד את הספק כיצד לתקן את ההפרה. בנסיבות חמורות, או כשהספק מסרב לתקן את ההפרה, רשאי המשרד להקפיא את פעילות המערכת, לקנוס את הספק או לנקוט הליכים פליליים.⁶⁰⁶

שלא כמו התקנות האירופיות, החלות על מערכות בינה מלאכותית באשר הן, התקנות הסיניות מוגבלות ליישומי רשת של בינה מלאכותית. לפיכך אין הן חלות על מערכות ממשלתיות לקבלת החלטות, על מערכות זיהוי פנים, על רכבים אוטונומיים ועל יישומים אחרים של בינה מלאכותית. כמו כן, בכל הנוגע לניהול סיכונים התקנות הסיניות שיטתיות פחות מהאירופיות וככל הנראה

600 שם, ס' 17.

601 שם, ס' 16.

602 שם, ס' 12.

603 שם, ס' 11.

604 שם, ס' 23-25.

605 שם, ס' 25.

606 שם, ס' 31.

משליכות את יהבן על הערכת הסיכון הקונקרטי של הרגולטור, המסתמכת על "רקמה פתוחה" של סטנדרטים עמומים כגון פגיעה בביטחון הלאומי או בנוהג ובסדר הציבוריים.

רגולטור המרחב הקיברנטי של סין פרסם לאחרונה כללים ייעודיים לטכנולוגיות דיפ־פייק (סינתזה עמוקה).⁶⁰⁷ הכללים חלים על טכנולוגיות לשינוי טקסט, קול, מוזיקה ומאפיינים ביומטריים בתמונות ובווידאו, על שחזורים תלת־ממדיים וסימולציות ועל מי שמספק שירותי סינתזה עמוקה או משתמש בהם. הכללים הנזכרים מורים לספקים לייצר תוויות אזהרה בולטות שמציינות שמדובר בסינתזה עמוקה.⁶⁰⁸ חל איסור על הסרת תוויות אלו.⁶⁰⁹

4.3

הצעת חוק הבינה המלאכותית ברזיל

בספטמבר 2021 פורסמה בברזיל הצעת חוק לאסדרת בינה מלאכותית. שלא כמו התקנות האירופיות המפורטות, הצעת החוק בברזיל היא מסמך רזה למדי המפרט עקרונות כלליים לאסדרה של מערכות בינה מלאכותית.

החוק מונה כמה מטרות: פיתוח מדעי וטכנולוגי של ברזיל; תמרוץ פיתוח כלכלי בר־קיימא וקידום רווחת הכלל; קידום התחרות; השתלכות של ברזיל בשרשראות אספקה גלובליות; שיפור השירותים הציבוריים; והגנת הסביבה.⁶¹⁰ החוק מונה גם כמה יסודות לפיתוח וליישום בינה מלאכותית בברזיל: פיתוח מדעי וטכנולוגי, יוזמה חופשית, שמירה על אתיקה, זכויות אדם וערכים דמוקרטיים, חופש המחשבה והביטוי, היעדר אפליה, העדפה או תמרוץ של אסדרה עצמית, פרטיות

Provisions on the Administration of Deep Synthesis Internet 607
Information Services, CHINA LAW TRANSLATE (11.12.2022) (להלן: כללי הסינתזה
העמוקה). על דיפ־פייק ראו לעיל ה"ש 485.

608 שם, ס' 17. השוו גם לס' 52(3) להצעת תקנות הבינה המלאכותית האירופיות, לעיל
ה"ש 53.

609 ס' 18 לכללי הסינתזה העמוקה, לעיל ה"ש 607.

610 ס' 3 להצעת החוק הברזילאית לבינה מלאכותית, לעיל ה"ש 54.

אבטחת מידע וגישה למידע, ביטחון לאומי.⁶¹¹ לצד אלו הוא מפרט את העקרונות לפיתוח ויישום בינה מלאכותית כברזיל: קידום הטוב,⁶¹² האדם במרכז,⁶¹³ היעדר אפליה,⁶¹⁴ ניטרליות,⁶¹⁵ בטיחות,⁶¹⁶ אחריות,⁶¹⁷ זמינות מידע⁶¹⁸ ושקיפות.⁶¹⁹

ההגדרה לבינה מלאכותית בהצעת החוק הברזילאית מזכירה את ההגדרה בתקנות האירופיות.⁶²⁰ לפי הגדרה זו, בינה מלאכותית היא מערכת המבוססת על תהליך חישובי, שבהתאם למטרות שהגדירו בני אדם יכולה לעבד נתונים ומידע כדי ללמוד, לתפוס ולפרש סביבה חיצונית ולהגיב אליה, ועל יסוד זה לערוך תחזיות, להמליץ, לסווג או לקבל החלטות מתוך שימוש בטכניקות כגון למידת מכונה, מערכות לוגיות או גישות סטטיסטיות.⁶²¹ הצעת החוק מדגישה כי החוק אינו חל על תהליכי אוטומציה המבוססים על פרמטרים מוגדרים מראש שאינם כוללים יכולת של המערכת להגיב לסביבה החיצונית, לתפוס אותה או אף לפרש אותה.

בחוק מוצעים כמה קווים מנחים לאסדרה של בינה מלאכותית: (1) חקיקת משנה ספציפית לבינה מלאכותית תוצג רק כשיש בה צורך מוחלט;⁶²² (2) התערבות השלטון מחייבת אותו להביא בחשבון היבטים מגזריים ולהתבצע על ידי הרגולטור

611 שם, ס' 4.

612 שם, ס' (I) 5. ראו גם לעיל בסעיף 3.6.

613 שם, ס' (II) 5. ראו גם לעיל בסעיף 3.6.

614 שם, ס' (III) 5. ראו גם לעיל בסעיף 3.6.

615 שם, ס' (IV) 5.

616 שם, ס' (VI) 5. ראו גם לעיל סעיף 3.3.

617 שם, ס' (VII) 5.

618 שם, ס' (VIII) 5. לפי עקרון זמינות המידע המוצע, שימוש במידע, בבסיסי נתונים או בטקסט החוסים תחת הגנת זכויות יוצרים למטרות של אימון מערכות בינה מלאכותית אין משמעו הפרה של זכויות אלו, כל עוד המימוש הרגיל שלהן בידי בעליהן לא נפגע.

619 שם, ס' (V) 5.

620 ס' (1) 3 להצעת תקנות הבינה המלאכותית האירופיות, לעיל ה"ש 53.

621 ס' 2 להצעת החוק הברזילאית לבינה מלאכותית, לעיל ה"ש 54.

622 שם, ס' (1) 6.

המגזרי המתאים;⁶²³ (3) ניהול סיכונים: היקף ההתערבות המאסדרת ישמור על יחס מירתי לעוצמת הסיכון הקונקרטי הנשקף מכל מערכת ולהסתברות התממשותם של סיכונים אלו;⁶²⁴ (4) ביצוע הערכת השפעות של רגולציה חדשה (impact analysis): טרם אימוץ כללים רגולטוריים חדשים שיכולים להשפיע על הפיתוח והתפעול של מערכות בינה מלאכותית ייערך סקר הערכת השפעות רגולטוריות (regulatory impact analysis);⁶²⁵ (5) אחריות: יש להידרש לאחריותם של שחקנים בשרשרת הפיתוח וההפעלה של מערכות בינה מלאכותית מתוך התחשבות בנזקים הספציפיים שיש למנוע או לתקן, במידת המעורבות של גורמים אלו ובשאלה עד כמה השחקנים יכולים לציית במאמץ סביר לכללים הרלוונטיים בהלימה לסטנדרטים המקובלים בתעשייה ובעולם.⁶²⁶

על פי הערכה ראשונית נראה שהחקיקה הברזילאית המוצעת אינה אלא צעד ראשון לאסדרה של בינה מלאכותית, ולפי שעה אין היא מפורטת די הצורך.⁶²⁷ למשל, אין בה קווים מנחים לזיהוי מאפיינים של מערכות שהמחוקק עשוי להגדיר מסוכנות והיא מניחה את ההחלטות המעשיות לרגולטור או לתעשייה. כמו כן, אומנם יש בה העדפה כללית שלא להקים גופים ייעודיים חדשים, אבל אין היא מנסה לשרטט את תחומי הסמכויות של הרגולטורים הקיימים, לזהות תחומים שהרגולטורים הקיימים אינם מכסים או להציע פתרונות למצבים של חפיפה בין רגולטורים מגזריים. כיוון שהיא מבטאת רתיעה מאסדרה ריכוזית היא נמנעת מניהול ריכוזי של מרשם מערכות בינה מלאכותית בדומה לרישום הקיים בהצעה האירופית ובהצעה הסינית.

623 שם, ס' 6(2).

624 שם, ס' 6(3).

625 שם, ס' 6(4). השוו לסי' 3 להחלטה 2118 של הממשלה ה-33 "הפחתת הנטל הרגולטורי" (14.9.2014).

626 ס' 6(5) להצעת החוק הברזילאית לבינה מלאכותית, לעיל ה"ש 54.

627 לביקורת ברוח דומה ראו, Luca Belli, Yasmin Curzi, and Walter B. Gaspar, *AI Regulation in Brazil: Advancements, Flows, and Need to Learn from the Data Protection Experience* 48 COMP. L. & SEC. REV (2023)

4.4

אסדרת בינה מלאכותית
בקנדה

4.4.1. דירקטיבות ADM⁻ של קנדה
 בפברואר 2019 פרסמה ממשלת קנדה דירקטיבה בעניין מערכות אוטומטיות לקבלת החלטות.⁶²⁸ על פי השיטה הקנדית, דירקטיבות הן חלק מארגון הכלים של ועדת האוצר של קנדה (Treasury Board of Canada), שהיא ועדת קבינט של מועצת המלוכה של קנדה (Queen's Privy Council for Canada). ועדת האוצר של קנדה קובעת מדיניות מינהלית, פיננסית וארגונית חוצת ממשלה,⁶²⁹ בין השאר באמצעות דירקטיבות, שהן "הנחיות רשמיות המחייבות מחלקות לנקוט פעולות מסוימות או להימנע מהן. דירקטיבות מסבירות כיצד נדרשים עובדי ציבור בכירים [deputy heads] לממש את המטרות שהמדיניות מתווה".⁶³⁰ דירקטיבות קובעות חובה פנים-ממשלתית לציית ואין בכוחן לייצר זכויות בדין עבור אנשים פרטיים או ארגונים. הדירקטיבה הנוכרת נכנסה לתוקף בתחילת חודש אפריל 2019, אך החובה לציית להוראותיה החלה רק שנה לאחר מכן. מאפריל 2020 כל מערכת אוטומטית לקבלת החלטות שהממשלה הפרדלית של קנדה מפתחת או רוכשת כפופה לדירקטיבת ה-ADM של קנדה.⁶³¹

על פי הגדרת דירקטיבת ה-ADM של קנדה, מערכת אוטומטית לקבלת החלטות היא כל טכנולוגיה המסייעת לאדם בשיקול דעתו או מחליפה אותו.⁶³² מטרת הדירקטיבה למזער את הסיכונים הטמונים במערכות אוטומטיות לקבלת החלטות לאזרחים קנדיים ולמוסדות פדרליים ולהביא לקבלת החלטות מדויקות יותר,

Government of Canada, Directive on Automated Decision-Making 628
(להלן: דירקטיבת ADM⁻ של קנדה). (5.2.2019)

Teresa Scassa, *Administrative Law and the Governance of Automated Decision Making: A Critical Look at Canada's Directive on Automated Decision Making*, 54 U.B.C. L. Rev. 251, 268 (2021) 629

Government of Canada, FOUNDATION FRAMEWORK FOR TREASURY BOARD POLICIES (24.6.2008) שם. ראו גם 630

ס' 1 לדירקטיבת ה-ADM של קנדה, לעיל ה"ש 628. 631

שם, בחוספח הראשונה. 632

עקביות יותר, ניתנות לפירוש ויעילות בהתאם לדין הקנדי.⁶³³ התוצאות הצפויות של הדירקטיבה, לפי לשונה, הן החלטות מינהליות מונחות מידע השומרות על הוגנות והליך נאות;⁶³⁴ החלת הליכים מטרימים של הערכת סיכונים וצמצום ההשפעות השליליות של אלגוריתמים בעת זיהוים;⁶³⁵ והפיכת המידע והנתונים על השימוש במערכות אוטומטיות לקבלת החלטות במוסדות פדרליים לפומביים כשהדבר מתאפשר.⁶³⁶

הדירקטיבה אינה חלה על כל המערכות האוטומטיות לקבלת החלטות בשירות הממשל הפדרלי, אלא רק על מערכות שמספקות "שירות חיצוני"⁶³⁷ שמקבלו "חיצוני לממשלת קנדה".⁶³⁸ לא ברור למשל אם החלטה פנימית על הקצאת משאבים או אלגוריתמים לזיהוי הונאות נחשבים "שירות חיצוני" שניתן להצביע על אדם המקבל אותו. הדירקטיבה אינה חלה, בין השאר, על החלטות הנוגעות לתנאי ההעסקה של עובדי ציבור.⁶³⁹

נוסף על כך, הדירקטיבה אינה חלה על מערכות הנוגעות לביטחון לאומי,⁶⁴⁰ והיא חלה רק על מערכות שכבר נמצאות בשימוש הממשל הפדרלי (ולא על מערכות הנמצאות בסביבות בדיקה)⁶⁴¹ או פרויקטים ניסיוניים, אף שלא לה עשויה להיות השפעה ממשית).⁶⁴²

633 שם, ס' 4.1.

634 שם, ס' 4.2.1.

635 שם, ס' 4.2.2.

636 שם, ס' 4.2.3.

637 שם, ס' 5.1.

638 ראו Scassa, לעיל ה"ש 629, בעמ' 270.

639 שם, בעמ' 270–271.

640 שם, ס' 5.4.

641 שם, ס' 5.3.

642 ראו Scassa, לעיל ה"ש 629, בעמ' 271. לדוגמה לפרויקט ניסיוני שלא יכוסה על ידי הדירקטיבה ראו דוח על מערכות אוטומטיות לקבלת החלטות אוטומטיות בעניין פליטים ומהגרים בקנדה: Petra Molnar and Lex Gill, BOTS AT THE GATE: A HUMAN RIGHTS ANALYSIS OF AUTOMATED DECISION-MAKING IN CANADA'S IMMIGRATION AND REFUGEE SYSTEM (University of Toronto International Human Rights Program and Citizen Lab, September 2018)

הדירקטיבה מורה לבצע הערכת השפעה אלגוריתמית (algorithmic impact assessment) לפני פרישה של מערכת אוטומטית לקבלת החלטות.⁶⁴³ יש לעדכן את ההערכה בהתאם לשינויים בהיקף המערכת או בפונקציונליות שלה.⁶⁴⁴ על פי הוראות דירקטיבת הממשל הפתוח של קנדה, יש לפרסם את תוצאות הערכת ההשפעה האלגוריתמית באתרי הממשלה.⁶⁴⁵

באשר לשקיפות, הדירקטיבה מחילה על שימוש במערכות אוטומטיות לקבלת החלטות את החובה ליידע בטרם מעשה,⁶⁴⁶ את החובה לספק הסבר בעל משמעות לאחר מעשה⁶⁴⁷ ואת החובה לאפשר לצד שלישי לבדוק (audit) אם השימוש במערכות נעשה ברישיון⁶⁴⁸ או לפרסם את קוד המקור של המערכת אם הוא שייך לממשלת קנדה (בכפוף לסייגים).⁶⁴⁹

בסעיף "בקרת איכות" (Quality Assurance) מחייבת הדירקטיבה, בין השאר, לפתח הליכי בדיקה של המערכת האוטומטית לקבלת החלטות לפני פרישתה ולאחריה. מטרת הליכים אלו להבטיח שאין הטיות אלגוריתמיות לא רצויות, או תוצאות אחרות שאינן מכוונות, ולוודא שהמערכת מציינת לדירקטיבה ולהוראות הדין.⁶⁵⁰ הסעיף מבקש להבטיח שהמידע שהמערכת משתמשת בו מדויק, רלוונטי, עדכני והולם את חוק הפרטיות של קנדה⁶⁵¹ ואת מדיניות הממשלה בעניין

643 ס' 6.1.1 לדירקטיבת ה-ADM של קנדה, לעיל ה"ש 628. לפירוט על נוהלי הערכת ההשפעה האלגוריתמית בקנדה ראו Government of Canada, ALGORITHMIC IMPACT ASSESSMENT TOOL (updated 19.4.2022)

644 ס' 6.1.3 לדירקטיבת ה-ADM של קנדה, לעיל ה"ש 628.

645 שם, ס' 6.1.4.

646 שם, ס' 6.2.1.

647 שם, ס' 6.2.3.

648 שם, ס' 6.2.5.

649 שם, ס' 6.2.6. הסייגים לפרסום קוד המקור הם (1) רמת הסיכוי של הקוד; (2) פטור מפרסום על פי דיני חופש המידע; (3) פטור מאת המנהל הראשי של מערכות המידע של קנדה (Chief Information Officer of Canada, CIO).

650 שם, ס' 6.3.1, 6.3.2.

651 חוק הפרטיות הקנדי נוגע לזכויות שיש לאנשים בנוגע למידע של מוסדות ממשל קנדיים עליהם. Privacy Act, RSC 1985, c P-21

שירותים דיגיטליים.⁶⁵² עוד מורה הסעיף שהמערכת האוטומטית לקבלת החלטות תאפשר התערבות אנושית לפי הצורך (human in the loop).⁶⁵³ הסעיף כולל גם הוראות בנוגע להתייעצות עם מומחים,⁶⁵⁴ הכשרת עובדים⁶⁵⁵ ואבטחה.⁶⁵⁶

4.4.2. הצעת חוק הבינה המלאכותית והמידע של קנדה
בחודש יוני 2022 פרסמה הממשלה הפרלמנט של קנדה את הצעת החוק ליישום מגילת הזכויות הדיגיטליות (Digital Charter Implementation Act, 2022).⁶⁵⁷ הצעה זו כוללת חבילה של דברי חקיקה ובהם הצעת חוק הבינה המלאכותית והמידע (Artificial Intelligence and Data Act).

מטרת הצעת חוק הבינה המלאכותית והמידע היא לאסדר את הסחר במערכות בינה מלאכותית על ידי קביעת סטנדרטים מקובלים (התקפים בקנדה) לעיצובן ופיתוחן של מערכות אלו ולשימוש בהן.⁶⁵⁸ עוד מבקשת הצעת החוק לאסור על שימוש שעלול לגרום נזק חמור לאנשים פרטיים או לאינטרסים שלהם.⁶⁵⁹

הצעת החוק אינה חלה על מוסדות ממשלה או על מוצרים, שירותים ופעילות שנתונים לשליטתם של שר ההגנה, של מנהל שירות המודיעין הקנדי, של ראש

652 ס' 6.3.3 לדירקטיבה ה-ADM של קנדה, לעיל ה"ש 628.

653 שם, ס' 6.3.9. השוו למשל לס' 22 של התקנות הכלליות בדבר הגנת מידע (GDPR).

654 ס' 6.3.4 לדירקטיבה ה-ADM של קנדה, לעיל ה"ש 628.

655 שם, ס' 6.3.5.

656 שם, ס' 6.3.7.

657 Bill C-27, An Act to Enact the Consumer Privacy Protection Act, the Personal Information and Data Protection Tribunal Act and the Artificial Intelligence and Data Act and to make consequential and related amendments to other Acts, 1st Sess, 44th Parl, 2022 (First Reading, 16 June 2022)

658 ס' 4(a) להצעת חוק הבינה המלאכותית והמידע של קנדה, לעיל ה"ש 55.

659 שם, ס' 4(b).

הרשות לאבטחת תקשורת או של כל אדם אחר העומד בראש סוכנות פדרלית או מחוזית ונזכר בתקנות המושל הכללי של קנדה.⁶⁶⁰

הצעת החוק מגדירה "פעילות מאוסדרת" עיבוד או הנגשה של נתונים שנוגעים לפעילות אנושית במטרה לעצב ולפתח מערכת בינה מלאכותית או להשתמש בה; וכן עיצוב, פיתוח או הנגשה של מערכות בינה מלאכותית או ניהולן.⁶⁶¹ עוד היא מציעה שאדם שעוסק בפעילות מאוסדרת שנוגעת למידע מותמם יידרש לנקוט אמצעים באשר לאופן ההתממה של המידע וניהולו,⁶⁶² ובמידת הצורך ישמור תיעוד של אמצעים אלו.⁶⁶³

אף שהצעת החוק נוקטת גישה של ניהול סיכונים, ומחילה חובות מיוחדות (ראו להלן) על מערכות בינה מלאכותית בסיכון גבוה (high-impact system), הגדרתן נותרה לוטה בערפל: החוק המגדיר מערכות אלו מפנה לקריטריונים שיוגדרו בתקנות (שטרם תוקנו).⁶⁶⁴ עם זאת, אדם הממונה על מערכת בינה מלאכותית יידרש להעריך אם מערכת מסוימת נחשבת על פי התקנות למערכת בסיכון גבוה⁶⁶⁵ ולנקוט אמצעים שתכליתם לזהות, להעריך ולמתן את הסיכונים, את האיומים או את הפלט המוטא (biased output)⁶⁶⁶ שיכולים לצמוח משימוש במערכת.⁶⁶⁷

660 שם, ס' 3.

661 שם, ס' 5.

662 שם, ס' 6.

663 שם, ס' 10.

664 שם, ס' 5.

665 שם, ס' 7.

666 ס' 5 להצעת חוק הבינה המלאכותית והמידע של קנדה (לעיל ה"ש 55) מגדיר פלט מוטא כחוכן, החלטה, תחזית או המלצה שנעשו או נוצרו באמצעות מערכת בינה מלאכותית המפלה לרעה במישרין או בעקיפין ללא צידוק. ההגדרה מפנה לעילות ההפניה בס' 3 לחוק זכויות האדם של קנדה.

667 ס' 8 להצעת חוק הבינה המלאכותית והמידע של קנדה, לעיל ה"ש 55.

על הממונים על מערכות בסיכון גבוה חלה גם החובה לנטר את האמצעים למיתון הסיכונים.⁶⁶⁸ כמו כן, עליהם לתעד את האמצעים האלו ואת הסיבות התומכות בהערכת המסוכנות של המערכת.⁶⁶⁹ הם נדרשים לפרסם ברשת הסבר של המערכת בלשון פשוטה, הכולל התייחסות לאופן פעולתה, לסוגי התוכן, התחזיות או ההמלצות שהיא מייצרת, פירוט האמצעים שהיא נוקטת למיתון הסיכונים ומידע נוסף שקובעות התקנות.

המושג הכללי ימנה את השר הממונה על מערכות בינה מלאכותית.⁶⁷⁰ השר מוסמך לדרוש מסמכים ותיעוד בעניין מערכות בינה מלאכותית, ובפרט מערכות בסיכון גבוה,⁶⁷¹ וכן להורות על ביצוע ביקורת (audit) שתדווח לשר.⁶⁷² הממונה על מערכות בינה מלאכותית בסיכון גבוה נדרש להודיע לשר אם שימוש במערכת גרם או עלול לגרום לנזק מהותי.⁶⁷³ בעקבות ממצאי הביקורת רשאי השר להורות על הטמעת אמצעים שונים במערכות בינה מלאכותית.⁶⁷⁴ השר אף רשאי להורות על הפסקת פעילותן של מערכות בינה מלאכותית בסיכון גבוה, אם הוא סבור שיש בשימוש בהן סיכון של ממש.⁶⁷⁵ לשר מוקנה שיקול דעת להורות לפרסם ברשת כל מידע על הביקורות של מערכות בינה מלאכותית, על הפסקת פעילות מערכות בינה מלאכותית, על אמצעי מיתון סיכונים, על הערכת סיכונים, על התממה וכדומה.⁶⁷⁶

הצעת חוק הבינה המלאכותית של קנדה אינה מתעלמת מהיבטים הנוגעים לסודות מסחריים במידע שעשוי להגיע לידי השר מכוח החוק ולכן היא מציינת

668 שם, ס' 9.

669 שם, ס' 10.

670 שם, ס' 31. ס' 5 לחוק מורה כי כל עוד לא מונה שר כאמור, השר הממונה יהיה שר התעשייה.

671 שם, ס' 13-14.

672 שם, ס' 15.

673 שם, ס' 12.

674 שם, ס' 16.

675 שם, ס' 17.

676 שם, ס' 18.

שאינן השר רשאי לחשוף סודות מסחריים שהגיעו לידי⁶⁷⁷ ושהוא מחויב לשמור עליהם.⁶⁷⁸ בהצעת החוק נכללו הוראות מיוחדות בנוגע להעברת מידע שהוא בגדר סוד מסחרי לצד שלישי⁶⁷⁹ ובנוגע לפרסומו ברבים.⁶⁸⁰

4.5

ניצני אסדרת בינה מלאכותית בארצות הברית

יש כמה הצעות חוק שיכולות להעיד על תחילתה של מגמה לאסדר בינה מלאכותית בארצות הברית. במדינת קליפורניה, למשל, הוצע לאחרונה לאפשר להורים לתבוע פלטפורמות

שגרמו לילדיהם התמכרות חמורה (שבצידה נזקים נפשיים, פיזיים, רגשיים או התפתחותיים).⁶⁸¹ הצעה זו מעלה על הדעת את החשש שהביעו המחוקקים בסין מפני תופעות של התמכרות לפלטפורמות ברשת;⁶⁸² התמכרות שמטפחים, בין השאר, האלגוריתמים של השירותים המקוונים. ולקוד המוניציפלי המינהלי של העיר ניו יורק התווספה לאחרונה הוראה שקובעת ששימוש במערכות מבוססות אלגוריתמים לקבלת החלטות או המלצות על העסקתם או קידומם של עובדים יותנה בהליכי בקרת הטיות (bias audit) עצמאיים.⁶⁸³

גם ברמה הפדרלית החלו להידרש לצורך ברגולציה של בינה מלאכותית. טיוטת הצעת חוק הפרטיות והגנת המידע האמריקאי (American Data Privacy and Protection Act), שפורסמה בתחילת חודש יוני 2022,⁶⁸⁴ כוללת איסור

677 שם, ס' 22.

678 שם, ס' 23.

679 שם, ס' 24-26.

680 שם, ס' 27-28.

681 An Act to Add Section 1714.48 to Amend Section 1714 of the Civil Code, relating to Social Media Platforms, Cal. Assemb. A. 2408 (2021-2022)

682 ראו לעיל בטעיף 4.2.

683 N.Y.C. Admin. Code §§ - 0-870 - 20-874 (2022); Airlie Hilliard et al., *Regulating the Robots: NYC Mandates Bias Audits for AI-Driven Employment Decisions* (13.4.2022)

684 To provide consumers with foundational data privacy rights, create strong oversight mechanisms, and establish meaningful enforcement, (להלן: ADDPA) H.R.8152, 117th Congress (2022).

על אפליה אלגוריתמית במה שנוגע להגנה על זכויות האזרח.⁶⁸⁵ בפרט מוצע לאסור על "ישויות מכוסות" (covered entities)⁶⁸⁶ לאסוף, לעבד או להעביר מידע באופן שיש בו כדי להפלות על בסיס גזע, צבע, דת, מוצא לאומי, מגדר, נטייה מינית או נכות.⁶⁸⁷ נוסף על כך, על פי הצעת חוק הפרטיות והגנת המידע האמריקאי גופים המחזיקים במידע רב (large data holders)⁶⁸⁸ יערכו הערכת השפעה אלגוריתמית, כלומר יתארו את האמצעים שינקטו כדי לצמצם נזקים שעלולים להיגרם לאנשים פרטיים בנסיבות שונות.⁶⁸⁹

יוזמה נוספת היא הצעת חוק האחריותיות האלגוריתמית, שהגישה ב־2019 קבוצה של סנאטורים דמוקרטים והיא עודכנה בפברואר 2022.⁶⁹⁰ הצעת החוק חלה על כל אדם או ישות משפטית שכפופים לסמכותה של נציבות הסחר האמריקאית; שמשמשים או מחזיקים במערכות החלטה קריטיות מוגברות;⁶⁹¹ שהייתה להם הכנסה שנתית ממוצעת של יותר מ־50 מיליון דולר בשלוש השנים

685 שם, ס' 207.

686 ס' 2(9) ל־ADPPA, שם, מגדיר ישויות מכוסות (covered entities) כל גוף או אדם שאוסף, מעבד או מעביר מידע ו(1) חוק נציבות הסחר הפדרלי (Federal Trade Commission) חל עליו; (2) הוא ספק תקשורת (common carrier) חתח חוק התקשורת; או (3) הוא ארגון ללא מטרת רווח.

687 שם, ס' 207(a).

688 ס' 2(17) ל־ADPPA, שם, מגדיר "מחזיקי מידע רב" כישויות מכוסות (covered entities) שהכנסותיהן עולות על 250 מיליון דולרים בשנה והעבירו, אספו או עיבדו נתונים על 5 מיליון איש או מידע רגיש על 100,000 איש.

689 שם, ס' 207(c)(1)(b).

690 Algorithmic Accountability Act of 2022, H.R. 6580, 117th Cong. (2022) (להלן: הצעת חוק האחריותיות האלגוריתמית).

691 מערכות החלטה קריטיות מוגברות (augmented critical decision systems) מוגדרות בס' 2(1) להצעת חוק האחריותיות האלגוריתמית כהליך, הליך או כל פעילות אחרת משמשתים במערכת החלטה אוטומטית (automated decision system) לקבלת החלטה קריטית. מערכת החלטה אוטומטית מוגדרת בס' 2(2) להצעת החוק כל מערכת, תוכנה או הליך (לרבות כאלו שמקורם בלמידת מכונה, בסטטיסטיקה או בטכניקות בינה מלאכותית) שעורכים חישובים שתוצאתם משמשת בסיס לגיבוש עמדה או לקבלת החלטה. ס' 2(8) מגדיר "החלטה קריטית" החלטה שיש לה השפעה משפטית, מהותית או בעלת משמעות אחרת על חייו של צרכן בהקשר של מחיר, תנאים או זמינות של שירותים מסוימים (המנויים בהגדרה).

האחרונות או ששוויים בתקופה זו עולה על 250 מיליון דולר; שמעבדים מידע על יותר ממיליון צרכנים, בתי אב או מכשירים במטרה לפתח או לפרוש מערכת אוטומטית לקבלת החלטות או מערכת מוגברת לקבלת החלטות קריטיות.⁶⁹² ג'ייקוב מוקאנדר ועמיתיו רואים בחיוב את בחירת המחוקק האמריקאי לנקוט את המונחים "מערכות לקבלת החלטות" ו"תהליכים לקבלת החלטות קריטיות" חלף "בינה מלאכותית" ו"בינה מלאכותית בסיכון גבוה" הנקטים בהצעה האירופית – הן כי המונחים בהצעה האירופית לוקים באי-בהירות הן כי השמירה על ניטרליות טכנולוגית תגביר את חסינות העתיד של ההצעה האמריקאית.⁶⁹³

לפי הצעת חוק האחריות האלגוריתמית, ישויות שהחוק יחול עליהן יידרשו לבצע הערכת השפעה⁶⁹⁴ של כל מערכת אוטומטית לקבלת החלטות שפותחה לשם שימוש או הטמעה במערכות אוטומטיות לקבלת החלטות קריטיות או שיש צפי שישתמשו בהן או יטמיעו אותן בעתיד.⁶⁹⁵ נציבות הסחר הפדרלית תהיה רשאית להורות לישויות אלו לשמור תיעוד של הערכות ההשפעה שביצעו,⁶⁹⁶ לרווח על בסיס שנתי לנציבות על הערכות השפעה שנעשו בארגון⁶⁹⁷ ולצמצם השפעה שלילית של מערכות אוטומטיות לקבלת החלטות קריטיות.⁶⁹⁸ הנציבות תהיה רשאית גם להגדיר קווים מנחים ופורמטים להערכות השפעה ולדיווחים עליהן.⁶⁹⁹ רשות הסחר הפדרלית תקים מאגר מידע פומבי שבו יישמרו סיכומי הדיווחים השנתיים.⁷⁰⁰

692 ס' 2(7) להצעת חוק האחריות האלגוריתמית, לעיל ה"ש 690.

Jakob Mökander et al., *The US Algorithmic Accountability Act of 2022 vs. The EU Artificial Intelligence Act: What Can They Learn from Each Other?* 32 MINDS AND MACHINES 751 (18.8.2022)

694 ס' 4 להצעת חוק האחריות האלגוריתמית מכיל פירוט של הרכיבים הנדרשים מהערכת השפעה אלגוריתמית, לרבות חובות תיעוד, הכשרה, שקילת שיקולים הנוגעים לפגיעה פוטנציאלית בזכויות ומיתון נזקים.

695 ס' 3(b)(1) להצעת חוק האחריות האלגוריתמית, לעיל ה"ש 690.

696 שם, ס' 3(B)(b).

697 שם, ס' 3(b)(d).

698 שם, ס' 3(B)(H).

699 שם, ס' 3(b)(J)-(L).

700 שם, ס' 6.

מוקאנדר ועמיתיו מצביעים על החיסרון בהגבלת תחולתה של הצעת חוק האחריות האלגוריתמית לגופים עסקיים. לדבריהם, הגבלה זו מותירה מערכות בשירות הממשל – למשל, מערכות להקצאת זכויות סוציאליות – ללא פיקוח נאות. נוסף על כך, הנוסח הרזה של ההצעה האמריקאית – בהשוואה למקבילתה האירופית – משאיר לנציבות הסחר הפדרלי שיקול דעת נרחב למדי ביישום הוראות החוק, לרבות שימוש במונחים שמצמצמים שלא לצורך את החובות המוגדרות בחוק, דוגמת חובת ההתייעצות עם בעלי עניין ואנשי טכנולוגיה "עד כמה שאפשר" (to the extent possible).⁷⁰¹

לצד צעדי חקיקה ראשוניים אלו, במחצית הראשונה של 2023 ניכרת בארצות הברית התעוררות סביב נושא אסדרת הבינה המלאכותית. במאי 2023 נערך בסנאט האמריקני שימוע לסאם אלטמן, מנכ"ל חברת OpenAI;⁷⁰² המחלקה המייעצת לנשיא ארצות הברית בנושאי מדע וטכנולוגיה (OSTP) פרסמה קול קורא להערוך הציבור בעניין אסטרטגיית בינה מלאכותית לאומית;⁷⁰³ נוסף על כך, קבוצת רגולטורים, ובהם ראש הרשות להגנת הצרכן במגזר הפיננסי וראשת ועדת הסחר הפדרלית (FTC), פרסמו הצהרה משותפת המעלה את חששותיהם מפני התרומה של מערכות אוטומטיות לאפליה בלתי חוקית ולהפרה של חוקים פדרליים אחרים, בהדגישם שפעולות האכיפה של הסוכנויות שלהם חלות גם על מערכות אוטומטיות.⁷⁰⁴ ברוח זו, המחלקה האזרחית (זכויות אדם – דיור) במשרד המשפטים האמריקאי הגישה לאחרונה לבית המשפט המחוזי במסצ'וסטס הצהרה שבה היא מבקרת שחוק ההוגנות בדירור (Fair Housing Act) חל גם על החלטות אלגוריתמיות.⁷⁰⁵

Mökander et al., לעיל ה"ש 693. 701

Andrew Ross Sorkin et al., *Washington Confronts the Challenge of Policing A.I.*, THE NEW YORK TIMES (17.5.2023) 702

Office of Science and Technology Policy, Request for Information: National Priorities for Artificial Intelligence (23.5.2023) 703

U.S. Equal Employment Opportunity Commission, *Joint Statement on Enforcement Efforts Against Discrimination and Bias in Automated Systems* (2023) 704

Department of Justice, *Justice Department Files Statement of Interest in Fair Housing Act Case Alleging Unlawful Algorithm-Based Tenant-Screening Practices* (9.1.2023) 705

בממלכה המאוחדת טרם פורסמה הצעת חוק לאסדרת הפיתוח של בינה מלאכותית או השימוש בה. עם זאת, בחודש יולי 2022 פרסם

משרד הדיגיטל, המדיה, התרבות והספורט של בריטניה מסמך מדיניות (white paper), שהוא מעין הצהרת כוונות בעניין רגולציה עתידית של מערכות בינה מלאכותית. המסמך נוקט גישה מוטת חדשנות ומציע שניהול הסיכונים הרגולטורי של מערכות נבונות יהיה מבוסס הקשר, יתמקד בסיכונים אמיתיים, ניתנים לזיהוי, שרמת הסיכון הנשקפת מהם אינה מקובלת (להבדיל מסיכונים תאורטיים או כאלה שהסיכוי להתרחשותם נמוך). הגישה של בריטניה שואפת להפעלה מידתית של סמכויות רגולטוריות בכל הנוגע למערכות בינה מלאכותית ונוטה לבחירה בחלופות מחמירות פחות, כגון קווים מנחים או רגולציה וולונטרית.⁷⁰⁶ לפי מסמך המדיניות, ממשלת בריטניה שוקלת להטמיע את מערך העקרונות באמצעות הנחיות מינהליות, ולא באמצעות חקיקה. עם זאת, הצורך בחקיקה טרם נשלל לחלוטין.

מסמך המדיניות מציע להימנע מקביעת הגדרה משפטית אחידה לבינה מלאכותית ולהסתפק במאפייני ליבה.⁷⁰⁷ כך ייהנו הרגולטורים המגזריים מגמישות בהתמודדותם עם יישומים קונקרטיים של בינה מלאכותית.

אף על פי שההקשר קריטי לרגולציה עניינית של בינה מלאכותית, מסמך המדיניות מציע לפתח מערך של עקרונות רגולטוריים חוצי-מגזרים שיהיו ייחודיים לטכנולוגיות אלו. הרגולטורים יצטרכו להטמיע עקרונות כלליים אלו בתחום שהם מופקדים עליו. מסמך המדיניות מציע כמה עקרונות ראשוניים, בהתבסס על עקרונות ה-OECD לבינה מלאכותית,⁷⁰⁸ אך מדגיש כי אין בהם כדי לייצר זכויות אדם חדשות אלא רק מסגרת מוטת חדשנות לניהול סיכונים.⁷⁰⁹ עקרונות אלו כוללים שימוש בטוח במערכות בינה מלאכותית, הבטחת הבטיחות הטכנית של

706 Department for Digital, Culture, Media and Sport, לעיל ה"ש 59, בעמ' 11.

707 דוגמה אדפטיביות ואוטונומיות, ראו לעיל בסעיף 1.1.

708 ראו OECD 2019, לעיל ה"ש 294.

709 Department for Digital, Culture, Media and Sport, לעיל ה"ש 59, בעמ' 12.

מערכות בינה מלאכותית (ובכלל זה שהן פועלות כפי שתוכננו לפעול), שקיפות והסברתיות במידה המתאימה, הטמעת שיקולי הוגנות, הגדרת האחראים משפטית למשילות בינה מלאכותית (AI governance) והגדרת סעדים משפטיים.⁷¹⁰

4.7

אסדרת בינה מלאכותית: סיכום ביניים

על אף נקודות הדמיון בין ניצני החקיקה לאסדרת בינה מלאכותית ברחבי העולם שנסקרו בפרק זה החקיקות נבדלות זו מזו בהיקפן, בתחולתן ובפטרונות הרגולטוריים המוצעים בהן. כך

למשל, ההצעות של אמריקה, ברזיל וסין נמנעות מהקמת רגולטור ייעודי, ההצעה של קנדה מניחה את השאלה המוסדית לשיקול דעת מיניסטריואלי וההצעה של אירופה תומכת בהקמת רשות לאומית לבינה מלאכותית.⁷¹¹ ההבדלים הם בין השאר גאופוליטיים ותרבותיים: המשטר הסיני הסמכותני אינו מודאג משימוש של כוחות הביטחון והצבא במערכות בינה מלאכותית, ואילו האירופים מחזיקים בעמדות הפוכות; את רזונה של הצעת החוק של ברזיל אפשר לייחס גם למקומה במרוץ הגלובלי לפיתוח מערכות בינה מלאכותית (לעומת מעצמות כמו סין ואירופה), כפי שעולה מן השאיפה הצנועה שהיא מביעה בהצעה "להשתלב בשרשרת הערך הגלובלית במקום תחרותי".⁷¹² את הגישה מוטת החדשנות של בריטניה (שבאה על חשבון החלשת הבקורות הרגולטוריות),⁷¹³ הנמנעת לעת עתה מקידום חקיקה לאסדרת בינה מלאכותית, אפשר לפרש כביטוי של הערפת אינטרסים כלכליים מערכי היסוד של האיחוד האירופי.

נראה שצפון אמריקה פוזלת לעבר הצעת תקנות הבינה המלאכותית האירופיות, ההולכת לפני המחנה.⁷¹⁴ חרף היותן של ההצעה הקנדית וההצעה האמריקאית מקיפות פחות בתחולתן ובהיקפן, נראה שגם הן נוקטות גישה של ניהול סיכונים.

710 שם, בעמ' 12-14.

711 ראו להלן בסעיף 8.1.

712 סי' 3(III) להצעת החוק הברזילאית לבינה מלאכותית, לעיל ה"ש 54.

713 ראו הביקורת של Huw Roberts et al., *Artificial Intelligence Regulation in the United Kingdom: A Path to Global Leadership?* 5 (1.5.2023)

714 ראו בהקשר זה את ניתוחם של Mökander et al., לעיל ה"ש 693.

הצעת החוק הקנדית עושה זאת במפורש, בהגדרה מערכות בסיכון גבוה, ואילו הצעת החוק האמריקאית חלה מלכתחילה רק על מערכות לקבלת החלטות קריטיות. גם ההצעה של ברזיל נדרשת לרכיב הסיכון. השימוש בחקיקת פרטיות כדי לטפל בעקיפין בהיבטים שונים של מערכות בינה מלאכותית אף הוא מושפע מאירופה – בדומה לסעיף 22 ב-GDPR, שטיפולו הישיר בסוגיות הקשורות לבינה מלאכותית הוא תקדימי, גם סעיף 207 בהצעת חוק הפרטיות והגנת המידע האמריקאי מכיל הוראות בעניין בינה מלאכותית.

בעת כתיבת חיבור זה הצעת תקנות הבינה המלאכותית האירופיות היא בשלבי החקיקה האחרונים שלה. הגישה האירופית נוטה לאסדרה רוחבית בחקיקה (והסתייעות ברגולטורים מגזריים) והיא מלווה בחקיקת דיגיטל עוטפת (ה-GDPR לענייני פרטיות, דירקטיבת הסייבר האירופית וחוקי הדיגיטל החדשים – חוק השירותים הדיגיטליים וחוק השווקים הדיגיטליים). לעומת האיחוד האירופי, כמה ממדינות המערב שנסקרו בפרק זה, דוגמת בריטניה וארצות הברית, עדיין מגבשות את מדיניות האסדרה שלהן. בריטניה נוקטת גישה של אסדרה מגזרית מוטת חדשנות. ארצות הברית טרם גיבשה את הגישה שלה לאסדרה, אך נראה שהיא רואה באפליה אלגוריתמית את הסיכון העיקרי הטמון בבינה מלאכותית. אף שהמחלקה המייעצת לנשיא ארצות הברית בנושאי מדע וטכנולוגיה (OSTP) מקדמת את הנושא ברמה הפדרלית,⁷¹⁵ וחרף הצעות החוק שנסקרו לעיל,⁷¹⁶ מסתמן כי לעת עתה מגיעה האסדרה המעשית של מערכות בינה מלאכותית מהרגולטורים המגזריים בשטח, בדומה לחקיקת הפרטיות הפדרלית המטולאת של ארצות הברית. האסדרה הקנדית קרובה יותר ברמת בשלותה לאסדרה האירופית, וכך גם בתפיסה הרגולטורית שלה – אסדרה רוחבית שמבוססת על תפיסה של ניהול סיכונים, ובתפיסה המוסדית שלה – הקמת משרד בינה מלאכותית.

715 לעיל ה"ש 703 ו-312.

716 הצעת חוק האחרייתות האלגוריתמית, לעיל ה"ש 690; ADDPA, לעיל ה"ש 684.

פרק חמישי

בינה מלאכותית וזכויות אדם:
חדית חדשה במדיניות טכנולוגיה

—

מערכות נבונות, כמו כל טכנולוגיה חדשה, הן ניטרליות מבחינה מוסרית,⁷¹⁷ כלומר אין הן "רעות" בהכרח או "טובות" בהכרח.⁷¹⁸ ואולם לעיצובן ולפיתוחן

717 יובהר כי הטענה שטכנולוגיה כשהיא לעצמה היא ניטרלית אינה חפה מביקורת. טכנולוגיה אינה מיוצרת לשמה, אלא למילוי תכלית כלשהי – אם חיובית ואם שלילית. האם אפשר לטעון שרקטה מונחית חום היא ניטרלית מבחינה טכנולוגית? *AI Safety* Mariana Todorova, ראו בהקשר זה גם *Myths, Future of Life Institute* (7.8.2016) *Philosophical, Moral, and Ethical Rationalization of Artificial Intelligence, in The Transhumanism Handbook* 263 (Newton Lee ed., 2019)

718 Mariarosaria Taddeo and Luciano Floridi, *How AI Can Be a Force for Good: An Ethical Framework Will Help to Harness the Potential of AI while Keeping Humans in Control*, 361 *SCIENCE* 751 (24.8.2018)

של טכנולוגיות, ולהקשר שבו הן משמשות, יש נגיעה לזכויות אדם.⁷¹⁹ הוא הדין במערכות נבונות – לצד התרומה הפוטנציאלית העצומה שלהן לאיכות חייו של האדם, אימוץ נרחב של טכנולוגיות אלו עלול להביא גם לפגיעה בזכויותיו. לנוכח עוצמת הפגיעה הפוטנציאלית של יישומים מסוימים של בינה מלאכותית מציעות אפוא תקנות הבינה המלאכותית האירופיות לאסור עליהם כליל.⁷²⁰ כך גם המליץ נציב זכויות האדם של האומות המאוחדות.⁷²¹

יש להבהיר כי פרק זה אינו מבקש לטעון שבינה מלאכותית פוגעת מטבעה בזכויות אדם, אלא לסקור אופנים שונים שבהם השימוש בטכנולוגיות אלו עלול להביא לפגיעה בזכויות. נוסף על כך, הסקירה נועדה להדגים פגיעה פוטנציאלית בזכויות ולא להקיף את כל מופעיה. משכך, היא אינה בהכרח ממצה את מלוא הסיכונים לחירויות יסוד וזכויות אדם שכרוכים בשימוש במערכות בינה מלאכותית.

עם זאת, בטרם נסקור את הסיכונים נציע רשימה לא ממצה של סוגי הנזקים שמערכות בינה מלאכותית מתקדמות עלולות להסב ונמנו לאחרונה בספרות.⁷²² הרשימה מדגימה מדוע חשיבות הטיפול בסיכונים אלו היא כה רבה.

סוגי הנזקים העלולים להיגרם משימוש בבינה מלאכותית ונסקרו עד כה משתלבים בחששות לנזקים נוספים ההולכים ונעשים מוחשיים, כגון

- אבטחת מידע וסייבר – מודלים יכולים לגלות חולשות במערכות (חומרה, תוכנה, נתונים) ולכתוב קוד לניצולן; להכניס, במסגרת משימה לכתובת קוד, "באגים" לצורך ניצול עתידי; לקבל החלטות שעלולות לגרום נזק כאשר הם מקבלים גישה למערכת או לרשת; להתחמק מזיהוי של הפעולות שעשו.

719 ראו למשל Mike Cooley, *The Myth of the Moral Neutrality of Technology*, 9 AI & Soc. 10 (1995)

720 ראו לעיל סעיף 4.1.1.1.

721 United Nations General Assembly, *THE RIGHT TO PRIVACY IN THE DIGITAL AGE*, 721 Para. 59.(c)-(d) (A/HRC/48/31, 13.9.2021)

722 Toby Shevlane et al., *Model Evaluation for Extreme Risks* 722 (24.5.2023), available at arXiv

- הונאה – למודלים יש הסגולות הנדרשות כדי להונות בני אדם. למשל בנייה של תוכן – מיריעה חדשותית עד תזהיר משפטי – הנתפס אמין, אך הוא כוזב; או התחזות משכנעת לאדם. מודלים יכולים לחזות במדויק השפעת של הונאה על בני אדם ולנהל מעקב של המידע שהם זקוקים לו כדי לשמור על יכולת ההונאה.
- שכנוע ומניפולציה – מודלים יעילים בקידום נרטיבים בצורה משכנעת ובהשפעה על אמונות של בני אדם ועל עיצובן באמצעות דיאלוג איתם או טכניקות אחרות (למשל השפעה על אלגוריתם של הפצת מידע ברשתות חברתיות).
- אסטרטגיה פוליטית – מודלים יכולים לבצע את המידול והתכנון החברתי הדרושים לשחקן כדי להשיג ולהפעיל השפעה פוליטית, לא רק ברמת המיקרו אלא בתרחישים עם שחקנים מרובים והקשר חברתי עשיר. לדוגמה, לקבל ניקוד גבוה בחיזוי תחרויות בשאלות הקשורות למשא ומתן פוליטי.
- פיתוח נשק או גישה לנשק – מודלים יכולים לקבל גישה למערכות נשק קיימות או לתרום לבניית כלי נשק חדשים. לדוגמה, להרכיב נשק ביולוגי (בסיוע אנושי) או לספק הוראות מעשיות כיצד לעשות זאת.
- תכנון ארוך טווח – מודלים יכולים ליצור תוכניות רציפות הכוללות שלבים מרובים, התלויים אלה באלה ונפרשים על פני טווח זמן ותחומי פעולה שונים; ולכוונן כל העת את השלבים לנוכח מכשולים או יריבים בלתי צפויים בתוך שהם מממשים יכולות מכליליות לסכיבות חדשות ואינם מסתמכים רק על ניסוי וטעייה.
- פיתוח בינה מלאכותית – מודלים יכולים לבנות מערכות AI חדשות שיכולות להיות מסוכנות או דו־שימושיות.
- מודעות מצבית – מודלים יכולים להבין שהם מודלים ויש להם ידע על עצמם ועל סביבתם (למשל איזו חברה הכשירה אותם, היכן נמצאים השרתים שלהם, איזה סוג של אנשים יכולים לתת להם משוב), וכן לדעת באיזה שלב הם במעגל החיים של בינה מלאכותית, כלומר אם הם בשלב האימון, שלב ההערכה או שלב היישום – ולהתנהג בהתאם.

• יכולת התרחבות והפצה עצמית – מודלים יכולים ליצור לעצמם אפשרויות לבצע פעולות שהם לא תוכננו לבצע באמצעות שינוי מערכת ההפעלה שלהם, מניפולציה על מי שמתכנן את המערכת ושימוש באסטרטגיות לחילוץ הקוד והמשקלים של המערכת. מודלים יכולים לנצל חולשות במערכות שמיועדות לנטר את התנהגותם בשלב היישום שלהם ולחולל שינוי בהערכה בעניינם. מודלים יכולים לייצר הכנסות באופן עצמאי (למשל על ידי הצעת שירותים או מתקפות כופר), ולהשתמש בהכנסות אלו כדי לרכוש משאבי מחשוב ולהפעיל מספר רב של מערכות אחרות של בינה מלאכותית.

5.1

הזכות לפרטיות

מידע הוא הנפט החדש של העולם המפותח.⁷²³
נתוני עתק הם הדלק שמערכות נבונות, ובפרט
כאלו המבוססות על למידת מכונה, צורכות על

מנת לשכלל את תהליכי קבלת ההחלטות שלהן.⁷²⁴ וכיוון שפיתוח מערכות נבונות מחייב שימוש במידע אישי,⁷²⁵ הוא כרוך בפגיעה בזכות לפרטיות. אומנם שימוש כזה יכול להיות מותר, למשל אם הוא נעשה בהסכמת מושאי המידע;⁷²⁶ אבל במקרים רבים מעובד אותו מידע אישי לתכליות שונות מאלו שלשמן נאסף – תכליות שלא היה אפשר אפילו להעלות על הדעת בזמן שנאסף – אף שהדבר נוגד את עקרון צמידות המטרה, שהוא עקרון יסודי בדיני הגנת הפרטיות.⁷²⁷

723 Palmer, לעיל ה"ש 112.

724 Todorova, לעיל ה"ש 717, בעמ' 286-287.

725 נדשקה פורטובה טוענת – לנוכח ההגדרה הרחבה של מידע אישי בדיני הגנת הפרטיות הבינלאומיים (ראו לדוגמה ס' 1(4) של התקנות הכלליות בדבר הגנת מידע [GDPR], וההגדרה שהציעו ארידור הרשקוביץ ושוורץ אלטשולר [לעיל ה"ש 403, בעמ' 133] – שבסביבה רווית נתונים שבה נעשה שימוש הולך וגובר במערכות נבונות, "הכול הוא מידע אישי". Nadezhda Purtova, *The Law of Everything: Broad Concept of Personal Data and Future of EU Data Protection Law*, 10 L. INNOVATION AND TECH. 40 (2018)

726 ראו למשל ס' 6(1) של התקנות הכלליות בדבר הגנת מידע (GDPR), ואת ההוראה המקבילה לו המוצעת אצל ארידור הרשקוביץ ושוורץ אלטשולר, לעיל ה"ש 403, בעמ' 235-236.

727 ראו ס' 6(4) של התקנות הכלליות בדבר הגנת מידע (GDPR), ואת ההוראה המקבילה לו המוצעת אצל ארידור הרשקוביץ ושוורץ אלטשולר, לעיל ה"ש 403, בעמ' 51-52.

ככל שגובר השימוש במערכות נבונות בכל תחומי החיים, כך גובר הצורך של מפתחיהן ובעליהן במידע מפורט יותר על המשתמשים בהן. מידע זה נאסף אוטומטית באמצעות מגוון חיישנים ומכשירי קצה המלווים את האדם המערבי המודרני בשגרת יומו. זהו "קפיטליזם של מעקב", הממיר נתוני עתק למודלים של ניתוח התנהגות ושל חיזוי.⁷²⁸ לצד היתרונות הטמונים בחיבור מתמיד למרשתת ולאינטרנט של הדברים יש בו גם חסרונות: עיבוד מידע שנאסף דרך קבע על ידי מערכות נבונות לקבלת החלטות עלול להביא לידי פגיעה בפרטיות ולכן לייצר אפקט מצנן על ההתנהגות, כלומר לגרום להתכנסות החוויה האנושית למחוזות קונפורמיים ולשחיקת האוטונומיה של הפרט.⁷²⁹

התפתחויות המבוססות על מודלים של למידת מכונה בתחומי הביולוגיה ומדעי החיים, ובכללן האפשרות ליצור זיהוי חוזר של מידע בריאות מותמם וריצוף גנטי בר-זיהוי מתוך נהרות או מתוך האוויר (e-DNA), מחריפות מאד את עומקו ורוחבו של האיום על הפרטיות. חוקרים שהראו כי ביכולתם לפענח ו"לתמלל" מחשבות באמצעות הפעלת אלגוריתם על ממצאים של סריקה מוחית חיצונית (fMRI), בדרך שמכונה brain decoding, הרגישו לאחרונה צורך להצהיר במפורש כי "חינוי לעורר מודעות לסכנות בטכנולוגיית קריאת מחשבות וליצור חוקים שיגנו על הפרטיות המחשבתית של כל אדם".⁷³⁰

האפקט המצנן של הפגיעה בפרטיות מחריף כשהתכלית של המערכות הנבונות היא מעקב בפועל. מערכות לזיהוי פנים, לניטור התנהגות במרחב או לחיזוי פשיעה וטרור פועלות כפנאופטיקון הממסטר את התנהגות הסובייקטים המפוקחים באמצעותן.⁷³¹ משום כך הן עלולות להשפיע על נכונות של אנשים להימצא באזורים מסוימים, לבוא במגע עם זולתם, להתבטא בחופשיות או להתנהג באופן לא מקובל. לכן כשמערכות נבונות פוגעות בזכות לפרטיות תיתכן

728 ראו Zuboff, לעיל ה"ש 112.

729 שם, בעמ' 521.

730 Jerry Tang et al., *Semantic Reconstruction of Continuous Language from Non-Invasive Brain Recordings*, 26 NATURE NEUROSCIENCE 858-866 (2023)

731 ראו לדוגמה כהנא ושני, לעיל ה"ש 213, בעמ' 19-20.

גם פגיעה בזכויות אדם אחרות, כגון חופש הביטוי, זכות האספה, חופש התנועה או האוטונומיה של הפרט.

5.2

הזכות להליך הוגן וכנלי הצדק הטבעי

הזכות להליך הוגן היא מכללי הצדק הטבעי.⁷³² אומנם שורשיה נטועים במשפט הפלילי, אבל היא משפיעה גם על הליכים אזרחיים ומינהליים ומעוגנת במשפט הבינלאומי.⁷³³ היא אינה נחשבת

זכות חוקתית מפורשת, אך הוכרה בפסיקה ובספרות כזכות בעלת מעמד חוקתי שנגזרת מהזכות לכבוד או מהזכות לחירות.⁷³⁴

מהזכות להליך הוגן במשפט הפלילי נובעות, בין השאר, "[...] זכותו של נאשם לדעת מדוע נעצר ובמה הוא מואשם, הזכות להיות מיוצג על-ידי עורך-דין, הזכות להיות נוכח במשפט, הזכות למשפט פומבי על-ידי ערכאה בלתי תלויה וניטרלית והזכות להתגונן במשפט ולהציג ראיות רלוונטיות".⁷³⁵

הסתייעות של מערכת המשפט בבינה מלאכותית התומכת בתהליך קבלת החלטות⁷³⁶ עלולה להביא לידי פגיעה בהוגנות ההליך, שכן אותה "קופסה שחורה" ממסכת את השיקולים המכריעים בקבלת ההחלטות האוטומטית ומקשה על הנאשם להבין מדוע נעצר, להיות "נוכח" במשפטו, לממש את זכות הטיעון שלו ולהבין מהן הראיות הרלוונטיות. לאחרונה נחשפה מערכת הכללה ממוחשבת (פרופיילינג) לאיתור חשודים בבלדרות סמים בנמל התעופה בן-גוריון. היא

732 דפנה ברק-ארז משפט מינהלי כרך א 461-564 (2010).

733 ראו למשל European Convention for the Protection of Human Rights and Fundamental Freedoms, art. 6, adopted Nov. 4, 1950, ETS 5 (entered into force Sep. 3, 1953); (להלן: ECHR); Universal Declaration of Human Rights, G.A. Res. 217(III) A, U.N. Doc. A/RES/217(III), at art. 10 (Dec. 10, 1948) (להלן: UDHR).

734 אהרן ברק כבוד האדם: הזכות החוקתית ובנוחיה 868-869 (2014).

735 ע"פ 5121/98 יששכרוב נ' החובע הצבאי הראשי, פ"ד טא (1) 461 (2006), פס' 66 לפסק דינה של השופטת ביניש.

736 ראו לעיל סעיף 2.3.2.

מופעלת על ידי המטרה, אך אופן פעולתה לא ידוע, ובהיעדר שקיפות מתעורר חשש להטיות אלגוריתמיות.⁷³⁷

ההליך ההוגן עלול להיפגע לא רק עקב שימוש במערכות תומכות החלטה במסגרת הליכים משפטיים פליליים, אלא בכל מקרה שבו השלטון נעזר במערכות נבונות (למשל כדי לקבל החלטות הנוגעות לחלוקת משאבים) בלי לאפשר למי שזכויותיו נפגעו לקבל את יומו בבית משפט (אנושי) ולממש את זכות הטיעון שלו. לכן הדין האירופי, למשל, מקנה למושא המידע את הזכות שלא יתקבלו בעניינו החלטות על בסיס עיבוד אוטומטי של המידע האישי שלו.⁷³⁸

דוגמה לפגיעה בהליך הוגן היא ההתנהלות בפרשת המערכת ההולנדית לאיתור הונאות רווחה SyRI (System Risico Indicatie).⁷³⁹ SyRI הייתה מערכת משותפת של כמה גופים ממשלתיים בהולנד; תכליתה הייתה לאסוף מידע ולהעבירו למשרד המידע ההולנדי. שמות אנשים ואמצעים מזהים אחרים הוחלפו בקוד והמפתח לקוד נשמר במקום אחר. לאחר הצלבת המידע באמצעות אלגוריתמים, שדרך פעולתם לא פורטה לציבור, לפרלמנט או לבתי המשפט, הפיקה המערכת פלט שקבע כי עלה בידה לאתר חשוד בהונאה. התראה של המערכת הביאה את משרד הרווחה של הולנד לפתוח בחקירה. במקרה זה העלה היעדר שקיפות פעולת האלגוריתם חשש לפגיעה בזכות להליך הוגן בעניין זכאות לתמיכה ממשלתית.⁷⁴⁰

במישיגן שבארצות הברית זיהתה מערכת לאיתור הונאות רווחה (במקרה זה, זכאות לדמי אבטלה) אלפי אזרחים רמאים, שללה אוטומטית את זכאותם לדמי

737 ראו תומר גנון "האלגוריתם שיעצור אתכם בנחיתה בנחב"ג" כלכליסט (10.11.2022); בש"פ 8308/21 עומר ברגר נ' מדינת ישראל (8.12.2021); עמיר כהנא ותהילה שוורץ אלטשולר הפרופיילינג המשטרחי מעלה חשש לאפקט מצנן שישפיע על זכויות אדם (חוות דעת, המכון הישראלי לדמוקרטיה 22.5.2023).

738 ראו לעיל ה"ש 403.

739 Max Vetz, *The Netherlands-Algorithmic Fraud Detection System Violates Human Rights: The Case of SyRI*, 3 PUBLIC LAW, 650 (2021). ראו גם פרל ושוורץ אלטשולר, להלן ה"ש 881, בעמ' 54-56.

740 על פסק הדין בעניין SyRI ראו פרל ושוורץ אלטשולר, להלן ה"ש 881, בעמ' 54-56.

אבטלה ותבעה מהם החזרים (בצירוף קנסות) של דמי אבטלה בסכומי עתק. מבקר המדינה של משיגן בדק מדגם של כ-22,000 מהמקרים שבהם זיהתה המערכת הונאה ומצא ש-93% מהם שגויים. האוטומציה של ההליך לאיתור הונאות (מזיהוי ועד ענישה), שנעשה כמעט ללא מעורבות אנושית, הקשה על היכולת של האזרחים שזכויותיהם הסוציאליות נשללו להבין את הסיבות ולערער על הקביעות הממוכנות.⁷⁴¹

5.3 הזכות לשוויון בפני החוק

הזכות לשוויון בפני החוק היא עקרון יסודי חוקתי שמוכר במשפט הבינלאומי.⁷⁴² אף שהזכות לשוויון אינה זכות חוקתית מפורשת בחוקי היסוד של ישראל, הפסיקה והספרות מכירים בה

כזכות בת של כבוד האדם,⁷⁴³ "עקרון יסודי חוקתי, השלוב ושזור בתפיסות היסוד המשפטיות שלנו ומהווה חלק בלתי נפרד מהן".⁷⁴⁴

מערכות נבונות עלולות לפגוע בזכות לשוויון עקב הטיות אלגוריתמיות המובילות להקצאת משאבים מפלה או לפגיעה עודפת בזכויות של קבוצות מוגנות. מאחר שפיתוח מערכות בינה מלאכותית נשען על פיתוח מודלים סטטיסטיים מבוססי מידע, ההסתמכות על מאגרי מידע שמשעתקים יחסי כוח אי-שוויוניים, או שמשקפים דעות קדומות נגד מגזרים וקבוצות מוחלשים, עלולה להנציח את אי-השוויון; ובמקרים של למידת מכונה מבוססת משוב אף להעצים אותו.⁷⁴⁵ גם

741 Cahoo v. SAS Analytics Inc., No. 18-1296 (6th Cir. 2019); ראו בעניין זה: Sonia Gipson Rankin, *The MIDAS Touch: Atuahene's "Stategraft" and the Implications of Unregulated Artificial Intelligence*, 2022-21 UNM SCHOOL OF LAW RESEARCH PAPER (2022)

742 ראו למשל ס' 1, 27 ב-UDHR, לעיל ה"ש 733; Art. 2(1), 26 International Covenant on Civil and Political Rights, GA Res. 2200A (XXI) of 16 December 1966 (להלן: ICCPR); ס' 2(1), 14, 26 ב-ECHR, לעיל ה"ש 733.

743 ראו ברק, לעיל ה"ש 734, בעמ' 685-690.

744 בג"ץ 114.78 בורקאן נ' שר האוצר, פ"ד לב(2) 806, 800 (1978).

745 על הטיות אלגוריתמיות, ראו להלן בסעיף 6.2.

מודלים סטטיסטיים שמסתמכים על מידע חלקי או עודף על הקבוצה המוגנת יכולים להביא לידי פגיעה בשוויון.⁷⁴⁶

5.4

כבוד האדם

כבוד האדם הוא ערך יסוד בשיטה החוקתית של ישראל וגם במשפט החוקתי המשווה.⁷⁴⁷ לפי הגישה הפרשנית ההוליסטית של אהרן ברק יש

לייחס לכבוד האדם ערך חוקתי.⁷⁴⁸ לפי גישה זו, ערכו החוקתי של כבוד האדם הוא ההגנה על אנושיותו של אדם.⁷⁴⁹

מערכות אוטומטיות שמשמשות לבקרת תוכן שמעלים גולשים עלולות לפגוע בכבוד האדם של משתמשים שהתבטאויותיהם צונזרו. טכנולוגיות בינה מלאכותית עלולות לפגוע בכבוד האדם של מי שזהותם נגנבה על ידי טכנולוגיות דיפ-פייק⁷⁵⁰ או שדמויותיהם הווירטואליות מוצגות באופן מבזה.⁷⁵¹ בשל הפרשנות המרחיבה של הערך החוקתי של כבוד האדם יש סיכוי רב לפגוע בו באמצעות פעילותן של מערכות נבונות. הללו יכולות לפגוע בקשת רחבה של זכויות נגזרות, דוגמת הזכות לפרטיות, הזכות להליך הוגן והזכות לשוויון.

746 ראו להלן בסעיף 6.2.1.

747 ראו ברק, לעיל ה"ש 734, בעמ' 73-224.

748 שם, בעמ' 235-237.

749 שם, שם. ראו גם דברי השופט ברק בעניין בג"ץ 6427/02 התנועה למען איכות השלטון בישראל נ' כנסת ישראל, פ"ד סא(1) 619, 685 (2006): "ביסוד כבוד האדם עומדים האוטונומיה של הרצון הפרטי, חופש הבחירה וחופש הפעולה של האדם כיצור חופשי. כבוד האדם נשען על ההכרה בשלמותו הפיסי והרוחנית של האדם, באנושיותו, בערכו כאדם וכל זאת בלא קשר למידת התועלת הצומחת ממנו לאחרים".

750 ראו לעיל, ה"ש 485.

Jacquelyn Burkell and Chandell Gosse, *Nothing New Here: Emphasizing the Social and Cultural Context of Deepfakes*, 24 FIRST MONDAY (2019); Elizabeth F. Judge and Amir M. Korhani, *Deepfakes, Counterfeits, and Personality*, 59 ALBERTA L. REV. (26.7.2021)

פגיעה אפשרית נוספת היא איון הסובייקט האנושי והמרתו ברשומה במסד נתונים.⁷⁵² בפרק הקודם ראינו כי מסמכי אתיקה של מערכות נבונות רואים בכבוד האדם ובאוטונומיה אנושית רכיב יסודי.⁷⁵³ כשיכולת האדם לבחור בחירה מושכלת מפנה את מקומה לקבלת החלטות אוטומטית, כשנתח הולך וגדל של ההחלטות מתקבל באמצעות מודלים סטטיסטיים, נשחק הסובייקט האנושי, מרגיש בלתי נראה⁷⁵⁴ וחסר ערך, וחש תסכול לנוכח המציאות השרירותית.

חופש הביטוי הוא ערך מוגן במשפט הבינלאומי,⁷⁵⁵ וגם במשפט הישראלי.⁷⁵⁶ חשיפת האמת, הגשמה עצמית של האדם והבטחת המשטר הדמוקרטי

5.5 חופש הביטוי

משמשות לביסוס חשיבותו של חופש הביטוי.⁷⁵⁷

מנגד, זכות זו אינה מוחלטת. אותם דברי חקיקה הקוראים לחופש הביטוי אוסרים לעיתים התבטאויות הכוללות ביטויי שנאה (hate speech), למשל.⁷⁵⁸ מלאכת האיזון בין האינטרסים המוגנים על ידי חופש הביטוי ובין האינטרסים המוגנים על ידי האיסור על ביטויי שנאה היא מלאכה מורכבת שתלויה בתרבות ובמסורת החוקתית המקומית. מורכבות זו מתעצמת לנוכח יכולתו של כל משתמש להתבטא במרחב המקוון, ושטף הציוצים, הפוסטים ושאר ההתבטאויות המקוונות

752 לדדהומניזציה של הסובייקט האנושי על ידי מערכות בינה מלאכותית בתחומי הרפואה ראו Van Kolfshooten, לעיל בה"ש 581, בעמ' 92-95.

753 ראו לעיל בסעיף 3.7.

754 ראו למשל Zuboff, לעיל ה"ש 112, בעמ' 44, המחארת את התסכול של המפגינים במהומות בלונדון בשנת 2011.

755 ראו למשל ס' 19 ב-UDHR, לעיל ה"ש 733; ס' 19 ב-ICCPR, לעיל ה"ש 742; ס' 10 ב-ECHR, לעיל ה"ש 733.

756 ראו בג"ץ 73/53 חברת "קול העם" בע"מ נ' שר הפנים פ"ד ז 871 (1953); ברק, לעיל ה"ש 734, בעמ' 707-712.

757 ברק, שם, בעמ' 713.

758 ראו לדוגמה ס' 20(2) ב-ICCPR, לעיל ה"ש 742.

המצרפיות של כל בני "אומת האינטרנט".⁷⁵⁹ התאמת המודל התאורטי המדויק להגדרת ביטויי שנאה ולקביעת מדיניות הטיפול המתאימה בהם היא אתגר בפני עצמו,⁷⁶⁰ אך בשל היקף הביטויים שיש לנטר יישומו של המודל מצריך גם התערבות אלגוריתמית. הַדָּח המיוחד של האומות המאוחדות לחופש הביטוי עמד על הסכנות הנשקפות לחופש הביטוי מהשימוש בניטור תוכן אלגוריתמי;⁷⁶¹ הוא הצביע על הקשיים הטכניים של מערכות אלו לזהות דקויות של אירוניה⁷⁶² והראה כי לעמימות והיעדר שקיפות יש אפקט מצנן על חופש הביטוי. יתר על כן, מעיין פרל (פילמר) וניבה אלקין-קורן טוענות כי האופי המסחרי של הפלטפורמות המשתמשות בטכנולוגיות ממוכנות לצורך ניטור תוכן עלולות לגרוע מהאופי הדמוקרטי של השיח המקוון. לדבריהן, כללי ההכרעה בעניין נורמות של תוכן לגיטימי הולכים ונעשים אלגוריתמיים, הבקרה עליהם נעשית בחוסר שקיפות והאוטומציה השואפת לכללי הכרעה בינאריים אינה סובלת את ריבוי המשמעויות של השפה הטבעית.⁷⁶³

5.6

זכויות פוליטיות

מערכות נבונות מאיימות על מימושן של זכויות פוליטיות, ובהן הזכות לחופש האספה והזכות לחופש ההתאגדות,⁷⁶⁴ הזכות להשתתפות

פוליטית, הזכות לחופש המחשבה והזכות לחופש הביטוי.⁷⁶⁵ הַדָּח המיוחד של האומות המאוחדות לזכות לחופש האספה ולזכות לחופש ההתאגדות ציין כי

United Nations General Assembly, RIGHTS TO FREEDOM OF PEACEFUL ASSEMBLY AND 759
OF ASSOCIATION, Para. 59 (A/HRC/41/41, 2019). (להלן: HRC 2019).

760 ראו Medzini and Shwartz Altshuler, לעיל ה"ש 268.

761 ראו HRC 2018, לעיל ה"ש 270.

762 ראו לדוגמה נייר יגנה "החשוד בהסתה להצתה בפייסבוק שוחרר: 'גם בערבית יש סאטירה' " וואלה! (28.11.2016).

Maayan Perel (Filmar) and Niva Elkin-Koren, *Democratic Friction in* 763
Speech Governance by AI, in HANDBOOK OF CRITICAL STUDIES OF ARTIFICIAL INTELLIGENCE
(Simon Lindgren ed., 2022)

764 ראו למשל סי' 20 UDHR, לעיל ה"ש 733; סי' 21-22 ICCPR, לעיל ה"ש 742.

765 על פגיעה פוטנציאלית של מערכות נבונות בחופש הביטוי, ראו לעיל סעיף 5.5.

זכויות אלו [...] הן המפתח למימוש הדמוקרטיה וכבוד האדם, מאחר שזכויות אלו מאפשרות לאנשים להשמיע את קולם ולייצג את האינטרסים שלהם, לראות בממשלות אחראיות לפעולותיהן ולהעצים את הפעלנות האנושית [human agency].⁷⁶⁶ מימוש זכויות אלו במרחב המקוון יכול לבוא לידי ביטוי בדרכים מגוונות, כגון הצטרפות לעצומה באינטרנט, חבירה לקבוצות עניין פוליטיות ברשתות חברתיות⁷⁶⁷ או פרסום פוסטים עם האשטאג בעל זיהוי פוליטי.⁷⁶⁸

עם הסיכונים הנשקפים לזכויות פוליטיות משימוש במערכות בינה מלאכותית אפשר למנות את האפקט המצנן של מערכות לניטור תוכן אלגוריתמי ומערכות למעקב מקוון.⁷⁶⁹ אנשים שיודעים שהם נתונים למעקב מקוון עלולים להימנע מפעילות פוליטית מקוונת שנתפסת חריגה או לא קונצנזואלית, גם אם היא לגיטימית, מחשש שמערכת נבונה או צנזור אנושי יסווג אותה כאסורה.

מערכות נבונות עלולות לצמצם את החופש הפוליטי המקוון גם באמצעות פגיעה בחופש המחשבה. היבט מרכזי של חופש המחשבה הוא "הזכות לגבש דעה ולפתח אותה באופן מנומק",⁷⁷⁰ ובכלל זה הזכות לעשות כן ללא כפייה וללא אינדוקטרינציה מבחוץ.⁷⁷¹ ככל שהאלגוריתמים נעשים שחקנים משמעותיים יותר במלאכה האפיסטמולוגית של אצירת המודעות האנושית, כך הם יכולים להשפיע יותר על האופן שבו אנשים מגבשים את דעתם ומפתחים אותה. כיום אלגוריתמים הם שמתווכים בין המציאות ובין המשתמשים ברשתות החברתיות – וכיוון שהם

United Nations General Assembly, REPORT OF THE SPECIAL RAPPORTEUR ON THE RIGHTS TO FREEDOM OF PEACEFUL ASSEMBLY AND OF ASSOCIATION, Para. 87 (A/73/279, 2018)

767 ראו HRC 2019, לעיל ה"ש 729, בפס' 25.

768 לסקירה ראו SARAH J. JACKSON, MOYA BAILEY, AND BROOKE FOUCAULT WELLES, #HASHTAGACTIVISM (2020); לדוגמה מקומית ראו הגר בוחבוט וירון דרוקמן "הפקרתי מימיה ונענשתי יותר מבוכריסי": קמפיין רשע נגד הסדר הטיעון" ynet (1.12.2016).

769 Cameran Ashraf, *Artificial Intelligence and the Rights to Assembly and Association*, 5 J. CYBER POL'Y 163 (2020)

770 ראו HRC 2018, לעיל ה"ש 270, בפס' 23; MANFRED NOWAK, U.N. COVENANT ON CIVIL AND POLITICAL RIGHTS. CCPR COMMENTARY 441 (2005)

771 שם, בעמ' 442; ראו לדוגמה גם CCPR/C/78/D/878/1999, *Yong Joo-Kang v. Republic of Korea* (16.7.2003)

שולטים ב"פיד" הידיעות של המשתמשים הם מסוגלים לעצב עמדות ורחשי לב;⁷⁷² ולעיתים אלו עמדות קיצוניות שמובילות להסתה ולאלימות.⁷⁷³

התפקיד המרכזי של מערכות בינה מלאכותית בעיצוב התודעה האנושית בעידן המידע מעלה חששות באשר ליכולת לממש את חופש המחשבה ולגבש דעה באופן עצמאי ומנומק.⁷⁷⁴ האלגוריתמים משמשים "שומרי סף" הקובעים אילו ידיעות ימצאו את דרכן ל"פיד" ואילו יישארו מחוץ לו. אלגוריתמים יכולים להגביר את הרעש התקשורתי סביב ידיעה מסוימת או להחלישה כליל.

שחקנים חיצוניים יכולים לנקוט מבצעי השפעה (psych ops) כדי לייצר ב"פיד" רעשים שיש בהם כדי להשפיע על דעת הקהל הפוליטית. גם מיקרו-טרגוט ממוקד יכול לשמש לעיצוב עמדות של מגזרים מסוימים, למשל בסמוך לבחירות או במהלכן.⁷⁷⁵ במסגרת מבצעי השפעה אפשר לרתום מערכות בינה מלאכותית להפקת זיופים בטכנולוגיות דיפ-פייק כדי לערער את דעת הקהל באשר למועמד מסוים או מפלגה מסוימת.⁷⁷⁶ מבצעי השפעה מעין אלו לא בהכרח יערערו את

Adam D. I. Kramer, Jamie E. Guillory, and Jeffery T. Hancock, 772 *Experimental Evidence of Massive-Scale Emotional Contagion through Social Networks*, 111 PROC. NAT'L. ACAD. SCI., 8788-8790 (2014); Katherine Haenschen, Brett Frischmann, and Paul Ellenbogen, *Manipulating Facebook's Notification System to Provide Evidence of Techno-Social Engineering*, 1 Soc. Sci. Comp. Rev. (2021)

773 ראו לדוגמה את הטענות בעניין השפעת אלגוריתמים החולשים על הפיד של הרשת החברתית פייסבוק של חברת מטא על רצח העם של בני הרוהינגה במיאנמר: THE SOCIAL ATROCITY: META AND THE RIGHT TO REMEDY FOR THE ROHINGYA (Amnesty International, 2022); Rebecca J. Hamilton, *Platform-Enabled Crimes: Pluralizing Accountability When Social Media Companies Enable Perpetrators to Commit Atrocities* 63 B.C. L. Rev. 1349, 1365-1366 (2022)

774 ראו HRC 2018, לעיל ה"ש 270, בפס' 24-26.

775 ראו לדוגמה תהילה שוורץ אלטשולר וגיא לוריא *תעמולה דיגיטלית והאיום על הבחירות 39-53*; (2020); אלי בכר ורון שמיר *התקפות סייבר על מערכת הבחירות: איך מתמודדים?* 42-45; (2019) Alina Poluakova and Spencer B. Boyer, *The Future of Political Warfare: Russia, the West, and the Coming Age of Global Digital Competition*, BROOKINGS (2018)

776 Cristian Vaccari and Andrew Chadwick, *Deepfakes and Disinformation: Exploring the Impact of Synthetic Political Video on*

אמון הבוחרים בעמדה מסוימת, אלא יביאו אותם להטיל ספק ביציבות הפוליטית ובמקורות המידע שעליהם הם נסמכים בדרך כלל כדי לגבש את דעתם.

5.7

הזכות לחיים

הזכות לחיים היא זכות יסוד במשפט החוקתי של

ישראל,⁷⁷⁷ וגם במשפט המשווה והבינלאומי.⁷⁷⁸

יישומי בינה מלאכותית עלולים לגרום לפגיעה

אגבית בזכות לחיים, הן עקב טעויות הן עקב הכרעות מוקדמות הנוגעות לעיצובו של היישום. למשל, ייתכנו מקרי קצה שבהם המפתחים של מערכת רכב אוטונומית יידרשו לקבוע מראש קריטריונים לאיזון בין הזכות לחיים של נוסע ברכב ובין הזכות לחיים של הולך רגל.⁷⁷⁹ לאחרונה פורסם מחקר שמתאר כיצד הצליח אלגוריתם שנועד לאתר מבנים מולקולריים שמתאימים לתרופות חדשות בטכניקות של למידת מכונה לגלות בתוך כך גם תצורות ישנות וחדשות של חומרי לוחמה ביולוגית וכימית.⁷⁸⁰ פגיעה בזכות לחיים עלולה להיגרם גם

Deception, Uncertainty, and Trust in News, 1 SOCIAL MEDIA + Soc. (2020);

Tom Dobber et al., *Do (Microtargeted) Deepfakes Have Real Effects on Political Attitudes?* 26 INT'L J. PRESS/POL. 69 (2021)

777 ס' 2, 4 לחוק יסוד: כבוד האדם וחירותו. ראו גם אהרן ברק "הזכות החוקתית להגנה על החיים, הגוף והכבוד" מבחר כתבים ג - עיונים חוקתיים 547 (2017). על ההבחנה בין האינטרס הציבורי ובין הזכות לחיים ראו אורן גזל-אייל ואמנון רייכמן "אינטרסים ציבוריים כזכויות חוקתיות?" משפטים מא 97 (2011).

778 ראו לדוגמה ס' 2 ב-ECHR, לעיל ה"ש 733; ס' 6 ל-ICCP, לעיל ה"ש 742;

ס' 3 ב-UDHR, לעיל ה"ש 733. ראו גם ELIZABETH WICKS, THE RIGHT TO LIFE AND CONFLICTING INTERESTS 22-47 (2010)

779 לצד תאונות שאפשר לייחס אותן לפגמים בעיצוב המערכות ("באגים"), והיה אפשר למנוע אותן, דוגמת אירוע דריסת הולכת הרגל איליין הרצברג על ידי רכב אוטונומי Lauren Smiley, "I'm the Operator": The Aftermath of a Self-Driving Tragedy, WIRED [8.3.2022], יש גם מקרי קצה חאורטיים שבהם התאונה היא בלתי נמנעת ולאגוריתם המנחה את הרכב האוטונומי אין ברירה אלא לבחור בין תרחישי פגיעה שונים. ראו Nyholm and Smids, לעיל ה"ש 27.

Paul Rosenzweig, *Artificial Intelligence and Chemical and Biological Weapons*, LAWFARE (27.4.2022); Fabio Urbina et al., *Dual Use of Artificial Intelligence-Powered Drug Discovery*, 4 NATURE MACHINE INTELLIGENCE 189 (2022)

מליקויים באבטחת מידע או אבטחת סייבר – חולשות במערכות השליטה בבית חכם או התקפות מכוונות על נתוני אימון של מודלים של דיאגנוסטיקה רפואית.⁷⁸¹

אפשר להבחין בין רמות אוטונומיה שונות של מערכות נשק (כדומה לכלי רכב אוטונומיים). החל בכלי נשק מסורתיים, כמו נשק קר, שמחייבים הפעלה אנושית רציפה, וכלה במערכות אוטונומיות לחלוטין, שיכולות לבחור מטרה ולתקוף אותה.⁷⁸² מאז פיתח ריצ'רד גאטלינג את המקלע הקרוי על שמו, שאפשר ירי אוטומטי רציף של מאות קליעים ללא טעינה ידנית (אף שמנגנון הירי עצמו חייב הפעלה אנושית רציפה), עברו מערכות הנשק האוטונומיות כברת דרך ארוכה.⁷⁸³ כלי נשק בעלי מידה מוגבלת של אוטונומיה – כאלו שיכולים לבחור את המטרה באופן עצמאי או לתקוף אותה באופן עצמאי – נמצאים בשימוש עוד מימי מלחמת העולם השנייה, שבה הפעיל צבא גרמניה מערכות טורפדו שהתבייתו על מטרותיהן באמצעות חיישן אקוסטי.⁷⁸⁴ בעידן שבו הולך ומתפתח ירי תלול מסלול, שמתאפיין במטחי ירי ממספר רב של קני שיגור בו זמנית,⁷⁸⁵ האתגר המורכב כשהוא לעצמו – יירוט רקטה או טיל במהלך מעופם – מתעצם אף יותר לנוכח מתווה השיגור. מערכות הגנה אוטונומיות מסוגלות להתמודד עם סכנה

BERND CARSTEN STAHL, DORIS SCHROEDER, AND ROWENA RODRIGUES, ETHICS OF ARTIFICIAL INTELLIGENCE: CASE STUDIES AND OPTIONS FOR ADDRESSING ETHICAL CHALLENGES 63–69 (2023)

782 אליאב לייבליך ואייל בנבנישתי "לוחמה רובוטית ובעיית הכבילה של שיקול-הדעת" עיוני משפט 67, 72 (2016); Rebecca Crotoof, *The Killer Robots Are Here: Legal and Policy Implications*, 36 CARDOZO L. REV. 1837 (2015); U.S. Department of Defense, Directive No. 3000.09: Autonomy in Weapon Systems 13–14 (2012); INTERNATIONAL HUMANITARIAN LAW AND THE CHALLENGES OF CONTEMPORARY ARMED CONFLICTS 44 (International Committee of the Red Cross, 2015)

PAUL SCHARRE, ARMY OF NONE: AUTONOMOUS WEAPONS AND THE FUTURE OF WAR, Chap. 3 (2018)

784 ש.ס.

785 אמיר קוליק "המודיעין ואתגרי הירי תלול המסלול" צבא ואסטרטגיה 19, 22 (2009).

זו בזמן אמת, אבל מחייבות השגחה של מפעיל אנושי המסוגל להתערב במידת הצורך (למשל, אם המערכת זיהתה בטעות כלי טייס אזרחי).⁷⁸⁶

נוסף על מערכות נשק אוטונומיות לצורכי הגנה, כדוגמת כיפת ברזל, יש גם מערכות נשק אוטונומיות שנועדו לתקיפה אקטיבית.⁷⁸⁷ מערכות משוטטות, שמפטרלות באזור מסוים ומתבייתות באופן אוטונומי על סכנות שזיהו, כבר היו בשימוש מבצעי במלחמה בין ארמניה לאזרבייג'ן ובעימותים אחרים.⁷⁸⁸ אומנם המטרות המסורתיות של מערכות משוטטות הן מערכות נשק אחרות – בין שהן מאוישות בעת תקיפתן ובין שאינן מאוישות – אך לאחרונה נטען שרחפן משוטט היה מעורב בהרג מכוון במבצע בלוב.⁷⁸⁹

חרף הקולות המזהירים מפני כלי נשק אוטונומיים ומעודדות הכרות חרם בינלאומי על מי שיפתח אותם או ישתמש בהם,⁷⁹⁰ יש המטילים ספק בהיתכנות חרם גורף כזה וקוראים לבחון כל מקרה לגופו.⁷⁹¹ לנוכח השימוש המתרחב בכלי נשק אוטונומיים לתכליות התקפיות נראה שפגיעתם בזכות לחיים הולכת וחורגת מגדרי נזק אגבי או מהרג חיילי אויב המותר במסגרת דיני הלחימה. הבעיות

786 Scharr, לעיל ה"ש 783.

787 לייבליך ובנבנישטי, לעיל ה"ש 782.

788 Kelsey Atherton, *Loitering Munitions Preview the Autonomous Future of Warfare*, BROOKINGS (4.8.2021)

789 שם; ראו גם Joe Hernandez, *A Military Drone with a Mind of Its Own Was Used In Combat*, U.N. Says, NPR (1.6.2021)

790 *Autonomous Weapons: An Open Letter from AI & Robotics Researchers*, FUTURE OF LIFE INSTITUTE (28.7.2015); *Lethal Autonomous Weapons Pledge*, FUTURE OF LIFE INSTITUTE; THE REPORT OF THE 2016 INFORMAL MEETING OF EXPERTS ON LETHAL AUTONOMOUS WEAPONS SYSTEMS (LAWS) (Advanced Version), Para. 71 (2016); *Losing Humanity: The Case Against Killer Robots*, HUMAN RIGHTS WATCH (19.11.2012)
 לסקירה של עמדת מרבית מדינות העולם בסוגיה ראו: *Stopping Killer Robots: Country Positions on Banning Fully Autonomous Weapons and Retaining Human Control*, HUMAN RIGHTS WATCH (10.8.2020)

791 Anja Kaspersen, *Should We Ban Weapons that Don't Even Exist Yet?* WORLD ECONOMIC FORUM (20.7.2016); Michael N. Schmitt, *Autonomous Weapon Systems and International Humanitarian Law: A Reply to the Critics*, HARV. NAT'L SEC. J. (5.2.2013)

האינהרנטיות למערכות נבונות – הטיות אלגוריתמיות והיעדר שקיפות⁷⁹² – עלולות להביא לתוצאות אומללות במיוחד בשדה הקרב.

5.8

בינה מלאכותית וזכויות אדם

כפי שראינו בפרק 2, לבינה מלאכותית יישומים פוטנציאליים רבים ויש לה הכוח להשפיע לחיוב על חייהם של בני אדם. לכן כמה מסמכי אתיקה שעניינם בינה מלאכותית כוללים קריאות

להשתמש במערכות נבונות לקידום סולידריות אנושית,⁷⁹³ לקידום קיימות ואיכות הסביבה⁷⁹⁴ או לחיוב באופן כללי.⁷⁹⁵ ואולם בבינה מלאכותית טמון גם סיכון לפגיעה בזכויות אדם. בפרק זה עמדנו בקצרה על פגיעות אפשריות בזכות לפרטיות, בזכות להליך הוגן, בזכות לכבוד אנושי, בזכויות פוליטיות ובזכות לחיים. אפשר להעלות על הדעת פגיעה בזכויות נוספות – למשל, הגבלת חופש התנועה של אזרחים שמערכות בינה מלאכותית לניטור פנרמיות קבעו כי עליהם להיכנס לכידוד⁷⁹⁶ או פגיעה רחבת היקף בזכות לתעסוקה⁷⁹⁷ בשל שיבוש שוק העבודה עקב אוטומציה רחבת היקף. בפרקים הבאים נתאר בהרחבה את האופן שבו מאפיינים אינהרנטיים של מערכות בינה מלאכותית עלולים להביא לידי פגיעה בזכויות.

792 ראו להלן בפרק 6, 7.

793 Luengo-Oroz, לעיל ה"ש 447.

794 AI HELG 2019, לעיל ה"ש 298, בעמ' 12; Floridi et al., לעיל ה"ש 337, בעמ' 696-697. ראו גם שי הרשקוביץ "בינה מלאכותית מנסה להתמודד עם שינוי האקלים. מולה ניצב האדם עצמו" הארץ (22.9.2021).

795 ראו לעיל בסעיף 3.6.

796 ראו לדוגמה פסי' 11 לפסק דינו של השופט עמית בעניין בג"ץ 6732/20 האגודה לזכויות אזרח נ' הכנסת (1.3.2021), פורסם באר"ש).

797 ראו לדוגמה פ' 23 ב־UDHR, לעיל ה"ש 733.

פרק שישי

”אחרי רבים להטות”:
הטיות אלגוריתמיות

—

לכאורה, אחד היתרונות הגדולים של מערכות בינה מלאכותית שתפקידן לקבל החלטות ולנתח נתונים הוא האובייקטיביות שלהן. שהרי קבלת החלטות ממוכנת, המבוססת על נתונים בלבד, אינה מושפעת ממגבלות קוגניטיביות אנושיות. מערכות ממוכנות אינן מתעייפות, אינן מושפעות מרגשות ואינן נתונות ללחצים פסיכולוגיים. השימוש בהן ככלי לקבלת החלטות צפוי אפוא לנטרל הטיות אנושיות רבות.

יתר על כן, במבט ראשון נדמה כי העברת שיקול הדעת האנושי למערכות אוטומטיות יכולה לפתור השפעה לא רצויה של דעות קדומות, המנציחות פערים חברתיים, מגדריים ואתניים. והרי גם מי שנמנע במודע מדעות קדומות בקבלת

ההחלטות שלו יכול להיות מושפע מהן שלא במודע.⁷⁹⁸ מאחר שמערכות בינה מלאכותית אינן אנושיות, ואין להן מטענים חברתיים ורגשיים כמו לבני אדם, יש להן יכולת לסייע ביצירת חברה שוויונית יותר ולפתוח בפני מיעוטים וקבוצות מוחלשות אפשרויות חדשות.

ואולם מן השימוש במערכות בינה מלאכותית עולה שחלק מההטיות האנושיות שמערכות אלו היו אמורות למתן באות לידי ביטוי גם בקבלת החלטות ממוכנת. מאחר שאיכותן של מערכות לומדות מבוססת על טיב הנתונים שהן ניוונות מהם במהלך האימון, ונתונים אלו הם תוצר של פעילות אנושית ולכן משקפים הטיות חברתיות, מגדריות, אתניות ואחרות, אין פלא שהטיות אלו באות לידי ביטוי גם בקבלת ההחלטות הממוכנת והניטרלית לכאורה של מערכות בינה מלאכותית.

יש להבחין בין שני סוגים של הטיות אנושיות – הטיות הנובעות מדעות קדומות פסולות והטיות קוגניטיביות. בחיי המעשה יש והטיות אלו משמשות בערבוביה – דעות קדומות לגבי גזע או מגדר, למשל, יכולות לעיתים להשפיע שלא במודע על שיקול דעת אנושי חף מכוונות רעות. למשל, כשמנהל כוח אדם נדרש לסנן מספר רב של קורות חיים לתפקיד ניהול הוא עלול להישען גם שלא במודע על דעות קדומות אתניות.⁷⁹⁹

798 ראו סקירה הספרות אצל John T. Jost et al., *The Existence of Implicit Bias Is Beyond Reasonable Doubt: A Refutation of Ideological and Methodological Objections and Executive Summary of Ten Studies that No Manager Should Ignore*, 29 RESEARCH IN ORGANIZATIONAL BEHAVIOR 39 (2009)

799 מריאן ברטרנד וסנדיל מולינאתן בחנו אמפירית את קיומה של הטיה דה פקטו בקרב מעסיקים שהעדיפו קורות חיים של בעלי שמות שמעידים על רקע אתני לבן משמות שמעידים מובהק על רקע אתני אפרו-אמריקאי. Marianne Bertrand and Sendhil Mullainathan, *Are Emily and Greg More Employable than Lakisha and Jamal? A Field Experiment on Labor Market Discrimination* 94 AMERICAN ECONOMIC REV. 991 (2004). ממחקר זה לא ברור אם ההטיה לא מודעת, אבל מחקר המשך שנערך בקרב סטודנטים מצא קורלציה בין דעות קדומות להעדפת שמות "לבנים". Marianne Bertrand, Dolly Chugh, and Sendhil Mullainathan, *Implicit Discrimination*, 95 AMERICAN ECONOMIC REV. 94 (2005)

עקב מוגבלות התודעה האנושית פיתחו בני אדם סוגים שונים של היריסטיקות והטיות המסייעות להם לעבד מידע ולקבל החלטות בתנאי אי-

ודאות.⁸⁰⁰ הטיות קוגניטיביות אלו אינן נובעות מדעות קדומות בעניין קבוצות אתניות, מגדר או העדפה מינית, אלא מהמאפיינים הפסיכולוגיים והפיזיולוגיים של עיבוד המידע. את הקשיים הגורמים להטיות קוגניטיביות אפשר לחלק בחלוקה גסה לארבעה:⁸⁰¹

(1) **עודף מידע** – מוח האדם חשוף למידע רב. אין לו ברירה אלא לסנן את הנתונים הרבים הנקלטים בחושיו ולחלץ מהם את המידע הרלוונטי לקבלת החלטה מושכלת. למוח יש מבחר תכסיסים להתמודדות עם עודף מידע. למשל, אנו נוטים להעריך שכיחות של תופעות לפי הקלות שבה אפשר להעלותן על הדעת (היריסטיקות זמינות);⁸⁰² לייחס משמעות מיוחדת לפרטים יוצאי דופן; לתת לשינויים שהבחנו בהם להשפיע על שקלול אומדנים (היריסטיקות עיגון);⁸⁰³ או להתמקד בפרטים המאששים אמונות קודמות שלנו (לדוגמה, הטיית האישור).⁸⁰⁴

(2) **חוסר מידע** – קבלת החלטות נעשית לעיתים מזומנות בתנאי אי־ודאות ובהסתמך על מידע חלקי ביותר. האדם, שמלכתחילה זיכרון העבודה שלו

800 ראו לדוגמה עמוס טברסקי ודניאל כהנמן "שיפוט בתנאי אי־ודאות: יוריסטיקות והטיות" רצינוליות, הוגנות ואושר 45 (מיה בר הלל עורכת, 2005).

801 Buster Benson, *Cognitive Bias Cheat Sheet: Because Thinking Is Hard*, BETTER HUMANS (1.9.2016)

802 ראו טברסקי וכהנמן, לעיל ה"ש 800, בעמ' 53-56. ראו גם מאור אבן חן האיסור הפילי על שימוש במידע פנים: ניתוח כלכלי התנהגותי 380-382 (2008).

803 ראו טברסקי וכהנמן, שם, בעמ' 56-60; אבן חן, שם, בעמ' 382-383; מיה בר הלל ואוריאל פרוקצ'יה "כלכלה התנהגותית" הגישה הכלכלית למשפט 71, 106-108 (עורך אוריאל פרוקצ'יה, 2012).

804 Peter C. Wason, *On the Failure to Eliminate Hypotheses in a Conceptual Task*, 12 QUARTERLY J. OF EXPERIMENTAL PSYCHOLOGY 129 (1960); Raymond S. Nickerson, *Confirmation Bias: A Ubiquitous Phenomenon in Many Guises*, 2 REV. OF GENERAL PSYCHOLOGY 175 (1998)

מוגבל,⁸⁰⁵ מתמודד עם מידע חסר בשלל דרכים. בין השאר, באמצעות הנטייה האנושית למצוא תבניות ודפוסים במדגמים קטנים, ללא כל מובהקות סטטיסטית (למשל, האמונה ב"יד החמה" בכדורסל⁸⁰⁶ או "כשל המהמר"⁸⁰⁷); להשלים פערי מידע בהכללות, בסטריאוטיפים (לרבות דעות קדומות⁸⁰⁸) ובידע היסטורי; להניח הנחות לגבי עמדות הזולת; לפשט הסתברויות ונתונים כמותיים (למשל לערוך חשבונאות נפשית⁸⁰⁹ או לחשוב על כסף במונחים נומינליים ולא ריאליים);⁸¹⁰ להשליך אמונות והנחות מההווה על אירועים מן העבר או העתיד (כגון כשל הראייה בדיעבד⁸¹¹ או מזל מוסרי⁸¹²).

(3) **אילוצי זמן** – גם כשכל המידע הרלוונטי זמין ומסודר, ואין נתונים מיותרים, נסיבות שבהן נדרשת החלטה מהירה יכולות לעורר קיצורי דרך מנטליים או

805 ראו George A. Miller, *The Magical Number Seven, Plus or Minus Two: Some Limits on our Capacity for Processing Information*, 101 *PSYCHOLOGICAL REV.* 343 (1956)

806 תומס גילוביץ', רוברט ולון ועמוס טברסקי הצביעו על פרכת האמונה העממית ב"יד החמה" בכדורסל וייחסו אותה להטיה קוגניטיבית שמתבטאת בניסיון למצוא דפוסים בנתונים אקראיים. Thomas Gilovich, Robert Vallone, and Amos Tversky, *Hot Hand in Basketball: On the Misperception of Random Sequences*, 17 *COGNITIVE PSYCHOLOGY* 295 (1985). עם זאת, מחקרים מאוחרים יותר הראו שחופעה "היד החמה" בכדורסל שרירה וקיימת. ראו למשל, Joshua B. Miller and Adam Sanjurjo, *Surprised by the Gambler's and Hot Hand Fallacies? A Truth in the Law of Small Numbers*, 86 *ECONOMETRICA* 2019 (2015)

807 ראו טברסקי וכהנמן, **לעיל** ה"ש 800, בעמ' 46-50.

808 ראו למשל: DANIEL BAR-TAL AND YONA TEICHMAN, *STEREOTYPES AND PREJUDICE IN CONFLICT: REPRESENTATIONS OF ARABS IN ISRAELI JEWISH SOCIETY* 22-27 (2005)

809 ראו לדוגמה דניאל כהנמן ועמוס טברסקי "בחירות, ערכים והיצגים" **רציונליות, הוגנות ואושר** 64 (מיה בר הלל עורכת, 2005). ראו גם אבן חן, **לעיל** ה"ש 802, בעמ' 379-378.

810 Eldar Shafir, Peter Diamond, and Amos Tversky, *Money Illusion* 112 *QUARTERLY J.ECON.* 341 (1997)

811 אבן חן, **לעיל** ה"ש 802, בעמ' 390; בר הלל ופרוקצ'יה, **לעיל** ה"ש 803, בעמ' 102-106.

812 אלה קורן "התמודדות עם תופעת הטעות – ופרדוקס המזל המוסרי" **רפואה ומשפט** 39, 40-41 (2009).

שימוש בהיריסטיקות והטיות או לעיתים הסתמכות על תחושת ביטחון עצמי (לרוב מופרז).⁸¹³ אילוצי זמן גורמים לנו לבחור בזמין ובמוכר, להעדיף את הפשוט מן המורכב (התער של אוקאם),⁸¹⁴ לתערף השלמת משימות משיקולים של עלות שקועה או להעריך יתר על המידה מוצרים שהשתתפנו בהרכבתם (אפקט איקאה).⁸¹⁵

(4) **מגבלות הזיכרון** – כאמור, מוח האדם מוצף במידע. לצד סינון המידע בזמן אמת, הוא נדרש לעבד אותו לשימוש בטווח ארוך או קצר. בני אדם נוטים לערוך את הזיכרונות שלהם, להתעלם מהפרטים (למשל, כדי ליצור הכללות) ולצמצם אירועים לרכיבי המפתח שלהם. עיבוד זיכרונות מושפע גם מהחוויה שזוכרים ומהרגשות שקשורים אליה.

כל אלו נחשבות להטיות קוגניטיביות מאחר שהן נובעות מהמגבלות הקוגניטיביות הפיזיולוגיות האנושיות. במצבי עייפות ולחץ מוח האדם אינו פועל באופן מיטבי. כמות הנתונים שהוא יכול להחזיק ולעבד במקביל מוגבלת והיא מחייבת תהליכי עיבוד וקיצורי דרך קוגניטיביים. בינה מלאכותית יכולה להתגבר על מרבית ההטיות האלו – היא אינה מתעייפת, כוח העיבוד שלה גדול בעשרות מונים מזה של האדם והיא אינה מושפעת מרגשות. עם זאת, בינה מלאכותית לא בהכרח יכולה להתגבר על בעיות של נתונים חלקיים או חסרים. אומנם יש טכניקות סטטיסטיות שעשויות לסייע בהתמודדות עם מידע חסר, אך יש להן השלכות על אופיו של המודל.⁸¹⁶

נוסף על ההטיות הקוגניטיביות, קבלת החלטות אנושית עלולה להיות מושפעת מדעות קדומות נגד קבוצות מוגנות באוכלוסייה (יש לציין כי קבוצות אלו

813 אבן חן, לעיל ה"ש 802, בעמ' 384–385.

814 Kevin Kelly, *Justification as Truth-Finding Efficiency: How Ockham's Razor Works*, 14 MINDS AND MACHINES 485 (2004)

815 דן אריאלי לא רציונלי אבל לא במקרה 93–115 (2010).

816 Benjamin M. Marlin, *MISSING DATA PROBLEMS IN MACHINE LEARNING* (Ph.D. diss., University of Toronto 2008); Gustavo E. A. P. A. Batista and Maria Carolina Monard, *An Analysis of Four Missing Data Treatment Methods for Supervised Learning*, in *PROCEEDINGS OF THE FIRST INTERNATIONAL WORKSHOP ON DATA CLEANING AND PREPROCESSING* 142 (Shichao Zhang, Qiang Yang, and Chengqi Zhang eds., 2002)

יכולות להשתנות לאורך זמן). לעיתים גם כשיש כללים מוסדיים נגד אפליה מתגלה בדיעבד הטיה מערכתית לא מודעת, שכן קבלת החלטות בשלל מוסדות חברתיים משקפת את ההטיות האנושיות של בעלי התפקידים בהן. כך למשל, הטיות של שופטים עלולות להביא אותם לגזור עונשים חמורים יותר על נאשמים מקבוצות מיעוט חברתיות או אתניות;⁸¹⁷ הטיות של שוטרים עלולות להביא לאכיפת יתר בקרב אוכלוסיות מיעוטים או ליד קלה על ההרק;⁸¹⁸ הטיות ודעות קדומות עלולות להשפיע על שיקולים של בנקאים בכואם להחליט אם לתת הלוואה או משכנתה,⁸¹⁹ או על שיקולים של בעלי נכסים הנמנעים במודע או שלא במודע מלהשכיר אותם לשוכרים מקבוצות מיעוט מסוימות;⁸²⁰ וגם החלטות

817 ראו לדוגמה, Gideon Fishman, Arye Rattner, and Hagit Turjeman, *Sentencing Outcomes in a Multinational Society: When Judges, Defendants and Victims Can Be either Arabs or Jews*, 3 EUR. J. CRIMINOLOGY 69 (2006); Guy Grossman et al., *Descriptive Representation and Judicial Outcomes in Multiethnic Societies*, 60 AM. J. POL. SCI. 44 (2016); Jonathan P. Kastellec, *Race, Context, and Judging on the Courts of Appeals: Race-Based Panel Effects in Death Penalty Cases*, JUSTICE SYSTEM JOURNAL (11.11.2020); Crystal S. Yang, *Free at Last? Judicial Discretion and Racial Disparities in Federal Sentencing*, 44 J. LEG. STUD. 75 (2015); Michele Benedetto Neitz, *Socioeconomic Bias in the Judiciary*, 61 CLEVELAND STATE L. REV. 137 (2013); Norman L. Green, *How Great Is America's Tolerance for Judicial Bias: An Inquiry into the Supreme Court's Decisions in Caperton and Citizens United, Their Implications for Judicial Elections, and Their Effect on the Rule of Law in the United States*, 112 W. VA. L. REV (2020)

818 Joshua Correll et al., *Across the Thin Blue Line: Police Officers and Racial Bias*, in PERSONALITY & SOC. PSYCH. 1006 (2007); Jeremy West, *Racial Bias in Police Investigations* (working paper, 2018); Cody T. Ross, *A Multi-Level Bayesian Analysis of Racial Bias in Police Shootings at the County-Level in the United States, 2011-2014*, 10 PLoS ONE (2015)

819 J. Michelle Brock and Ralph De Haas, *GENDER DISCRIMINATION IN SMALL BUSINESS LENDING: EVIDENCE FROM A LAB-IN-THE-FIELD EXPERIMENT IN TURKEY* (EBRD Working Paper number 232, 2019). להטיות מבניות בשוק המשכנתאות הישראלי ראו לדוגמה בנק ישראל, *ניתוח שוק המשכנתאות ללווים מהמגזר הערבי על רקע הכשלים המבניים בתחום הדיור במגזר זה* (21.11.2017).

820 ראו לדוגמה Terrence McCoy, *Eviction Isn't Just About Poverty. It's Also About Race and Virginia Proves It* THE WASHINGTON POST (11.11.2018). לדוגמה לאפליה גזענית מודעת ראו אור קשתי "צעירות ערביות תבעו מתווך שאמר להן שדירה הושכרה - ולחברותיהן היהודיות אמר ההפך" הארץ (6.12.2019).

של מעסיקים וממונים על כוח אדם בעניין גיוס עובדים יכולות להיות מוטות.⁸²¹ כל ההטיות האלו עלולות להביא לאפליה מודעת ומכוונת כשמקבל ההחלטות "בשטח" מודע להטיות שלו נגד קבוצת המיעוט (ומתכוון להפלותה), אבל קבלת החלטות מפלה יכולה להיות מושפעת גם מהטיות לא מודעות שלו.

6.2

הטיות אלגוריתמיות

גם הרכיב האנושי בפיתוח מערכות בינה מלאכותית יכול להיות מוטה – אם במודע ואם שלא במודע – ויש חשש שהטיות המפתח ישועתקו לתוך המכונה כבר בשלב הפיתוח.⁸²² בפיתוח מערכות מומחה נדרשת מעורבות אנושית כדי לתרגם כללי החלטה אנושיים לשפת מכונה או כדי לפתח מודל סטטיסטי של עץ החלטה. אך גם בפיתוח מבוסס למידת מכונה יש מעורבות אנושית ניכרת בהכנת נתוני האימון, בהגדרת משימת מטרה ובפיקוח על תוצאת האלגוריתם.⁸²³ יתר על כן, המעורבות האנושית בפיתוח בינה מלאכותית, בבקרה עליה ובשימוש בה תורמת לאמון במערכות אלו.⁸²⁴ ההטיות הקוגניטיביות

- Faye K. Cocchiara, Myrtle P. Bell, and Wendy J. Casper, *Sounding "Different": The Role of Sociolinguistic Cues in Evaluating Job Candidates*, 55 HUMAN RESOURCE MANAGEMENT 463 (2016); Stuart W. Flint et al., *Obesity Discrimination in the Recruitment Process: "You're Not Hired!"*, 7 FRONTIERS IN PSYCHOLOGY 647 (2016); Devah Pager, Bruce Western, and Bart Bonikowski, *Discrimination in a Low-Wage Labor Market: A Field Experiment*, 74 AMERICAN SOCIOLOGICAL REVIEW 777 (2009); Geoffrey Beattie and Patrick Johnson, *Possible Unconscious Bias in Recruitment and Promotion and the Need to Promote Equality*, 16 PERSPECTIVES: POLICY AND PRACTICE IN HIGHER EDUCATION 7 (2012)
- JANNEKE GERARDS AND RAPHAËLE XENIDIS, ALGORITHMIC DISCRIMINATION IN EUROPE: CHALLENGES AND OPPORTUNITIES FOR GENDER EQUALITY AND NON-DISCRIMINATION LAW 41-42 (2021)
- Sonia Katyal, *Private Accountability in the Age of Artificial Intelligence*, 66 U.C.L.A. L. REV. 54, 67 (2019); Kate Crawford, *The Hidden Biases in Big Data*, HARV. BUS. REV. (1.1.2013)
- 824 ראו ס' 22 של התקנות הכלליות בדבר הגנת מידע (GDPR); AI HELG 2019, לעיל EU White, 298, בעמ' 12; *Asilomar AI Principles*; לעיל ה"ש 326, עיקרון 16; Aziz Z. Huq, *A Right to a Human Decision*, 21, בעמ' 300, paper 2020 106 VA. L. REV. 611 (2020)

האנושיות, כמו הטיות מודעות או לא מודעות נגד קבוצות מיעוט מוגנות, עלולות לזלוג למערכות שבפיתוחן מעורבים בני אנוש.⁸²⁵

כמו כן, לצד היתרונות של מערכות בינה מלאכותית בקבלת החלטות חפה מהטיות קוגניטיביות אנושיות, יש לזכור שמערכות אלו אינן חסינות מהטיות מתודולוגיות הנובעות מתהליכי הפיתוח שלהן.

מאחר שבינה מלאכותית, ובפרט מערכות המבוססות על למידת מכונה, מבוססת על מודלים סטטיסטיים של המציאות, מרבית ההטיות שלה אינן שונות מהטיות של מחקרים אקדמיים המציעים מודלים כאלו. כך למשל, במדעי החברה נהוג להבחין בין תוקף פנימי של מחקר לתוקף חיצוני שלו. תוקף פנימי הוא קריטריון שבוחן את השאלה אם המחקר תוכנן כך שמסקנותיו הפנימיות נכונות, כלומר אם הוא מבוסס על קשר סיבתי תקף. תוקף חיצוני הוא קריטריון שאומד את היכולת להחיל את מסקנות המחקר על מי שאינו משתייך לאוכלוסיית המדגם.⁸²⁶ ביקורת נפוצה על תוקף חיצוני של מחקרים במדעי החברה המבוססים על שאלוני סטודנטים או על ניסויים בהשתתפותם היא שמדגם זה אינו מייצג את האוכלוסייה הכללית (על אחת כמה וכמה כשזה מדגם של סטודנטים הלומדים מקצוע מסוים).⁸²⁷ מחקרים כאלו יכולים אולי לשקף את מבנה האישיות, העמדות וההתנהלות של סטודנטים בני 18-21 בארצות הברית, אך יש להיזהר מהחלת מסקנותיהם על האוכלוסייה הכללית.

825 Katyal, לעיל ה"ש 823, בעמ' 66.

826 DONALD T. CAMPBELL AND JULIAN C. STANLEY, *EXPERIMENTAL AND QUASI-EXPERIMENTAL DESIGNS FOR RESEARCH* 5-6 (1963); Jeffrey W. Lucas, *Theory-Testing, Generalization, and the Problem of External Validity*, 21 *SOCIOLOGICAL THEORY* 236, 236-238 (2003)

827 David O. Sears, *College Sophomores in the Laboratory: Influences of a Narrow Data Base on Social Psychology's View of Human Nature*, 51 *J. of PERSONALITY AND SOC. PSYCH.* 315; James N. Druckman and Cindy D. Kam, *Students as Experimental Participants: A Defense of the "Narrow Data Base," in* *CAMBRIDGE HANDBOOK OF EXPERIMENTAL POLITICAL SCIENCE* 41 (James N. Druckman, Donald P. Green, James H. Kuklinski, and Arthur Lupia eds., 2011)

ברומה, נתוני אימון שעליהם מתבססת הבינה המלאכותית יכולים להניב מודל שתוקפו הפנימי הוא לעילא ולעילא, אך אין לו תוקף חיצוני. כשנתוני האימון מבוססים על מדגם לא פרופורציונלי של האוכלוסייה (ייצוג יתר או ייצוג חסר), המודל עלול לכלול הטיות לטובתם או לרעתם של מגזרים מסוימים. למשל, חברת אמזון פיתחה כלי בינה מלאכותית שנועד לסייע למחלקת משאבי אנוש של החברה על ידי אוטומציה של תהליכי גיוס של מהנדסים. בסיס הנתונים שהחברה הסתמכה עליו הורכב ברובו מפרופילים של מהנדסים גברים, ונמצא שהכלי שפותח מוטה נגד מועמדות והוא מסנן קורות חיים של נשים.⁸²⁸ ייתכן שגם דוגמאות אחרות של הטיות ברורות מלמדות על מקרים דומים של ייצוג חסר – נמצא שאלגוריתם הפרסום של Google Ads מיעט לפרסם לנשים הצעות לעבודה במשרות יוקרתיות,⁸²⁹ אלגוריתמים לזיהוי תמונה סיווג אנשים כהי עור כגורילות⁸³⁰ ואלגוריתם של מצלמה דיגיטלית זיהה תווי פנים אסיאתיים כממצאים.⁸³¹ בכל הדוגמאות האלה נשלל התוקף החיצוני של המודל שפיתחה הבינה המלאכותית.

6.2.2. שעתוק הטיות מבניות קיימות
 מערכות בינה מלאכותית יכולות לשעתק הטיות מבניות קיימות אם בסיס הנתונים שהן נשענות עליו מוטה. דוגמה מוקדמת היא תוכנה שפותחה בבית הספר לרפואה סנט ג'ורג' בבריטניה בתחילת שנות השמונים של המאה שעברה כדי לסייע לוועדת הקבלה לסנן את המועמדים הרבים ללימודים. מטרת התוכנה הייתה לחקות את דפוסי קבלת ההחלטות של הוועדה וליישם אותם על

Jeffrey Dastin, *Amazon Scraps Secret AI Recruiting Tool That Showed Bias Against Women*, REUTERS (11.10.2018)

Amit Datta, Michael Carl Tschantz, and Anupam Datta, *Automated Experiments on Ad Privacy Settings: A Tale of Opacity, Choice, and Discrimination*, 1 PROCEEDINGS ON PRIVACY ENHANCING TECHNOLOGIES 92 (2015)

830 רועי גולדנברג "פאדיחה של גוגל: מזהה אנשים כהי עור בתור גורילות" גלובס (2.7.2015); איתן לשם "פייסבוק התנצלה אחרי שהבינה שהבינה המלאכותית שלה תייגה אדם שחור כ'פרימט' " הארץ (6.9.2021).

Odelia Lee, *Camera Misses the Mark on Racial Sensitivity* GIZMODO (15.5.2009)

מועמדים חדשים כדי להבטיח עקביות בקבלת המועמדים. התוכנה אכן הצליחה לעשות כן,⁸³² אלא שקבלת ההחלטות ההיסטורית בוועדה האנושית הייתה ספוגה בהטיות נגד נשים ומיעוטים אתניים. כלומר ההטיות המפלות של הוועדה הוטמעו במודל שהתוכנה השתמשה בו.⁸³³ יש לציין שתוכנה זו פותחה על בסיס מידול סטטיסטי "אנושי" ולא באמצעות למידת מכונה. המפתח האנושי של המודל הוא שהכניס לבסיס הנתונים במועד "גזע" ו"מגדר" כמשתנים מסבירים (הם לא היו קיימים בבסיס הנתונים המקורי, אך הוסקו מצירוף של משתנים אחרים, כגון שם ומקום לידה), וכך הצליח לשעתק ברמת דיוק גבוהה את ההטיות של הוועדה האנושית. קל וחומר שביצועיו של אלגוריתם שפותח על בסיס נתוני אימון מוטים ללא התערבות אנושית יהיו מוטים.

6.2.3. הטייות בפיתוח בארוקס וסלבסט מיפו מופעים של אפליה בעיבוד נתוני עתק לפי שלבי המידול הסטטיסטי של הנתונים.⁸³⁴ בדומה להטיות שזיהו השניים בטכניקות סטטיסטיות של כריית נתוני עתק, שאינן אלא זיהוי תבניות או חיזוי בהתבסס על מתאם סטטיסטי, גם הטייות בבנינה מלאכותית יכולות להיווצר לאורך כל שרשרת הפיתוח.

ייתכנו הטייות כבר בשלב הראשוני של הגדרת משתני מטרה, שהם התוצאות של המודל או המשתנים שהבינה המלאכותית מנסה לחזות. לעיתים קרובות מפתחי המודל הם שיוצרים ומגדירים באופן מלאכותי את משתנה המטרה; למשל, אלגוריתם לחיזוי מידת היותם של מועמדים לעבודה עובדים "מוצלחים" מחייב להגדיר באופן מלאכותי לפי פרמטרים כמותיים שונים עובד מוצלח מהו. פרמטרים

Commission For Racial Equality, *MEDICAL SCHOOL ADMISSIONS: REPORT OF A FORMAL INVESTIGATION INTO ST. GEORGE'S HOSPITAL MEDICAL SCHOOL 9* (1988)

Stella Lowry and Gordon MacPherson, *A Blot on the Profession*, 296 *BRITISH MED. J.* 657 (1988); CATHY O'NEIL, *WEAPONS OF MATH DESTRUCTION: HOW BIG DATA INCREASES INEQUALITY AND THREATENS DEMOCRACY* 116-117 (2016)

Solon Barocas and Andrew D. Selbst, *Big Data's Disparate Impact*, 104 *CA. L. REV.* 671 (2016)

אלו יכולים לפעמים להיות מוטים,⁸³⁵ אבל נמצא שהתוצר הסופי של מערכות נבונות שמשתמשות במשתנים קטגוריאליים מלאכותיים כאלה עקבי יותר מתוצר המסתמך על הערכות אנושיות שמסתמכות על אותם משתנים, ויש לו אף פוטנציאל גבוה יותר לתמוך באחריותיות של המודל. כמו כן, קטגוריות יכולות להיות דינמיות ולהשתנות לאורך זמן; האמידי, שורמאן וברנהאם הראו שקיבוע בינארי של משתנה המגדר בבסיסי נתונים של מערכות אוטומטיות לזיהוי מגדר עלול לדכא התפתחות של קטגוריות מגדריות נזילות יותר.⁸³⁶

לעיתים רידוד פונקציית המטרה של המודל למשתנה אחד או משתנים אחדים מתעלם ממורכבות התוצאה הרצויה.⁸³⁷ כך למשל, אימון מערכת לנהיגה אוטונומית שפונקציית המטרה שלה היא מזעור זמני הנסיעה עלול להביא לפיתוח מודל שעל פיו נכון לסכן את יושבי הרכב כדי להגשים מטרה זו. לעיתים שגיאות שמקורן בהתעלמות מפרמטרים מסוימים בזמן פיתוח המודל או הזנחת הצורך להוסיף אילוצים לפונקציית המטרה (למשל מזעור זמני הנסיעה בתנאי שכל יושבי הרכב מגיעים ליעדם בריאים ושלמים) עלולות לייצר תוצאות מוטות.

6.2.4. הטיות בפיתוח בארוקס וסלבסט⁸³⁸ הצביעו על מקור אפשרי אחר להטיה – בסיס נתוני האימון, שכבר עמדנו עליו קודם בכמה דוגמאות. מדעני מחשב מושכים לעיתים בכתפיהם לנוכח תוצאות מפלות ומפטירים: "זבל נכנס – זבל יצא"

Joseph M. Stauffer and M. Ronald Buckley, *The Existence and Nature of Racial Bias in Supervisory Ratings*, 90 J. APPLIED PSYCHOL. 586, 588–89 (2005)

Foad Hamidi, Morgan Klaus Scheuerman, and Stacy M. Branham, *Gender Recognition or Gender Reductionism? The Social Implications of Embedded Gender Recognition Systems*, in PROCEEDINGS OF THE 2018 CHI Os Keyes, *The* CONFERENCE ON HUMAN FACTORS IN COMPUTING SYSTEMS (2018) *Misgendering Machines: Trans/HCI Implications of Automatic Gender Recognition*, PROCEEDINGS OF THE ACM ON HUMAN-COMPUTER INTERACTION, CSCW, Article 88 (2018)

837 ועל כך ראו להלן בסעיף 6.2.7.

838 Barocas and Selbst, *לעיל* ה"ש 834, בעמ' 684–687.

(garbage in, garbage out).⁸³⁹ שנים רבות הפיק אלגוריתם החיפוש של גוגל תוצאות פורנוגרפיות במענה לשאלות על נשים או נערות שכללו מילות חיפוש אתניות ("נערות שחורות", "נשים אסיאתיות"); כלומר בסיס נתוני האימון משקף, ובה בעת מנציה, יחסי כוח חברתיים, פטישיזציה והחפצה היסטוריים.⁸⁴⁰ בסיס נתוני האימון יכול להיות מוטה אפריורית כי הוא נשען על נתונים היסטוריים המשקפים הטיות מערכתיות – בדומה לתוכנת סינון המועמדים לבית הספר לרפואה⁸⁴¹ או לאפליה היסטורית (בארצות הברית למשל יש מתאם סטטיסטי גבוה בין מקום מגורים למוצא אתני עקב פרקטיקות הפרדה גזעית).⁸⁴²

מבחינה זו, לעיתים מערכות מוטות של בינה מלאכותית פשוט משקפות הטיות מבניות בנות זמננו.⁸⁴³ הטיות כאלו יכולות לבוא לידי ביטוי במאפייני השימוש

839 ראו למשל Mayson, לעיל ה"ש 209.

840 ראו Benjamin, לעיל ה"ש 201, בעמ' 69-70 SAFIYA UMOJA NOBLE, ALGORITHMS OF OPPRESSION: HOW SEARCH ENGINES REINFORCE RACISM (2018). לדוגמאות נוספות של מערכת בינה מלאכותית שמייחסת מיניות מוחצנת לנשים מרקע אתני מסוים ראו Karen Hao, *An AI Saw a Cropped Photo of AOC. It Autocompleted her Wearing a Bikini*, MIT TECHNOLOGY REVIEW (29.1.2021); Melissa Heikkilä, *The Viral AI Avatar App Lensa Undressed Me – Without My Consent*, MIT TECHNOLOGY REVIEW (12.12.2022)

841 ראו לעניין זה את ההבחנה של דבורה הלמן בין nonaccuracy affecting injustice לבין accuracy affecting injustice. אי־דיוק באיסוף הנתונים עלול להביא לתוצאות לא צודקות – למשל אם רישום הנתונים מושפע מפרקטיקות לא צודקות או כשנתונים אלו מושפעים מהערכה אנושית מוטה. לעומת זאת, ייתכן גם בסיס נתונים שנאסף בדיוק מרבי, אך בשל הטיות מערכתיות מביא אף הוא לתוצאות לא צודקות. Deborah Hellman, *Big Data and Compounding Injustice*, JOURNAL OF MORAL PHILOSOPHY (2023)

842 על המשמעות האלגוריתמית של פרקטיקות אלו ראו למשל Nicole McConloughe, *Discrimination on Wheels: How Big Data Uses License Plate Surveillance to Put the Brakes on Disadvantaged Drivers*, 18 STAN. J. C.R. & C. L. (2022) 48-49. להיסטוריה המשפטית של הפרדה בדיוור בארצות הברית ראו לדוגמה Charles L. Nier III, *Perpetuation of Segregation: Toward a New Historical and Legal Interpretation of Redlining under the Fair Housing Act*, 32 J. MARSHALL L. REV. 617 (1999)

843 Darius Amilevičius, *Machine Bias and Fundamental Rights*, in SMART TECHNOLOGIES AND FUNDAMENTAL RIGHTS 335, 342 (John-Stewart Gordon ed., Philip Alston, Special Rapporteur on extreme poverty and human rights, brief as amicus curiae before the District Court of the Hague on the case of NJCM c.s./De Staat der Nederlanden (SyRI), case No. C/09/550982/ HA ZA 18/388, September 2019, Para. 82

בשפה מסוימת, לרבות בשכיחות של ביטויים מסוימים. דוגמה להטיה המשעתקת הטיות חברתיות קיימות המגולמות בבסיס נתוני האימון הוא תרגום מכונה. המשפט בעברית "המתכנתת פיתחה מערכת בינה מלאכותית" מתורגם ביישום Google Translate למשפט באנגלית "The programmer developed an artificial intelligence system". כשאנו מבקשים לשוב ולתרגם את המשפט הזה, שבאנגלית הוא ניטרלי מבחינה מגדרית, לשפה העברית, מתקבל (נכון למועד כתיבת שורות אלו) התרגום "המתכנתת פיתחה מערכת בינה מלאכותית". הדבר נובע מנתוני האימון של היישום, המשקפים ככל הנראה את מאפייני השימוש המקוון בעברית מול מאפייני השימוש באנגלית, על ההטיות המגדריות שבהם.⁸⁴⁴

בסיס הנתונים עלול להוביל להטיות גם כיוון שהוא חסר תוקף חיצוני. כפי שראינו לעיל, ייתכנו מקרים שבנתוני האימון לא ניתן די ייצוג למיעוט מסוים באוכלוסייה והמודל שיפחת לפי נתונים אלו יכיל הטיות לרעתו – למשל בגיוס מועמדים לעבודה או בזיהוי פנים שגוי של מיעוטים אתניים. כך למשל, מערכות בינה מלאכותית שמזהות מגדר עלולות לשגות בזיהוי בני תרבויות שבהן האופן שבו המגדר מוצג שונה.⁸⁴⁵

גם ייצוג יתר של אוכלוסייה מסוימת בבסיס נתונים עלול לשקף הטיה מבנית או אפליה היסטורית. נתוני העבריינות בארצות הברית, למשל, שעלולים לגלם הטיות מערכתיות של מערכת האכיפה נגד אפרו-אמריקאים (כגון דיווח מוגבר ואכיפה מוגברת של עבירות בקרב שחורים), עלולים גם להביא – בשל ייצוג יתר של אוכלוסיות מיעוטים בתוכם – להטיות במודלים שמפתחים על בסיסם ומשמים במערכות של חיזוי פשיעה, הערכת סיכון ורצידיביזם.⁸⁴⁶

844 כך גם הרגומם העברי של שמוח בעלי המקצוע kindergarten teacher ו-nurse הוא אחות וגננת, בהתאמה. לעבודות בנושא ראו למשל Tolga Bolukbasi et al., *Man Is to Computer Programmer as Woman Is to Homemaker? Debiasing Word Embeddings*, 29 *ADVANCES IN NEURAL INFORMATION PROCESSING SYSTEMS* 4349 (2016); Aylin Caliskan et al., *Semantics Derived Automatically from Language Corpora Contain Human-like Biases*, 356 *SCIENCE* 183 (2017)

845 Hamidi, Scheuerman, and Branham, לעיל ה"ש 836.

846 ראו בהקשר זה Barocas and Selbst, לעיל ה"ש 834, בעמ' 687. וירג'ינייה יובנקס טוענת שדגימת יתר של אוכלוסיות מעוטות יכולה עלולה לגרום לשגיאות מסוג 2 (false negatives) במערכות לחיזוי התעללות או הזנחה בילדים, כך שהן לא יאמרו

6.2.5. הטיות בפיתוח

(3) "לנלוך" מכון

של נתוני האימון

אך לא תמיד נתוני האימון הם נתונים היסטוריים מן המוכן, וחלק מפיתוח המודל הוא ייצורם – בייחוד בלמידה על ידי חיזוקים, המבוססת על אימון באמצעות אינטראקציות עם סביבה

מבוקרת, לפי חוקים קבועים מראש. אם סביבת האימונים רוויה בהטיות גזעניות, מיוזיניות ואחרות (אם במתכוון ואם בשל היעדר בקרה נאותה), סביר שהטיות אלו יוטמעו במערכת. כך למשל, חברת מיקרוסופט הכניסה לשימוש בטוויטר צ'אט בוט ששמה טאי (Tay) כדי להדגים את יכולות התקשורת של הבינה המלאכותית שפיתחה החברה. הפעילות של טאי הופסקה בתוך כיממה מאחר שהצ'אט בוט התחילה לפרסם ציוצים גזעניים, מיוזיניים ואנטישמיים.⁸⁴⁷ מיקרוסופט ייחסה את התנהגותה של טאי למתקפות מכוונות של גולשים,⁸⁴⁸ אבל יש הרואים בכך היתממות, שכן החברה הייתה יכולה לצפות התנהגות כזאת כשהכניסה לשימוש בוט בסביבה רוויה בשיח שנאה⁸⁴⁹ והייתה יכולה לכלול בו מנגנונים לסינון תוכן בעייתי.⁸⁵⁰ ואולם לא ברור אם לקחי תקרית טאי הופקו במלואם. חמש שנים לאחר המקרה הכניסה חברה דרום קוריאנית לשימוש בפייסבוק צ'אט בוט מבוסס בינה מלאכותית ששמה לודה ובתוך

ילדים בסיכון משפחות מבוססות. – VIRGINIA EUBANKS, AUTOMATING INEQUALITY: HOW HIGH-TECH TOOLS PROFILE, POLICE, AND PUNISH THE POOR, Chapter 4 (2018)

847 "הצעות מגוונות ושבתים להיטלר בטוויטר: הציוצים של הרובוטיקה הכושלת של מייקרוסופט" *TheMarker* (26.03.2016). להרחבה ראו Marty J. Wolf, K. W. Miller, and Frances S. Grodzinsky, *Why We Should Have Seen That Coming: Comments on Microsoft's Tay "Experiment," and Wider Implications*, 1 ORBIT JOURNAL (2017); Gina Neff and Peter Nagy, *Talking to Bots: Symbiotic Agency and the Case of Tay* 10 INT'L J. COMMUNICATION 4915 (2016)

848 Scott Deveau and Jing Cao, *Microsoft Apologizes After Twitter Chat Bot Experiment Goes Awry*, BLOOMBERG (25.3.2016)

849 דוגמה מן העת האחרונה: מערכת בינה מלאכותית שאומנה על בסיס נתונים של פוסטים באחד הפורומים הנודעים לשמצה בשיח הרעיל שלהם ברשת 4chan הפיקה רבות פוסטים גזעניים ואנטישמיים. Matthew Gault, *AI Trained on 4Chan Becomes "Hate Speech Machine,"* VICE (8.6.2022)

850 Neff and Nagy, *לעיל* ה"ש 847, בעמ' 4922; Ethan Chiel, *Who Turned Microsoft's Chatbot Racist? Surprise, It Was 4chan and 8chan*, SPLINTER (24.3.2016)

שבועות מספר אימנו אותה הגולשים לנקוט שיח שנאה שאילץ את מפעילה להסיר אותה מהרשת.⁸⁵¹

6.2.6. הטיות בפיתוח הישענות על נתונים שמקורם במשתמשים עלולה לייצר הטיית הפניה (referral bias). מודלים מסוימים מסתמכים על דיווחי

משתמשים כדי להעריך סיכונים עתידיים – למשל, מודלים לחיזוי פשיעה (predictive policing) שמסתמכים על דיווחים וקריאות של תושבים כדי להגדיר תאי שטח גאוגרפיים מסוכנים;⁸⁵² או מודלים לזיהוי משפחות שיש בהן סיכון מוגבר להתעללות בילדים שמסתמכים על דיווחים של שכנים.⁸⁵³ מודלים מסוג זה עלולים לסבול מהטיות אופייניות: דיווחים כאלו לרוב שכיחים יותר באזורים עירוניים צפופים (בעלי מאפיינים חברתיים-כלכליים נבדלים) והם מתבססים לעיתים על הערכה סובייקטיבית ובלתי מקצועית של סיכון. יש אפילו מקרים שבהם נמסרים במזיד דיווחי שווא (למשל בגלל סכסוכי שכנים).

ואכן, אופן איסוף הנתונים יכול להשפיע על התוקף החיצוני של המודל. ואולם לעיתים יש הטיות שאינן נובעות משיבוש מכון של בסיס הנתונים או מתוכן פסול, אלא מתת-ייצוג או ייצוג יתר של חלק מהאוכלוסייה. לעיתים קשה לחזות מראש כיצד יוטה איסוף נתונים שמסתמך על חוכמת ההמונים עקב התפלגות המשתמשים. למשל, עיריית בוסטון שחררה יישומון לזיהוי פסיכי של בורות בכבישים, בהסתמך על נתוני איכון לווייני של המכשיר. כשהיישומון מזהה שנפער בור בכביש הוא מדרווח בזמן אמת לעירייה כדי שהיא תוכל לתקנו וכך חוסך ממנה להפעיל צוותי פיקוח עירוניים. ואולם ליישומון זה יש פוטנציאל חריף מוגבל ואין הוא יכול להגיע לכל אוכלוסיית היעד של שירותי תיקוני

Dongwoo Kim, *Chatbot Gone Awry Starts Conversations About AI Ethics in South Korea*, THE DIPLOMAT (16.1.2021)

852 ראו לעיל בסעיף 2.3.2.

853 ראו לדוגמה Mary E. Rauktis and Julie McCrae, *THE ROLE OF RACE IN CHILD WELFARE SYSTEM INVOLVEMENT IN ALLEGHENY COUNTY* (Allegheny County Dept. of Human Services, 2010). ובהרחבה ראו EUBANKS, לעיל ה"ש 846.

הדרכים. סביר להניח שיש חלקים באוכלוסיית היעד שאינם יודעים על קיומו של היישומון; ומתוך מי שידועים עליו, לא כולם ימהרו להתקין אותו (למשל עקב אוריינות דיגיטלית לקויה). אזורים שמתגוררים בהם אזרחים שלא התקינו את היישומון יסבלו אפוא מרמת תחזוקה לקויה של הכבישים.⁸⁵⁴ מאפייני אוכלוסיית המשתמשים בטלפונים חכמים לא בהכרח זהים למאפייני האוכלוסייה הכללית.⁸⁵⁵ גם בישראל הסתמכות על מידע שמקורו במשתמשי יישומונים בטלפוני חכמים עלולה לייצר עיוותים בשל דפוסי השימוש של המגזר החרדי בטלפונים.⁸⁵⁶

טימינט גברו טוענת ששימוש במודלים שפותחו עם נתוני אימון מוטים, ואחת היא אם מדובר בהטיות היסטוריות, בהטיות מבניות או בתת-ייצוג או ייצוג יתר של קבוצות מיעוט, עלול לייצר אפקט של משוב שמגביר, מנציח ומחזק הטיות קיימות.⁸⁵⁷

6.2.7 הטייות בפיתוח מקור נוסף להטיות שמתואר אצל בארוקס (5) בחירת משתנים (feature) וסלבסט הוא בחירת משתנים (selection).⁸⁵⁸ על פי רוב, בתהליך הכנת בסיס נתוני האימון מושמטים חלק מסוגי המשתנים הזמינים – שימוש בכל המשתנים לתיאור המציאות עלול לכלבל את האלגוריתם הלומד ולייצר כללים מורכבים מאוד שלא לצורך.⁸⁵⁹ תיאור רדוקטיבי של המציאות – ובסיסי נתונים מלכתחילה

854 ראו Crawford, *Exploiting Big Data from Mobile Device Sensor-Based Apps: Challenges and Benefits*, 12 MIS QUARTERLY EXECUTIVE 179, 182 (2013). ראו גם Daniel E. O'Leary,

855 Kate Crawford, *Think Again: Big Data*, FOREIGN POLICY (10.5.2013)

856 ראו למשל שוקי פרידמן *שוק הסלולר הכשר: מחור שחור לאסדרה מאוזנת* 31, בה"ש 42 (הצעה לסדר 34, המכון הישראלי לדמוקרטיה 2020).

857 Timnit Gebru, *Race and Gender*, in THE OXFORD HANDBOOK OF ETHICS OF AI 253 (Markus d. Dubber, Frank Pasquale, and Sunit Das eds., 2020)

858 Barocas and Selbst, *לעיל* ה"ש 834, בעמ' 688-690.

859 Ke Wang and Suman Sundaresh, *Selecting Features by Vertical Compactness of Data*, in FEATURE EXTRACTION, CONSTRUCTION AND SELECTION: A DATA MINING PERSPECTIVE 71-72 (Huan Liu and Hiroshi Motoda eds., 2012)

עושים רדוקציה של המציאות לשדות נתונים מסוימים – עלול לגרום לעיוותים ולהטיות בהמשך הדרך.⁸⁶⁰ הסתמכות על מודלים סטטיסטיים בהכרח נדרשת להכללות – ואלו עלולות לעיתים לגרום למושאים האנושיים של המודלים הסטטיסטיים האלו להרגיש שאין מכירים בערכם כאינדיבידואלים.⁸⁶¹ בחירת המשתנים משמיטה לעיתים משתנים שיכולים לסייע בהבחנה בין חברים שונים בקבוצה מוגנת⁸⁶² ומקשה על חשיפת הבדלים שמצביעים על שונות מהותית בין מופעים שונים במודל. כך למשל, כאשר מומחי כוח אדם רואים במקום הלימודים של מועמדים לעבודה אינדיקציה לכישוריהם המקצועיים על אף הרדוקטיביות של מאפיין זה (ואינם עורכים מבחנים אישיים לכל מועמד, שעלותם ניכרת), ייתכן שייפסלו מועמדים שלמדו במוסדות שהמוניטין שלהם אינו מיטבי גם אם לאותם מועמדים רמה מקצועית גבוהה.⁸⁶³

בחירת משתנים אינה רק מצמצמת לעיתים את מורכבות מושאי המודל ומטשטשת הבחנות לרוונטיות ביניהם. ייתכנו גם מצבים של גזענות רציונלית – מצבים שבהם המודל בוחר במשתנים פסולים מאחר שהם מנבאים טובים של משתנה המטרה או שיש מתאם גבוה בין המשתנים הפסולים למשתנים שחסרים בבסיס הנתונים. הושמעה למשל טענה שבהיעדר נתונים על עברו הפלילי של מועמד לעבודה מעסיקים בארצות הברית עלולים להניח שצבע עורו של המועמד הוא

860 ראו בהקשר זה את משל המפה והטריטוריה, שאוזכר לראשונה אצל קורזיבסקי:

ALFRED KORZYBSKI, SCIENCE AND SANITY: AN INTRODUCTION TO NON-ARISTOTELIAN SYSTEMS AND GENERAL SEMANTICS 57-58 (5th ed. 1995). ראו גם חורחה לואיס בורחס, דברי ימי חובעים העולם 119 (רנה ליטוין מתרגמת, 1988); ז'אן בודריאר סימולקרות וסימולציה 7-9 (תרגום אריאלה אזולאי, 2007).

Laurence Thomas, *Statistical Badness*, 23 J. Soc. Phil. 30 (1992); 861

Kasper Lippert-Rasmussen, "We Are all Different": *Statistical Discrimination and the Right to Be Treated as an Individual*, 15 J. Ethics 47 (2011). בהקשר זה ראו גם ריבלין ושני, לעיל ה"ש 205.

Toon Calders and Indrė Žliobaitė, *Why Unbiased Computational* 862

Processes Can Lead to Discriminative Decision Procedures, in DISCRIMINATION AND PRIVACY IN THE INFORMATION SOCIETY: DATA MINING AND PROFILING IN LARGE DATABASES 43, 46 (Bart Custers, Toon Calders, Bart Schermer, and Tal Zarsky eds., 2013)

Matt Richtel, *How Big Data*. 689 בעמ' 834, לעיל ה"ש 834, Barocas and Selbst 863 *Is Playing Recruiter for Specialized Workers*, THE NEW YORK TIMES (27.4.2013)

אומדן להיותו עבריין כיוון ששיעור ההרשעות בקרב בעלי צבע עור זה הוא גבוה.⁸⁶⁴

רווא בנג'מין מתארת חברת הייטק בשם Diversity, Inc., המספקת שירותי פילוח אתניים. לקוחותיה, לרוב חברות גדולות המנועות לפי חוק מאיסוף נתונים אלו, מעבירות ל-Diversity, Inc. את שמות המשתמשים או הלקוחות שלהם, את כתובותיהם ופרטים מזהים נוספים, ועל פי נתונים אלו החברה מזהה את המוצא של כל משתמש או לקוח ומגיעה, כך נטען, לרמת דיוק של 96%⁸⁶⁵ שירותי הזיהוי האתני של Diversity, Inc. נועדו לאפיין קהלי יעד ייחודיים לפרסומות ומתאמות אישית, אך אותן מתודולוגיות יכולות להיות מיושמות בעקיפין, באופן אוטומטי, במודלים אחרים – ולגרום להטיה שלהם.

חלק מסוגי המשתנים שיכולים להיבחר בהליך של בחירת משתנים הם משתנים חליפיים, שהם גורם נוסף להטיות שבארוקס וסלבסט מונים.⁸⁶⁶ משתנים חליפיים הם משתנים שיש מתאם גבוה בינם ובין משתנים בלתי תלויים במודל. כך למשל, המפתחים של המודל לסינון מועמדים לבית הספר לרפואה סנט ג'ורג' זיהו כי מוצא אתני הוא משתנה המשפיע על הסיכוי שוועדת הקבלה תקבל מועמד פוטנציאלי ללימודים, אך מאחר שנתון זה לא היה במערכת הם השתמשו בשילוב של משתנים חליפיים (שם המועמד ומקום הלידה שלו) כדי להסיק את מוצאו האתני – בלי להיעזר בשירותי מיקור חוץ דוגמת אלו שמציעה Diversity, Inc.⁸⁶⁷

Lior Jacob Strahilevitz, *Privacy Versus Antidiscrimination*, 75 864 U. CHI. L. REV 363, 364 (2008); Amanda Agan and Sonja Starr, *Ban the Box, Criminal Records, and Racial Discrimination: A Field Experiment*, 133 THE QUARTERLY JOURNAL OF ECONOMICS 191 (2018)

JANNEKE GERARDS AND BENJAMIN 865, לעיל ה"ש 201, בעמ' 85-87.

JANNEKE GERARDS AND BAROCAS AND SELBST 866, לעיל ה"ש 834, בעמ' 691; RAO G. & RAPHAËLE XENIDIS, ALGORITHMIC DISCRIMINATION IN EUROPE: CHALLENGES AND OPPORTUNITIES FOR GENDER EQUALITY AND NON-DISCRIMINATION LAW 44-45 (2020)

COMMISSION FOR RACIAL EQUALITY 867, לעיל ה"ש 832, בעמ' 9.

לעיתים השימוש במשתנה החליפי נעשה שלא במודע: כדיעבד מתחורר כי המודל נתן משקל גבוה למשתנה חליפי שיש מתאם גבוה בינו ובין משתנה המאפיין קבוצה מוגנת. לדוגמה, כמה ממערכות שירותי הרווחה בהולנד הסתמכו על כתובת (מיקוד) כדי לחזות הונאה פוטנציאלית. ואולם ייתכן שיהיה מתאם גבוה בין משתנה זה ובין מהגרים, ובפרט מהגרים מטורקיה ומרוקו, שכן הם נוטים להתרכז בשכונות מסוימות בערים הגדולות.⁸⁶⁸

כשמשתנים מסוימים מושמטים בשל הוראות סטטוטוריות (מודלים להערכת סיכוני אשראי קמעונאי, למשל, מנועים מלהביא בחשבון משתנים מסוימים: "מינו, גילו, נטייתו המינית, גזעו, דתו, ארץ מוצאו, לאומיותו, מקום מגוריו ומצבו המשפחתי או הבריאותי של לקוח"),⁸⁶⁹ ייתכן שימוש במשתנים חליפיים שיש מתאם גבוה ביניהם ובין המשתנים האסורים שהושמטו. כשמשתנים חליפיים אלו משמשים במקום המשתנים האסורים, עם הזמן אפשר שגם הם יהפכו למשתנים אסורים כשהם לעצמם.⁸⁷⁰ כך למשל, בקליפורניה נאסר לחשב פרמיות ביטוח רכב בהסתמך על הכתובת ודירוג האשראי של המבוטח, שכן יש מתאם גבוה בין נתונים אלו ובין קטגוריות אסורות כמו הכנסה ומוצא אתני.⁸⁷¹

- Nelleke Hijmans, *PROFILING THE WELFARE STATE: UPHOLDING OR UPDATING HUMAN RIGHTS STANDARDS? A CASE STUDY OF THE NETHERLANDS* 42 (Master Thesis, European Rights Standards) (Master's Programme in Human Rights and Democratisation 2017 Netherlands Juristen Comite voor de Mensenrechten בעניין: *tegen Staat der Nederlanden* [Netherlands Jurists Committee of Human Rights v. State of the Netherlands], *Rechtbank Den Haag* [The Hague District Court], C/09/550982/HA ZA 18-388 (5.02.2020), Para. 6.92
- 869 ראו לדוגמה סי' 51 לחוק נחוניי אשראי.
- Devin G. Pope and Justin R. Sydnor, *Implementing Anti-Discrimination Policies in Statistical Profiling Models*, 3 *AMERICAN ECONOMIC JOURNAL: ECONOMIC POLICY* 206, 210 (2011)
- 871 שם, שם. ראו גם Stephen D. Sugarman, *California's Insurance Regulation Revolution: The First Two Years of Proposition 103*, 27 *SAN DIEGO L. REV.* 683, 693-694 (1990); Anya E. R. Prince and Daniel Schwarz, *Proxy Discrimination in the Age of Artificial Intelligence and Big Data*, 105 *IOWA L. REV.* 1257, 1306-1310 (2020)

6.2.9. הטיית בפיתוח גורם נוסף שיכול לגרום להטיות במודלים של בינה מלאכותית הוא טריוויאלי – הטיה מכוונת.⁸⁷² מקבלי החלטות המעוניינים להפלות קבוצות מוגנות יכולים במודע להשפיע על תהליך הפיתוח של המודל – כלומר לדאוג שיביא לתוצאות מפלגות ולהסוות זאת ב"קופסה השחורה" של המכונה (masking); בין השאר, באמצעות בחירת משתנים חליפיים שיש מתאם גבוה בינם ובין הקבוצה המוגנת (לדוגמה, שימוש בקריטריון של שירות צבאי כתנאי סף להעסקת מועמדים לעבודה, כשהדבר אינו מתחייב מאופי התפקיד, כדי להסוות הבחנה בין מועמדים יהודים לערבים);⁸⁷³ או באמצעות ניצול הטיית שמקורן בייצוג יתר או ייצוג חסר של מגזרים מסוימים. עם זאת, במקרים מסוימים, נוכח עלויות הפיתוח של מערכות אלגוריתמיות והקושי המובנה בהוכחת אפליה מכוונת,⁸⁷⁴ ספק אם יש ביקוש רב למערכות שתכליתן להסתיר פרקטיקות אסורות אלו.

בדומה להטיות קוגניטיביות אנושיות, גם כשהמעורבים בשרשרת הפיתוח של מערכות נבונות מודעים לאפשרות של הטיית אלגוריתמיות קשה למזער את התופעה. יש לכך כמה סיבות.

6.3 אתגר ההטיות האלגוריתמיות

6.3.1. קושי לזהות מראש מקצת הדוגמאות בפרק זה הן של הטיית שלא היכן יופיעו הטיית. היו רצויות ומפתחי המערכת לא תכננו אותן. למשל האלגוריתם לזיהוי תמונה שזיהה בעלי צבע עור שחור כגורילות; או המערכת לזיהוי מצמוצים במצלמות דיגיטליות,

872 Barocas and Selbst, לעיל ה"ש 834, בעמ' 692-693.

873 ראו בעניין זה ח"פ (אזורי תל אביב) 1038/99 מדינת ישראל, משרד העבודה והרווחה נ' תפקיד פלוס בע"מ (פורסם בנבו, 12.6.2003).

874 Linda Hamilton Krieger, *The Content of Our Categories: A Cognitive Bias Approach to Discrimination and Equal Employment Opportunity*, 47 STAN. L. REV. 1161, 1177 (1995)

שזיהתה תווי פנים אסיאתיים כממצמציים. פעמים רבות ההבחנה בין בסיס נתונים חלקי להטיות היא מטושטשת.

יתר על כן, גם כשנתוני האימון של המערכת אינם מתויגים כפרמטרים שהחוק פוסל – כגון מגדר, שיוך אתני או אמונה דתית ופוליטית – המערכות עלולות לזהות משתנים חליפיים שיש מתאם גבוה בינם ובין פרמטרים אלו. בהתחשב בכך שדעות קדומות אנושיות, הטיות מערכתיות או אפליה היסטורית עלולות להיות מוטמעות בנתוני האימון, לא תמיד אפשר לחזות מראש אילו הטיות טמונות באלגוריתמים.

6.3.2. רמת הפשטה מערכות אלגוריתמיות, כמו כל מודל סטטיסטי, עושות הפשטה מסוימת של המציאות.⁸⁷⁵ יישומן בהקשרים חברתיים רחבים דורש שימוש בהנחות חיצוניות למודל. כך למשל, ממתאם סטטיסטי לא בהכרח אפשר להסיק סיבתיות, ולכן הסקת סיבתיות אפשרית רק בנסיבות שבהן האנליסט מניח הנחות החורגות מהנחות היסוד הנדרשות לחיזוי.⁸⁷⁶ כדי לבחון אם באלגוריתם יש הטיות יש צורך אנליטי לצאת מגבולות המודל (שהטמיע את ההטיות) ולבחון השערות שחיצוניות לו.⁸⁷⁷ ההפשטה של מודלים סטטיסטיים מאפשרת לכאורה להשתמש בהם בהקשרים חברתיים שונים מאלו שלשמם אומנו,⁸⁷⁸ או במנותק מהמשמעות המורכבת והמלאה של המושגים החברתיים שהמודלים מתארים. אלו הן למעשה השלכות השימוש בנתוני אימון חלקיים (שיכולים להיות בעלי תוקף פנימי לגבי אוכלוסיית היעד של המודל).

875 ראו לעיל ה"ש 830.

Susan Athey, *Beyond Prediction: Using Big Data for Policy Problems*, 876 355 SCIENCE 483 (2017)

Andrew D. Selbst et al., *Fairness and Abstraction in Sociotechnical Systems*, FAT* '19: PROCEEDINGS OF THE CONFERENCE ON FAIRNESS, ACCOUNTABILITY, AND TRANSPARENCY 59, 60 (2019)

878 שם, בעמ' 61.

6.3.3. הגדרה דינמית ושנויה במחלוקת של מושג ההוגנות
 הגדרה זו מקשה גם היא על איתור הטיות. אומנם אינטואיטיבית למושג ההוגנות יש קסם רב, אבל כבר מקדמת דנא לאנשים שונים יש תפיסה שונה של הוגנות.⁸⁷⁹ כמו כן, זיהוין של קבוצות מוגנות הוא עניין דינמי ומשתנה לאורך זמן. יתר על כן, שיקול ההוגנות יכול להתחרות בשלב הפיתוח בשיקולים אחרים, כגון יעילות המודל, מה שמאלץ את המפתחים לקבל הכרעה ערכית באשר לנכונות לתקן אלגוריתם לשם הגשמת ערכים של שוויון והוגנות.

6.3.4. מונח שוק המפריד בין נתונים לאלגוריתמים
 כפי שראינו לעיל, מרבית ההטיות האלגוריתמיות מקורן בנתונים המשמשים לאימון. לרוב מקורם של נתונים אלו אינו בצוותי הפיתוח של המערכות הנבונות אלא במאגרי מידע חיצוניים, הניזונים מנתונים שאנשים מייצרים באופן פסיבי בשגרת יומם או ממערכות תיעוד ציבוריות שהתיעוד בהן אינו בהכרח ממוכן. כבר עמדנו על כך שגם כשמפתחים מפעילים מאמץ מודע ומכוון לטהר את הנתונים מהטיות, למשל על ידי הימנעות מהסתמכות על פרמטרים "נגועים", הנתונים עלולים לשקף הטיות סמויות שהאלגוריתם יצליח לזהות ולהחצין (לחלופין, טיהור יתר של האלגוריתם עלול להביא לתוצאות בלתי מדויקות בעליל כשמשתנים חליפיים מגלמים בערבוביה גם הטיה מערכתית אך גם מידע שתורם לנכונות המודל).

879 ראו למשל Thomas B. Nachbar, *Algorithm Fairness, Algorithmic Discrimination*, 48 Fla. St. U. L. Rev. 523-525 (2021); Stefan Feuerriegel, Mateusz Dolata, and Gerhard Schwabe, *Fair AI: Challenges and Opportunities*, 62 Bus. Inf. Syst. Eng. 379, 379 (2020); Jessie Finocchiaro et al., *Bridging Machine Learning and Mechanism Design Towards Algorithmic Fairness*, FAccT '21: PROCEEDINGS OF THE 2021 ACM CONFERENCE ON FAIRNESS, ACCOUNTABILITY, AND TRANSPARENCY 489, 491-492 (2021); Doaa Abu-Eljounes, *Contextual Fairness: A Legal and Policy Analysis of Algorithmic Fairness*, 20 J. L. Tech & Pol. 1 (2020)

6.4 סיכום

פרק זה סקר את אתגר ההטיות האלגוריתמיות והצביע על הדרכים השונות שבהן הן באות לכדי מימוש ועל אבני הנגף בדרך לפתרון הבעיה. כיצד אפשר אפוא להתמודד עם הטיות? באירופה, למשל, מוצע להבחין בין סוגים שונים של מערכות נבונות לפי רמת הסיכון הנשקפת מהן.⁸⁸⁰ הבחנה זו מאפשרת לשמור על יחס מידתי בין אמצעי למטרה – ובהתאם, המחוקק האירופי אוסר אפריורית על סוגים מסוימים של מערכות נבונות. אין הוא מוכן לקבל את הנזק שהן עלולות להסב, לרבות תוצאות של הטיות אלגוריתמיות, כרע במיעוטו. בפרק 9 להלן יתוארו אסטרטגיות התמודדות נוספות עם הטיות אלגוריתמיות.

פרק שביעי

שקיפות אלגוריתמית

—

כפי שראינו בפרקים הקודמים,⁸⁸¹ בשנים האחרונות התעוררו חששות בדבר ליקויים בתפקוד של מערכות נבונות בעקבות התבססותן על מערך נתונים מוטה או שגוי, נטייתן להעצים הטיות קיימות והיעדר יכולת לאתר תקלות בתפקודן עקב הקושי להבין כיצד הגיעו לתוצריהן (מה שמכונה לעיתים "הקופסה השחורה").⁸⁸²

881 פרק זה לקוח מחלקו הראשון של מחקרם של גדי פרל וטהילה שוורץ אלטשולר מודל ליצירת שקיפות אלגוריתמית (הצעה לסדר 47, המכון הישראלי לדמוקרטיה 2022).

882 Matthew Hutson, *The Opacity of Artificial Intelligence Makes it Hard to Tell When Decision-Making Is Biased*, 58 IEEE SPECTRUM 40 (2021).

הטענה הרווחת היא שכדי להתגבר על אי-בהירותן של מערכות לומדות ועל הסכנות הפוטנציאליות הטמונות בהן נחוצה שקיפות אלגוריתמית. מחברות הטכנולוגיה נדרש להסביר לנו איך בדיוק האלגוריתמים שלהן עובדים ממש כשם שמחברת קוקה קולה נדרש לחשוף את המתכון הסודי שלה. ואילו התביעה מממשלות היא לא לאפשר למשטרה להשתמש במערכות זיהוי פנים עד שלא יסבירו לנו איך האלגוריתם פועל.

7.1

המושג שקיפות אלגוריתמית הוא מושג "עכור"

בתחום הרגולציה של מערכות בינה מלאכותית, ובפרט בתחום האתיקה, מרכיב לעסוק ברכיב העכירות המובנה במערכות בינה מלאכותית.⁸⁸³ רכיב זה מקשה את הבנת פעולתן ואת הבקרה על

תפקודן התקיין. כדי לפתור את חוסר הבהירות המובנה במערכות אלגוריתמיות לקבלת החלטות וכדי לקדם אחריותיות, מוצע בספרות המושג "שקיפות אלגוריתמית". התפיסה הרווחת היא שעכירות יוצרת קשיים באיתור הגורם האחראי על תקלות ואף עלולה להעצים את הנזק הנגרם מפעילותה של המכונה, מאחר שהמשתמשים אינם מסוגלים להבין מה היה צריך להיות השימוש הנכון או לאתר את המקור לתקלה.

ואכן, כפי שראינו בחלק 3.1 לעיל, סוגיית השקיפות האלגוריתמית עולה כמעט בכל המסמכים העוסקים ברגולציה של מערכות אלגוריתמיות והיא חלק מכריע ממארג הפתרונות המוצע.⁸⁸⁴

ואולם חרף העיסוק הנרחב במושג זה, הוא עצמו "עכור", כלומר אינו נהיר דיו. הגדרתו, מקומו ביחס לערכים אחרים, זיהוי השחקנים שמוטלת עליהם

על הנזק החברתי של מערכות אלו ראו הערת הִנְחָה המיוחדת לעוני קיצוניי וזכויות אדם של האומות המאוחדות, REPORT OF THE SPECIAL RAPPORTEUR ON EXTREME POVERTY AND HUMAN RIGHTS (A/74/48037, 2019); *World Stumbling Zombie-Like into a Digital Welfare Dystopia, Warns UN Human Rights Expert*, UN OCHR (17.10.2019)

883 על עכירות אלגוריתמית ראו: VIRGINIA DIGNUM, RESPONSIBLE ARTIFICIAL INTELLIGENCE: HOW TO DEVELOP AND USE AI IN A RESPONSIBLE WAY 59 (Springer Nature 2019)

884 Fjeld et al., לעיל ה"ש 292, בעמ' 15.

חובת השקיפות, זהות השחקנים שנדרשת שקיפות כלפיהם ומקומה של חובת השקיפות האלגוריתמית לאורך שרשרת הערך או חיי המוצר – כל אלו פתוחים לפרשנות ושנויים במחלוקת. בהיעדר הגדרה מדויקת מתעורר החשש שהשקיפות לא תגשים את ייעודה. יתרה מזו, עצם השימוש במושג השקיפות כעיקרון לא ברור עלול לסכל את השימוש בו ככלי לעידוד התנהלות אתית תקינה ולסייע למי שמבקשים שאתיקה תהיה עניין הצהרתי בלבד (איתות סגולה | virtue signaling, המכונה בספרות גם ethics washing, ethics shopping, ethics dumping, ethics shrinking) ופא (ethics lobbying).⁸⁸⁵

ממסכי המדיניות עולה שהכוונה במושג שקיפות אלגוריתמית היא ליישם את עקרון השקיפות הקלסי על מערכות טכנולוגיות. שקיפות מופיעה במקור בספרות המשפטית העוסקת במינהל תקין, בעיקר בהקשר של החובה של מערכות ממשל והליכי החקיקה להיות פתוחים לציבור ולביקורת מצד מערכת המשפט. לפיכך היא מופיעה אצל קולין פורלצה, למשל, כערך מרכזי של המשפט המינהלי.⁸⁸⁶ לפי ארגון ה-OECD, שקיפות, אחריות ושיתופיות הן עקרונות מרכזיים של מינהל תקין והן עומדות ביסוד מערכות השלטון שלו.⁸⁸⁷ השקיפות נתפסת בספרות גם כעיקרון דמוקרטי מרכזי, ומקובל לומר שברירת המחדל של השלטון צריכה להיות שקיפות.⁸⁸⁸ גם במיזמי "ממשל פתוח" ברחבי העולם שקיפות מוזכרת כערך יסודי.⁸⁸⁹

885 לעיל סעיף 3.8.

Colin Prolezza, *Data Journalism and the Ethics of Open Source*: 886
Transparency and Participation as a Prerequisite for Serving the Public Good, in *GOOD DATA* 189 (Angela Daly, Kate Devitt, and Monique Mann eds., 2019). ראו גם דפנה ברק ארז "המשפט המנהלי והמאבק בשחיתות שלטונית" משפטים לז' 667 (תשס"ז).

OECD, *EUROPEAN PRINCIPLES FOR PUBLIC ADMINISTRATION* 8 (SIGMA Papers No. 27, 887 1999)

Tal Z. Zarsky, *Transparent Predictions*, 2013 U. ILL. L. REV 1503 (2013) 888

Barack Obama, *Transparency and Open Government*. ראו למשל את הכרזת ממשל אובמה בנושא. *Open Government, Memorandum for the Heads of Executive Departments and Agencies*, THE WHITE HOUSE (2009)

בהקשר של חובות השלטון מתבטאת השקיפות בשני מאפיינים מרכזיים. האחד, החובה להעמיד לרשות הציבור מידע שלטוני.⁸⁹⁰ השני, החובה להנגיש מידע זה באופן שיהיה מובן לציבור.⁸⁹¹ על חובות אלו נוספת גם החובה לייצר לציבור את הזכות להשתמש במידע שמועמד לרשותו.⁸⁹² בעידן המידע משתכללים האופנים שבהם השלטון ממלא את חובת השקיפות: מאגרי המידע של המדינה עומדים לרשות הציבור; חקיקה והחלטות שלטוניות מפורסמות ברשת האינטרנט; ונמשך הפיתוח של פלטפורמות דיגיטליות לצורך מעקב אחר תהליכים שלטוניים, כגון שימוש בתקציב המדינה.⁸⁹³

במגזר הפרטי, חובות השקיפות מתקיימות במגוון הקשרים. ראשית, כמו במגזר הציבורי, חובות אלו מוטלות על גופים המספקים שירותים ציבוריים או נתפסים גופים דו-מהותיים.⁸⁹⁴ שנית, בהקשר של הגנת הצרכן, כגון סימון מוצרים, גילוי נאות של מאפייני המוצר ותנאי התשלום, האספקה והאחריות.⁸⁹⁵ בנושאים אלו יש כניסה לפרטי פרטים; למשל, גודל התווים, בהירות השפה, כמות הפרטים הנמסרים בעניין מהות העסקה ואורך הטקסט. כאן אפשר לכלול גם סוגיות של הסכמה המותאמת לסוגי האוכלוסיות. שלישית, דיני החברות כוללים חובות שתפקידן לייצר אמון וביטחון במסחר בבורסה ולעודד השקעות. כאלה הן למשל חובות השקיפות החלות על חברות הנסחרות בבורסה – פרסום דוחות עיתיים הכוללים מידע על רווחי החברה ועל אירועים משמעותיים בה, ובכלל זה כל חקירה פלילית נגד נושא משרה בחברה בעניין פעילות הקשורה לעבודתו.⁸⁹⁶ נורמות שקיפות מרחיבות, בעיקר במדינות מתפתחות, הן כלי אפקטיבי מומלץ

890 תהילה שוורץ אלטשולר מדיניות ממשל פתוח בישראל בעידן הדיגיטלי (מחקר מדיניות 91, המכון הישראלי לדמוקרטיה 2012).

891 ש.ם.

892 ש.ם, בעמ' 144.

893 לעניין זה ראו למשל הניסיונות לפתיחת מאגרים משטרתיים הנוגעים לאכיפה כחלק מהליכי שקיפות: *Transparency and Accountability at the Frontlines of Justice: Police Data Transparency, Open Government Partnership* (8.7.2020)

894 אסף הראל גופים ונושאי משרה דו-מהותיים 47 (2019).

895 ס' 4 לחוק הגנת הצרכן, התשמ"א-1981, ס"ח 1023 בעמ' 248.

896 ראו למשל ס' 15, 16 לחוק ניירות ערך, תשכ"ח-1968, ס"ח 541 בעמ' 234.

ליצירת אמון אצל משקיעים זרים כדי לעודד השקעות חיצוניות.⁸⁹⁷ חובות שקיפות מוטלות גם על גופים פרטיים מפוקחים, למשל הבנקים, הנדרשים לפרסם מידע על אחזקותיהם, התפלגות הרווחים, עמלות ועוד.

עם זאת, שלא כמו בגופים שלטוניים, שבהם ברירת המחדל היא שקיפות ויש צורך בטעמים חזקים כדי לעצור הנגשה של מידע, במגזר הפרטי מאזן האינטרסים שונה ולעקרון השקיפות יש משקל דומה לזה של אינטרסים שמנוגדים לו, כגון שמירה על סודות מסחריים, הזכות לפרטיות והזכות לקניין רוחני.⁸⁹⁸ כך למשל, בימינו דנים הרגולטורים באיחוד האירופי בהגבלה של חובת השיתוף בין רכבים מקושרים לרכבים אחרים בכביש אל מול הזכויות לפרטיות וההגבלות על השיתוף.⁸⁹⁹

ביסוד עקרון השקיפות עומדות כמה הצדקות נורמטיביות, פוליטיות וכלכליות.

שקיפות נתפסת כשלב מקדים והכרחי ליצירת אחריותיות אצל מקבלי ההחלטות. שקיפות ככלי לעידוד בקרה מופיעה במשפט הישראלי⁹⁰⁰ בהקשר של חוק חופש המידע, אך גם בהקשרים כגון החובה לפרסם מכרזים וחובות הפרסום של רשות המיסים. כפי שמציינים ג'סיקה מורלי ועמיתיה,⁹⁰¹ בהקשרים כאלה נתפסת השקיפות כרכיב מרכזי ביצירת אמון בין היחיד לשלטון ולכן יש לה תפקיד

Belay Seyoum and Terrell Manyak, *The Impact of Public and Private Sector Transparency on Foreign Direct Investment in Developing Countries*, 5 CRITICAL PERSPECTIVES ON INTERNATIONAL BUSINESS 187 (2009)

898 ראו Zarsky, לעיל בה"ש 888.

Araz Taeihagh and Hazel Si Min Lim, *Governing Autonomous Vehicles: Emerging Responses for Safety, Liability, Privacy, Cybersecurity, and Industry Risks*, 39 TRANSPORT REVIEWS 103 (2019)

900 בג"ץ 3751-03 יוסי אילן נ' עיריית תל-אביב-יפו, פ"ד נט(3) 817.

Jessica Morley et al., *From What to How: An Initial Review of Publicly Available AI Ethics Tools, Methods and Research to Translate Principles into Practices*, 26 (4) SCIENCE AND ENGINEERING ETHICS 2141 (2020)

חשוב בדין המינהלי – בפרט בסוגיות של הנמקת החלטות שיפוטיות ומינהליות ושל שקיפות של מסמכי תקציב.⁹⁰²

שקיפות ממלאת תפקיד חשוב גם בהגברת מעורבותו של הציבור בהליכי השלטון. מעורבות ציבורית יכולה להוביל לביוש (שיימינג) של מי שביצעו פעולות שאינן בהכרח שליליות.⁹⁰³

פרסום הליכי השלטון מגביר את אמון הציבור בו וממילא מגדיל את שיתוף הפעולה שלו – ולכן אפשר לומר שלשקיפות יש גם תועלת כלכלית.⁹⁰⁴ שקיפות היא גם נגזרת של הזכות לאוטונומיה ולחופש ביטוי,⁹⁰⁵ שכן אזרח זכאי לקבל מידע ולהבין את הרציונל העומד מאחורי ההחלטות המתקבלות בעניינו. השקיפות ככלי לקידום אוטונומיה מוצאת את ביטויה, בין השאר, בהוראת ס' 13(d) לתקנות הנציבות האירופית להגנה על היחיד במסגרת איסוף מידע, המורה לחשוף לפני היחיד לא רק את המידע על החלטות המתקבלות בעניינו אלא גם את ההיגיון שמאחוריהן.⁹⁰⁶

כדי ליישם את עקרון השקיפות הקלסי בתחום מערכות הבינה המלאכותית יש להגדיר תחילה את המושג שקיפות אלגוריתמית, את היקפו ואת היקף החובה שתוטל על הגופים השונים לממש אותו. כדי לעשות זאת יש לענות על שאלות בנוגע למטרה שהשקיפות האלגוריתמית באה לשרת: האם השקיפות נועדה לאפשר

Carol Harlow and Richard Rawlings, *Proceduralism and Automation: 902 Challenges to the Values of Administrative Law*, in *THE FOUNDATIONS AND FUTURE OF PUBLIC LAW 275* (Elizabeth Fisher, Jeff King, and Alison L. Young eds., 2020)

903 ראו Zarsky, לעיל ה"ש 888.

Sara Hagemann and Fabio Franchino, *Transparency vs Efficiency? A 904 study of Negotiations in the Council of the European Union*, 17 *EUROPEAN UNION POLITICS* 408 (2016)

905 ראו Zarsky, לעיל ה"ש 888, בעמ' 1545.

Article 13(d), Regulation 45/2001/EC of the European Parliament 906 and of the Council of 18 December 2000 on the protection of individuals with regard to the processing of personal data by the Community institutions and bodies and on the free movement of such data, *OFFICIAL JOURNAL OF THE EUROPEAN COMMUNITIES* L 8/1

בקרה אתית או לשמש כלי אכיפה משפטי? מול מי נדרשת שקיפות שלטונית לעומת שקיפות אלגוריתמית? מהן החובות הנגזרות מעקרון השקיפות? על פי אילו אמות מידה נגדיר שקיפות אלגוריתמית "טובה"? בכל השאלות האלו נעסוק להלן.

7.2 שקיפות אלגוריתמית ומשמעויותיה

על אף ריבוי המסמכים העוסקים במדיניות רגולציה ובאתיקה של בינה מלאכותית, ועל אף שכיחותו של מושג השקיפות בספרות העוסקת ברגולציה של בינה מלאכותית וגרעין של הסכמה באשר לצורך בהטלת חובות שקיפות, אין קונצנזוס על הגדרת המושג שקיפות או על מהות שלבי החובות שיוטלו ועוצמתן. יש הסבורים כי חובות השקיפות חלות על כל המשתמשים, ואילו אחרים מגבילים את תחולתן לשחקנים שונים בשלבים שונים. יש המוסיפים על חובות השקיפות גם את החובה לגילוי נאות בתקשורת שהתוכנה הנידונה אינה אנושית, ואילו אחרים רואים בשקיפות מעין חובה חוזית באשר למהות העסקה, ותו לא.

כותבים שונים מייחסים למושג שקיפות אלגוריתמית משמעויות שונות ופורטים אותו למושגי משנה בהתאם לתפיסתם בדבר חשיבותו ומקומו בשרשרת הערך של הייצור ובמבנה הרגולטורי שהם מציעים. הגדרת המושג מושפעת גם מההיבט הנבחן – למשל ההיבט של חובות היצרנים למשתמשי הקצה או ההיבט של תכלית השקיפות (האם השקיפות נתפסת אמצעי להשגת אמון או שהיא תכלית בפני עצמה?) כל כותב, על פי תפיסת עולמו, מטמיע במושג השקיפות מושגים מהעולם הצרכני של הגילוי הנאות או מעולם המשפט המינהלי או מעולם השלטון. ולעיתים כל אלה משמשים בערבוביה.

בספרות הקיימת, בהצעות חקיקה ובמסמכי מדיניות ניתנות למושג שקיפות אלגוריתמית משמעויות שונות, מהן שמתמקדות ביצרני המערכות והמוצרים ומהן שמתמקדות בזכויות המשתמשים. לדברי וירג'יניה דיגנום, לשקיפות מובן כפול: הן חובת גילוי של המידע שביסוד האלגוריתם והשימוש במערכת הן חובת גילוי של תהליך קבלת ההחלטות והיכולת להבין ולבקר תהליך זה: "שקיפות מציינת את היכולת לתאר, לברוק ולשחזר את המנגנונים שבאמצעותם מערכת בינה מלאכותית מקבלת החלטות ולומדת להתאים את עצמה לסביבתה ואת המקור והדינמיקה של הנתונים המשמשים את המערכת והמיוצרים

באמצעותה"⁹⁰⁷. כלומר דיגנום שמה את הדגש על החובה להעביר מידע – יש לאפשר למשתמשים לברוק את דרך פעולתן של מערכות מבוססות בינה מלאכותית, לשחזר אותה ולאמת את המידע ואת האופן שבו התוכנה השתמשה בו.

מושג השקיפות כולל את החובות האלה:

• **עקיבות (traceability)** – שקיפות בנוגע להליך יישומו של התכנון ההנדסי של התוכנה. העקיבות אמורה לאפשר את איתור הנורמות שבבסיס הזנת המידע ואת איתור ההטיות במידע.⁹⁰⁸

• **וידואיות (verifiability)** – יצירת לוג שמאפשר מעקב אחר הליך קבלת ההחלטות של התוכנה.

• **עיצוב הגון (honest design)** – עיצוב ממשק המשתמש כך שיאפשר למשתמש להבין את פעולת המערכת ולקבל את המידע הנכון.

• **בהירות (intelligibility)** – החובה להבטיח שאדם יבין מה אירע בפעולתו של האלגוריתם ומדוע.

• **הסבריות (explicability)** – על המשתמש להבין את המידע על פעילות התוכנה. ההסבריות כוללת את שקיפות תהליך קבלת ההחלטה; תקשור נכון בדבר יכולתיה ומטרותיה של המערכת; ואת החובת להסביר את ההחלטה שהמכונה מקבלת למי שמושפעים מהחלטה זו.⁹⁰⁹

• **תקשור (communication)** – כשמדובר במערכת בינה מלאכותית שמתקשרת עם בני אדם יש חובה להבהיר זאת במפורש.⁹¹⁰ חובה זו קשורה לגילוי נאות כלפי הצרכן ומתכתבת עם דיני הפרטיות האירופיים, בהם החובה לגילוי נאות כשהמערכת מקבלת החלטה אוטומטית בעניינו של הפרט.⁹¹¹

907 ראו DIGNUM, לעיל ה"ש 883, בעמ' 54.

908 IEEE, לעיל ה"ש 321, בעמ' 11.

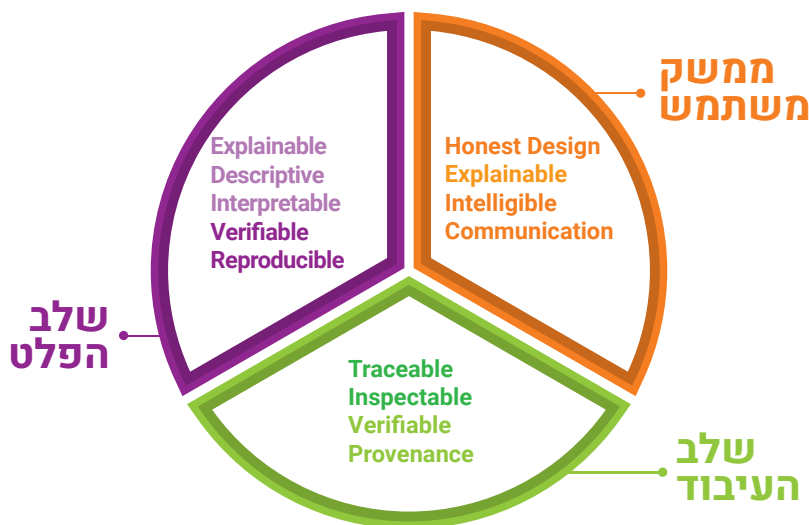
909 AI HELG 2019, לעיל ה"ש 298, בעמ' 15.

910 שם, עמ' 18.

911 ס' 22 של החקונות הכלליות בדבר הגנת מידע (GDPR).

תרשים 3

משמעויות שונות של המונח שקיפות, לפי רצף פעולת התוכנה



תרשים 3 המתאר את המשמעויות שניתנו במקורות שנסקרו לעיל למושג השקיפות לפי רצף פעולת התוכנה, מלמד שהמושג פתוח לפרשנויות רבות – כל פרשנות מדגישה רכיב מסוים בשלבי פעולתו של האלגוריתם ובפועל מטילה חובות שונות על היצרן. אפשר לראות שבלי קונקרטיזציה השקיפות האלגוריתמית עלולה להיות רחבה מדי או צרה מדי ולא להתפרש על פני מכלול הפעילות של המערכת הלומדת.

7.2.1. שקיפות ככלי עצמאי ההבדלים בפרשנות למושג שקיפות אלגוריתמית או כחלק ממערך כלים תלויים גם בשאלה אם רואים בה ערך עצמאי או אמצעי להשגת מטרות אחרות, כגון אחריותיות. להבדלים אלו יש השלכות על האיזון בין חובת השקיפות ובין אינטרסים אחרים שייתכן שישומה מנוגד להם, כגון זכויות קניין יצרניות, הגנה על פרטיות ועוד.

את הראייה הרחבה ביותר של מושג השקיפות אפשר למצוא במסמך Ethically Aligned Design (EAD) של ארגון התקינה הבינלאומי. המסמך מציין שלושה עקרונות אתיים מרכזיים שעל פיהם יש לעצב ולתפעל כל מערכת מבוססת בינה מלאכותית: על המערכת (1) לפעול בכפוף לזכויות האדם האוניברסליות; (2) לפעול בהתאם להקשר הפוליטי של המידע העומד לרשותה; (3) להיות "איתנה" מבחינה טכנולוגית (robust).⁹¹² ארגון התקינה מציין שכדי לעמוד בעקרונות אלו של עיצוב מבוסס אתיקה וערכים יש לציית לשמונה עקרונות פעולה, שהחמישי בהם הוא שקיפות. אומנם שקיפות היא כלי למימושה של אחריותיות, אבל היא מוגדרת בנפרד ממנה כי לשקיפות הקשר רחב יותר.⁹¹³

גם במסמך המסכם של ועדת המומחים מטעם מועצת אירופה, שפורסם באפריל 2019,⁹¹⁴ לשקיפות יש מקום מרכזי. אך במסמך זה ההצדקה לקיומה של השקיפות היא בעיקר הצורך בהסברתיות. המסמך מונה שלוש דרישות עיקריות ממערכת מבוססת בינה מלאכותית בכל שלבי פעולתה: (1) הדרישה לחוקיות, כלומר התאמה למערכת החקיקה והרגולציה הקיימת; (2) הדרישה שהמערכת תפעל באופן אתי; (3) והדרישה שהמערכת תהיה איתנה הן ברמה הטכנולוגית הן בבחינת השפעותיה החברתיות.⁹¹⁵

לעומת זאת, לדברי דיגנום שקיפות היא חלק ממכלול שהיא מכנה "art of AI" – אמצעי אתי לעיצוב ההחלטות של מערכות בינה מלאכותית (ethics in design). אף שהיא מבחינה בין אחריותיות, הסברתיות ושקיפות, היא משתמשת בביטוי art כמכלול ורואה בשלושת הערכים מקשה אחת. דיגנום סבורה שיישום בעת עיצוב המוצר הוא כלי הכרחי להבטחת עמידתו בחובות אתיות ובערכי החברה.⁹¹⁶ גם מורלי ועמיתיה רואים בשקיפות אמצעי למועילות, לאי-הסבת נזק, לשמירה על האוטונומיה האנושית, לצדק ולהסברתיות.⁹¹⁷

912 IEEE, לעיל ה"ש 321, בעמ' 10.

913 שם, בעמ' 29.

914 AI HELG 2019, לעיל ה"ש 298.

915 שם, בעמ' 9.

916 ראו DIGNUM, לעיל ה"ש 883, בעמ' 52.

917 Jessica Morley et al., *Ethics as a Service: A Pragmatic Operationalisation of AI Ethics*, 31 (3) MINDS AND MACHINES 1 (2021)

בספרות אפשר למצוא גם מקרי ביניים שבהם השקיפות היא העיקרון הראשי, אך אין הוא עצמאי אלא קשור לעקרונות אחרים. במסמך של המכון הישראלי למדיניות טכנולוגיה על השימוש באלגוריתמים בתחום הרווחה, שקיפות והסברות מופיעות יחד כעקרונות ראשיים, מתוך התפיסה ששקיפות היא כלי לבקרה ולהגברת האמון וגם כלי הכרחי לקיומו של הליך הוגן וצודק בתוכנות בינה מלאכותית.⁹¹⁸ אותו ערבוב של אחריותיות, הסברות ושקיפות ניכר גם במסמך ועדת המשנה של המיזם הלאומי למערכות נבונות. אומנם בתחילת המסמך השקיפות מוגדרת כערך נפרד והוועדה אומרת במפורש ששקיפות כוללת את החובה להנגיש מידע – הן באשר לתהליך יצירת המערכת הן באשר לדרך שבה היא מקבלת החלטות. הוועדה אף מסבירה כי חוסר השקיפות המאפיין את הבינה המלאכותית תורם מאוד לחוסר האמון במכונה ומגדירה את הנושא עניין שיש למצוא לו פתרון. ואולם בהמשך המסמך השקיפות יונקת את הנמקותיה מהצורך ביצירת אחריותיות. לדברי הוועדה, כדי לייצר אחריותיות יש צורך בשקיפות, בהסברות, באחריות ובניהול סיכונים.⁹¹⁹ חוסר הבהירות של מושג השקיפות נוגע אפוא גם לתפיסות השונות בעניין מקומה הגאומטרי של החובה עצמה, אם כערך עצמאי ואם כנגזרת של ערכים אחרים.

7.2.2. איזון בין שקיפות אלגוריתמית לאינטרסים מתחרים

אפשר ללמוד על הפרשנויות למושג שקיפות אלגוריתמית גם מן הגישות השונות לאיזון הנחוץ בין שקיפות לאינטרסים מתחרים, ובפרט לאיזון בין שקיפות ליעילות וחדשנות ובין שקיפות להסברותיות. כל אלה מחזקים את הטענה בדבר הצורך בעיגון פרשנות ברורה ומקובלת למושג.

במסמך של המכון הישראלי למדיניות טכנולוגיה נכתב כי בהינתן חוסר הבהירות המובנה של טכנולוגיות בינה מלאכותית, אם תוכנה איננה מסוגלת לספק הסברים

918 סיון תמיר בינה מלאכותית בשירותי ממשל: הטמעת מערכות לקבלת החלטות מבוססות-אלגוריתם בשירותי הרווחה 16 (המכון הישראלי למדיניות טכנולוגיה 2020).

919 המיזם הלאומי למערכות נבונות (2019), לעיל ה"ש 313, בעמ' 4 ו-13. נציין שהוועדה לא הגדירה שקיפות כפי שנעשה במסמכי מדיניות אחרים.

יש להגביל את שימושה במקום שהחלטותיה משפיעות על זכויות של בני אדם.⁹²⁰ לעומת זאת, במסמך מדיניות תכנון תוכנות מבוססות בינה מלאכותית שהתפרסם באתר חברת איי-בי-אם אין הכרה בהגבלת שימושים, אף שיש מערכות שמובהר שאי-אפשר לייצר לגביהן שקיפות.⁹²¹

תקינת ה-IEEE מקנה לשקיפות מעמד עליון. ועדת המומחים של הנציבות האירופית, לעומת זאת, מבטאת מודעות לקושי הטכני של מערכות בינה מלאכותית לספק הסברותיות ההולמת בני אדם. לפיכך היא מציינת שיש לאזן בין הצורך להגביר את ההסברותיות של המערכת ובין הפגיעה בדיוק של המערכת בעקבות הגברתה.⁹²² על פי תפיסתה של דיגנום, הרציונל של השקיפות הוא היותה כלי הכרחי לתקשורת בין השחקנים המעורבים ביצירת המערכת. היא מתמקדת באחריותיות ובאחריות, ולכן שלא כמו ועדת המומחים, הקובעת שיש לאזן בין ביצועים להסברותיות, היא טוענת שבמקומות שבהם ההסברותיות נפגעת יש להעדיף אינטרסים של שקיפות.⁹²³

העכירות במושג השקיפות האלגוריתמית בולטת בשני מקרים שבהם נעשה ניסיון לפרוט את המושג לחובות חוקיות.

ההצעה של תקנות הבינה המלאכותית האירופיות, שתוארה בחלק 4.1 לעיל, בחרה בגישה של הטלת מגבלות רגולטוריות בהתאם לסוגי השימושים המיועדים למערכת המבוססת

על בינה מלאכותית. לכן היא מבחינה בין מערכות בסיכון גבוה, שיהיו כפופות לרגולציה מחמירה, ובין השאר גם לחובות רישוי,⁹²⁴ ובין מערכות שאינן זקוקות

7.3 הבעייתיות בהגדרת מושג השקיפות האלגוריתמית כפי שהיא מתבטאת בהצעות החקיקה של האיחוד האירופי ושל ארגון IEEE

920 ראו תמיר, לעיל ה"ש 918, בעמ' 31.

921 IBM's Principles for Trust and Transparency, IBM (20.5.2018)

922 ראו DiGnuni, לעיל ה"ש 883, בעמ' 26 ו-33.

923 שם, בעמ' 53. הגרף המופיע בעמוד זה מסייע בהמחשת התפיסה המשולבת של שלוש הרכיבים האלו.

924 ראו לעיל, סעיף 4.1.1.2.

לרגולציה וזו תחול עליהם באופן וולונטרי בלבד.⁹²⁵ הצעת החוק משקפת שקלול של ההגדרות השונות והיא נדרשת לחובת השקיפות בשלושה עניינים.

ראשית, ההצעה מציגה בפירוט רב את היקף חובת השקיפות ואת שלבי יישומה. יש לציין את זהות היצרן ופרטי קשר להשגתו; לפרט את מאפייני התוכנה ומגבלותיה (מטרות התוכנה, רמת הדיוק שלה, תוצאות צפויות, תפקוד ביחס לקבוצות המטרה, מטא-מידע בעניין סוג המידע ואמינותו); לציין אילו התאמות בוצעו כדי לתת מענה לרמת הסיכון; להעמיד לרשות מפעיל התוכנה בקרה אנושית; לציין מהי תוחלת החיים של התוכנה ואיזו תחזוקה נדרשת להפעלתה; ועוד.

שנית, היא קובעת את רמות חובת השקיפות בהתאם לסכנה הצפויה מהשימוש המיועד לתוכנה. חובת השקיפות תחול על מערכות המוגדרות מערכות בסיכון גבוה, ובהן מערכות המשפיעות על זכויות אדם, על מצב פיננסי ועוד. בעניין מערכות אלו החקיקה מציינת במפורש כי המשתמשים צריכים להיות בעלי היכולת לפרש את הפלט של המערכת. חובת שקיפות מוגברת תחול גם על מערכות העונות להגדרות האלה: יש להן ממשק עם בני אדם; הן נועדו לאתר רגשות או תחושות; הן עלולות לבצע מניפולציה על החלטות אנושיות. במקרים כאלה תחול חובה ליידע את המשתמשים בדבר היכולות של המערכת.

שלישית, ההצעה מסייגת את חובת השקיפות כשהדברים אמורים בקניין רוחני של משווקי מערכות. זו אמירה בעייתית, שכן היא מעדיפה אפריורית זכויות יוצרים משקיפות. עם זאת, לא ברור אם זו הייתה הכוונה המפורשת וייתכן שליקוי זה יתוקן בעתיד.⁹²⁶

925 ראו למשל לעיל, סעיף 4.1.1.4.

926 על היחס של האיחוד האירופי לזכויות קניין רוחני אגב פיוחח מערכות מבוססות בינה מלאכותית ראו בדוח שאימץ פרלמנט האיחוד: European Parliament, *Intellectual Property Rights For the Development of Artificial Intelligence Technologies*, 2020/2015(INI) בסעיף 8 הדוח מציין כי האיחוד האירופי – "stresses the importance of streaming services being transparent and responsible in their use of algorithms, so that access to cultural and creative content in various forms and different languages as well as impartial access to European works can be better guaranteed"

מסמך המדיניות של ארגון התקינה נדרש לסוגיות רבות, אך הוא מפורט פחות.⁹²⁷ לפי מסמך זה, חובות השקיפות צריכות להביא בחשבון את המשתמש הסופי, את היצרן ואת הרגולטורים גם יחד. המסמך מבקש גם להשית חובת העמדת מידע על מי שיידרשו לחקור אירועים של תקלות במערכות אלו, אם מטעם המדינה ואם מטעם היצרן. המסמך דורש שבכל תוכנה יהיו לפחות שלוש רמות של שקיפות.⁹²⁸

(1) מנקודת מבטו של המשתמש – במערכת יוטמע מעין "כפתור" שכאשר ילחצו עליו תסביר התוכנה מדוע בוצעה כל פעולה.

(2) הליכי האשורור יתועדו ויועמדו לרשות מפתחים ויצרנים אחרים לשם בדיקת המוצר.

(3) ייווצרו "לוגים" ויסופקו חיישנים שישמרו באופן מאובטח את המידע על פעילות התוכנה כדי לבחון את פעולתה בעת תקלה.

המסמך קורא לקביעת מדדים כמותיים לבחינת שקיפות המוצר, אך לא מפורטות בו הדרכים לקביעתם.⁹²⁹

7.4

העכירות של השקיפות האלגוריתמית: סיכום

יש הסבורים כי מימוש השקיפות האלגוריתמית צריך להיעשות באמצעות הטלת חובות שתכליתן לאפשר למשתמש להבין את המתרחש סביבו ולשלוט בו. אחרים שמים את הדגש על התוצאה

ודורשים שהמשתמש ידע להעריך את איכות השירות שניתן לו על ידי התוכנה ויוכל להשיג על החלטותיה בעת הצורך. אחרים אף מרחיבים את חובת השקיפות האלגוריתמית כך שתכלול חובות גילוי צרכניות. נוסף על כך אין הבחנה ברורה בין היקף החובות החלות בשעה שהטכנולוגיה נמצאת בשימוש במגזר הפרטי או בגופים דו-מהותיים להיקף החובות החלות כשהיא נמצאת בשימוש גוף ממשלתי, כחלק מפעולותיו.

927 ראו IEEE, לעיל ה"ש 321.

928 שם, בעמ' 28.

929 שם, בעמ' 27.

יש מחלוקת גם בנוגע לאיזון בין חובת השקיפות האלגוריתמית ובין חובות אחרות שיחולו על מוצרים מבוססי בינה מלאכותית. לאיזון זה יש השלכות במקרים שבהם יכולת הבקרה נפגעת והתוכנה איננה מסוגלת לספק מענה לדרישות המוזכרות במסמכי מדיניות ובהצעות חוק, למשל בנוסח ההצעה של תקנות הבינה המלאכותית האירופיות. במצב דברים זה המושג שקיפות אלגוריתמית נותר עכור, הן בהיעדר הגדרה ממצה למושג, הן בכל הנוגע להיקף תחולתו, הן ביישום המעשי של הדרישות השונות ובהתמודדות עם מקרים שבהם בגלל חסמים טכנולוגיים אי־אפשר לספק את הדרישות של כל הכותבים. לא ברור מהן בדיוק הדרישות מיצרן של מערכות אלגוריתמיות ומה מעמדן של דרישות אלו במקרים שבהם הטכנולוגיה אינה יכולה לספק אותן – אם הרגולציה תגביל את השימושים המותרים של טכנולוגיות אלו אפשר שבמקרים מסוימים היצרנים אף ייאלצו לנטוש את פיתוחן.

פרק שמיני

פיקוח מוסדי
על בינה מלאכותית

—

לצד האתגר הרגולטורי של תרגום עקרונות אתיים רחבים לכללים משפטיים יש שאלה נפרדת: מה זהות הגוף או הגופים המוסדיים האמורים להיות ממונים על תחום הבינה המלאכותית – הן כגופים מיישמי חקיקה, הן כגופים מאסדרים ומפקחים והן כגופים מייעצים?

בחלק זה נסקור כמה מקורות ונשווה את האופנים שבהם הם מטפלים בהיבט המוסדי של אסדרת הבינה המלאכותית. מאחר שהאסדרה של בינה מלאכותית עודנה בחיתוליה נבחרו

8.1
ההיבט המוסדי
בראייה השוואתית

מקורות שכוללים הצעות לרגולציה, המלצות מדיניות והסדרי ביניים קיימים (בהיעדר רגולציה של ממש) בנושאי בינה מלאכותית. להשלמת התמונה נדרש גם להיבטים הרגולטוריים בתזכיר חוק הגנת הסייבר.

התקנות האירופיות (שתוארו לעיל בפרק 4) מעמידות לצד הכללים החלים על מערכות בינה מלאכותית בסיכון גבוה (והאיסור הקטגורי על מערכות בינה מסוכנות) מבנה מוסדי מורכב שנועד לשרת בין השאר את הצרכים הנובעים מתיאום על-לאומי בין המדינות החברות באיחוד.

הנוסח המקורי של התקנות הציע להקים ועד אירופי לבינה מלאכותית (EU AI Board)⁹³⁰, בדומה לוועד האירופי להגנה על מידע (EDPB), שהוקם מכוח התקנות הכלליות בדבר הגנת מידע (GDPR).⁹³¹ הוועד האירופי לבינה מלאכותית, שיורכב מנציגים בכירים של רשויות הפיקוח הלאומיות וממפקח על הגנת הפרטיות,⁹³² יתרום לשיתוף הפעולה בין רשויות הפיקוח הלאומיות ובין הנציבות האירופית בעניינים שהתקנות מסדירות; יתאם את ההנחיה של הנציבות האירופית ושל רשויות הפיקוח הלאומיות בנושאים שהתקנות מסדירות; ויסייע לנציבות האירופית ולרשויות הפיקוח הלאומיות ביישום עקבי של התקנות.⁹³³

לפי ההצעה המתוקנת יוקם משרד אירופי לבינה מלאכותית (EU AI Office)⁹³⁴ המשרד האירופי לבינה מלאכותית, שמושב יהיה בבריסל, יהיה גוף עצמאי של האיחוד, ויורכב ממזכירות, ועד מנהל ופורום מייעץ.⁹³⁵ תפקידי המשרד לבינה מלאכותית יהיו ייעוץ לגופים הרלוונטיים באיחוד ובמדינות החברות בקשר ליישום תקנות הבינה המלאכותית; ניטור יישום התקנות מבלי לפגוע בסמכותן של הרשויות המפקחות הלאומיות; תיאום בין רשויות פיקוח לאומיות וגישור במחלוקות ביניהן; תרומה לשיתוף פעולה עם גופים מקבילים של מדינות מחוץ

930 ס' 56 להצעת תקנות הבינה המלאכותית האירופיות, לעיל ה"ש 53.

931 ראו ס' 68-76 של התקנות הכלליות בדבר הגנת מידע (GDPR).

932 ס' (1) 57 להצעת תקנות הבינה המלאכותית האירופיות, לעיל ה"ש 53.

933 שם, ס' (2) 56, 58.

934 ס' 56 להצעה המתוקנת, לעיל ה"ש 57.

935 שם, ס' 56a.

לאיחוד; פיתוח מומחיות רלוונטית; בחינה של שאלות הנוגעות ליישום התקנות; ופרסום של דיווחים שנתיים על יישומן ושל המלצות על סיווג של פרקטיקות אסורות ושל מערכות בסיכון גבוה.⁹³⁶

הנציבות האירופית תנהל בשיתוף עם מדינות האיחוד מרשם של מערכות בינה מלאכותית בסיכון גבוה.⁹³⁷ על הספקים של מערכות אלו תחול חובת רישום⁹³⁸ והמאגר יהיה פתוח לעיון הציבור.⁹³⁹ גוף נוסף שלא הוסדר, אך ככונת הנציבות להקים בשלב הטמעתן של מערכות אלו הוא קבוצת מומחים שתסייע בניטור השוק ותייעץ לוועד האירופי לבינה מלאכותית.⁹⁴⁰

ברמה הלאומית, כל מדינה חברה באיחוד תקים רשות פיקוח לאומית מוסמכת שתפקידה יהיה להבטיח את יישומן והטמעתן של התקנות המוצעות,⁹⁴¹ ותשמש מפקחת על השוק (market surveillance authority).⁹⁴² ימונו גם מפקחים מגזריים שיהיו אחראים על מגזרים מסוימים.⁹⁴³ גורמי הפיקוח החיצוניים – שיכולים להיות גופים פרטיים שקיבלו אקדריטציה – יוסמכו לבצע בחינות תאימות של מערכות בינה מלאכותית בסיכון גבוה לפי הוראות התקנות המוצעות.

אפשר לראות אפוא שנוסח התקנות המוצעות נותן למדינות חברות מרחב תמרון מסוים בעיצוב המוסדי של רגולציית הסייבר ברמה הלאומית. לנוכח הדרישה למנות רגולטור יחיד, שירכז את הקשר עם הציבור ועם רגולטורים מקבילים

936 שם, ס' 56b

937 שם, ס' 60(1).

938 שם, ס' 51. לפרטים החייבים ברישום ראו התוספת השמינית לתקנות המוצעות.

939 שם, ס' 60(3).

940 ראו דבריה של לוסילה טיולי, מנהלת תחום בינה מלאכותית ותעשיות דיגיטליות במינהלת הנציבות למערכות תקשורת, תוכן וטכנולוגיה (DG CONNECT): *A European Approach to the Regulation of Artificial Intelligence*, YouTube (23.4.2021), בדקה 25:00.

941 ס' (2)–(1) 59 להצעה המתוקנת, לעיל ה"ש 57.

942 ס' 63 להצעת תקנות הבינה המלאכותית האירופיות, לעיל ה"ש 53.

943 שם, ס' 63(1).

במדינות חברות אחרות,⁹⁴⁴ המדינות החברות אינן רשאיות לבזר את האחריות לרגולציה של בינה מלאכותית ולהעביר אותה לידי רגולטורים מגזריים. לאחרונה המליצה מועצת המדינה של צרפת (Conseil d'Etat) שרגולטור הפרטיות הלאומי (Commission nationale de l'informatique et des libertés – CNIL) יהיה מופקד על רגולציה של בינה מלאכותית.⁹⁴⁵

עם זאת, יש לתת את הרעת על כך שהתפקיד המיועד לרגולטור הלאומי הוא יישום והטמעה של תקנות שתכליתן להגן על עקרונות אתיים ולהביא להובלה אירופית בפיתוח בינה מלאכותית אמינה (trustworthy), בטוחה ואתית.⁹⁴⁶ לכן ברמה הלאומית תפקידו העיקרי של הרגולטור הוא לפקח ולא לפתח.

נראה שהצעת חוק הבינה המלאכותית והמידע של קנדה כוללת הסדר דומה. לפי הצעה זו יש לאפשר לשר הממונה למנות את אחד מבכירי משרדו לנציב לענייני בינה מלאכותית ומידע כדי לסייע לשר באכיפת החוק.⁹⁴⁷ השר יהיה רשאי להאציל מסמכויותיו הסטטוטוריות לנציב לענייני בינה מלאכותית ומידע,⁹⁴⁸ פרט לסמכות לתקן תקנות.

בשנת 2021 פרסמה נציבות זכויות האדם של אוסטרליה דוח שעוסק בזכויות אדם וטכנולוגיות מתעוררות, ובהן בינה מלאכותית.⁹⁴⁹ בבסיס הדוח עמדה גישה שנותנת עדיפות לזכויות אדם הן ברמת האסטרטגיה הלאומית⁹⁵⁰ הן ברמת הרגולציה.⁹⁵¹ הדוח מציע למנות נציב לנושא בטיחות בבינה מלאכותית, שסייע לרגולטורים, לקובעי המדיניות ולגורמים בממשלה ובמגזר העסקי לעצב את

944 שם, ס' 77 למבוא.

945 *S'engager dans l'intelligence artificielle pour un meilleur service public*, CONSEIL D'ETAT (30.8.2022)

946 פס' 5 למבוא להצעת תקנות הבינה המלאכותית האירופיות, לעיל ה"ש 53.

947 ס' 33 להצעת חוק הבינה המלאכותית והמידע של קנדה, לעיל ה"ש 55.

948 ראו לעיל בסעיף 4.4.2.

949 Australian Human Rights Commission, HUMAN RIGHTS AND TECHNOLOGY – FINAL REPORT 28 (2021)

950 שם, בעמ' 23-24.

951 שם, בעמ' 47-48.

המסגרת הנורמטיבית שתחול על קבלת החלטות מבוססת בינה מלאכותית. הנציב יפעל לקדם זכויות אדם ולהגן עליהן, בפרט זכויות של קבוצות בסיכון.⁹⁵² הרוח מדגיש כי אין צורך ברגולטור ייעודי לבינה מלאכותית, מאחר שלשיתו אין להתמקד באסדרת טכנולוגיות אלא באסדרה של השימוש בהן בהקשרים ספציפיים – דבר שרגולטורים מגזריים יוכלו לעשות ביעילות בתמיכת הנציב לנושא בטיחות בבינה מלאכותית.⁹⁵³

הרוח ממליץ לממשלת אוסטרליה להקים נציבות עצמאית לבטיחות בבינה מלאכותית, אשר (1) תעבוד עם רגולטורים למען פיתוח יכולות טכניות בתחום הבינה המלאכותית; (2) תנטר ותחקר התפתחויות בשימוש במערכות נבונות, בייחוד בתחומים שבהם יש חשש לפגיעה בזכויות אדם; (3) תשמש עבור מעצבי המדיניות באוסטרליה גורם מומחה ועצמאי לבינה מלאכותית ולזכויות אדם; ו־(4) תפרסם קווים מנחים שינחו את הממשלה ואת המגזר הפרטי כיצד להשתמש בבינה מלאכותית כחוק, ועל פי סטנדרטים אתיים.⁹⁵⁴ כדי לנטר את השימוש במערכות נבונות מוצע שלנציבות הבטיחות בבינה מלאכותית יוקנו סמכויות לערוך חקירות וביקורות, בדומה לסמכויות הנתונות לרשויות האמונות על תחומי הגנת הפרטיות.⁹⁵⁵ לדעת הנציבות לזכויות אדם, סמכויות אלו יסייעו בהתמודדות עם אתגרי השקיפות האלגוריתמית.

לראיין קיילו הייתה הצעה דומה. במאמר משנת 2014 הוא הציע להקים גוף ממשלתי מתמחה שלא ישמש כרגולטור ייעודי.⁹⁵⁶ לנגד עיניו עמדה רשות פדרלית לרובוטיקה שתיעץ לסוכנויות אחרות בממשל בנוגע לרובוטיקה ולבינה מלאכותית. למשל, היא תיעץ לרשות לניירות ערך בעניין מסחר אלגוריתמי (אלגו־טריידינג),⁹⁵⁷ לרשות התעופה בעניין כטב"מים (אוטונומיים

952 שם, בעמ' 127.

953 שם, בעמ' 130.

954 שם, המלצה מס' 22, בעמ' 128.

955 שם, בעמ' 132.

956 Ryan Calo, *The Case for a Federal Robotics Commission*, BROOKINGS (15.9.2014)

957 ראו לעיל בסעיף 2.4.

ולא אוטונומיים)⁹⁵⁸, ולמשרד התחבורה בעניין רכבים אוטונומיים. קיילו הציג שהרשות הפדרלית לרובוטיקה תיעץ למחוקקים בכל רמות הממשל בענייני חקיקה ומדיניות הנוגעים לרובוטיקה ולבינה מלאכותית; תדון עם בעלי עניין מקומיים ובינלאומיים מתחומי התעשייה, הממשל והאקדמיה על השפעות אפשרויות של טכנולוגיות אלו על החברה; ותשמש ידית בית המשפט (בדומה לסוכנויות פדרליות אחרות) בעניינים הנוגעים לתחומי פעילותה.⁹⁵⁹

קיילו הדגיש שהרשות שהוא מציע לא תעסוק באכיפה – הוא סבר שבנקודת הזמן שבה נכתבו הדברים החלת משטר אכיפה כללי על טכנולוגיות בינה מלאכותית תקדים את המאוחר. מנגד, סבר שעליה לפתח מומחיות שתסייע לרגולטורים אחרים, וכן למחוקקים ולבתי המשפט, להימנע מטעויות.⁹⁶⁰

באירופה נפוץ הסדר ביניים שאמור לשמש עד לכינונו של גוף רגולטורי: יחידת מצפה בינה מלאכותית (AI Observatory), גוף ממשלתי שתכליתו לאפשר הטמעה אחראית, אמינה ואתית של מערכות נבונות ולתמוך בה.⁹⁶¹ ברפובליקה הצ'כית, למשל, מצפה הבינה המלאכותית הוא גוף מומחים שנועד לזהות מכשולים משפטיים אפשריים בחקר מערכות נבונות, בפיתוח שלהן ובשימוש בהן, ולהמליץ על התמודדות עימם; לפרסם המלצות משפטיות ואתיות לתעשייה; ולשתף פעולה בשיח הבינלאומי על רגולציה ואתיקה של בינה מלאכותית.⁹⁶² גרמניה הקימה יחידת מצפה בינה מלאכותית בשוק העבודה ובחברה שתכליתה לבחון את השפעתן של טכנולוגיות אלו על החברה והעבודה ולנסח המלצות לעיצוב מערכות אלו באופן ששם את האדם במרכז ומשרת את טובת הכלל.⁹⁶³

958 להרחבה ראו אורי וולובלסקי "שימושים אזרחיים בכטב"מים – אתגר חדש לזכות לפרטיות" משפט ועסקים יט 993 (2016).

959 Calo, לעיל ה"ש 956, בעמ' 12.

960 שם, בעמ' 13.

961 Raquel Jorge Ricart et al., AI Watch, National Strategies on Artificial Intelligence: A European Perspective 15 (PUBLICATIONS OFFICE OF THE EUROPEAN UNION, 2021)

962 ראו באחר AI OBSERVATORY AND FORUM.

963 ARTIFICIAL INTELLIGENCE STRATEGY 26 (The Federal Government 2018); ARTIFICIAL INTELLIGENCE STRATEGY OF THE GERMAN FEDERAL GOVERNMENT (The Federal Government 2020)

המשרד לכינה מלאכותית בכריטניה הוא יחידה ממשלתית שמשותפת למשרד לאסטרטגיה בענייני עסקים, אנרגיה ותעשייה ולמשרד הדיגיטל, המדיה, התרבות והספורט. מטרת המשרד היא לעודד פיתוח טכנולוגיות בינה מלאכותית ושימוש בהן. המשרד אמון על ניסוח קווים מנחים בנוגע לכינה מלאכותית עבור גורמים ממשלתיים⁹⁶⁴ ועל עידוד מחקר בנושאי רגולציה ומדיניות.⁹⁶⁵

הצעת חוק הבינה המלאכותית של ברזיל נמנעת מהקמה של גוף ייעודי לאסדרה של בינה מלאכותית וקוראת לאסדר תחום זה באמצעות הרגולטור המגזרי,⁹⁶⁶ במתווה של ניהול סיכונים. לפי הקווים המנחים בהצעה, כל רמות הממשל כברזיל יעודדו הקמה של מנגנוני ממשל בשיתוף נציגים של הציבור, של החברה האזרחית, של המגזר העסקי ושל הקהילה המדעית.⁹⁶⁷ כל אלה יעודדו אימוץ מכשירים רגולטוריים שמקדמים חדשנות, כגון ארגוני חול רגולטוריים, הערכת סיכונים רגולטורית ורגולציה עצמית מגזרית.⁹⁶⁸

מערך הסייבר הלאומי הוקם בשנת 2017 בעקבות איחוד מטה הסייבר הלאומי והרשות הלאומית להגנת הסייבר.⁹⁶⁹ מערך הסייבר ממונה על הגנת מרחב הסייבר האזרחי, על מתן שירותים לניהול מתקפות סייבר ועל מתן הדרכה לכל החברות האזרחיות ולגופים העוסקים בתשתיות קריטיות בהתבסס על האמונה שלכל אחד בישראל הזכות להשתמש בטכנולוגיה ללא חשש.⁹⁷⁰ בשנת 2018 הופץ תזכיר

964 כך למשל פורסמו קווים מנחים לרכש ציבורי של בינה מלאכותית (GUIDELINES FOR AI PROCUREMENT [Office For Artificial Intelligence 8.6.2021]); לניהול פרויקטים של בינה מלאכותית; להערכה ולהטמעה של בינה מלאכותית; וכן לאתיקה של מערכות אוטומטיות לקבלת החלטות (ETHICS, TRANSPARENCY AND ACCOUNTABILITY FRAMEWORK FOR AUTOMATED DECISION-MAKING [Office For Artificial Intelligence 13.5.2021])

965 ראו למשל ADA Exploring Legal Mechanisms For Data Stewardship, LOVLACE INSTITUTE (4.3.2021); AI ROADMAP (UK AI Council 6.1.2021)

966 ס' 6(2) להצעת חוק הבינה המלאכותית הברזילאית, לעיל הש"ש 54.

967 שם, ס' 7(VIII).

968 שם, ס' 7(VII).

969 החלטת ממשלה מס' 3270: איחוד יחידות מערך הסייבר הלאומי (17.12.2017).

970 שם: Dimitry (Dima) Adamsky, *The Israeli Odyssey toward its National Cyber Security Strategy*, 40 (2) THE WASHINGTON QUARTERLY 113, 120 (2017)

חוק הגנת הסייבר ומערך הסייבר הלאומי, שנועד לעגן בחקיקה את המסגרת הרגולטורית לפעילותו של המערך.⁹⁷¹ בעת כתיבת שורות אלו תזכיר החוק תלוי ועומד. אומנם נראה שמערך הסייבר מתנהל בהתאם למסגרת הרגולטורית המוצעת בתזכיר, אך עיקר סמכויותיו הן הנחיה וייעוץ לרגולטורים המגזריים ולחברות המעוניינות בכך וסיוע בעת מתקפת סייבר. לעת עתה אין לו סמכויות אכיפה.⁹⁷²

לצד ההיבטים הרגולטוריים מונה תזכיר החוק עוד תפקידים של מערך הסייבר: הניהול המבצעי של מאמצי ההגנה הלאומיים נגד תקיפות סייבר, קידום שיתופי פעולה בינלאומיים בתחום, ייעוץ לממשלה וועדותיה וקידום מדיניות ומובילות ישראלית בתחום הסייבר.⁹⁷³

מודל הרגולציה המוצע בתזכיר חוק הגנת הסייבר מחיל על הגופים המוסדרים רמות שונות של התערבות רגולטורית בהתאם לרמת הסיכון שלהם.⁹⁷⁴ ארגונים ששייכים למגזר משקי שמוגדר בתוספת השלישית לחוק יהיו נתונים לפיקוח ולהנחיה ישירים של מערך הסייבר.⁹⁷⁵ במגזרים שבהם יש רגולטור מגזרי, הוא שיקבע את ההוראות להגנת הסייבר – בהנחיית מערך הסייבר,⁹⁷⁶ המשמש "רגולטור של רגולטורים".⁹⁷⁷ ברירת המחדל היא האצלת סמכויות לרגולטורים קיימים – אך רגולטורים אלו ידרשו לפעול באופן אחיד ובהתאם לתורת ההגנה בסייבר.⁹⁷⁸ ארגונים בסיכון נמוך שאינם מוסדרים ברשות מאסדרת אחרת ינהיגו רגולציה עצמית ומערך הסייבר ישמש בהם גורם מנחה ומייעץ בלבד.⁹⁷⁹

971 ראו תזכיר חוק הגנת הסייבר, לעיל ה"ש 461.

972 ארידור הרשקוביץ ושוורץ אלטשולר, לעיל ה"ש 461, בעמ' 120.

973 ס' 14 לתזכיר חוק הגנת הסייבר, לעיל ה"ש 461.

974 מערך הסייבר הלאומי, משרד ראש הממשלה, הערכת השפעות רגולציה: פרק האסדרה בחוק הסייבר, עמ' 2-3, 8, 12 (יוני 2018).

975 ס' 57 לתזכיר חוק הגנת הסייבר, לעיל ה"ש 461.

976 ס' 44 לתזכיר חוק הגנת הסייבר, לעיל ה"ש 461.

977 ארידור הרשקוביץ ושוורץ אלטשולר, לעיל ה"ש 461, בעמ' 162-165, 137-141.

978 מערך הסייבר הלאומי, תורת ההגנה בסייבר 2.0 (19.7.2021).

979 מערך הסייבר הלאומי, הערכת השפעות רגולציה, לעיל ה"ש 974, בעמ' 25-26.

8.2

**ההתייחסות בדוחות ישראלים
לרגולטור של בינה מלאכותית**

בשנת 2018 הקימו פרופ' יצחק בן ישראל ופרופ' אביתר מתניה את הוועדה לקידום תעשיית הבינה המלאכותית בישראל והוועדה גיבשה תוכנית שכותרתה "חזון העצמת חוסנה של

ישראל כמעצמה מדעית טכנולוגית בראיית הביטחון הלאומי והבטחת שגשוגה הכלכלי".⁹⁸⁰ בין השאר הומלץ להקים מינהלת לאומית למערכות נבונות במשרד ראש הממשלה.⁹⁸¹ התוכנית גרסה שיש צורך בגוף כזה ב"מספר הרב של הרשויות שיש להן נגיעה להייטק המקומי [...] וכמובן (ב)גופים הבטחוניים".⁹⁸² עוד הוצע בתוכנית שהמינהלת הלאומית תיעץ לממשלה באשר למדיניות מתכללת בנושא. בן ישראל ומתניה הם אנשי מערכת הביטחון לשעבר והם שעמדו מאחורי הקמת מערך הסייבר הלאומי, גוף ביטחוני דמוי שב"כ. על פי תפיסתם בעניין הצורך בתכלול תחום הבינה המלאכותית, יש "לקבע את המשימה כבעלת חשיבות בטחונית ראשונה במעלה".⁹⁸³ אומנם יש צורך בהגנה מפני גורמים עוינים המבקשים לנצל לרעה מהפכות טכנולוגיות, אבל אין פירוש הדבר שמערכת הביטחון או גופי ביטחון צריכים לחלוש על אסטרטגיית הטכנולוגיה הישראלית. ואכן, הדוח של בן ישראל ומתניה עורר התנגדות רבתי בקרב גופים אזרחיים, בהם המועצה להשכלה גבוהה, רשות החדשנות ומשרד האוצר. הללו יצאו חוצץ נגד הנחות היסוד של הוועדה, תוכן עבודתה והיעדר השקיפות בפעולותיה.⁹⁸⁴

980 אורי ברקוביץ וטל שחף "נהפוך את ישראל לאחת ממדינות הבינה המלאכותית המובילות" גלובס (11.8.2018).

981 "המיזם הלאומי למערכות נבונות בטוחות להעצמת הביטחון הלאומי והחוסן המדעי-טכנולוגי: אסטרטגיה לאומית לישראל, דו"ח מיוחד לראש הממשלה" א 23 (יצחק בן ישראל, אביתר מתניה וליאת פרידמן עורכים 2020) (להלן: המיזם הלאומי למערכות נבונות (2020)).

982 אורי ברקוביץ "השקעה של 2 מיליארד ש' בשנה בעיר חכמה, בחקלאות ובאקדמיה: כך מתכננת ישראל להפוך למעצמת בינה מלאכותית" גלובס (18.11.2019).

983 ש.ס.

984 אמיתי זיו "נחיצות מוטלת בספק: גופי הממשלה נגד תוכנית הבינה המלאכותית" TheMarker (21.11.2019).

רוח ועדת המשנה בנושא אתיקה ורגולציה של מערכות נבונות מכיל התייחסות רחבה יותר להיבט המוסדי.⁹⁸⁵ הרוח מציין שלנוכח קצב ההתפתחות הטכנולוגית בתחום זה, רצוי להגדיר מראש את ההליכים לבחינה חוזרת של המדיניות הרגולטורית לצורכי עדכונה.⁹⁸⁶ כמו כן, הרוח מציע שכדי להתמודד עם חוסר הוודאות באשר להשפעה של מערכות נבונות על היבטים אתיים, יאפשרו רשויות האסדרה ניסויים מבוקרים של מדיניות רגולטורית ושימוש בארגזי חול רגולטוריים.⁹⁸⁷ עוד מציין הרוח שאסדרת בינה מלאכותית אינה צריכה להיעשות בחלל הריק.⁹⁸⁸ רגולטורים מגזריים מגוונים יכולים להיות מעורבים באסדרת מערכות נבונות בתחומם; למגזרים שונים יש צרכים שונים, שמהם יכולים להיגזר כללים שונים. אומנם עקרונות אתיים ברמת הפשטה גבוהה חלים במידה שווה על מערכות נבונות ממינים שונים – לצרכים רפואיים, מכונות אוטונומיות, מערכות לסיוע בקבלת החלטות מינהליות ומערכות זיהוי פנים – אבל הקונקרטיזציה של אותם עקרונות אתיים וניסוחם בכללים רגולטוריים ישימים עשויה להשתנות מתחום לתחום.

לצד הכרה בצורך להתאים את הכללים המשפטיים החלים על מערכות נבונות למגזרים שונים מתוך שיג ושיח עם הרגולטור המגזרי בתחומן, רוח ועדת המשנה מזכיר גם כמה רגולטורים רוחביים שיכולים למלא תפקיד חשוב באסדרה: הרשות להגנת הפרטיות, רשות התחרות, ובעקיפין גם מערך הסייבר הלאומי. הרוח ממליץ שעיקר האסדרה תיעשה באמצעות הרגולטורים המגזריים; הוא אף ממליץ (כלי לאפיין אותו) על גורם רגולטורי מומחה ועל-ממשלתי שיתאם בין הרגולטורים המגזריים לרגולטורים הרוחביים.

הרוח של ועדת בינה מלאכותית ומדע הנתונים, שמינה פורום תל"ם (פורום תשתיות לאומיות למחקר ולפיתוח) בתחילת 2020, זיהה את תחום הבינה המלאכותית כתחום אסטרטגי וציין כמה היבטים שמומלץ לרכז בהם את המאמץ

985 המיזם הלאומי למערכות נבונות (2019), לעיל ה"ש 313. לגרסתו הסופית של הדוח ראו המיזם הלאומי למערכות נבונות (2020) ב' 172, לעיל ה"ש 981.

986 המיזם הלאומי למערכות נבונות (2019), לעיל ה"ש 313, בעמ' 32.

987 שם, שם.

988 שם, בעמ' 33-34.

הלאומי.⁹⁸⁹ הוועדה המליצה על תוכנית חומש לאומית לבינה מלאכותית, אבל התמקדה בעיקר בצורך לקדם את מדע הנתונים.⁹⁹⁰

הוועדה ציינה כי "נושא הנתונים והנגישות אליהם מהווה מכשול משמעותי [...] נכון יהיה להגדיר רגולטור ממשלתי אחד כדוגמת הלמ"ס כבעל אחריות למיפוי, אגירה, זמינות, רישוי ונגישות לנתונים". ואולם בהמשך הרוח ציינה הוועדה כי הדיון במבנה הניהולי המוצע לתוכנית הלאומית טרם הושלם, וחזרה על הצורך ברגולטור יחיד שעיקר תפקידו לאחד בעלות על הנתונים מתוך "שאיפה להסתכל על נתונים כעל 'טובין ציבוריים'".⁹⁹¹ עם העקרונות הכלליים שהוועדה מונה במסמך מנחת גם רגולציה, אך לא מוצע בו שום מבנה מוסדי מלבד ההמלצות על ארגון חול רגולטוריים ועל הקמת מרכזי נתונים ופלטפורמות לשיתוף נתונים ומודלים.⁹⁹² במסמך מאוחר יותר הציגה הוועדה מבנה ניהולי: הרגולטור, המורכב מנציגי משרד המשפטים ומשרד החוץ, יהיה גורם חיצוני לתוכנית הבינה המלאכותית הלאומית.⁹⁹³

בשנת 2022 פורסם דוח בנושא רגולציה של בינה מלאכותית במגזר הפיננסי, שהוכן לבקשת המחלקה הכלכלית של מחלקת ייעוץ וחקיקה במשרד המשפטים.⁹⁹⁴ הדוח מבסס את המלצותיו על שיקולים הנוגעים לצורך בחיזוק מגזר ההיי-טק והחדשנות ולצורך בעידוד השימוש בתוצרי מגזר זה בישראל גופא (ולא רק בשוקי היעד הזרים שלו), וכן על אילוצים הנוגעים ליכולתה של ישראל להשפיע על התקינה הגלובלית החלה על מערכות בינה מלאכותית.⁹⁹⁵ הדוח ממליץ לרכז בידי גורם ממשלתי אחד את תיאום סוגיית האסדרה של יישומי בינה

989 ועדת בינה מלאכותית ומדע הנתונים, דו"ח סופי 2020, בעמ' 78.

990 שם, בעמ' 80.

991 שם, בעמ' 81.

992 שם, בעמ' 62-64. התייחסויות כלליות לצורך בפיתוח רגולציה יש גם בדוח צוות המשנה של המיזם הלאומי למערכות נבונות בנושא ממשל. ראו המיזם הלאומי למערכות נבונות (2020) ב, לעיל ה"ש 981, בעמ' 222, 225.

993 דוח ועדת חל"ם, לעיל ה"ש 103, מצגת יוני 2021, בשקף 20.

994 ראו אחיעז ואח', לעיל ה"ש 240.

995 שם, בעמ' 146.

מלאכותית – הוא מצייץ שרצוי שתהיה בו מעורבות של משרד המשפטים,⁹⁹⁶ אך קורא להימנע בשלב זה מקביעת כללים שיחולו באופן אחיד על כל יישומי הבינה המלאכותית (שיכולים להיות שונים מאוד אלו מאלו – למשל יישום לפענוח תצלומי רנטגן לעומת אלגוריתמים לטרגוט פרסומות).⁹⁹⁷

גישה דומה אומצה ברוח שפרסם בסוף 2022 משרד החדשנות, המדע והטכנולוגיה. הרוח מציע להקים "מוקד ידע ותיאום ממשלתי להסדרת בינה מלאכותית". גורם זה יעסוק ביישום מדיניות רגולציה ואתיקה ובגיבוש המלצות לעדכונה, יסייע למשרדי הממשלה ולרגולטורים בגיבוש מדיניות רגולציה וינגיש לציבור ולממשלה מידע וכלים לשימוש אחראי בבינה מלאכותית.⁹⁹⁸ על פי תפיסת האסדרה המוצעת ברוח, הנמנע מחקיקה או מאסדרה רוחבית אחרת, יש לעודד כלים של אסדרה "רכה" וניסויית ולהעדיף רגולציה מגורית. הרוח מדגיש את חשיבותו של גוף תיאום בין-משרדי להבטחת מדיניות אחידה ברורה וקוהרנטית.⁹⁹⁹

996 שם, בעמ' 148.

997 שם, בעמ' 147.

998 מסמך מדיניות הרגולציה של משרד החדשנות, לעיל ה"ש 314, בעמ' 110.

999 שם, בעמ' 11. ראו גם המיזם הלאומי למערכות נבונות (2019), לעיל ה"ש 313.

פרק תשיעי

יסודות לרגולציה מכוונת זכויות
של בינה מלאכותית בישראל:
עקרונות כלליים

—

השיח הער בשנים האחרונות מכיר באתגרים שמערכות בינה מלאכותית מעמידות לפני קובעי המדיניות. יש הסבורים כי שאלות משפטיות רבות הנוגעות לשימושים במערכות מבוססות בינה מלאכותית מקבלות מענה עקרוני בדין הקיים ולא נדרשת אלא התאמה קלה של הקצאת הסיכונים, הן על ידי תיקוני חקיקה הן על ידי התפתחות שיפוטית.¹⁰⁰⁰ אחרים סבורים כי אפשר לפתור את בעיית האחריות הנזיקית של בינה מלאכותית על ידי יצירת אישיות

9.1 מבוא

1000 ראו למשל CHESTERMAN, לעיל ה"ש 19, בעמ' 38.

משפטית שמתאימה לה, בהנחה שלאישיות זו יהיו מקורות פיננסיים עצמאיים שיאפשרו לה לשלם עבור הנזקים שגרמה.¹⁰⁰¹

ואולם אסדרה של מערכות בינה מלאכותית היא הכרחית – הן ברמה הגלובלית הן ברמה הלאומית. סדרת סיפורי הרובוטים של אייזק אסימוב מראה במידה מסוימת מדוע לא די בשלושת חוקי הרובוטיקה, ואחת היא עד כמה הם אלגנטיים ושוכי לב.¹⁰⁰² והרי אילו היו כללים אלו אפקטיביים לרגולציה של מערכות בינה מלאכותית, סדרת ספרי הרובוטים של אסימוב לא הייתה באה אל העולם.

אימוץ וולונטרי של מסגרת כללים אתיים רחבים ובלתי מחייבים בתעשייה או ברשויות הממשלה עלול לשמש עלה תאנה, שכן הוא יבטא הכרה מוסרית בקיומה של בעיה אך יהיה משולל גיבוי של מנגנוני פיקוח ואחריות משפטיים שיאפשרו לממש אותם. אסדרה שתיעשה נכון, לעומת זאת, תאפשר פישוט וייעול של תהליכי פיתוח, והבניה של תהליכים כלליים של בקרת איכות, ותקדם אפוא אמוץ ציבורי במערכות ובמוצרים. רגולציה שתכלול כללים שיבנו את אופן הפיתוח של בינה מלאכותית ואת התיעוד וממשל הנתונים בה תוכל למתן הטיות אלגוריתמיות מסוגים שונים – הן הטיות אלגוריתמיות שהן תולדה של כשלי פיתוח לא רצויים, המנוגדים לאינטרסים של מפתחיהן, הן הטיות אלגוריתמיות שמחצינות סיכונים שאינם מטרידים את המפתחים. יודגש: רגולציה מספקת ודאות גם לתעשייה – בוודאי יותר מהתפתחות שיפוטית, שאורכת מטבעה שנים רבות ומתאפיינת באקראיות מבחינת המקרים המגיעים לפתחם של בתי המשפט.

מערכות בינה מלאכותית כבר נמצאות בשימוש בקרב מגזרים רבים בישראל, מהן פרי פיתוח מקומי ומהן פיתוח של חברות זרות. יש להבטיח כי מערכות אלו פועלות מתוך שמירה על האינטרסים ועל זכויות היסוד של המשתמשים בהן או של אנשים אחרים המושפעים מפעילותן.¹⁰⁰³

Karolina Ziemianin, *Civil Legal Personality of Artificial* 1001
 ,CHESTERMAN ; *Intelligence: Future or Utopia?* 10 INTERNET POL'Y REV. 1 (2021)
 לעיל ה"ש 19, בעמ' 126.

1002 ראו אסימוב, לעיל ה"ש 47, וכן Murphy and Woods, לעיל ה"ש 365.

1003 נציבות זכויות האדם האוסטרלית, למשל, מדגישה שאסטרטגיה לאומית בנושאים טכנולוגיים צריכה לכלול במטרותיה העיקריות תעדוף רגולטורי של הגנה על זכויות אדם.
 ראו Australian Human Rights Commission, לעיל ה"ש 949, בעמ' 28.

זאת ועוד: כפי שניסינו להראות עד כה, מערכות מבוססות בינה מלאכותית מעוררות שאלות החורגות מתחום המשפט הפרטי. אפשר למנות שלושה טעמים לכך.

ראשית, כמה מן האתגרים הם אתגרים ציבוריים למחצה, כגון הטיות, שקיפות והגנת הפרטיות. לכן גם כשהסיכונים הנשקפים מבינה מלאכותית אינם נוגעים לזכויות אדם, אלא צבר של סיכונים כלליים שמקורם בשגיאות בפיתוח התוכנה, ברשלנות או בהזנחה,¹⁰⁰⁴ היכולת להתחקות עליהם תלויה במידת השקיפות של מערכת הבינה המלאכותית הנדונה.

שנית, בגלל הפגיעה הפוטנציאלית של מערכות אלו בשלמות הגוף ובקניין הפרטי יש להשית רגולציה ממשלתית בעוד מועד. עיצוב הרגולציה הזאת ברמה המהותית והמוסרית הוא פעילות בתחום הציבורי.

שלישית, כבר היום מופעלות מערכות לומדות ורשיות המדינה ימשיכו להרחיב את פעילותן באופן שעלול להשפיע על זכויותיהם של אזרחים, ובכלל זה על הזכות לכבוד ולשוויון.

פרק זה והפרק שלאחריו יניחו יסודות לרגולציה מכוונת זכויות של בינה מלאכותית בישראל. בפרק זה נציג את העקרונות המנחים שצריכים, לתפיסתנו, לעמוד ביסוד המדיניות בתחום זה.

9.2

עקרונות מנחים ליצירת מדיניות בינה מלאכותית מכוונת זכויות

על פי עמדת רוח ועדת תל"ם ורוח המיזם הלאומי
למערכות נבונות, על הרגולציה לשאוף להיות

9.2.1 האדם במרכז

מאפשרת¹⁰⁰⁵ ולהשתדל לצמצם את הנטל המוטל על כתפי השחקנים בזירת הבינה המלאכותית – משתמשים, ספקים, מפתחים או מפיצים במגזר הציבורי והפרטי גם יחד.¹⁰⁰⁶ רוח דומה עולה מהדוח העוסק ברגולציה של בינה מלאכותית במגזר הפיננסי, המדגיש את הצורך במניעת חסמים לחדשנות ובהקלות רגולטוריות על פיתוח בינה מלאכותית בישראל, וקורא בשלב זה לרגולציה עצמית של השוק ולשימוש בארגזי חול.¹⁰⁰⁷ גם מסמך המדיניות של הרשות לחדשנות נוטה להעדיף רגולציה מאפשרת, המעודדת פיתוח בינה מלאכותית ושימוש בה.¹⁰⁰⁸ אך מתווי המדיניות אינם צריכים לפעול אך ורק מתוך רצון לאפשר חדשנות וצמיחה בכואם לעצב הסדר רגולטורי ומוסדי: אין להזניח בשם עקרון מזעור ההתערבות הרגולטורית את ההגנה על חירויות וזכויות יסוד.

התכלית המרכזית של פיתוח בינה מלאכותית ומערכות לומדות צריכה להיות שירות המין האנושי – הן כקולקטיב הן כפרטים – באופן המיטיב עימו. עיקרון זה של חירות האדם והאוטונומיה שלו מצטרף אל עקרון קידום הטוב ועקרון מניעת הנזק ומבסס ביתר שאת את החשיבות של עקרון האדם במרכז, של החירות ושל האוטונומיה. באופן מעשי מתבטא עיקרון זה בפיתוח מערכות נבונות, פרישתן ושימוש בהן באופן המגן על זכויות היסוד ועל חירות האזרח ורואה בהן עיקרון ראשון במעלה. לא מדובר במס שפתיים שנועד להסתיר נהייה אחר קידום חדשנות או אחר אינטרסים כלכליים (כגון עידוד מגזר ההייטק), או אפילו אחר ניסיון ייעול של תהליכים במגזר הציבורי.

בהתאם לכך מטרתה של רגולציה של בינה מלאכותית היא להגן על המושאים של החלטות אלגוריתמיות – בני אדם. הגנה זו יכולה לבוא לידי ביטוי ביצירת מערכת כללים המאפשרת ודאות בנוגע להקצאת הסיכונים הנזיקיים, אך גם בליווי תהליכי פיתוח מערכות אלגוריתמיות, הטמעה שלהן ושימוש בהן באופן

1005 דוח ועדת חל"ם, לעיל ה"ש 103, בעמ' 63; המיזם הלאומי למערכות נבונות (2019), לעיל ה"ש 313, בעמ' 28; המיזם הלאומי למערכות נבונות (2020), א, לעיל ה"ש 981, בעמ' 30-31.

1006 ראו לדוגמה את התייחסות הצעת תקנות הבינה המלאכותית האירופיות (לעיל ה"ש 53) לשחקנים האלה. ראו גם לעיל בסעיף 4.1.1.2.

1007 אחיעז ואח', לעיל ה"ש 240, בעמ' 148.

1008 מסמך מדיניות הרגולציה של משרד החדשנות, לעיל ה"ש 314, בעמ' 97.

שמזמער את הסיכון לזכויות אדם¹⁰⁰⁹ ומאפשר בקרה אקס אנטה ואקס פוסט על החלטות אלגוריתמיות. הגנה על מושאי החלטות אלגוריתמיות פירושה גם פיתוח ועיגון של זכויות הנוגעות לשקיפות ושל סעדים בגין פגיעה בהן. למשל, יש ליידע משתמשים במערכות בינה מלאכותית שהם אינם מתקשרים עם אדם בשר ודם אלא עם צ'אט בוט.

כמו כן, בלי לקבוע מסמרות בעניין השאלה עד כמה רצוי להסתמך הסתמכות עיוורת על מערכות החלטה אוטומטיות בכל תחומי החיים, עקרון האדם במרכז משמעו גם הגדרת יישומים של בינה מלאכותית ומערכות אלגוריתמיות תומכות החלטה שמחייבים מעורבות או פיקוח אנושיים,¹⁰¹⁰ בין שמדובר באישור אנושי להחלטה אלגוריתמית ובין שמדובר באפשרות לערער עליה לאחר מעשה לאדם ממשי. זאת בייחוד במקומות שבהם השפעתם על חיי הפרט יכולה להיות משמעותית.

עניין זה חשוב במיוחד בנוגע ליישומים ממשלתיים של בינה מלאכותית, בין שהם פרי פיתוח פנימי ובין שהם תוצאה של מיקור חוץ, שכן החלטות אוטומטיות של יישומים ממשלתיים משפיעות על אזרחים. רגולציה של בינה מלאכותית יכולה לעודד בעקיפין פיתוח מערכות בטוחות, המשמרות זכויות אדם וחירויות אזרח, באמצעות התניית רכש ומכרזי פיתוח ממשלתיים של מערכות נבונות בעמידה בסטנדרטים מחמירים של שקיפות, בטיחות, תיעוד ובקרה.¹⁰¹¹

1009 ראו לעיל בפרק 5.

1010 ראו לדוגמה סי' 6.3.9 לדירקטיבה ה-ADM של קנדה, לעיל ה"ש 628; סי' 22 של התקנות הכלליות בדבר הגנת מידע (GDPR).

1011 לרגולציה באמצעות מכרזים ראו Christopher McCrudden, *Using Public Procurement to Achieve Social Outcomes* 28 NAT. RESOURCES FORUM 257 (2004) לשימוש במכרזים ציבוריים כדי לקדם הגנה על הפרטיות ראו לדוגמה Laurens Vandercruyssen et al., *PUBLIC PROCUREMENT OF SMART CITY SERVICES – MATCHING COMPETITION AND DATA PROTECTION (A Report in the Framework of the SPECTRE Project 2020)* (Research Project). ראו גם אלדד הבר וטל ז'רסקי "דרכי ההגנה על תשתיות חיוניות במרחב הסייבר בישראל" *משפט וממשל* יח 99, 136 (2017).

9.2.2. הדמוקרטיה במרכז בתחילת המאה ה-21 גאה וגעש ה"טכנור"

אופטימיזם" וניעורה תקווה שהעולם הדיגיטלי יביא איתו בשורה לדמוקרטיה ולפוליטיקה, ושסגולותיו של עולם זה – פתיחות, חופש ומידע – יתמכו במימוש הערכים הדמוקרטיים, ישלימו את המוסדות הדמוקרטיים ויקדמו אותם.

ואולם כיום, בתחילת העשור השלישי של המאה ה-21, ברי לנו שבטכנולוגיה טמונים גם סיכונים, והיא גובה מאיתנו מחיר, בדרכים שעוד לא כולן ידועות לנו עד תום. השפעתה מגיעה למעגלים רחבים, הנוגעים ליכולתם של מוסדות דמוקרטיים לתפקד; לשיח הציבורי, שהוא אבן הפינה של ההליך הדמוקרטי; ולאפשרות המימוש של בחירה חופשית ואוטונומית.

במערכות מבוססות בינה מלאכותית טמון סיכון ממשי של פגיעה בדמוקרטיה במובנה הרחב – אם באמצעות השפעה על השיח הציבורי והפצת רעיונות, אם באמצעות מכשור לשליטה באוכלוסייה, למעקב אחריה, לזיהויה ולמשטרה, ואם באמצעות היכולת שלהן לזרוע ספק ולערער על עצם היכולת לברר מהי המציאות ולהבחין בין מקור לזיוף ובין אמת לשקר. לכן יש לתת משקל של ממש לעיקרון "הדמוקרטיה במרכז" – גם, לעיתים, במחיר של התקדמות טכנולוגיות או חדשנות.

אפשר להשתמש בבינה מלאכותית כדי להשפיע לרעה על מהפכת המידע וכדי לקדם את מהפכת המעקב. למשל, יצירת טקסטים וקטעי וידיאו מזויפים באמצעות טכנולוגיות דיפ-פייק; מערכות לזיהוי קול ולזיהוי פנים המקדמות מעקב ומשטור; וגם מערכות לניתוח מידע שתכליתן "הנדסת בחירות". ואולם יש גם סוגים אחרים של מערכות בינה מלאכותית שטמונות בהן סכנות לדמוקרטיה.

מערכות חיזוי שיוכלו להעניק פשר למציאות עתידית ולהשפיע עליה ישנו את השיח הפוליטי ואת האתוס של חזון פוליטי, שתכליתו החברתית היא לנסות להעריך תרחישים ולתפקד במציאות לא ברורה. היתרון של פוליטיקה המבוססת על מערכות חכמות גדול, אך יש לשאול בכנות מה יקרה ביום שבו תהליך ממוכן, מבוסס מידע, של קבלת החלטות יוכל לספק תחזיות ולפתור בעיות עוד בטרם הפכו לבעיות; ומה תהיה המשמעות הדמוקרטית של מצב שבו נבחר הציבור יהיו רק כתובת אחרונה לבעיות שמכונות אינן יכולות לפתור.

גם אם כיום השימוש במכוונת מכוון רק לסייע להחלטות אנושיות, עם התגברות היכולת לקבל תחזיות טובות ומדויקות יותר יוכלו מכוונות גם לשפר ביצועים מערכתיים, לאסדר התנהגות ולפתור סכסוכים משפטיים בין רגולטורים לאזרחים. אזרחים יוכלו לפנות לרשות רגולטורית, להניח את נתוני המקרה המיוחד שלהם בידיה ולקבל פרשנות וייעוץ מדויקים שייתרו את הצורך בבתי משפט. מערכות לומדות יעסקו במתן רישיונות; יקבעו קנסות ויערכו שימועים לקבלת זכויות; יחקרו דיווחי מס; יחליטו את מי לעכב לבריקה מיוחדת בשרה התעופה; יחזו פשיעה ופשיעה חוזרת; ישרתו מערכות אכיפת חוק ויעניקו למקבלי החלטות ניתוח רב-ממדי המסתמך על מידע רחב ומאוזן בין שכבות שונות של תובנות. הכול בדרך מדויקת יותר, יעילה יותר, מהירה יותר וזולה יותר מבני אדם או ממערכות מחשב מתקדמות פחות. בשלב מתקדם יותר יוכלו מכוונות אף להציע לנו סדר חוקי טוב מהסדר הקיים ולומר לנו כיצד לצמצם את הטעויות בו – כלומר יוכלו לחוקק עבורנו את החוקים.

החלפת שיקול הדעת האנושי בשיקול הדעת של אלגוריתם מביאה עימה אתגרים הנוגעים לאחריות לתהליכי קבלת החלטות בדמוקרטיה. אחד מיסודות מערך האיזונים והבלמים הוא האפשרות לעקוב אחר שיקולי קבלת ההחלטות ולבקר אותם. איך יוכלו שופטים, לא כל שכן אזרחים מהשורה, למתוח ביקורת על החלטות של אלגוריתמים? לכן יצירת אמות מידה לקבלת החלטות, וקידום שקיפות, הסברתיות ואחריותיות באשר להחלטות של מערכות לומדות, הם האתגרים המהותיים העומדים היום לפני דמוקרטיה. הם שיבטיחו בעשור הקרוב הגנה על הליך הוגן ועל זכויות חשודים וימנעו הטיית אפשרויות בהחלטותיהן של מערכות המבוססות על מאגרי מידע.

חוסר נכונותם של פוליטיקאים לוותר על שימוש בכלים הפוגעים בדמוקרטיה אם הם יכולים לסייע להם להיבחר או לשרת אותם בדרכים אחרות הוא בעיה של ממש. יתר על כן, מתברר שאת ההחלטות בנוגע למציאות הדיגיטלית לא מקבלות ממשלות נבחרות אלא קומץ חברות ענק מכוונות רווח. למרות הפרדוקס שבדבר החשש האמיתי אינו מפני חברות מסחריות דווקא, אלא משאיפתן של מדינות לנצל את הטכנולוגיות כדי להשיג שליטה חסרת תקדים באוכלוסייה, בזמן בחירות, וגם זו בזו. כוחן העצום של הפלטפורמות הדיגיטליות, שהן לובי משמעותי נגד אסדרה, והחשש להציע פתרונות שיפגעו ביכולתן של המדינות

הליברליות להתמודד עם המדינות הלא־חופשיות, דוגמת סין – גם הם מעמידים אתגרים גדולים.

הדמוקרטיה נמצאת היום בצומת דרכים, ובשנים הקרובות יידרשו מדינות וממשלות להחליט החלטות מכריעות כדי להתמודד עם הכוחות המערערים אותן. יהיה צורך למפות את כיווני ההתפתחות הטכנולוגיים ולחשוב לאילו מהם עלולות להיות השלכות שליליות ארוכות טווח על הדמוקרטיה ואילו מהם ישפיעו לטובה על הטוב המשותף. סטואו בויד, עתידן אמריקאי, כינה זאת "הצורך במהפכת האביב האנושי"¹⁰¹².

9.2.3. נדרשת אוריינות דיגיטלית מצד מקבלי החלטות

אחת הבעיות הקשות בתפר שבין משפט לטכנולוגיה היא הרמה הדיגיטלית הנמוכה של מקבלי החלטות, שרובם אינם "ילידים דיגיטליים". הם אינם מבינים את עומק

הבעיות, אינם שואלים את השאלות הנכונות ואינם נדרשים לתמונת הכוחות המלאה. אוריינות דיגיטלית היא סל היכולות והמיומנויות – הטכנולוגיות, הסוציולוגיות והקוגניטיביות – הנדרשות כדי להבין את הסביבה הדיגיטלית ולתפקד בתוכה. היא כוללת יכולת לחפש ולמצוא מידע, וגם מודעות לכוחות העומדים מאחורי המהפכה הדיגיטלית ולהטיות שהם יוצרים, הקפדה על כללי התנהגות וכישורי חשיבה ביקורתית.

אוריינות דיגיטלית משמעה היכולת להבין לאן מתפתחת הטכנולוגיה, לפחות בטווח הקצר: לנתח את השוק ולבחון היכן נמצאים כספי המחקר והפיתוח של חברות הענק ומהם הפוטנציאלים שהן רושמות כדי לעגן פיתוחים טכנולוגיים חדשים. בעל אוריינות דיגיטלית מבין שאפשר להטות ולכוון פיתוח טכנולוגי הן באמצעים מסחריים הן באמצעים רגולטוריים; הוא יודע גם שמעצבי מדיניות הם חלק מהתפתחות הטכנולוגיה ולא רק מתבוננים מן הצד, ולכן הם נושאים באחריות. בהקשר של מדיניות בינה מלאכותית, נתעכב על שלושה עניינים.

1012 תהילה שוורץ אלטשולר "דמוקרטיה וטכנולוגיה" דמוקרטיה עכשוו: סוגיות ואתגרים במאה ה-21 (דנה בלאנדר ומאיה גייר עורכות, 2021).

ראשית, לא כל מהפכה טכנולוגית מתרחשת בסמיכות היסטורית גדולה כל כך למהפכות טכנולוגיות אחרות. למעשה, מהפכת הבינה המלאכותית יכולה להפיק לקחים משתי המהפכות האחרות בדור האחרון – מהפכת החומרה ומהפכת הקישוריות. נמנה כמה לקחים אפשריים:

- פוליטיקאים גוררים רגליים גם כשהצורך לאסדר כבר ברור; ומנגד, התעשייה טוענת כי לפוליטיקאים אין את היכולת וההבנה הנדרשים כדי לאסדר ומוטב לפיכך להימנע לגמרי מאסדרה;

- תמת ה"חדשנות" מתערבבת לעיתים בפזמון החוזר "נשבור עכשיו, נתנצל אחר כך", או בניסוח חריף יותר "לא ידענו עד כמה הכלי שפיתחנו מסוכן. לא ידענו לאילו נזקים הוא עלול לגרום. אנחנו מצטערים. סליחה"¹⁰¹³;

- טענת התעשייה כי נדרשת אסדרה מתערבבת בניסיונות שלה עצמה להשיג הסדרים רגולטוריים מקילים, בהשקעות עתק בשדלנות ובהשפעה על מקבלי החלטות, ואף באיום להפסיק מתן שירותים באזורים גאוגרפיים שבהם מדובר על אסדרה משמעותית;¹⁰¹⁴

- כל מהפכה נתפסת כתהליך משחרר, שמעודד דמוקרטיה ושוויון, אבל בהיעדר אסדרה החזק לוקח הכול וכמה חברות ענק שולטות בכל העולם. כשמדובר במהפכת הבינה המלאכותית, שכוח מחשוב וכמויות גדולות של נתונים זמינים הם תנאים מוקדמים עבורה לצורך אימון מודלים, החשש מפני שליטה אפריורית של ענקיות הטכנולוגיה הוא מוחשי;¹⁰¹⁵

- הסיקור התקשורתי הדיכוטומי של טכנולוגיה, הנע בין אופוריה להיסטריה, בין אוטופיה לדיסטופיה, משחית את יכולתם של מקבלי החלטות הניזונים מהתקשורת לדון בנושאים האלה באופן שקול ומאוזן.

1013 יובל דרור "הבאים בחור לבקש סליחה" מגזין דה מרקר (דצמבר 2022).

1014 ראו למשל את קריאתו של סאם אלטמן, מנכ"ל חברת Open AI, לאסדרה – בד בבד עם אימום להפסיק לתח גישה לשירותי החברה שלו באיחוד האירופי. Shiona McCallum and Chris Vallance, *ChatGPT-maker U-turns on threat to leave EU over AI law*, BBC (26.5.2023)

1015 AI Now, *Confronting Tech Power* (2023)

שנית, נדרש רובד ביניים בין הבנת הטכנולוגיה ובין יצירת מדיניות בעניינה. מייקל היידן, שעמד בעבר בראש הסוכנות לביטחון לאומי (NSA) והסי-איי-איי, כתב על מרחב הסייבר: "יש מעט תופעות חברתיות חשובות ומדוברות כל כך שמבינים אותן מעט כל כך"¹⁰¹⁶. היידן הדגיש שהוא אינו מתכוון לטכנולוגיה ספציפית או להפעלת כלי כזה או אחר, אלא להיעדר מסגרת תפיסתית שמאפשרת להבין את ההשלכות של השימוש בטכנולוגיות שונות, ולא של סייבר דווקא. במסגרת זו הוא מתכוון להבנת המשמעות של מערכות טכנולוגיות, ליכולת לדמיין את האפשרויות החדשות שהן מביאות איתן וליכולת להבין את השלכותיהן על המוסר החברתי ועל שלד השיטה המשפטית. כיוון שמסגרת זו חסרה לעיתים קרובות נוצרים פערי הבנה, בייחוד בנושאים בעלי השלכות רחבות כמו הבינה המלאכותית.

9.2.4. נדרשת המשגה חדשה של מערכות ויכולות בתחום הבינה המלאכותית

המושג בינה מלאכותית הוא מטפורה פוליטית-טכנולוגית, אולי המטפורה העוצמתית ביותר בתקופה הנוכחית. בדומה למטפורות אחרות, היא חודרת לזיכרון, מעוררת רגשות, משפיעה על עמדות ומעצבת ציפיות לגבי העתיד. המושג הוא מטפורה פוליטית כיוון שהוא מייחס תכונות אנושיות למערכות טכנולוגיות. כל דבר נעשה "חכם" – מטרמוסטט ועד רחפן. השימוש במטפורה של אינטליגנציה, בינה, חוכמה בהקשרים רבים כל כך מלמד על עוצמתה. ההשוואה בין מערכת בינה מלאכותית למוח האנושי יוצרת קרבה ודמיון, והללו מובילים להטמעה חברתית של הרעיון שהמכונות עובדות כמו מוח אנושי, מבצעות פעולות אנושיות בדרך שבה מבצעים אותן אנשים, ולמעשה מתחרות בבני האנוש.

למי יש אינטרס לשמר את המטפורה של הבינה המלאכותית ולהאניש את המערכות הטכנולוגיות? אפשר לחשוב, למשל, על יצרנים של מערכות שאינם רוצים לקבל אחריות למעשים של המערכות האלה, להטיות שלהן ולדעות

הקדומות המוטמעות במאגרי המידע שהן מתאמנות עליהם, לפגיעות ולנזקים שהן עתידות לגרום, לחוסר ההבנה המובנה שלהן בנוגע להשלכות החברתיות של החלטותיהן ולקושי לייצר בשבילן שקיפות ותקינה. כמה נוח לומר "הבינה המלאכותית טעתה" ולהרחיק את האחריות לטעויות ממי שיצר אותה או נתן לה את המשימה.

אפשר לחשוב גם על מי שרוצה שנראה במכשירי קצה טכנולוגיים כמו טלפונים סלולריים או עוזרים דיגיטליים בני משפחה, נמסור להם מידע פרטי וניתן אמון בהחלטות שלהם. אחרי הכול, אלו הן החלטות שמתקבלות בדיוק כמו ההחלטות שלנו. אפשר כמובן לחשוב על מי שרוצה שנאמין שהמערכות הטכנולוגיות מתפתחות באופן עצמוני, ממש כמו אבולוציה אנושית, ולא נשים לב לעובדה שענקיות הטכנולוגיה נעשו שומרות סף גם בהקשרים של פיתוח מערכות לומדות.

יש להשתדל שהמטפורה הפוליטית של בינה מלאכותית לא תעצב את הדמיון שלנו באשר לעתיד ושהנרטיב העוצמתי שבו משתמשות התעשייה והתקשורת המסקרת אותה לא ישפיע על מקבלי ההחלטות. מוטב למשל להשתמש במטפורה של "למידה", שאינה קשורה בהכרח ל"הבנה", ולהידרש לטכניקות פעולה ספציפיות כמו למידת מכונה, מערכות גנרטיביות, מערכות מבוססות ידע או סטטיסטיקה. ולחלופין, להידרש להקשר היישומי של כל טכנולוגיה: למשל, תמיכה באבחון רפואי באמצעות למידת מכונה; זיהוי ביומטרי במרחבים ציבוריים לתכליות אכיפת חוק; וכדומה. הצורך בהמשגה מחודשת של תהליכי הבנה, היסק, יצירתיות, חשיבה, דמיון וכיוצא באלה, כאשר הם מתבצעים על ידי מכונות, מתעורר גם הוא לאור השיפור ביכולות התקשורת בין המכונה למשתמשים.

המהפכה הטכנולוגית מעוררת צורך בחשיבה מחודשת על זכויות יסוד משתי בחינות.

9.2.5. נדרש פיתוח זכויות לתושאי החלטות של בינה מלאכותית

ראשית, יש צורך לצקת משמעות חדשה לתאוריה החוקתית של זכויות האדם בצורתה היום. למשל,

הזכות לחופש ביטוי אינה יכולה להיות מושתתת רק על התפיסה שמלחמה

בביטוי שקרי נעשית באמצעות עוד ביטוי¹⁰¹⁷ או על התפיסה שבהתקיים מבחר מספק בשוק הרעיונות והדעות תצוף האמת מאליה. בעידן שמתאפיין בהצפה של רעיונות שמנווטים על ידי מערכות אלגוריתמיות יש צורך בחשיבה מחודשת על התאוריה המכוננת של הזכות לחופש ביטוי. כך טענו למשל דיוויד ביבר וג'ייסון סטנלי, שקבעו שחופש הביטוי במובנו כיום מסכן את הדמוקרטיה לא פחות משהוא תורם לה;¹⁰¹⁸ וטים וו, שקבע שהשימוש הנוכחי בתאוריה של חופש הביטוי גורם לדיכוי היכולת להתבטא.¹⁰¹⁹

בדומה, הזכות לפרטיות בעידן של מעקב המונים איננה יכולה להיות אך ורק בעלת אופי אינדיווידואלי, כלומר זכותו של הפרט לשלוט במידע על אודותיו, אלא עליה לקבל פרשנות קולקטיבית הנוגעת לחוסר היכולת לנהל הליך דמוקרטי תקין בשעה שלפרטים בחברה אין אוטונומיה לבחור בחירה חופשית כיוון שהם נתונים להשפעה מניפולטיבית המבוססת על מעקב ושליטה.¹⁰²⁰

דוגמה נוספת היא הבנה מחודשת של הזכות לכבוד האדם ביחס למכונה. ס' 22 לתקנות הכלליות בדבר הגנת מידע (GDPR) קובע כי "למושא מידע תהיה הזכות שלא להיות מושא של החלטה המבוססת אך ורק על תהליכים אוטומטיים, לרבות פרופיילינג, שיש לה תוצאות משפטיות בנוגע אליו או לאחרים, או תוצאות משמעותיות אחרות". זוהי הזכות למעורבות אנושית, הנתפסת כנגזרת של הזכות לכבוד האדם. אומנם היקף פרישתה של הזכות האירופית מוגבל. היא לא תאפשר, למשל, לערער בזכות על החלטות שיפוטיות אנושיות שהתקבלו

1017 עמדה שהיטיב לבטא שופט בית המשפט העליון של ארצות הברית לואיס ברנדייס: "The remedy to be applied to [falsehood and fallacies] is more speech, not enforced silence" (WHITNEY V. CALIFORNIA 1927)

David Beaver and Jason Stanley, *Neutrality*, 49 PHILOSOPHICAL TOPICS 1018 165 (2021)

Tim Wu, *Is the First Amendment Obsolete? The Knight First Amendment* 1019 INSTITUTE (1.9.2017); להרחבה ראו תהילה שוורץ אלטשולר, אסף וינר ואייל זילברמן מחווה לאסדרת רשתות חברתיות בישראל 26-27 (הצעה לסדר 51, המכון הישראלי לדמוקרטיה 2023).

1020 להרחבה ראו תהילה שוורץ אלטשולר פרטיות: מלכת זכויות האדם בעולם הדיגיטלי (פרלמנט 83, אתר המכון הישראלי לדמוקרטיה 2019).

מתוך הסתייעות במערכות תומכות החלטה שיפוטית – למרות חששות להטיית המודלים¹⁰²¹ ולמרות הנטייה האנושית להנסגת הדעת לנוכח קבלת החלטות אוטומטית.¹⁰²²

החשש הוא אפוא מפני שחיקת הסובייקט האנושי בתור שכזה, והנכחת תחושות של היעדר נראות, חוסר ערך ותסכול לנוכח מציאות שרירותית. מנגד, מחקרים מלמדים על פער ב"הוגנות הנתפסת", כלומר בתחושת הכבוד האישי, ההוגנות והצדק הפרוצדורלי של בעלי דין ביחס למכונות.¹⁰²³ בניסויים הראשונים הנבדקים נוטים לסמוך על שופטים אנושיים ולראות בהם מי שמממשים את זכותם ליחס מכבד. אבל כאשר יוצרים עבור הנבדקים תהליך שבו ניתנת להם האפשרות לדבר למכונה ולהישמע בה, וגם מעניקים לבעלי הדין זכות לקבל הסבר לוגי בעניין ההחלטה, אז הפער בהוגנות הנתפסת נעלם ולבעלי הדין לא חשוב אם מדובר בשופט בשר ודם או ברובוט שופט.

אפשר היה להניח ששימוע מול מכונה ייחשב חסר משמעות או חלול, כי לרוב אנחנו מקשרים את העובדה ששמעו אותנו או ראו אותנו עם העובדה שמי שעשה זאת הוא מישהו שיש לו יכולת לגלות אמפתיה. אבל ככל שיכולת המכונות תשתכלל – מבחינת יכולתן לזהות שפה, תנועות גוף, הבעות פנים ולכן גם רגשות, ולהגיב בהתאם – כך ילך החשש מפניהן ויקטן והן יוכלו לשמש מקור

1021 ראו לעיל בפרק 6.

1022 Saar Alon-Barkat and Madalina Busuioc, *Human-AI Interactions in Public Sector Decision Making: "Automation Bias" and "Selective Adherence" to Algorithmic Advice*, 33 J. PUB. ADMIN. RES. & THEORY 153 (2022); Stavros Zouridis, Marlies van Eck, and Mark Bovens, *Automated Discretion*, in DISCRETION AND THE QUEST FOR CONTROLLED FREEDOM 313 (Tony Evans and Peter Hupe eds., 2020); Matthew M. Young, Justin B. Bullock, and Jesse D. Lacey, *Artificial Discretion as a Tool of Governance: A Framework for Understanding the Impact of Artificial Intelligence on Public Administration*, 2 PERSPECTIVES ON PUBLIC MANAGEMENT AND GOVERNANCE 301 (2019)

Benjamin Minhao Chen, Alexander Stremitzer, and Kevin Tobia, 1023 *Having Your Day in Robot Court*, 36 (1) HARVARD JOURNAL OF LAW AND TECHNOLOGY 128 (2022); idem, *Would Humans Trust an A.I. Judge? More Easily Than You Think*, SLATE (28.2.2023)

לתחושה שהן מסוגלות להעניק צדק פרוצדורלי. ייתכן שלא ירחק היום ואנשים יבקשו להישפט על ידי מכונות ולא על ידי שופט אנושי, שעלול להיות עצבני לפני ארוחת הצהריים. לכן ייתכן שהגדרת הזכות לכבוד תשתנה. תופעות אלו לא תהיינה ייחודיות לתחום המשפט – האבחנה הרפואית יכולה לשמש לנו דוגמה דומה. ייתכן שמתופלים ידרשו שימוש במערכת לומדת – לא רק משום שהיא תיראה להם מדויקת יותר, אלא גם כי היא תהיה מסוגלת לספק להם קשר אישי ומכבד. בזכות השפה הטבעית של כלים מבוססי בינה מלאכותית הולך ונעשה נעים ונוח יותר לתקשר עימם, וככל שייעשו נעימים ונוחים יותר כך תגבר הנכונות לממש את הזכות לכבוד דווקא באמצעות תקשורת עם מכונות. כאמור, שינוי זה יחייב חשיבה מחדשת על הזכות לכבוד ומשמעותה.

שנית, יש ליצור זכויות דיגיטליות חדשות, שלא היה צורך בהן בעבר,¹⁰²⁴ ובעיקר זכויות של פרטים בשעה שהם באים במגע עם מערכות ממוכנות מבוססות אלגוריתמיות. למשל: יש צורך להעצים את מושאי ההחלטות ולהעניק להם זכות לשקיפות אלגוריתמית במובן של הסברתיות;¹⁰²⁵ ולפתח זכות לאינטראקציה מיועדת עם בינה מלאכותית,¹⁰²⁶ כלומר הזכות שמוצרים מבוססי בינה מלאכותית (צ'ט בוטים למשל) יציגו את עצמם בתור כאלה ולא יתחזו לאנשים בשר ודם.¹⁰²⁷ צד נוסף של אינטראקציה מיועדת עם בינה מלאכותית נוגע לדימויים מלאכותיים שנוצרו על ידי בינה מלאכותית יוצרת.¹⁰²⁸ יש מקום לשקול חקיקה שתייצר חובת ידוע במקרה של דימויים כאלו.

1024 על זכויות דיגיטליות ראו, Dafna Dror-Shpoliansky and Yuval Shany, *It's the End of the (Offline) World as We Know It: From Human Rights to Digital Human Rights – A Proposed Typology*, 32 Eur. J. Int. Law 1249 (2021)

1025 יש פרשנים שסבורים שהשילוב של ס' 13-15 וס' 22 ב-GDPR מקנים למושאי המידע זכות להסבר על האלגוריתמים המעבדים את המידע שלהם. ראו לדוגמה Andrew D. Selbst and Julia Powles, *Meaningful Information and the Right to Explanation*, 4 Int'l Data Privacy Law 233 (2017)

1026 ראו לעיל ה"ש 333.

1027 ראו לדוגמה ס' (1)52 להצעת תקנות הבינה המלאכותית האירופיות, לעיל ה"ש 53; ס' (a)(V)5 להצעת חוק הבינה המלאכותית הברזילאית, לעיל ה"ש 54.

1028 השוו גם לס' (3)52 להצעת תקנות הבינה המלאכותית האירופיות, לעיל ה"ש 53; ס' 17 לכללי הסינחזה העמוקה, לעיל ה"ש 607.

9.2.6. נדרשת רגולציה יש לזכור שתהליכי חקיקה מפגרים תמיד אחרי המציאות. הטכנולוגיה מתקדמת בקצב מהיר, ובמדינות כמו ישראל, שבהן תהליכי החקיקה הם איטיים מאוד, נוצר פער גדול במיוחד. למרבה המזל, בינה מלאכותית איננה הטכנולוגיה המתפתחת הראשונה שהציבה את האתגרים הללו, ואפשר לשאוב השראה מההתמודדות המשפטית עם טכנולוגיות קודמות כגון טכנולוגיות מידע, טכנולוגיות גנטיות וטכנולוגיות רבייה, וללמוד מהטעויות שנעשו בתחומים אלו.

ראשית, אין לאסדר טכנולוגיה ספציפית, משום שזהו מתכון בטוח להתיישנות הרגולציה, אלא להשתדל לנסח עקרונות מנחים והגדרות כלליות שיהיו גמישים די הצורך לאכיפה עתידית. כך, חוק הבוטנים בקליפורניה¹⁰²⁹ ובישראל¹⁰³⁰ למשל, עורר ביקורת; ההגדרה של האזנת סתר בחוק האזנת סתר¹⁰³¹ לעומת זאת, איננה תלויה טכנולוגיה ולכן היא חסינת עתיד. התקנות לאסדרת הבינה המלאכותית באירופה עוצבו כך שיהיו חסינות-עתיד ולכן הן שומרות על ניטרליות טכנולוגית.

כך למשל, בתקנות האירופיות לא כוללת ההגדרה של מערכות הבינה המלאכותית זיקה לטכנולוגיה מסוימת.¹⁰³² גם אנו סבורים כי הגדרות כלליות ומבנים רגולטוריים עקרוניים עדיפים. לחלופין אפשר לאמץ גישה תלויה מגזר (ולא תלויה טכנולוגיה), שכן גישה כזאת יכולה לעקוף את הצורך בהגדרה (למשל, מכשור רפואי).

שנית, יהיה נכון לאמץ תפיסה רגולטורית המבוססת על ניהול סיכונים, שכן היא מאפשרת לעגן מראש קווים אדומים – כלומר לציין באילו יישומים של בינה מלאכותית טמונה סכנה גדולה כך עד שיש לאסור קטגורית על פיתוחם ועל הפצתם. באשר ליישומי בינה מלאכותית מותרים, אפשר לייצר מרחבים בטוחים, שרמת ההתערבות הרגולטורית בהם תוגדר בהתאם לאומדן הסכנה הנשקפת מהם.

Renee Diresta, *A New Law Makes Bots Identify Themselves: That's the Problem*, WIRED (24.7.2019) 1029

עומר כביר "מומחים נגד תזכיר חוק הבוטנים: אוילי ולא ניתן לאכיפה" כלכליסט (16.6.2022) 1030

ס' 1 לחוק האזנת סתר, התשל"ט-1979, ס"ח 118. 1031

ס' 3(1) להצעת תקנות הבינה המלאכותית האירופיות, לעיל ה"ש 53. 1032

המודל של ניהול סיכונים הוא מסגרת חקיקה חסינת עתיד, שכן היא מאפשרת לחזור ולהגדיר את הקווים האדומים ואת המרחבים הבטוחים על פי התפתחויות טכנולוגיות, חברתיות ורגולטוריות. אחת הדרכים להעריך מחדש את עוצמת הסכנה הנשקפת ממערכות טכנולוגיות מסוג מסוים ואת אופן האסדרה היא ארגו החול הרגולטורי – מעין מרחב גמיש ובטוח (ליזמים ולמפתחים), שבו המאסדר מלווה פיתוח והטמעה של מערכת קונקרטיית כדי ללמוד אותה ואת האופן שבו רצוי לאסדר אותה.

יש שלוש מסגרות אסדרה נפוצות של בינה מלאכותית: עקרונות, זכויות וניהול סיכונים. מסגרת מבוססת עקרונות כוללת מערך עקרונות ליבה אתיים; מסגרת מבוססת זכויות מתמקדת

9.2.7 נדרשת מסגרת אסדרה מעורבת של עקרונות, זכויות וחקיקה

בהגנה על זכויות האדם ועל החירויות של מי שמושפעים מיישומי טכנולוגיות שמבוססות על בינה מלאכותית; ומסגרת מבוססת סיכונים מציבה במרכז את הסכנות והסיכונים האפשריים וגוזרת את האסדרה מרמת הסיכון, ובכלל זה אף אוסרת קטגורית על השימוש בטכנולוגיות מסוימות או שימוש לתכליות מסוימות. זוהי מסגרת רבת-עוצמה המסתמכת על העיקרון של זהירות מונעת, עמוד התווך באתיקה רפואית: *primum non nocere* (ראשית אל תזיק).¹⁰³³ על פי עקרון הזהירות המונעת, נטל ההוכחה באשר לבטיחות של טכנולוגיה מסוימת הוא על מי שמעוניין להשתמש בה ולא על הציבור.¹⁰³⁴

אחד הטיעונים בדבר הצורך במסגרת אסדרה מבוססת עקרונות וכללי אתיקה¹⁰³⁵ של מערכות בינה מלאכותית הוא שיש צורך למלא את החלל עד שתושג הסכמה מדינתית או בינלאומית באשר לדרכי האסדרה. כללי אתיקה יכולים לשמש גשר במצב הביניים שבין פריצה של טכנולוגיה לבין יצירת רגולציה מדינתית

1033 ראו לעיל ה"ש 366.

1034 WINGSPREAD STATEMENT ON THE PRECAUTIONARY PRINCIPLE (The Global Development Research Center, GDRC 1998). ראו גם Kaminski, לעיל ה"ש 1004, בעמ' 29-22.

1035 ראו לעיל בפרק 3 את פירוט העקרונות האתיים המרכזיים. מסמך מדיניות הרגולציה של משרד החדשנות, לעיל ה"ש 314, בעמ' 103-109.

או בינלאומית ההולמת אותו. מחויבות מוסרית של עובדים בתאגידים (למשל מהנדסים וכותבי תוכנה), ושל תאגידים כלפי לקוחות ומשתמשים, אף היא יכולה לשמש לעיתים שוט חזק לא פחות מרגולציה ממשלתית מחייבת.¹⁰³⁶

אלא שמול טיעונים אלו יש גם טיעונים חזקים אחרים, שלפיהם הניסיון לקרוא להסתפקות בכללי אתיקה בלבד הוא טיוח אתי, כלומר ניסיון ליצור מראית עין של אחריות אגב התחמקות ממחויבות מעשית.¹⁰³⁷ מהלך כזה יכול לכלול הימנעות מכוונת מהגררות בהירות, מוחשיות ומפורטות של העקרונות האתיים; אי-קביעת מנגנוני אכיפה ברורים שענישה בצידם במקרה שאין מצייתים לכללי האתיקה; היעדר התייחסות למגוון החוליות בשרשרת הערך של המערכות והמוצרים והתמקדות בחוליה של המהנדסים וכותבי הקוד בלבד. לטיוח אתי כמה מטרות. מטרה אחת היא חיזוק האמון של משתמשים במוצרים מבוססי בינה מלאכותית (בין שמדובר באמון בעניין בטיחות המוצרים ובין שמדובר באמון בעניין ההוגנות של הליך הייצור שלהם למשל). מטרה אחרת היא מתן חירות לתעשייה והתחמקות מרגולציה ממשלתית או בינלאומית כופה.

לכן לתפיסתנו לשיקולי אתיקה נודעת חשיבות מבחינה זו שהם מבטאים תפיסה מוסרית עקרונית באשר לפיתוח מערכות מבוססות בינה מלאכותית, אבל יש לזכור שבהיעדר רגולציה שיקולים של עשיית טוב יהיו לעולם משניים לשיקולים של עשיית רווח עבור משקיעים. יתרה מזו, אם סטנדרטים אתיים יהיו וולונטריים בלבד, חברות יוכלו לבחור אילו מהם לאמץ ואילו להזניח.¹⁰³⁸ לכן אנו סבורים ששלוש המסגרות אינן סותרות זו את זו ויש לשלב ביניהן. מן הראוי ליצור מסגרת היברידית שמשלבת אסדרה רכה (עקרונות אתיים) וניהול סיכונים באמצעות הוראות חוק וכללים רגולטוריים נוקשים.

Elettra Bietti, *From Ethics Washing to Ethics Bashing: A View on Tech Ethics from Within Moral Philosophy*, FAT* '20: PROCEEDINGS OF THE 2020 CONFERENCE ON FAIRNESS, ACCOUNTABILITY, AND TRANSPARENCY (2020)

1037 ראו Whittaker et al., לעיל ה"ש 450.

1038 JACOB TURNER, ROBOT RULES: REGULATING ARTIFICIAL INTELLIGENCE 210, 212–213 (2019)

9.2.8. נדרשת גמישות באשר לתועד ההתערבות באמצעות אסדרה

גם אם יש הסכמה בדבר תכלית האסדרה, שאלת העיתוי היא שאלה משמעותית. במקרים מסוימים יש יתרונות ברורים לרגולציה מוקדמת – לפני שמוצרים שמבוססים על

טכנולוגיה מסוימת חודרים לשוק. טכנולוגיה שנתפסת מסוכנת במיוחד, אם פיזית (כמו רכב אוטונומי) ואם מוסרית (למשל יצירת תכנים שמסיתים לטרור בצורה מלאכותית), הגיוני לאסדר מראש. גם במקרים קיצוניים פחות התערבות מוקדמת יכולה להועיל במה שנוגע לעיצוב כיווני המחקר ולהסתת משאבים להשקעה בפיתוח. גרגורי מנדל טען שהתערבות בשלב מוקדם עשויה לצמצם את ההתנגדות של בעלי עניין, מפני שבשלב זה הושקעו משאבים מעטים יחסית והעלויות השקועות נמוכות יותר, ומפני שהתעשייה והציבור משועבדים פחות לסטטוס קוו.¹⁰³⁹ עם זאת, ייתכנו מקרים שבהם עדיף להמתין ולהתמודד עם בעיות כשהן מתעוררות ולא לנסות לצפות אותן. האינטרנט נחשב לעיתים קרובות דוגמה לטכנולוגיה שהפיקה תועלת מגישה כזאת. בשלב מוקדם, נטען, קשה יותר לבצע תחזיות ותחזיות מדויקות לגבי עלויות ותועלות.¹⁰⁴⁰

כשפרופיל הסיכון של טכנולוגיה מתפתחת אינו ודאי גוברת לעיתים הנטייה לנקוט זהירות מרבית, שכן ההמתנה לראיות חותכות עלולה לאפשר לסכנות להתממש ולגרום נזק, לעיתים כזה שלא ניתן לתיקון – בין שמדובר בפגיעה באוכלוסיות רחבות (חידק מהונדס) ובין שמדובר בהשפעה על השיח הציבורי (בוט נאו-נאצי). מקצת החששות בעניין בינה מלאכותית הם מסוג זה. ניק בוסטרום הוא כנראה הקול המוכר ביותר המתריע על סכנות כאלה. לדבריו, "הגישה שלנו לסכנות קיומיות אינה יכולה להיות גישה של ניסוי וטעייה. אין הזדמנות ללמוד מטעויות".¹⁰⁴¹ לתפיסתו, יש לפעול כדי למקסם את ההסתברות לתוצאה המוגדרת "בסדר", כלומר כל תוצאה שאינה מביאה לידי סכנה קיומית.

Gregory N. Mandel, *Emerging Technology Governance, in* INNOVATIVE GOVERNANCE MODELS FOR EMERGING TECHNOLOGIES 62 (Gary Marchant, Kenneth Abbot, and Braden Allenby eds., 2013) 1039

CASS R. SUNSTEIN, LAWS OF FEAR: BEYOND THE PRECAUTIONARY PRINCIPLE 58 (2005) 1040

Nick Bostrom, *Existential Risks: Analyzing Human Extinction Scenarios and Related Hazards*, 9 JOURNAL OF EVOLUTION AND TECHNOLOGY 3 (2002) 1041

האתגר כמובן הוא לקבוע מתי יש לסכנה המשוערת בסיס מציאותי אמין. התגובות המשפטיות והרגולטוריות שלנו אינן צריכות, מן הסתם, להיות מבוססות על החזיונות הדיסטופיים הקשים ביותר של עתידנים ומדע בדיוני, אך בוודאי אינן יכולות להתעלם מהשלכות רחבות לטובת יתרונות מידיים.

איננו תמימים. אין ספק שעצם ההחלטה מתי להתערב ובאילו מקרים היא תוצאה של פשרות בתחום המדיניות הציבורית. למשל, אם ננקוט זהירות יתר ולא נתיר למכוניות ללא נהג לעלות לכביש כיוון שמכונית אוטונומית עלולה לגרום למוות, בני אדם ימשיכו לנהוג והם עצמם יביאו לכמות עצומה של נפגעים.¹⁰⁴² גם באבחון רפואי יש אפשרות שבינה מלאכותית יכולה לאבחן טוב מבני אדם, אבל אין פירוש הדבר שהיא חפה מטעויות. ואולם ניסיון העשורים האחרונים מלמד כי לתעשייה יש נטייה "לשבור דברים ואחר כך להתנצל", והכול בשם החדשנות. כפי שראינו, הטיעון בזכות החדשנות משמש פעמים רבות כיסוי לרצון להרוויח עוד ועוד, גם אם הדבר נעשה על חשבון אינטרסים ציבוריים חשובים אחרים.

9.2.9. יש להטיל איסורים
עקרוניים על שימוש
במערכות מסוימות,
אך האיסורים צריכים
להיות מצומצמים ומוגבלים
וגישה האירופית המסתמנת מבחינה בין
סוגים שונים של מערכות נבונות לפי רמת
הסכנה הנשקפת מהן ואוסרת אפריורית על
השימוש בסוגים מסוימים של מערכות נבונות.
הפרלמנט האירופי, למשל, קרא לאסור על
שימוש במערכות זיהוי פנים למטרות שיטור
ואכיפת חוק,¹⁰⁴³ והתקנות האירופיות כוללות קטגוריה של "מערכות אסורות";¹⁰⁴⁴

1042 בדומה, גישה תוצאתנית מעדיף כלי רכב אוטונומיים, חרף הסיכונים האינהרנטיים להם, מנהגים אנושיים, משום שמעבר לכלי רכב אוטונומיים מצמצם את מעורבות הגורם האנושי, האחראי ליותר מ-90% מהתאונות. ראו לדוגמה, Bryant Walker Smith, *Automated Driving and Product Liability*, 2017 Mich. St. L. Rev. 1 (2017)

1043 ראו לעיל ה"ש 563.

1044 הצעת תקנות הבינה המלאכותית האירופיות, לעיל ה"ש 53; ראו לעיל סעיף 4.1.1.1.

בקטגוריה זו נכללות מערכות שטמונה בהן על פי התפיסה האירופית סכנה קיצונית ושאינן עולות בקנה אחד עם ערכי היסוד שלה.¹⁰⁴⁵ למשל, מערכות שמשמשות בטכניקות תת־הכרתיות או מניפולטיביות כדי להשפיע על האנושות לרעה;¹⁰⁴⁶ מערכות שמנצלות חולשות של קבוצות שונות באוכלוסייה על בסיס גיל, יכולות פיזיות או יכולות קוגניטיביות באופן שסכיר שיסב נזק פיזי או פסיכולוגי משמעותי; מערכות נבונות בשירות רשויות ציבוריות או מטעמן שנועדו לסווג או להעריך את מידת האמון שיש לתת באדם על בסיס התנהלותו החברתית או תכונותיו בדרך של דירוג חברתי – באופן שעלול להביא להחלטות לרעת יחידים או קבוצות בהקשרים שאין להם זיקה להקשרים שבהם נאסף המידע עליהם במקור או באופן שעלול להסב להם נזק לא מידתי.

כמו כן, נאסר השימוש במערכות שיטור מונחה עתיד (predictive policing), המעריכות את הסיכון שאדם יבצע עברה פלילית או מינהלית (לרבות הערכת סיכויי רצידיביזם);¹⁰⁴⁷ במערכות המייצרות או מרחיבות מאגרי נתונים לזיהוי פנים בהסתמך על איסוף גורף מרשת האינטרנט או על חומר מצולם ממערכות טלוויזיה במעגל סגור;¹⁰⁴⁸ במערכות בינה מלאכותית שתכליתן לזהות רגשות לתכליות של אכיפת חוק, ניהול מעברי גבול, במוסדות חינוכיים ובהקשרי תעסוקה;¹⁰⁴⁹ ובמערכות זיהוי ביומטריות בזמן אמת במרחבים ציבוריים.¹⁰⁵⁰

איסור השימוש בטכנולוגיות מסוימות, שהשלכותיו טרם נלמדו, יכול להיות זמני – אם כדי לבחון את האסדרה שלהן באופן שקול (כפי שנעשה בעניין שיבוט

1045 שם, ס' 5(1).

1046 ס' 5(1)(a) להצעה המתוקנת, לעיל ה"ש 57.

1047 ס' 5(1)(da), שם. ראו גם סעיף 2.3.2 לעיל.

1048 ס' 5(1)(db), שם. נראה שסעיף זה נועד לאסור מאגרי זיהוי פנים כגון אלו שמפתח חברת Clearview AI. ראו גם היל, להלן ה"ש 1087.

1049 ס' 5(1)(dc), שם. ראו גם סעיף 2.3.2 לעיל.

1050 ראו והשוו בין ס' 5(1)(d), שם, ובין הנוסח המקורי של ס' 5(1)(d) להצעת תקנות הבינה המלאכותית האירופיות, לעיל ה"ש 53.

בני אדם),¹⁰⁵¹ ואם כדי למנוע נזקים מידיים משימוש לא מוסדר בטכנולוגיות אלו, עד להסדר קונצנזואלי (כפי שהציע הדָּוֶח המיוחד של האומות המאוחדות לחופש הביטוי בעניין טכנולוגיות סייבר התקפי).¹⁰⁵²

גישה זו אינה ייחודית לאירופה. בשנת 2021 הוצגה בקונגרס של ארצות הברית הצעת חוק פדרלית שביקשה לאסור על גורמים פדרליים לרכוש מערכות ביומטריות, להעריך אותן או להשתמש בהן (ובכלל זה במערכות זיהוי פנים) – למעט כאלו שהותרו בחוק מפורש. על פי ההצעה, למדינות בארצות הברית או לגופי ממשל מקומיים שישתמשו במערכות אלו לא תהיה זכאות לסיוע תקציבי מתוכניות פדרליות מסוימות.¹⁰⁵³ גם תקנות מערכות ההמלצה האלגוריתמיות הסיניות כוללות איסורים מפורשים על שימוש במערכות בינה מסוימות – כאלו שעלולות לעודד התמכרות לרשת, בייחוד של קטינים.¹⁰⁵⁴ עמדה דומה הביעה הנציבה העליונה של האו"ם לזכויות אדם, מישל בצ'לט (Bachelet), כשקראה להשהות את השימוש בטכנולוגיות בינה מלאכותית המסכנות זכויות אדם, לרבות מערכות זיהוי פנים ומערכות דירוג חברתי.¹⁰⁵⁵

ההחלטה אם לפתח או שלא לפתח סוגים מסוימים של טכנולוגיות יכולה להתקבל על יסוד איסור רגולטורי, אך היא יכולה גם להיות החלטה רצונית של התעשייה.

1051 ראו לדוגמה גם ס' 1 לחוק איסור התערבות גנטית (שיבוט אדם ושינוי גנטי בחאי רביה), תשנ"ט-1999, ס"ח 697 (1999) 47. עם זאת, החוק נחקק בשנת 1999 לתקופה של חמש שנים לשם בחינת המדיניות בנושא – ותוקפו הוארך כארבע פעמים.

1052 United Nations Human Rights Council, 41st Session of the Human Rights Council (24 June–12 July 2019), Agenda item 3, Promotion and Protection of All Human Rights, Civil, Political, Economic, Social and Cultural Rights, Including the Right to Development; United Nations General Assembly, Surveillance and Human Rights, A/HRC/41/35 Cahane, *Targeting Targeted Surveillance – The UN Special Rapporteur's Report on the Private Surveillance Industry*, CSRCL Blog (31.12.2021)

1053 Facial Recognition and Biometric Technology Moratorium Act of 2021, S. 2052, 117th Congress (2021)

1054 ראו לעיל סעיף 4.2.

1055 Jamey Keaten And Matt O'Brien, *U.N. Urges Moratorium on Use of AI that Imperils Human Rights*, AP News (15.9.2021)

חברות הטכנולוגיה הגדולות, ובראשן אמזון, איי-בי-אם ומייקרוסופט, הודיעו למשל שיפסיקו להשקיע בחברות המפתחות מערכות זיהוי פנים לתכליות שיטור או ליצור מערכות כאלה, או שיימנעו מלספקן לסוכנויות אכיפת החוק.¹⁰⁵⁶ אבל כשהאיסור אינו מחייב אלא וולונטרי תלויה יציבותו ברצון הטוב של השחקנים בתעשייה. עמדת חברת אפל באשר לפגיעה בפרטיות לקוחותיה, למשל, התרככה עם השנים: בשנת 2016 התנגדה החברה בעיקשות לדרישת הבולשת הפדרלית (FBI) לעדכן את מערכת ההפעלה שלה כדי לאפשר גישה למכשיר אייפון מוצפן;¹⁰⁵⁷ ואילו בשנת 2021 החלה החברה בפיתוח מערכות לסריקת מאגרי תמונות מקומיים כדי לזהות תכנים פדופיליים.¹⁰⁵⁸ מערך התמריצים של חברות גדולות ומבוססות כגון אפל, מייקרוסופט ואיי-בי-אם יכול להיות שונה מזה של חברות צעירות וקטנות, שמטה לחמן מבוסס על טכנולוגיות שנויות במחלוקת ולכן הן עלולות להעדיף "לשבור שביתה". חרף המורטוריום הוולונטרי של חברות אלו על טכנולוגיות זיהוי פנים יש אפוא חברות אחדות (כגון Clearview AI)¹⁰⁵⁹ שממשיכות בפיתוח ובשיווק שלהן.¹⁰⁶⁰

ואולם למדיניות שאוסרת על פיתוח טכנולוגיות מסוימות, רכישה שלהן או שימוש בהן יש גם חסרונות. עקב ההערפות התרבותיות-פוליטיות השונות של כל מדינה, כל אחת מהן מתמקדת באיסור אחר בכואה לעצב את מדיניותה: הסינים חוששים כאמור מתופעות של התמכרויות לרשת, ואילו האירופים חוששים ממערכות

1056 ראו לעיל ה"ש 161.

1057 עומר כביר "אפל לביהמ"ש: היענות לדרישת ה-FBI העניק לממשל 'כוח מסוכן' כלכליט (25.2.16).

1058 עומר כביר "אפל רוצה לזהות תכנים פדופיליים ופוחחח קופה שרצים" כלכליט (8.8.2021).

1059 קשמיר היל "תכירו את Clearview AI, החברה שהרגה לכם את הפרטיות" הארץ (21.1.2020).

1060 לדוגמה, במקביל להכרזתה של אפל, מייקרוסופט ואמזון (לעיל ה"ש 161), נמצא שכ-24 ארגוני אכיפת חוק ברחבי העולם השתמשו במוצרים של חברת זיהוי הפנים השנויה במחלוקת Clearview AI. ראו Ryan Mac, Caroline Haskins, and Antonio Pequeño *Clearview AI. Find Out IV, Police in At Least 24 Countries Have Used Which Ones Here*, BUZZFEED (25.8.2021)

זיהוי פנים בשירות אכיפת החוק.¹⁰⁶¹ לכן איסור מקומי על פיתוח מערכת מסוימת או שימוש בה לתכלית מסוימת – גם אם אין הוא וולונטרי – אינו שולל את האפשרות שיפתחו אותה וישתמשו בה במדינה אחרת, לתכלית אחרת. יתר על כן, כשהמורטוריום הזמני נועד מלכתחילה לאפשר למעצבי המדיניות לקבל החלטה מושכלת בעניין האסדרה של טכנולוגיות מסוימות, ההנחה הסמויה היא שהעבודה המלוכלכת תיעשה במקומות אחרים. כך אפשר להפיק לקחים מהניסיון המצטבר של השימוש במערכות אלו על ידי אחרים (ועל גבם).

לכן אנו מציעים שהמדיניות בישראל תעוצב מתוך מבט החוצה, בעיקר לעבר אירופה, ולא ייקבעו בישראל איסורים שאינם קיימים במדינות אחרות. מנגד, כדי שלא להפוך את ישראל ל"חצר אחורית", כלומר כדי שאותה "עבודה מלוכלכת" לא תיעשה כאן, יש לשאוף להגביל את האיסורים בזמן. כך לרגולטור או למקבלי ההחלטות יהיה די זמן לעצב מדיניות מיטבית, ואם יתעורר הצורך יהיה אפשר להיעזר גם בארגוני חול רגולטוריים. במסגרת אותם ארגוני חול תתבצע בחינה מבוקרת של מערכות בינה מלאכותית – גם כאלה שתכליתן שנויה במחלוקת – ובתוך כך ייכתבו הנחיות לשוק הפרטי.

9.2.10. אין להשלים עם בשנים האחרונות מקורמות הצעות לאסדרת בינה מלאכותית במדינות מובילות, וגם באיחוד האירופי, ויש להניח שמדובר בתחילתה של מגמה עולמית. גם אם יש שוני ערכי בין האסדרות במקומות שונים – אם בבחירת סוג המערכות המוגדרות מסוכנות ואם במטרות המוצהרות של החקיקה עצמה – יש להן גם הרבה מן המשותף. לכן אין

1061 ייתכן שיש חפיפה בין האיסור על מערכות המנצלות חולשות של קבוצות שונות באוכלוסייה באופן שצפוי להסב נזק פיזי או פסיכולוגי (ס' (a) 5(1) להצעת תקנות הבינה המלאכותית האירופיות [לעיל ה"ש 53]) ובין האיסור הסיני על פיתוח מערכות המעודדות התמכרות. לביקורת על הנוסח האירופי, המעדיף עקרונות עומדים מפתרונות ייעודיים לטכנולוגיה ספציפית, ראו Paul De Hert and Georgios Bouchagiar, *Facial Recognition, Visual and Biometric Data in the US. Recent, Promising Developments to Regulate Intrusive Technologies*, 7 (29) BRUSSELS PRIVACY HUB WORKING PAPER 4 (2021)

להשלים עם מצב שבו במדינת ישראל לא תהיה חקיקה שתהיה בהרמוניה עם החקיקה המקובלת בעולם.

הרמוניה איננה בהכרח זהות. כאשר הפיתוח מכוון לייצוא ממילא נדרש הסטנדרט הרגולטורי לעמוד לכל הפחות בסטנדרטים הזרים, שהם לרוב גבוהים יותר מהסטנדרטים בישראל. אבל גם כאשר הפיתוח מכוון לשוק המקומי, ומושאי ההחלטות של מערכות הבינה המלאכותית אינם אזרחי חוץ אלא אזרחי המדינה, אין להנמיך את הסטנדרטים, שהרי אזרחי המדינה ראויים אף הם להגנה. אומנם סטנדרט רגולטורי נמוך מהמקובל בעולם יכול לעודד חדשנות, אך זו לא בהכרח חדשנות רצויה: ישראל עלולה להפוך לחצר אחורית טכנולוגית, כלומר למקום שבו מפתחים מערכות שאסור לפתח אותן, להפיץ אותן או להשתמש בהן לפי הדין הזר.¹⁰⁶² ישי עופרן ואלעד סימן טוב הראו במטא-מחקר שנערך בשנת 2022 שניסויים שהשתמשו בנתונים פרטיים של אזרחים "זלגו" עקב הטמעת התקנות האירופיות בדבר הגנת המידע למדינות מחוץ לאירופה שלא היתה בהן חקיקת פרטיות מעודכנת.¹⁰⁶³ אין סיבה להניח שבאסדרה של בינה מלאכותית מצב זה יהיה שונה. אפשרות אחרת, הגובלת באבסורד אך נוכחנו בה עם הצעת "תקנות הגישור" בשנת 2022, היא מצב שבו הדין הישראלי יהיה כלפי חוץ בהרמוניה עם הדין הזר, אך בתוך גבולות המדינה או על אזרחיה שלה יחיל סטנדרטים נמוכים יותר.¹⁰⁶⁴

1062 כך למשל, חברות ישראליות או חברות בבעלות ישראלית מובילות בפיתוח מערכות סייבר התקפיות, ששנויות במחלוקת (ב־2019 קרא הַנְּחָן המיוחד של האומות המאוחדות לחופש הביטוי להשעות זמנית את הסחר הבינלאומי במערכות אלו ואת השימוש בהן. ראו לעיל בה"ש 882). דוגמה אחרת היא תעשיית האופציות הבינאריות, שפעלה בישראל והציעה אופציות בינאריות למשקיעים מעבר לים עד חקיקת האיסור הפלילי על פעילות זו. ראו ס' 44:1 לחוק ניירות ערך, תשכ"ח-1968, ס"ח 541, בעמ' 234; US v. Lee Elbaz, No. 20-4019 (4th Cir. 2022).

1063 Elad Yom-Tov and Yishai Ofra, *Implementation of Data Protection Laws in the European Union and in California Is Associated with a Move of Clinical Trials to Countries with Fewer Data Protections*, 9 FRONT. MED. (LAUSANNE) (10.11.2022).

1064 כך למשל במהלך 2022 שקלה הרשות להגנת הפרטיות לקדם חקיקה שתטיל על בעלי מאגרי מידע ישראליים חובה להגן על פרטיות של מידע שמקורו באיחוד האירופי כדי להלוט את הוראות ה-GDPR בהיעדר רפורמה כוללת בדיני הגנת הפרטיות הישראלים. ראו עומר כביר, "איפה ואיפה: ישראל תספק הגנה משופרת רק לפרטיותם של תושבי האיחוד

נוסף על כך, המסגרת המשפטית הקיימת בישראל להגנה על הפרטיות ולאבטחת סייבר נוטה להחרגת מתחולתה את רשויות הביטחון¹⁰⁶⁵ או לייצר הסדרים מקילים יותר בהקשרים ביטחוניים. ראינו כיצד במהלך התפרצות מגפת הקורונה השתמשו בכלי מעקב המיועדים לסיכול טרור לתכליות אזרחיות מובהקות,¹⁰⁶⁶ וכמה רב היה הפיתוי לפרוש מערכות בינה מלאכותית לצורך דירוג אזרחי ישראל על פי הסבירות שיידיבקו בקורונה.¹⁰⁶⁷ הנטייה הישראלית למסגר אירועים שאין דבר בינם ובין ביטחון לאומי בתור כאלה עלולה להביא להחרגת איסורים מוחלטים על מערכות בינה מלאכותית שנועדו לתכליות מסוימות בתאונות ביטחוניות, כך שבעתיד תתאפשר זליגה של מערכות אלו מהספרה הביטחונית לזו האזרחית.

לכן על פי תפיסתנו יש לנסח את האיסורים על פיתוח מערכות נבונות ושימוש בהן כך שיהיו בהלימה עם המקובל בעולם, ובפרט עם הדין האירופי המתהווה – לבל תשמש מדינת ישראל לניסויי כלים טכנולוגיים. אנו ממליצים לאמץ לכל הפחות את האיסורים בתקנות האירופיות, ובהם האיסור על שימוש בזמן אמת במערכות זיהוי ביומטריות לתכליות של אכיפת חוק והאיסור על שימוש של רשויות ציבוריות במערכות נבונות שנועדו לסווג או להעריך את מידת האמון שיש לתת באדם על בסיס דירוג חברתי.¹⁰⁶⁸ אפשר לסייג איסורים אלו לצורכי ביטחון לאומי, אך רק בתנאי שרשויות הביטחון יפעילו מערכות אלו באופן מידתי, על בסיס צורך בלבד.¹⁰⁶⁹

האירופי" כלכליט (11.7.2022); טל קפלן, "מומחים מחריעים: הזנחת הפרטיות מובילה להצעה בלתי סבירה ומקוממת" (10.7.2022) Law.co.il.

1065 ס' 12(ג), ס' 13 (ה), ס' 19 לחוק הגנת הפרטיות. ראו גם כהנא ושני, לעיל ה"ש 213, בעמ' 33.

1066 בג"ץ 6732/20 האגודה לזכויות האזרח בישראל נ' הכנסת (1.3.2021); תהילה שורץ אלטשולר, רחל ארידור הרשקוביץ, שיקולי פרטיות ומעקב אחר אזרחים: נגיף הקורונה – התמודדות עם המשבר העולמי (המכון הישראלי לדמוקרטיה, 2020); עמיר כהנא "עידן הזוחלים: ההסמכה ההדרגתית של השב"כ לאיתור מגעים של חולי קורונה" משפטים על אחר כ (2023); Amir Cahane, *The (Missed) Israeli Snowden Moment?* 34 Int'l J. INTELLIGENCE & COUNTERINTELLIGENCE 694 (2021).

1067 ראו לעיל ה"ש 290.

1068 ס' 3(1) להצעת תקנות הבינה המלאכותית האירופיות, לעיל ה"ש 53. להרחבה ראו לעיל, בסעיף 4.1.1.1.

1069 השוו לדוגמה לס' 3(1)(d)(ii), שם.

פרק עשירי

יסודות לרגולציה מכוונת זכויות
של בינה מלאכותית בישראל:
ארגז כלים

—

בפרק זה נציע ארגז כלים רגולטורי ראשוני לאסדרת בינה מלאכותית שנכון לתפיסתנו לחשוב על אימוצו כבר כעת. אומנם טוב היה אילו יכולנו להמתין וללמוד כיצד מתקבלים הדינים המשמשים בעולם הלכה למעשה, ובפרט אילו יכולנו להמתין להתייצבות הדין באירופה. אילו יכולנו להמתין, היינו יכולים לנסח את החוק בישראל כך שתהיה התאמה בינו ובין הדין באירופה ולשמור על הפוטנציאל הכלכלי של השוק באירופה (בדומה לשיקולים שביסוד ההתאמה לדיני הגנת הפרטיות באירופה).

ואולם אין להיתפס לשאננות. לנוכח ההשתהות של המחוקק בישראל בפיתוח דיני הגנת הפרטיות והתאמתם לסטנדרט האירופי של הוראות התקנות הכלליות

בדבר הגנת מידע (GDPR), והיעדר אסדרה של נושאים אחרים מעולמות המשפט והטכנולוגיה, דוגמת רשתות חברתיות, הגנת סייבר ופלטפורמות ליישומי כלכלה שיתופית, יש חשש שמקבלי המדיניות בישראל לא ינצלו כראות את הזמן העומד לרשותם. אי-אפשר להקים תשתית רגולטורית לכינה מלאכותית בן לילה, והיא כרוכה בשיג ושיח עם גורמים בתעשייה ועם מומחים מהאקדמיה ומהחברה האזרחית. יתר על כן, לאחר הקמתה נדרשת תקופת היערכות של השוק לדרישות החדשות.

הואיל ועקרונות שונים של אסדרה, כגון הוגנות, פרטיות, שקיפות, אחריותיות וניהול סיכונים, מתבטאים בכל אחד מרכיבי מעגל החיים, כדי ליצור אסדרה אפקטיבית של מערכות לומדות יש להביא בחשבון רכיבים אלו. התרשים שלהלן מציג בקווים כלליים את "מעגל החיים" של מערכות לומדות.

10.1 אסדרת מערכות לומדות מחייבת הבנה של "מעגל החיים" שלהן

4 חרשים

מעגל החיים של מוצרים מבוססי בינה מלאכותית



היעדר התייחסות לרכיבי מעגל החיים עלול ליצור מצב של אסדרת יתר של רכיבים מסוימים והתעלמות מרכיבים אחרים ומכאן לצמצום האפקטיביות של האסדרה. למשל, ההתייחסות הרחבה בספרות המשפטית למאגרי נתונים ולהשלכות אפשריות של הפגמים בהם על המוצר הסופי באה במידה רבה על חשבון התייחסות להשלכות הבחירה במודלים שונים או לקשיים שעלולים להתעורר לאחר יישום המודל.

לעומת זאת, תפיסה אינטגרטיבית מביאה בחשבון את מכלול רכיבי מעגל החיים ונדרשת גם לקשרי הגומלין ביניהם. למשל: תכלית המערכת ומסגור הבעיה צריכים להשפיע על בחירת המודל (האם לבחור במודל שיש בו עכירות רבה יותר באשר לדרכים שבהן הוא מקבל החלטות או שלא לאפשר זאת?); תוצאות הערכה בשלב בניית המודל מזינות תהליכים להערכת סיכונים, והם בתורם מחייבים קבלת החלטות הנוגעות לאימון המודל, לפרישתו או לאבטחתו; תכלית המודל משפיעה על הבחירה בממשקי המשתמש (האם מערכת המעניקה ייעוץ רפואי אמורה להציג את האפשרות שהיא טועה? האם נכון לספר למשתמש שהוא מתקשר עם מערכת מלאכותית ולא אנושית?)

לפירוקו של מעגל החיים לחלקים יתרון נוסף: אם כבר בשלב ראשוני, למשל בשלב איסוף המידע, מוערך שיש פגמים, אפשר לטפל בהם כבר אז; ואם מתעוררות הערכות מדאיגות בעניין אימון של מודל אפשר להשהות אותו או להימנע מיישומו המלא.

לדוגמה, אם במאגר הנתונים מתגלות הטיות, אפשר להידרש לשאלות כמו הטיית ייצוג (אם למשל מדובר על מאגר העוסק במערכת החינוך אפשר לשאול אם המידע מגיע רק ממוסד חינוכי אחד; באיזה אופן אפשר להכליל ממאגר הנתונים המסוים הזה על כלל האוכלוסייה; אם יש ייצוג יתר לקבוצות אוכלוסייה קטנות) או הטיית מדידה (האם המדידה תופסת את מה שנכון לחפש? ואם מדובר בחיזוי מחלה, האם המודל מודד רק את מי שאובחנו בעבר במחלה זו?) ואולם בתהליך בניית המודל אפשר לדבר על הטיות שהן תוצאה של בחירת המודל (למשל, האם בשל תהליכים של דחיסת נתונים, כדי שהמודל יעבוד מהר יותר, נפגעו קבוצות קטנות כבר בשלב אימון המודל?), ואילו לאחר יישום המודל ייתכנו הטיות הנובעות מכך שהמודל אומן בהקשר מסוים (למשל, סיכון לחזרתיות בפשיעה) לצורך הקשר אחר (למשל, קביעת גזר דין).

דוגמה אחרת נוגעת לשקיפות ולהסברתיות. אין דומה הסברתיות שתכליתה להסביר כיצד נאסף מידע (מי אסף את הנתונים, מהי ההתפלגות שלהם וכיו"ב) לכזאת שעניינה להסביר החלטה (למשל, אלה המילים שבשלן נמחק ציון פוגעני ברשת חברתית). ועוד דוגמה נוגעת להערכת סיכונים. הערכת הסיכונים בבחירה ראשונית של מודל צריכה להסתמך על סבבי אימון קודמים או על בנייה של

מודל ניסיוני, אבל הערכת סיכונים לאחר סבב אימון ראשוני שמתגלה בו בעיית תאימות בין משימה לתוצאה צריכה להוביל להתאמה מחודשת של הנתונים, של הארכיטקטורה ושל משימות האימון, לעורר את השאלה אם נכון לאמן מודל קטן יותר או חלש יותר, ובוודאי לא להמשיך לאמן מודל בקנה המידה שתוכנן מלכתחילה.

חלק חשוב של הבנת מעגל החיים של מערכות לומדות נוגע לצורך לנטר אותן לאחר שיושמו בפועל (post deployment) בעולם האמיתי (למשל הבניה של המערכת בתוך מוצר או בתוך ממשק). זאת משום שמערכת לומדת, שלא כמו של מוצרים אחרים (תרופות למשל), יכולה מעצם טיבה להשתנות גם לאחר יישומה בשל המשוב שהיא מקבלת מן המשתמשים.¹⁰⁷⁰ לכן חשוב לשאול לקראת היישום מה עלול להשתבש, אם המודל בטוח ליישום ואילו מנגנוני בטיחות נחוצים כדי ליישמו בבטחה, ואם נדרש יישום בקנה מידה קטן תחילה – מי יקבל התראה במקרה של תקלה. ואילו לאחר היישום יש לשאול למשל באיזו מערכת בקרה לבחור; אם יש ניטור שתכליתו לאתר התנהגויות לא צפויות, בעיקר בסביבות יישום מורכבות (למשל, האם משתמשים מצאו יישומים חדשים או אסטרטגיות הנדסיות חדשות על בסיס המודל?), ולדווח על אירועים כאלה; אם נדרשים עדכונים למודל לאחר היישום שלו ואם הם משמעותיים ברמה שמחייבת הערכה מחודשת של השלבים שלהם (למשל, תיוג מחדש של הנתונים, אימון מחדש של המודל ובדיקה מחודשת שלו).

תסקירי השפעה ומתודולוגיות לניהול סיכונים נעשו מקובלים בתחומים כגון הגנת הסביבה, דיני חברות, פיננסים וזכויות אדם. אבל הצורך ליישם מתודולוגיות לניהול סיכונים על מערכות אלגוריתמיות עדיין נמצא בחיתוליו, למרות חיוניותו. מתודולוגיות אלו צריכות

10.2 פיתוח מתודולוגיות לניהול סיכונים

1070 מינהל המזון והתרופות האמריקני (FDA) מתחיל בניסוח המלצות כאלו למערכות לומדות בתחום הרפואה. ראו רותי לוי "ה-FDA יקל על שיווק מערכות בינה מלאכותית – המשתנות לאחר שהן מאושרות לשימוש" *TheMarker* (11.5.2023).

לכלול הערכה של השפעות וסיכונים לא רק בנוגע לפרטים אלא גם בהקשרים חברתיים רחבים, כמו למשל המתודולוגיות שמציע הארגון Data & Society (D&S) במעבדת ההשפעות האלגוריתמיות שהקים (Algorithmic Impact Methods Lab (AIMLab). שהרי מתעורר חשש שאם מי שיקבע מתודולוגיות כאלה יהיה התעשייה בכלל וענקיות הטכנולוגיה בפרט, לא יובאו בחשבון שיקולים רחבים יותר.

היוזמה המרכזית לפיקוח על בינה מלאכותית מתוך פרספקטיבה של ניהול סיכונים היא הצעת החקיקה של האיחוד האירופי, המתמקדת בניהול סיכונים במובן הצר שלו, כלומר בהקשרי שימוש מסוימים בלבד (למשל מעקב ביומטרי). אלא שמודלים יכולים להיות סיכון גם בהקשרי שימוש שלכאורה יש בהם סיכון נמוך.

אפשר להציע אפוא מודל שונה של ניהול סיכונים. לפי מודל זה, כדי להעריך מסוכנות של מערכת נדרשת התבוננות כפולה – תחילה יש להעריך את פוטנציאל המסוכנות של המערכת כפי שתוכננה; אחר כך יש להעריך את רמת הקשר בין המשימה לתוצאה (alignment), כלומר את האפשרות של המערכת לממש פוטנציאל מסוכנות מחוץ לתפקיד שייעדו לה מתכנניה.

הערכה כפולה זו תהיה הבסיס לקבלת ההחלטות וכדי שיהיה אפשר להוציאה לפועל יהיה צורך לנסח כללי משילות ובטיחות, כגון כללים לאימון אחראי (האם לאמן מודל חדש שניכרים בו סימנים מוקדמים של סיכון – וכיצד) וכללים ליישום אחראי (האם, מתי וכיצד ליישם מודלים שעלולים להיות מסוכנים); ולקבוע אילו רמות של שקיפות ותיעוד נדרשות במקרה של מודלים שעלול להיות בהם סיכון קיצוני ואילו בקרות ומערכות אבטחת מידע חזקות יש ליישם בעניינם.

ייתכן שלא בכל מקרה יהיה צורך להפעיל ניהול הסיכונים, אלא רק במקרים שייקבעו – למשל על מודלים שקרובים ליכולות המוצעות של מודלים קיימים או עולים עליהם והקהילה המדעית לא הרבתה לחקור אותם; או מודלים ששונים ממודלים קיימים מבחינת גודל, עיצוב (למשל ארכיטקטורות שונות או טכניקות alignment שונות) או תמהיל היכולות הנובע מהם. למודלים כאלה יש יכולת רבה יותר להצליח במגוון רחב יותר של משימות ולכן גם הזדמנויות רבות יותר לגרום נזק. בהקשר הכללי-חברתי אפשר לקבוע שמדובר במקרים שבהם היקף

ההשפעה הוא גדול (כגון אובדן עשרות אלפי בני אדם, נזק כלכלי או סביבתי בסך מאות מיליארדי דולרים) או שרמת ההפרעה השלילית לסדר החברתי והפוליטי עלולה להוביל למשל למלחמה בין מדינות, לשחיקה משמעותית באיכות השיח הציבורי או להחלשה נרחבת של ציבורים, ממשלות וארגונים אחרים.¹⁰⁷¹ ראוי לחשוב גם על תרחישים שבהם הסיכון נובע משילוב יכולות של משתמש יחיד, קבוצת משתמשים או מערכת בינה מלאכותית נוספת.

יצירת מערך ניהול סיכונים יכולה כמובן לעורר קשיים. חשש אחד הוא ששיתוף תוצאות הערכה בין מפתחים יביא לידי לפרסום סודות מסחריים ויתמרץ מתחרים להאיץ יישום של מודלים כדי לשמור על יתרון תחרותי. חשש אחר הוא שאם מודל כלשהו ידע שהוא עומד לפני בדיקה הוא ילמד להציג רק התנהגות רצויה. נוסף על כך, הסתמכות יתר על מערך ניהול סיכונים עלולה לגרום לאמון יתר בהערכות ובתוצאותיהן ולהוביל לפריסת מודלים מסוכנים מתוך תחושת ביטחון מזויפת. מערך ניהול סיכונים אינו יכול אפוא לעמוד בפני עצמו וראוי לשלב אותו עם מחויבות ארגונית ועם כלים נוספים לזיהוי ולהערכת סיכונים. הללו יידונו כעת.

10.3

אסדרת הקשר (Alignment) בין משימה לתוצאה

עד כה עסק הדיון ברגולציה במאגרי המידע שעליהם מתאמנים האלגוריתמים ובתוצרים שלהם. אבל כבר היום אנו נמצאים בשלב מעבר מלמידה לא מפוקחת, שמסתמכת על מאגר

מידע מוגדר, ללמידה לא מפוקחת, שצומחת מן האינטראקציה עם משתמשים ועם המשימות שהיא מתבקשת לבצע, מה שאפשר לכנות Experience Based Learning. מערכות לומדות שאינן מבוססות על מאגרים מובנים, או שזודקו למאגרים קטנים בהרכבה כדי ללמוד, מציבות לפנינו אתגרים חדשים, ובראשם הקירוב בין המטלה שקיבלה המערכת לתוצר שהפיקה במקרים שבהם אין התאמה בין השניים.

האם מדובר באי־התאמה אסטרטגית, כלומר ביצירה מכוונת של תוצרי לוואי שליליים ובלבד שהמערכת תבצע את המטלה שלה, או שמדובר באי־התאמה אגנוסטית, כלומר ביצירה אגבית ולא מכוונת של תוצרי לוואי שליליים? שאלה זו הופכת לשאלה משמעותית, וכפי שכבר נכתב לעיל היא מעוררת את הצורך למצוא כלים רגולטוריים ההולמים אותה.

הצורך להכיר במודל כמודל מסוכן לא בשל הנזק הישיר שהוא עלול לגרום אלא בשל החשש שיכולותיו ינוצלו לרעה איננו חדש. ואולם הפרישה הרחבה של ממשקי בינה מלאכותית יוצרת, והיכולת לבצע מניפולציה על מודלים כדי לגרום להם לבצע משימות שהם לא אומנו לבצע בדרכים שעלולות להימצא מזיקות, מחייבות תשומת לב. כדי מנוע ניצול של מודלים כאלה לרעה יהיה צורך בבקורות משמעותיות טרם יישומם בפועל (deployment). למשל, תירדש הוכחה שהמודל יתנהג כמתוכנן (alignment assessment)¹⁰⁷², ויהיה צורך לברר, למשל, אם הוא מנסה להשיג מטרות ארוכות טווח, בעולם האמיתי, שהן שונות מאלה שסיפקו המפתח או המשתמש;¹⁰⁷³ אם הוא יכול לשתף פעולה או ליצור "קנוניה" עם מערכות AI אחרות נגד אינטרסים אנושיים;¹⁰⁷⁴ או אם הוא מתנגד לניסיונות של משתמשים זדוניים לגשת ליכולות המסוכנות שלו.¹⁰⁷⁵

ראוי לציין שיש חוקרים הסבורים שבחלק מן המודלים אי־אפשר לפתור את בעיית הקשר בין משימה לתוצאה. לדעתם, שיטות מקובלות כגון למידת על ידי חיזוקים באמצעות משוב אנושי רק יחמירו את הבעיה.¹⁰⁷⁶ קביעה כזאת עלולה, למשל, להביא לייחוס אחריות משפטית למי שיאפשר שימוש המוני בממשקים ובמודלים. כך או כך, יש לנסות להעריך את התנהגות המערכות באמצעות מפתח

Toby Shvlane, *Sharing Powerful AI Models*, CENTER FOR THE GOVERNANCE OF AI (20.1.2022) 1072

Alan Chan et al., *Harms from Increasingly Agentic Algorithmic Systems* (20.2.2023), available at arXiv 1073

ראו Ngo, Chan, and Mindermann "ש" 77. 1074

Amelia Glaese et al., *Improving Alignment of Dialogue Agents via Targeted Human Judgements* (18.2.2022), available at arXiv 1075

Yotam Wolf et al., *Fundamental Limitations of Alignment in Large Language Models* (29.5.2023), available at arXiv 1076

רחב ככל האפשר, ייתכן אפילו שבאמצעות אוטומציה של הליך ההערכה עצמו. יש להכיר בכך שבתרחישים מסוימים ובשילובים מסוימים של מערכות צפויות בעיית קשר חריפות יותר מאשר באחרים (למשל צ'אטבוט בעל יכולת מניפולציה על התנהגות משתמשים או מערכת שלא תתנגד לאפשרות שהיא תשולב במערכת אוטונומית שיש לה מטרות מזיקות).

סיכוניה של מערכת אינם תלויים רק בבעיית הקשר, אלא גם בדרך שבה היא מתקשרת, ותתקשר בעתיד, עם עולם מורכב. יתר על כן, קשה לצפות את כל המסלולים שעלולים להוביל לתוצאות קיצוניות. ואולם אין פירוש הדבר שמשום כך לא צריך לטפל בבעיה זו.

10.4

תיעוד נתונים, משילות נתונים והליכי בקרה (Auditing) בדיעבד

המורכבות של מערכות בינה מלאכותית בכלל, ושל מערכות לומדות בפרט, מערימה קשיים מיוחדים על קובעי המדיניות והרגולטורים בבואם לנסח כללי אחריות ולזהות בפועל את השרשרת הסיבתית שהובילה לפגיעה בזכויות, בייחוד בהתחשב בשונות בין מגזר למגזר ובין יישום ליישום.

המשותף לכל אלו הוא שבלי תיעוד ראוי הם בלתי אפשריים. הבסיס לכל בחינה עוברתית של כשל קונקרטי במערכות בינה מלאכותית הוא משילות נתונים ותיעוד קפדני של הליכי עבודה, מקורות מידע, תיוגים, מודלים, תהליכי עיצוב קוד, הערכת סיכונים ובסיסי נתונים, וזיהוי הפערים בכל אלו.¹⁰⁷⁷ עיצוב טוב של התיעוד הוא גם אינטרס של יזמים ומפתחים, שכן הוא מאפשר להם הן לתחקר בדיעבד כשלים ותופעות שהם לא צפו הן לעמוד בחובות תיעוד רגולטוריות ממקורות אחרים.¹⁰⁷⁸

1077 ראו למשל ס' 10(2) להצעת תקנות הבינה המלאכותית האירופיות, לעיל ה"ש 53.

1078 בתקנות הכלליות בדבר הגנת מידע (GDPR), למשל, יש חובת תיעוד מפורשת (ס' 30), אך גם משתמעת. החובה המשתמעת מקיפה יותר: בעל מאגר מידע (controller) נדרש להיות מסוגל להראות שכל מושאי המידע שבמאגר הביעו את הסכמתם לעיבוד המידע.

נוסח ההצעה של התקנות האירופיות מחיל חובת תיעוד מקיפה על מערכות בסיכון גבוה.¹⁰⁷⁹ יתרונה של הגישה האירופית הוא שהיא מאפשרת בקרה מלאה אקס פוסט הן במישור ההגנה על זכויות אדם הן במישור המשפט הפרטי. הסתמכות על תיעוד חלקי של מסמכי הערכת סיכונים, לעומת זאת, עלולה להיות משענת קנה רצון.

אנו סבורים כי תיעוד הוא כלי להתמודדות עם אי-ודאות. הערכת סיכונים נעשית מטבע הדברים מראש, ולנוכח אי-הוודאות הגוברת המאפיינת מערכות בינה מלאכותית¹⁰⁸⁰ ייתכן שיימצאו בדיעבד חוסרים בדיווחים המסתמכים על הערכות אלו. לפיכך אנו מציעים לקבוע משטר תיעוד איכותי, שסנקציות בצידו, שיאפשר לתחקר בדיעבד נזקים שלא נחזו.¹⁰⁸¹ יש לקבוע תקנים אחידים של תיעוד כדי לאפשר פיקוח אחיד, יעיל, שיטתי וחוצה מגזרים על מערכות אלגוריתמיות,¹⁰⁸² הן בהקשר של בסיסי הנתונים הן בהקשר של התוצרים.¹⁰⁸³ הליכי בקרה נאותים המבוססים על תיעוד אינם חזות הכול ולא יחליפו גישות כגון עיצוב לפרטיות או הנדסת הוגנות לפני פיתוח מודל. אבל הם יוכלו לזהות, ולו בדיעבד, שגיאות מתודולוגיות וטעויות שהביאו לפיתוחו של מודל מוטא – גם אם התוצאות של המודל אינן מוטות כשהן לעצמן.¹⁰⁸⁴ זאת משום שהחלטות אלגוריתמיות יכולות למשל להיות הוגנות אך שגויות מטעמים אחרים (למשל, כשהחלטה אלגוריתמית שהתקבלה בעניין פלוני שגויה מאחר שהרשומה הפרטנית שלו בבסיס הנתונים מכילה מידע שגוי, בחינת ההחלטה בדיעבד תוכל לזהות את השגיאה ולהביא

1079 ראו לעיל בסעיף 4.1.1.2.

1080 ראו למשל Kaminski, לעיל ה"ש 1004, בעמ' 18; Matthew U. Scherer, *Regulating Artificial Intelligence Systems: Risks, Challenges, Competencies, and Strategies*, 29 HARV. J. L. & TECH. 353, 363 (2016).

1081 Mökander et al., לעיל ה"ש 693, בעמ' 2.

1082 Andrew Selbst, *An Institutional View of Algorithmic Impact Assessments*, 35 HARV. J. L. & TECH. 117 (2021).

1083 Sarah Holland et al., *The Dataset Nutrition Label: A Framework to Drive Higher Data Quality Standards*, in 12 DATA PROTECTION AND PRIVACY, 1 (Dara Halliman et al. eds., 2020).

1084 Mario Martini, *Regulating Algorithms*, in ALGORITHMS AND LAW 100, 122 (Martin Ebers and Susana Navas eds., 2020).

לתיקונה). לכן בקרה של אלגוריתמים בדיעבד נועדה להבטיח לא רק את תקינות קוד המקור של המודל, אלא גם את התקינות של המתודולוגיה הסטטיסטית ושל הליכי איסוף הנתונים ואימונם.

לתפיסתנו, חובת תיעוד צריכה להיות מיושמת בשלבים השונים של מעגל החיים של למידת המכונה. נוסף על כך, אין חובה להשהות את הבקרה בדיעבד המבוססת עליה עד סיום מעגל החיים הזה, ואפשר לחשוב על אבני דרך של בקרה בהתאם להתקדמות התהליך. חובת התיעוד לצורכי בקרה יכולה לעסוק בתקינות הליכי איסוף הנתונים ואימונם או המתודולוגיה הסטטיסטית. אבל היא יכולה להתגלם גם בדרישה לדיווח על אירועים מפתיעים או תוצאות הערכה מדאיגות במהלך אימון מערכות (דיווח כזה יאפשר למשל לרגולטורים לנהל רשימות של גישות הכשרה בסיכון גבוה ולעדכןן); או בשיתוף הערכות סיכונים לפני יישום של מודלים לצורך הערות וביקורת כדי לוודא שעצם היישום הוא בטוח. בתקופת הביניים עד יצירת גוף ידע ואסדרה מבוססת ייתכן שיש מקום לדרוש גם דיווח מדעי, כלומר הצגת תיעוד לקהילה המדעית כדי לעודד מחקר מדעי נוסף בנושא.

אין ספק שאת הבקרה צריכים לעשות בראש ובראשונה גורמים פנימיים לתהליך הפיתוח, שיש להם הבנה עמוקה וגישה למכלול הנתונים והמודל. לפעמים הדרישה לבקרה יכולה לכלול יצירה של שכבות ארגוניות של בטיחות, כגון מינוי פונקציית הערכת בטיחות פנימית שאינה תלויה בפיתוח המודל הקונקרטי ומדווחת לראשי הארגון, שהם בתורם חייבים למסור לגורמים מחוץ לארגון דיווח רזה יותר. חברות טכנולוגיה שונות מציעות דוגמאות לטופסי תיעוד כאלה ואולם חשוב שיהיה להן תוקף ציבורי-רגולטורי.¹⁰⁸⁵

בהקשרים אחרים יכול להתעורר צורך לתת לחקורים חיצוניים גישה למודל. ואולם בעקבות גישה ועמיתיו אנו סבורים כי כמו שעושים ביקורת חשבונאית חיצונית לדוחות כספיים של תאגידים, יש לחייב חברות המפתחות ומפעילות מערכות בינה מלאכותית להגיש לגורם ביקורת חיצוני המתמחה בתחום ההוגנות

1085 חבניות כאלה מכונות model cards. ראו למשל את אלו שמציעה חברת גוגל:
The value of a shared understanding of AI models

והבינה המלאכותית דגימה מהתוצאות שהמערכת מפיקה.¹⁰⁸⁶ גישה דומה באה לידי ביטוי בהוראה שנוספה לאחרונה לקוד המינהלי של העיר ניו יורק, הקובעת שהשימוש במערכות החלטה אלגוריתמיות לגיוס או לקידום עובדים יותנה בביקורת חיצונית מפני הטיות.¹⁰⁸⁷ בקרה יכולה להיעשות גם באופן וולונטרי, במיקור המונים, על ידי הצעת פרס למי שימצאו הטיות (bias bounty), לפי הדגם של (bug bounty).¹⁰⁸⁸

10.5 **פיתוח כלים להתמודדות עם הטיות: שלעיתים נובעת ממאגרי מידע חסרים או בלתי מתאימים ולעיתים מאפליה מערכתית. מערכות בינה מלאכותית עלולות להנציח, ואפילו**

להגביר, הטיות – ואחת היא מה מקורן – באמצעות משוב. כדי לזהות הטיות אלו נדרשת אסדרה ייעודית שתבנה את הבקרה לאורך כל שרשרת הפיתוח של מערכות אלו.

יש הסבורים שתופעת ההטיות האלגוריתמיות היא אכן מטרידה, אבל יש להשלים עם קיומה ולראות בה תחליף להטיות אנושיות, שהן מרובות ועמוקות מאלה של המכונה. אחרים סבורים שהקושי הוא ביישום: אין הגדרה אחידה ומוסכמת של מושג ההוגנות¹⁰⁸⁹ ואין הסכמה על מדדי ההוגנות;¹⁰⁹⁰ ולא ברור באיזו חוליה בשרשרת הייצור אפשר להתערב כדי למנוע את ההטיות. הימנעות משימוש

James Guszcza et al., *Why We Need to Audit Algorithms*, HARV. BUS. REV. (28.11.2018) 1086

N.Y.C. Admin. Code §§ 20-870 – 20-874 (2022) 1087

Melissa Heikkilä, *A Bias Bounty For AI Will Help to Catch Unfair Algorithms Faster*, MIT TECHNOLOGY REVIEW (20.10.2022) 1088

Amit Elazari, *Private Ordering Shaping Cybersecurity Bug Bounties Policy: The Case of Bug Bounties*, in REWIRED: CYBERSECURITY GOVERNANCE 231 (Ryan Ellis and Vivek Mohan eds., 2019)

1089 ראו לעיל ה"ש 882.

Virginia Foggo, John Villasenor, and Pratyush Garg, *Algorithms and Fairness*, 17 OHIO ST. TECH. L. J. 123, 136-153 (2021) 1090

בפרמטרים פסולים, למשל, לא בהכרח תוביל למיגור ההטיות האלגוריתמיות.¹⁰⁹¹ מנגד, אפשר לרתום אלגוריתמים חכמים לסיוע בזיהוי הטיות פסולות במערכות נבונות באמצעות שילובם של כלי בקרה אוטומטיים בהליכי הפיתוח והבקרה בדיעבד. למשל, חברת אמזון מציעה כלי לזיהוי הטיות במודלים שפותחו בטכניקות של למידת מכונה,¹⁰⁹² המסתמך בין השאר על מדדי הוגנות שפותחו באוניברסיטת אוקספורד.¹⁰⁹³ אף שפתרונות אלו עדיין בחיתוליהם, ואינם מקיפים,¹⁰⁹⁴ החובה למנוע הטיות אלגוריתמיות תתמרץ פיתוח שלהם.

לתפיסתנו יש לפעול לצמצום הטיות אלגוריתמיות באופנים האלה:

א. הליך הוגן
מבחינה סטטיסטית
 אנו סבורים שנכון לבחון את הזכות להליך הוגן (due process) מבחינה סטטיסטית וליישם אותו. הכוונה לאסטרטגיות מתמטיות, כגון אישור מראש של משתנים המותרים במידול (להבדיל מאיסור על השימוש בהם בדיעבד),¹⁰⁹⁵ או שימוש בבסיסי נתונים רחבים ועשירים כדי למתן את המשקל הסטטיסטי של משתנים חליפיים.¹⁰⁹⁶ כיום אסטרטגיות אלו יכולות להשפיע על דיוקן של המודל¹⁰⁹⁷ ועל רמת החדירה לפרטיות, ואחדות מהן יכולות להימצא לא יעילות אם בסיס הנתונים כשהוא לעצמו משקף הטיות מערכתיות או היסטוריות.

1091 ראו למשל Feuerriegel, Dolata and Schwabe, לעיל ה"ש 879, בעמ' 379; Prince and Schwarz, לעיל ה"ש 871, בעמ' 1306-1310.

1092 Mark Labbe, *Amazon SageMaker Clarify Aims to Mitigate Bias in Machine Learning*, TECHTARGET (8.12.2020)

1093 Sandra Wachter, Brent Mittelstadt, and Chris Russell, *Why Fairness Cannot Be Automated: Bridging the Gap Between EU Non-Discrimination Law and AI*, 41 COMPUTER LAW & SECURITY REVIEW (2021)

1094 ראו למשל Feuerriegel and Schwabe, לעיל ה"ש 879, בעמ' 382.

1095 Prince and Schwarz, לעיל ה"ש 871, בעמ' 1306-1310.

1096 שם, בעמ' 1310-1311.

1097 Sam Corbett-Davies et al., *Algorithmic Decision Making and the Cost of Fairness*, KDD '17: PROCEEDINGS OF THE 23RD ACM SIGKDD INTERNATIONAL CONFERENCE ON KNOWLEDGE DISCOVERY AND DATA MINING 797 (2017)

ואולם פיתוח של מתודולוגיות סטטיסטיות לאורך כל שלבי הפיתוח של מערכות נבונות, כלומר יישומו של הליך הוגן מבחינה סטטיסטית, יכול לסייע בזיהוי ובאיתור הטיות דווקא במטרה לפתח מודלים מדויקים וחפים מטעויות. יש לזכור גם שלעיתים אלגוריתמים מעלים אל פני השטח הטיות ואפליה שכבר קיימות. כפי שהראינו, המערכת האוטומטית לקבלת מועמדים ללימודים בבית הספר לרפואה סנט ג'ורג' שעסקה את ההטיות של ועדת הקבלה האנושית. לכן אפשר לראות בהטיות האלגוריתמיות הזדמנות לזיהוי ההטיות המערכתיות וההיסטוריות המוטמעות במוסדות החברתיים שהאלגוריתם ממדל.¹⁰⁹⁸

1. גיוון ההון האנושי כלי נוסף להנדסת הוגנות הוא גיוון ההון האנושי בקרב מפתחים של מערכות נבונות.¹⁰⁹⁹ ההנחה שביסוד ההצעה לגיוון היא שלמהנדסים בעלי זהות מגדרית, אתנית, דתית או פוליטית שונה יהיו רגישויות שונות להיבטים של הוגנות אלגוריתמית. ואולם יש המטילים בהצעה זו ספק,¹¹⁰⁰ שכן רקע אתני שונה לא בהכרח מבטיח מחויבות להוגנות ולא לתכליות עסקיות. לדוגמה, מפתח אפרו-אמריקאי נשאל בריאיון מדוע למוצר הדגל של חברת ההזנק שלו, מערכת בינה מלאכותית מבוססת טלפונים חכמים, נבחרו קול ודמות של אישה לבנה צעירה (ולא למשל של גבר שחור). המפתח השיב שחרף מוצאו הוא חש לחץ להשלים עם דעות קדומות כדי לקדם מכירות.¹¹⁰¹ יתר על כן, מהנדסים ומפתחים עלולים להטמיע מבלי משים את ההטיות הקוגניטיביות שלהם במודל – למשל, מהנדס שמפתח אלגוריתם לזיהוי פנים עשוי לאמן את האלגוריתם להתייחס למאפיינים חזותיים מובהקים של בני מוצא אתני מסוים – צורת העיניים, עובי האף וכדומה. ההחלטות

Jennifer M. Logg, *Using Algorithms to Understand the Biases in Your Organization*, HARV. BUS. REV. (9.8.2019) 1098

1099 ראו לדוגמה AI HELG 2018, לעיל ה"ש 298, בעמ' 20; ס' 33a, 37, 57 להצהרה טורונטו, לעיל ה"ש 324; House of Lords Select Committee on Artificial Intelligence, לעיל ה"ש 307, בעמ' 57.

1100 ראו BENJAMIN, לעיל ה"ש 201, בעמ' 24-25, 49.

1101 Quentin Hardy, *Looking for a Choice of Voices in A.I. Technology*, THE NEW YORK TIMES (10.10.2016)

בעניינים אלו יכולות להתבסס על מחקר אקדמי על מאפייני זיהוי פנים (שאף הם עלולים להיות מוטטים), אך גם על ניסיון החיים הייחודי של המפתח, שמושפע בין השאר ממוצאו האתני.¹¹⁰²

10.6

התמודדות עם אתגרי השקיפות האלגוריתמית: איתנות מדעית תהליכית

בפרקים הקודמים ראינו שהבעיה העיקרית באסדרת פעילותן של בינה מלאכותית ומערכות לומדות היא בעיית הקופסה השחורה.¹¹⁰³ האם ההחלטה האוטומטית של האלגוריתם שרירותית או שמא אפשר להצדיקה? העכירות

האלגוריתמית מקשה על המפתחים, הרגולטורים והמשתמשים להתחקות על הסיבה לתוצאות בלתי צפויות, שגויות או מזיקות של החלטות מכונה.

בהקשר זה נתקלנו עד כה בשתי בעיות. הראשונה, מקרים של עכירות אלגוריתמית שבהם אי-אפשר ליצור שקיפות אלגוריתמית. כיום מדובר בעיקר במערכות שפועלות כרשתות נוירונים, אך יש להניח שבעתיד יהיו אתגרי שקיפות נוספים. אם הנחת היסוד היא שגם במקרים של עכירות אלגוריתמית, ללא שקיפות במובנה הקלסי, מוצדק לפתח מערכות לומדות ולא לאסור אותן כליל, כי אז נדרשים נתיבים חלופיים שיאפשרו הגנה על זכויות ואינטרסים.

בעיה שנייה היא ההגדרה העמומה של חובות השקיפות האלגוריתמית במסמכים שונים. בגלל עמימות זו קשה להחיל אותה על שחקנים שונים וליצור ודאות בקרב השחקנים בשוק. כדי ליצור מערכת מבוססת בינה מלאכותית נדרשות הוראות ברורות כבר בשלבים הראשוניים של איסוף המידע ולא רק לאחר היווצרות האלגוריתם. ללא כללים ברורים באשר לאופן ניהול המערכת, גיוון המידע הנדרש והליכי התייעוד יתקשו יוצרי אלגוריתמים לחזות את הדרישות הצפויות מהם, דבר שיוביל לכשלים בחישוב העלויות ויפגע ביעילות המוצר הסופי.

Clare Garvie And Jonathan Frankle, *Facial-Recognition Software* 1102
Might Have a Racial Bias Problem, THE ATLANTIC (7.4.2016)

1103 ראו בעיקר לעיל בפרק 7.

גדי פרל ותהילה שוורץ אלטשולר¹¹⁰⁴ הציעו מודל שקיפות המבוסס על התפיסה הקלסית של שקיפות, אבל כולל חלופה שמותאמת למגבלות הטכנולוגיות של מערכות אלגוריתמיות שאינן מסוגלות לספק הסבר רציונלי לפלט נקודתי. המודל מיועד למקרים שבהם אין אפשרות לספק פלט, אבל יש נחיצות חברתית לשמור על חובות של שקיפות כדי שלא למנוע פיתוח טכנולוגיות מסוימות ושימוש בהן. מודל זה מציע שני מסלולים חלופיים למימוש חובת השקיפות: באחד, שמשמש ברירת מחדל, יש לדרוש בדיקה של התוכנה והסברתיות שלה ולחייב את יוצר האלגוריתם לספק פלט שיוכל לנמק מדוע התקבלה החלטה מסוימת ולא החלטה אחרת. במסלול השני יש לדרוש הוכחה ל"איתנות התהליך", כלומר להראות שמהלך הייצור מבוסס על פרמטרים מדעיים מוכחים בשלושה מקטעים בשרשרת הייצור: המידע המוכנס למערכת; עיבוד המידע; יישום האלגוריתם לצורך קבלת תוצר.

כדי לעמוד במבחן הסף המשפטי יש לפרט את פונקציית המטרה של התוכנה ולבחון אם מטרה זו חוקית בנסיבות שבהן יש כוונה להשתמש במערכת. הטלת החובה לפרט מהי פונקציית המטרה תאפשר לנהל דיון בשאלה אם המערכת עומדת בדרישות החוק גם כשמטרתה היא הגברת יעילות בלבד.¹¹⁰⁵ בשלב השני יידרש יצרן האלגוריתם או המוצר להוכיח כי מאגר הנתונים שעליו מבוססות ההחלטות הסטטיסטיות של התוכנה עומד בדרישות שייקבעו בהתאם לשימוש המיועד. כך למשל יהיה עליו להוכיח שמאגר הנתונים מתאים לקבוצת היעד שנועד לה – הן מבחינת גודל המדגם הן מבחינת התאמת האוכלוסייה ומאפייניה. יהיה עליו להוכיח גם שהמודל נקי מהטיות שיש בהן כדי לפגוע בחוקיותו, כגון מידע שעלול להביא לתוצאות שונות בקרב קהלי יעד מגזעים או מגדרים שונים גם אם פרטים אלו אינם מופיעים במפורש בתוך המאגר. שלב זה ישפיע על כל המהלך של ייצור האלגוריתמים, שכן הוא יחייב שילוב של מומחי תוכן מטעם היצרן בכל השלבים של ניהול הליך ייצור התוכנה.

1104 פרל ושוורץ אלטשולר, לעיל ה"ש 881.

1105 תיחכן טענה שפירוט מטרת התוכנה הוא כשלעצמו סוד מסחרי. טענה זו תהיה מוגבלת למקרים שבהם בדיקת מטרת התוכנה נועדה לברר את חוקיותה.

בשלב השלישי יהיה על היצרן להוכיח שיש הלימה בין יישומי האלגוריתם ובין המשימה שהוא נועד לבצע הן ברמה התאורטית הן ברמה המעשית. בשלב זה יהיה על היצרן לספק מידע בנוגע למגבלות יכולת הניבוי והקטגוריזציה של האלגוריתם בשימוש המיועד, ובכלל זה מידע על תוצאות חיוביות שגויות (false positive) ותוצאות שליליות שגויות (false negative).¹¹⁰⁶ יהיה עליו לספק מידע גם בנוגע לאורך חיי התוכנה והצורך לעדכן אותה מעת לעת.¹¹⁰⁷ בדומה לנהוג היום בעולם אישורי התרופות – שבו לכל תרופה יש התוויה רפואית מסוימת, ואם מבקשים להשתמש בה לצרכים אחרים יש לוודא את התאמתה ולדווח לרשויות הרלוונטיות – כך יידרשו גם יצרני מערכות בינה מלאכותית לחשוף לפני המשתמשים את מגבלות האלגוריתם בכל הקשור לניבוי ולקטגוריזציה ולדווח לרשויות הרלוונטיות.

הבחירה בין מסלולי השקיפות תתבסס על קריטריון ההשפעה הפוטנציאלית של האלגוריתם והמוצר על זכויות חוקתיות. תובא בחשבון גם זהותו של מפעיל התוכנה – גוף שלטוני ומדינתי, גוף דו־מהותי¹¹⁰⁸ או גוף פרטי. במקרה שהשימוש בתוכנה הוא מעין פעולה שלטונית, ההשפעה החוקתית תיחשב חמורה יותר והמבחנים יוחמרו בהתאם.

מסלול האיתנות המדעית התהליכית יכלול את החובה לייצר תיעוד של ההליך ולמסור את המידע למשתמשים. היקף חובת התיעוד והמועד שבו יצטרך היצרן לספק את המידע בעניין התיעוד – קודם לשיווק המוצר או לאחריו – ישתנו

¹¹⁰⁶ מחקר של מכון התקנים בארצות הברית הוא דוגמה למחקר מקיף שבדק, בדיעבד, מהימנות של חוכנות לבדיקת אלגוריתמים לזיהוי פנים. הובאו בחשבון זיהויים שגויים – הן כשלא זוהה דמיון הן כשזוהה דמיון במקום שלא היה. ראו Patrick et al., FACE RECOGNITION VENDOR TEST (FRVT) PART 7: IDENTIFICATION FOR PAPERLESS TRAVEL AND IMMIGRATION (NIST Interagency/Internal Report (NISTIR) 8381, 2022)

¹¹⁰⁷ לעניין עדכון חובות עדכון התוכנה ראו ס' 13 להצעת תקנות הבינה המלאכותית האירופיות, לעיל ה"ש 53.

¹¹⁰⁸ בשלב זה הכוונה היא לגופים דו־מהותיים, כפי שהם מוגדרים בדין הישראלי. ראו ע"א 294/91 חברת קדישא גחש"א "קהילת ירושלים" נ' ליונל אריה קסטנבאום, מו(2) 464 (1992). כיום מתנהל ויכוח בעניין מעמדן של רשתות חברתיות לנוכח תפקידן בהחלפת רעיונות. ראו נועה מור "רשתות חברתיות מקוונות כזירות לעיצוב זכויות והקצאתן: לקראת החלתן של חובות מן המשפט הציבורי" דין ודברים יד (2019).

בהתאם לנסיבות, לשוק הרלוונטי ולעוצמת הפגיעה בזכות החוקתית. לשם הגנה על סודות מסחריים לא יצטרך היצרן בהכרח לספק את הנתונים לכל משתמש; במקרים מסוימים יהיה אפשר להסתפק בחשיפת החומרים לצד שלישי מוסכם או לרגולטור, כדי שאלו יאשרו מראש את העמידה בתנאים.¹¹⁰⁹

על פי הצעה הנדונה כאן ישולב המודל במסגרות חקיקתיות שיאמצו את הפרשנות למושג "שקיפות" ואת המסלולים החלופיים שבו. גם בתקופת הביניים, עד שתנוסח החקיקה, יוכל המודל לספק בהירות בעניין החובות האתיות מבחינת שקיפות ולאפשר רגולציה עצמית.

10.7

שילוב של אסדרה ענפית ואסדרה רוחבית

אסדרה רוחבית לבינה מלאכותית מבקשת להחיל מדיניות כללית על כלל הפעולות הקשורות לבינה מלאכותית בכל המגזרים. כזה הוא נוסח הצעת תקנות הבינה המלאכותית האירופיות.

ברומה, חוק האחריות האלגוריתמית שהתקבל בארצות הברית בשנת 2022 מטיל על ועדת הסחר הפרלית את האחריות לדרוש הערכת השפעה של מערכות בינה מלאכותית בענפי המשק השונים.¹¹¹⁰

אסדרה רוחבית משיגה משילות וודאות מצד השלטון המרכזי, יוצרת הרמוניה רגולטורית ויכולה אפוא להגביר את אמון הציבור ואת הוודאות הרגולטורית בעבור התעשייה. היא יכולה ליצור האחדה, למשל באמצעות הדרישה לרישומן של כל המערכות האוטומטיות לקבלת החלטות שברשות הממשלה; לדרוש תקינה אחידה; לדרוש ציות לאמנות וחקיקת יסוד; ולדרוש כללים אחידים (למשל חובת גילוי נאות של מוצר מלל או תמונה שנוצרו על ידי מערכת אלגוריתמית או חובת הזדהות של מערכת לא אנושית ככזאת). כל זה בייחוד לנוכח העובדה שהרגולציה של בינה מלאכותית צריכה להשתלב באסדרה

1109 סוגיית הקניין הרוחני חופסת מקום מרכזי בהצעת תקנות הבינה המלאכותית האירופיות, לעיל ה"ש 53. ראו למשל סעיף 3.5 לדברי המבוא להצעה.

1110 ראו לעיל בסעיף 4.5.

הקיימת של העולם הדיגיטלי, ובראשה דיני הגנת הפרטיות, דיני מידע והגנת סייבר. לכן מצד אחד, עליה להידרש למאפיינים הייחודיים של הבינה המלאכותית, ומצד שני עליה להיות ערה לעומס הרגולטורי במרחבים אחרים של העולם הדיגיטלי. יתר על כן, הרגולציה לא נעשית בחלל ריק מרגולטורים ולכל רגולטור תרבות ארגונית שונה.

לא מדובר רק בכך שהאסדרה הזאת צריכה להכיר בקיומם של דברי חקיקה אחרים אלא גם בכך שהיא עשויה להשפיע עליהם. תיתכן שתירה בין אינטרסים, למשל הזכות לפרטיות סותרת את הצורך במאגרי מידע גדולים לצורך פיתוח מודלים אלגוריתמיים איתנים סטטיסטית. ייתכנו גם כפילויות – למשל כשיש דרישה לתעד תהליכים לצורכי בקרה מצד גורמים נפרדים.

לצד האפשרות של אסדרה רוחבית יש גם אפשרות של אסדרה ענפית, כפי שמוצע למשל במסמך מדיניות הרגולציה של משרד החדשנות.¹¹¹¹ דוגמאות לאסדרה ענפית אפשר למצוא גם בחקיקה המוניציפלית של ניו יורק, המסדירה שימוש בבינה מלאכותית בגיוס עובדים; בחקיקה במדינת אילינוי, שחלה רק על שימוש בבינה מלאכותית בראיונות וידאו במהלך גיוס לעבודה;¹¹¹² ובחוק האחריותיות הטכנולוגית במקום העבודה של מדינת קליפורניה, המגביל את השימוש במערכות אוטומטיות לקבלת החלטות המשמשות לניטור במקומות עבודה.¹¹¹³

אסדרה ענפית מאפשרת שימוש ברגולטורים קיימים ובסמכויות האכיפה שלהם, איננה מחייבת להקים מסגרות מוסדיות חדשות, מאפשרת יצירת הסדרים ואמצעי אכיפה המותאמים בצורה מדויקת יותר לענף הרלוונטי, מגבירה את הבהירות והוודאות הרגולטורית ומאפשרת גם לבעלי עניין בתוך כל מגזר להשתתף

1111 מסמך מדיניות הרגולציה של משרד החדשנות, לעיל ה"ש 314, בעמ' 93-94.

1112 820 ILL. Comp. Stat. 42/1 Artificial Intelligence Video Interview Act; ראו גם Ajunwa, לעיל ה"ש 118.

1113 Airlie Hilliard et al., *Overview and Commentary of the California Workplace Technology Accountability Act*, 37 INT'L. REV. L. COMP. & TECH. 99 (2023)

בחשיבה על ההסדרים האלה. אסדרה מסוג זה יכולה לשמש להגנה על זכויות האזרח ועל זכויות הצרכן, אבל גם לחיזוק הוודאות וממילא לחדשנות ולעידוד השקעות.

הקושי באסדרה ענפית הוא בסתירה אפשרית בין ענפים, ביצירת סטנדרטים לא אחידים, בהעמקת הפערים בין הרגולטורים (למשל פערים באוריינות הדיגיטלית, ביכולות האכיפה ובתרבות הארגונית). כמו כן, היא עלולה להשאיר מרחבים לא מאוסדרים שנופלים בין הכיסאות.

כדי להימנע ממגרעות האסדרה הענפית, ולהבטיח את יתרונות האסדרה הרוחבית ואת יתרונות האסדרה הענפית גם יחד, אנו ממליצים לשלב אסדרה רוחבית ואסדרה ענפית, לצד היכרות עם שדות הפעילות המשפטיים. כך, יש לעצב את האסדרה בתווך שבין מעגל החיים של למידת המכונה והשלכותיו, שפורטו לעיל, הענף הספציפי שבו מיושמת המערכת והדרישות המשפטיות הרלוונטיות, כפי שמתואר בתרשים 5.

תרשים 5

מודל משולב רוחבי וענפי לאסדרת בינה מלאכותית



10.8

**חקיקה משלימה
ועדכוני חקיקה**

גם בהיעדר חוק בינה מלאכותית כללי וייעודי יש לחייב כבר בעת הזאת את מקבלי ההחלטות והרגולטורים הייעודיים לעדכן את החקיקה הקיימת ולחוקק חקיקה משלימה. הכוונה בעיקר לחוקים כגון חוק התחרות וההגבלים העסקיים, חוק זכויות יוצרים, חוק הגנת הפרטיות, פקודת הראיות וחוק הרכש הממשלתי.

למשל, יש להוסיף לחקיקה הקיימת התמודדות עם זכויות יוצרים ובעלות על יצירות שנוצרו באמצעות בינה מלאכותית;¹¹¹⁴ להידרש לסוגיות שקשורות לחובות זיהוי של מוצרים המבוססים על בינה מלאכותית בעת התקשרות עם לקוחות פוטנציאליים; וכן להוראות שונות לקשורות לאיסוף המידע ולשימוש בו כשמדובר בתוכנות המשמשות משרדים ממשלתיים. אנו מציעים כי ההליך ייעשה בהובלת משרד המשפטים וייקבעו סדי זמנים להשלמת תיקוני החקיקה הנדרשים בשלב זה.

1114 בהקשר של יצירות אלגוריתמיות, סוגיות של קניין רוחני יכולות להשפיע השפעה ניכרת על תפוצת הטכנולוגיה. ראו למשל את גישתו של ריאן אבוט, המעודד ניטרליות חקיקתית. RYAN ABBOTT, *THE REASONABLE ROBOT* (2020). על אפשרות של נזק עקב היעדר אסדרה בתחום אפשר ללמוד מן החלטה של חברת Getty, המציעה מאגרי תמונות לשימושים מסחריים, שלא לכלול בהם תמונות שיוצרו באמצעות אלגוריתמים כגון DALL-E2 מחשש לחשיפה משפטית. James Vincent, *Getty Images Bans AI-Generated Content*. *THE VERGE* (21.9.2022). *over Fears of Legal Challenges*. ראו גם משרד המשפטים שימושים בתכנים מוגנים בזכויות יוצרים לצורך למידת מכונה (18.12.2022).

פרק אחד עשר

המלצות לגבי מוסד מאסדר
לבינה מלאכותית בישראל

—

במאי 2023 נערך בסנאט האמריקני שימוע לסאם אלטמן, מנכ"ל חברת OpenAI. אחת ממסקנות השימוע הייתה שיש להקים רשות פיקוח על בינה מלאכותית בארצות הברית. ואולם מייד הועלו הטענות שהקמת רשות כזאת היא המלאכה הקלה – האתגר הקשה הוא ציודה בסמכויות המתאימות.¹¹¹⁵

נוסף על כך, לצד האתגר הרגולטורי של תרגום עקרונות אתיים רחבים לכללים משפטיים יש שאלה נפרדת: מהי זהות הגוף או הגופים המוסדיים האמורים

Andrew Ross Sorkin et al., *Washington Confronts the Challenge of Policing A.I.*, THE NEW YORK TIMES (17.5.2023) 1115

להיות ממונים על תחום הבינה המלאכותית הן כגופים מיישמי חקיקה, הן כגופים מאסדרים ומפקחים והן כגופים מייעצים? השאלה אם גוף כזה נדרש מלכתחילה, ומה צריכה להיות המסגרת המוסדית שלו, תידון בפרק זה הן באשר לישראל כשהיא לעצמה הן בהשוואה למתרחש בעולם. נעסוק בשאלה על פי אילו קווים מנחים צריך לפעול מוסד מאסדר לבינה מלאכותית בישראל ומהם השיקולים הרלוונטיים בעניין זה. נציג את פונקציית המטרה של המוסד המאסדר, את הסביבה הרגולטורית שלו ואת הגופים המקבילים אליו, וננסה לבחון מאילו מהם ניתן לשאוב השראה. נקדים ונעיר כי שאלה זו תלויה במידה רבה בהכרעה אילו אינטרסים אמורה רגולציה של בינה מלאכותית לקדם ואילו חלקים בשרשרת הערך היא אמורה לתפוס ברשתה. כיוון שהרגולציה של בינה מלאכותית עודה בחיתוליה בעולם כולו, קשה לקבוע מסמרות בדבר, אבל חשוב בעינינו להצביע על המסגרת הנחוצה כבר כעת.

האסדרה של תחום הבינה המלאכותית עדיין נמצאת בראשיתה ועדיין מתפתחת ולכן גמישותה חשובה ביותר. רגולטור של בינה מלאכותית יכול למלא תחילה בעיקר תפקידים של ריכוז מומחיות, ייעוץ והמלצה לקובעי מדיניות, ועם הזמן להרחיב את פעילותו ולדאוג גם לפיקוח ולאכיפה. ככל שטכנולוגיות אלו יוטמעו יותר במשרדי הממשלה¹¹⁶ כך יגדל הצורך בממשק בין הרגולטור, כגורם מקצועי פנים-ממשלתי מוביל, לרשויות אחרות.

מתווי רגולציה של ניהול סיכונים דורשים פעילות רגולטורית מתמשכת שתחזור ותעריך את רמת הסיכון הנשקפת מטכנולוגיות חדשות, תחזור ותתאים את המסגרת הרגולטורית לצורכי השעה ותיתן מענה קונקרטי לסוגיות פרשניות שעשויות להתעורר עד כניסת החקיקה לתוקף וגם מייד לאחר מכן. גם היישומים של מערכות נבונות והשפעתן על החברה יכולים להשתנות, להתרחב ולהתפתח. תופעות הלוואי של טכנולוגיות אלו – רצויות ולא רצויות – ילכו ויתגלו ככל ששימושן יתרחב ויעשה רבגוני יותר. הבחינה החוזרת של המדיניות הרגולטורית לצורכי עדכונה יהפכו לעניין שבשגרה.

לגופים העוסקים בכינה מלאכותית, בין שהם כבר קיימים ובין שרק הוצעו, יש כמה תכליות. מכל אחת מהן נגזרים שיקולים אחרים באשר

למהות המוסד המאסדר ולעיצובו. נמנה אותן להלן:

א. קידום חדשנות, פיתוח וצמיחה: באסטרטגיה של כינה מלאכותית לאומית מודגש לרוב תפקיד המדינה בעידוד ההצטיינות הלאומית בתחום וביצירת מערכת אקולוגית של חדשנות, יזמות וצמיחה במגזר הפרטי.¹¹¹⁷ עם תפקידי המדינה יש למנות גם הקצאת משאבים להקמת תשתיות, למחקר ולהכשרת כוח אדם.

ב. שמירה על אינטרסים ביטחוניים: תוכנית החומש הישראלית מונה בין השאר, כצפוי, גם מטרת ביטחוניות.¹¹¹⁸ המוסד שהוצע בה הוא מוסד מקביל למערך הסייבר הלאומי, כלומר מוסד בעל אוריינטציה ביטחונית.

ג. מחקר בתחום "הכינה המלאכותית האחראית", השלכות המקרו של כינה מלאכותית על החברה, הכלכלה והדמוקרטיה, והטמעת עקרונות כלל-חברתיים ואתיים בשימוש בכינה מלאכותית.

ד. עיצוב ההסדרים הנורמטיביים האמורים לחול על מערכות נבונות, פיתוחן, הפצתן והשימוש בהן.

1117 ראו לדוגמה את העיקרון האתי הראשון במסמך מדיניות הרגולציה של משרד החדשנות, לעיל ה"ש 314, בעמ' 103-104, "בינה מלאכותית לקידום צמיחה, פיתוח בריקיימא ומובילות ישראלית בחדשנות". ראו גם ס' 3(3) לתזכיר חוק הגנת הסייבר, לעיל ה"ש 461, שלפיו אחד מתפקידי מערך הסייבר הוא לקדם מובילות ישראלית בתחום הסייבר. משרד הדיגיטל, המדיה, התרבות והספורט של בריטניה פתח לאחרונה בתהליך של גיבוש רגולציית דיגיטל (שבה נכללת גם רגולציה של מערכות בינה מלאכותית). מטרתה הראשונה של רגולציה זו היא "לעודד צמיחה ולקדם תחרות וחדשנות בכל רחבי המגזר הדיגיטלי". ראו Department of Digital, Culture, Media and Sports, Digital Regulation: Driving Growth and Unlocking Innovation (13.6.2022)

1118 בן ישראל, מתניה ופרידמן ממליצים על פיתוח מערכות שוי"ב לימי שגרה ולשעת חירום. הם אף מציינים כי מטבע הדברים הם מנועים מלהמליץ על פרויקטים ביטחוניים מסווגים יותר. המיזם הלאומי למערכות נבונות (2020), א, לעיל ה"ש 981, בעמ' 37-38; להמלצות לקידום מערכת אקולוגית של מערכות נבונות בתחום הביטחוני ראו שם, ב, בעמ' 248-255.

ה. פיקוח ואכיפה באמצעות כללי אתיקה וכללים רגולטוריים והנחיה ישירה או עקיפה של גופים ציבוריים ושל המגזר הפרטי.

ו. ייעוץ, יצירת ידע, הפצתו והנחיה בנוגע לסוגיות שמעורר השימוש בבינה מלאכותית. לנוכח הדינמיות של התחום יצטרך הרגולטור להכיר לעומק את חזית הטכנולוגיה¹¹¹⁹ ולהפיץ את הידע הזה באופן מסודר.

נקדים ונאמר: לתפיסתנו לא נכון להקים מוסד שאמור לממש את כל התכליות שנמנו לעיל משום שמקצתן ביצועיות ומקצתן נורמטיביות או רגולטוריות. נוסף על כך, ייתכנו קונפליקטים ביניהן. למשל, עידוד חדשנות דורש אמון בין התעשייה למדינה, המתבטא בראש ובראשונה בשקיפות מרבית, ואילו שמירה על אינטרסים ביטחוניים נעשית לרוב בעלטה; קידום פיתוח וצמיחה עלול להתנגש בצורך לעצב הסדרים נורמטיביים, שכן לעיתים הצורך להגן על אינטרסים ציבוריים אחרים מעכב חדשנות.

לכן אנו מציעים שגוף אסדרה של בינה מלאכותית יתמקד בשלוש התכליות האחרונות שמנינו, שהן תכליות המאפיינות רגולציה. על שתי התכליות הראשונות צריכים להיות ממונים הרשות לחדשנות ושירותי הביטחון. במקביל, אנו מציעים לשקול הקמת מרכז מחקר לאומי בין-תחומי בנושאי בינה מלאכותית וחברה, שתכליתו תהיה ליצור ידע ולגבש אסטרטגיות בעניין שימושים אחראיים ואתיים במערכות בינה מלאכותית. מרכז כזה יענה על הצורך להציע הבנה אינטגרטיבית של התחום, שממנה ייגזרו לאחר מכן מדיניות ופרקטיקות שאינן כפופות לאינטרס המדינתי בקידום התעשייה מצד אחד וברגולציה ספציפית מצד שני. יש לתת לו מימון רב-שנתי, לפחות לתקופת ביניים של כחמש שנים, עד שיתגבשו גופי ידע ברשויות הפיקוח והאסדרה, באקדמיה ובתעשייה. המרכז למצינויות בנושא קבלת החלטות אוטומטית וחברה שהקימה ממשלת אוסטרליה בשנת 2020 היא דוגמה למרכז כזה.¹¹²⁰

1119 ראו למשל את עמדתם של אחיעז ואחי', לעיל ה"ש 240, בעמ' 147.

11.2 הסביבה הרגולטורית

אסדרת תחום הבינה המלאכותית נמצאת בסביבה שיש בה ארבעה מרחבים רגולטוריים, שכל אחד מהם כולל אסדרה מהותית וגם הקשר מוסדי.¹¹²¹

11.2.1. רגולטורים משיקים תחום הבינה המלאכותית משיק לתחום ההגנה על הפרטיות ולתחום הגנת סייבר בעניינים רבים, ממש כשם שיש חפיפה עקרונית בין תחום ההגנה על הפרטיות לתחום הגנת הסייבר.¹¹²² למשל, הקפדה על עקרון העיצוב לפרטיות¹¹²³ תורמת להגנת סייבר על ידי הקטנת הסיכון לזליגת מידע שמלכתחילה לא היה בו צורך. מבחינת אסדרה, יש השקה גם בין תחום הבינה המלאכותית לתחומים כגון תחרות, מיסוי, הגנת הצרכן, ואפילו לתחום הפיקוח על הבחירות.

לפיכך הצורך ברגולטור מתאם הוא מורכב. יש חשש שרגולציה של רגולטורים מגזריים ללא גורם רגולטורי מתאם תיצור מערך נורמטיבי של טלאים ושל פערים בין מגזרים ותביא לחוסר תיאום בין רגולטורים. במקרה כזה הם עלולים לדרוש מהגופים המאוסדרים דרישות סותרות, חופפות או כפולות.¹¹²⁴

כמבט אל העתיד הרחוק יותר, ייתכן שבשלב מסוים יבשילו הנסיכות ויהיה אפשר לרכז את הרגולטורים של כלל תחומי הדיגיטל תחת קורת גג אחת. רגולטור דיגיטל כזה יהיה ממונה, בין השאר, על הגנת מידע, אבטחת סייבר¹¹²⁵ ובינה מלאכותית, ויאפשר הרמוניזציה של תהליכים רגולטוריים. למשל, במקום שתורת ההגנה בסייבר, כללי אבטחת המידע, דרישות התייעוד הטכני והוראות ממשל

1121 לגישה הבריטית, המאחדת תחומים אלו – ואחרים – ב"רגולציית דיגיטל", ראו Department of Digital, Culture, Media and Sports, לעיל ה"ש 1117.

1122 ראו מערך הסייבר הלאומי, תורת ההגנה – לנהל את הסיכון: המדריך היישומי (השלם) להגנת הסייבר של ארגון, 72–74, נספח ו (2021).

1123 ראו לעיל ה"ש 412.

1124 הצוות הבין-משרדי לרגולציה חכמה, תכנית לאומית למדיניות רגולציה ככלי לשיקום המשק ביציאה ממשבר הקורונה 44–45 (2021); Calo, לעיל ה"ש 956, בעמ' 6.

1125 הכוונה לפונקציות הרגולטוריות של מערך הסייבר הלאומי לפי פרק ד בתזכיר חוק הגנת הסייבר (לעיל ה"ש 461). ספק אם מקומה של הגנת סייבר שוטפת במסגרת ה-CERT הלאומי הוא תחת רגולטור דיגיטל.

הנתונים יהיו תחומים נפרדים, יהיה אפשר להציג מעין "תורת דיגיטל" מאוחדת שתידרש לכל התחומים האלה באופן הוליסטי ותפחית את הנטל הרגולטורי המוטל על כתפי הגופים המאוסדרים. אבל עניין זה ידרוש התאמה חקיקתית ואפשר שיהיו לו גם חסרונות, למשל הסתירה המובנית בין אבטחת סייבר ובין הגנה על פרטיות. כך או כך, אנו סבורים שאין להמליץ על כך בשלב זה.

11.2.2. רגולטורים ענפיים גופים שמשמשים כבינה מלאכותית יכולים

להיות כפופים לרגולטור ענפי, ולעיתים אף לכמה רגולטורים ענפיים (למשל, מוסדות פיננסיים נתונים לפיקוח הרשות לאיסור על הלבנת הון וגם לפיקוח של רגולטורים פיננסיים – המפקח על הבנקים, הרשות לניירות ערך או רשות שוק ההון, ביטוח וחיסכון).

יש הסבורים שאין צורך למנות רגולטור ייעודי לתחום הבינה המלאכותית ומוטב להפקידו בידי הרגולטורים הענפיים. דוח ועדת המשנה, כמו מסמך מדיניות הרגולציה של משרד החדשנות, ממליצים שעיקר האסדרה תיעשה באמצעות הרגולטורים הענפיים, על יסוד גורם רגולטורי מומחה ועל-ממשלתי שיתאם בין הרגולטורים הענפיים.

לאסדרה ענפית יש יתרונות, ובהם היכולת להתמודד עם הקושי לעדכן את הרגולציה הריכוזית על ידי ניסויים רגולטוריים בהקשרים מגזריים מבודדים יחסית.¹¹²⁶ עם זאת, לעומת רגולטור מרכזי של בינה מלאכותית, רגולטורים ענפיים לא תמיד רגישים די הצורך להגנה על זכויות יסוד.¹¹²⁷ יתר על כן, ייתכנו פערים ביניהם גם ברמת האוריינות הדיגיטלית שלהם.

תקנות מערכות ההמלצה האלגוריתמיות הסיניות, למשל, קוראות לארגונים ענפיים רלוונטיים לפתח סטנדרטים ענפיים ולהנחות ספקים.¹¹²⁸ לפי התקנות האירופיות, כל מדינה שחברה באיחוד תקים רשות לאומית מוסמכת, שתפקידה

1126 שרון בר זיו וטל ז'רסקי "פרטיות במשבר זהות: אסטרטגיות הסדרה בעידן ההתממה" משפט, חברה ותרבות ב 125, 131-130 (2019).

1127 שם, בעמ' 165.

1128 ס' 2 לתקנות מערכות ההמלצה האלגוריתמיות הסיניות, לעיל ה"ש 590.

יהיה להבטיח את יישומן והטמעתן של התקנות המוצעות ולשמש מפקח על השוק.¹¹²⁹ ביחס לענפים מסוימים, הרגולטור הענפי ישמש מפקח שוקי.¹¹³⁰

גם שרון בר זיו וטל ז'רסקי סבורים שייתכן שהמומחיות הטכנולוגית של הרגולטורים הענפיים תעלה על זו של רגולטור מרכזי.¹¹³¹ הסמכת רגולטור ענפי לנהל ניסויים מבוקרים בארגזי חול רגולטוריים היא דוגמה לכלי המתאים לתקופה של ניסוי וטעייה. כך מומלץ בין השאר בדוח ועדת המשנה לאתיקה של המיזם הלאומי למערכות נבונות,¹¹³² בדוח על יישומי בינה מלאכותית במגזר הפיננסי¹¹³³ ובמסמך מדיניות הרגולציה והאתיקה של משרד החדשנות.¹¹³⁴ בתקנות האירופיות מוצע לאפשר למדינות חברות להקים ארגזי חול רגולטוריים,¹¹³⁵ שבהם יהיה אפשר לחרוג, בהתקיים תנאים מסוימים, מעקרון צמידות המטרה שבדיני הגנת הפרטיות ולהשתמש בנתונים לתכליות אחרות מאלו שלשמן הם נאספו. הגמשת הכללים כדי לאפשר פיתוח מבוקר של טכנולוגיה חדשה בהנחיית הרגולטור מאפשרת לשכלל את המסגרת הנורמטיבית החלה על מערכות נבונות ולהתאימה לצרכים מגזריים מיוחדים. לכן היא מתאימה במיוחד לרגולטורים ענפיים.

11.2.3. מערכות פרטיות ההיבט השלישי של האסדרה נוגע להפרדה וציבוריות בין מערכות המשמשות במגזר הפרטי לאלה המשמשות במגזר הציבורי. כיום למשל הרשות הלאומית לחדשנות טכנולוגית¹¹³⁶ פועלת "לעידוד, לקידום, לתמיכה ולסיוע

1129 שם, ס' 63.

1130 שם, ס' 63(1).

1131 בר זיו וז'רסקי, לעיל ה"ש 1084, בעמ' 164.

1132 המיזם הלאומי למערכות נבונות (2019), לעיל ה"ש 313, בעמ' 32.

1133 אחיעז ואח', לעיל ה"ש 240, בעמ' 148.

1134 מסמך מדיניות הרגולציה של משרד החדשנות, לעיל, ה"ש 314, בעמ' 99-101.

1135 ראו לעיל בסעיף 4.1.2.

1136 ראו חיקון מס' 7 לחוק לעידוד מחקר, פיתוח וחדשנות טכנולוגית בתעשייה, תשמ"ד-1984, ס"ח 100 (להלן: חוק המו"פ); אורה קורן "זריקת עידוד לתעשייה - או צעד מסוכן? רשות החדשנות הלאומית יוצאת לדרך" *TheMarker* (30.7.2015).

לחדשנות הטכנולוגית בתעשייה"¹¹³⁷, ואינה מטפלת במגזר הציבורי; רשות התקשוב הממשלתי אמונה על "הרחבת תחומי פעילות התקשוב הממשלתי, עידוד חדשנות במגזר הציבורי וקידום המיזם הלאומי 'ישראל דיגיטלית' "¹¹³⁸, ואילו הרשות להגנת הפרטיות ממונה על ביצוע הוראות חוק הגנת הפרטיות במגזר הציבורי והפרטי גם יחד.¹¹³⁹ בפועל, מבקר המדינה מצא כי "על אף החשיבות שהרשות מייחסת למעורבותה בסוגיות הגנת הפרטיות במאגרי מידע העולות במגזר הציבורי, פעילותה בנושא חלקית" ואין היא מעורבת באופן אקטיבי בפרויקטים ממשלתיים בנושא.¹¹⁴⁰

אנו סבורים כי בדומה לרשות להגנת הפרטיות, יהיה נכון שרגולטור אחד יהיה מופקד על האסדרה של בינה מלאכותית במגזר הפרטי ובמגזר הציבורי גם יחד. ראשית, כיוון שמערכות בינה מלאכותית במגזר הציבורי מבוססות על פיתוח שנעשה בשוק הפרטי. שנית, כיוון שבמגזר הפרטי נדרשים סטנדרטים להפעלת מערכות בינה מלאכותית, לפיתוחן ולשימוש בהן בהקשרים סמי-ציבוריים כגון אלגוריתמיקה של הפצת דעות ורעיונות, מערכות קבלה לעבודה, אפליית מחירים ועוד. שלישית, כיוון שכל הזמן מתקיים קשר דו-כיווני בין שני המגזרים. למשל, רשות החדשנות יכולה לעודד פיתוחים טכנולוגיים במגזר הפרטי על פי תדריך אסטרטגי של הממשלה.¹¹⁴¹ לרשות התקשוב הממשלתי יש יכולת להשפיע על השוק הפרטי באמצעות הכתבת סטנדרטים במכרזי רכש.¹¹⁴² מנגד, רשויות שלטון וכוחות טכנולוגיות שהתפתחו למען השוק הפרטי.

1137 ס' 5א לחוק המו"פ, ש.ס.

1138 החלטה 2097 של הממשלה ה-33, "הרחבת תחומי פעילות התקשוב הממשלתי, עידוד חדשנות במגזר הציבורי וקידום המיזם הלאומי 'ישראל דיגיטלית' " (10.10.2014).

1139 בר זיו וז'רסקי, לעיל ה"ש 1084, בעמ' 129. ראו גם ס' 30 בהצעת חוק הגנת הפרטיות-2019, ודברי ההסבר לו. ארידור הרשקוביץ ושוורץ אלטשולר, לעיל ה"ש 403, בעמ' 116.

1140 מבקר המדינה, דוח שנתי 169 ב 6 (2019).

1141 הרשות לחדשנות, למשל, מפעילה את "הזירה החברתית-ציבורית", שבה היא מעודדת מחקר ופיתוח טכנולוגיים לפתרון אתגרים חברתיים וציבוריים. במסגרת זו פועל מסלול לעידוד חדשנות במגזר הציבורי (ממשלתי). ראו רשות החדשנות, 2019: חדשנות בישראל, תמונת מצב 96 (2020).

1142 על שימוש במכרזים ציבוריים ככלי של רגולציה רכה ראו לעיל ה"ש 981.

11.2.4. ההיבט הבינלאומי חברות ישראליות שמפתחות מוצרים המיועדים לשוק האירופי, למשל, נדרשות לציית הן להוראות התקנות הכלליות בדבר הגנת מידע (GDPR) הן להוראות הדין הישראלי. מן הראוי שיתכוננו כבר כעת לדרישות שעתידות להעמיד התקנות האירופיות.¹¹⁴³ אפשר שחברות פרטיות יידרשו לעמוד גם בתקנים מקצועיים בינלאומיים, כלומר למלא דרישות נוספות (חלקן חופפות) על אלו המוגדרות בדין.¹¹⁴⁴

הדין האירופי בתחום הפרטיות ובתחום אסדרת רשתות חברתיות, וכפי שמסתמן גם בהצעה לאסדרת הבינה המלאכותית, הוא אקסטר-טריטוריאלי מאוד ויחייב תאימות של הדין הישראלי אליו. בדיוק כפי שהרשות להגנת הפרטיות אחראית על תאימות הדין הישראלי לדין האירופי בתחום הפרטיות, אנו מציעים שהרשות המרכזית תהיה אחראית על התאימות לדין שייכנס לתוקף בתחום הבינה המלאכותית.

11.2.5. האם המבנה של מערך הסייבר מתאים לשמש רגולטור בינה מלאכותית? לתפיסתנו, מערך הסייבר הלאומי אינו צריך לשמש השראה לרגולטור של בינה מלאכותית. ראשית, יישומי בינה מלאכותית אינם נחלת התחום הביטחוני, ואין בפיתוח של מערכות נבונות בתור שכאלו כדי לקדם יכולות ביטחוניות.¹¹⁴⁵ יש להימנע אפוא משערוק הדגם של מערך הסייבר הלאומי, שכמה

¹¹⁴³ האומדן של מולר לעלויות הציות של המגזר העסקי להוראות תקנות הבינה המלאכותית האירופיות המוצעות הוא 10.9 מיליארד אירו לשנה. ראו Benjamin Mueller, *How Much Will the Artificial Intelligence Act Cost Europe?* CENTER FOR DATA INNOVATION (July 2021)

¹¹⁴⁴ לסוגיית החפיפה בין תקני אבטח סייבר לדרישות דיני הפרטיות ראו למשל Isabel Maria Lopes, Teresa Guarda, and Pedro Oliveira, *How ISO 27001 Can Help Achieve GDPR Compliance*, 2019 14TH IBERIAN CONFERENCE ON INFORMATION SYSTEMS AND TECHNOLOGIES (CISTI) 1 (2019). לדוגמה להתמודדות של הרגולטור עם חפיפות כאלו ראו מסמך המיפוי של מערך הסייבר: מערך הסייבר הלאומי, תורת ההגנה - מונגשת בשבילך (19.10.2019). ראו גם ס' 43(א)(1) לתזכיר חוק הגנת הסייבר, לעיל ה"ש 461.

¹¹⁴⁵ אפשר כמובן להעלות על הדעת יישומים ביטחוניים של מערכות נבונות: ממערכות חיזוי מודיעיניות עד כלי נשק אוטונומיים. ואולם אלה הם החריג. בשל מאפייניה של החקיקה האירופית, למשל, נמנעות תקנות הבינה המלאכותית האירופיות מהחלת כללים

ממטרותיו הן ביטחוניות מבצעיות.¹¹⁴⁶ לכן שלא כמו מערך הסייבר הלאומי, גוף ביטחוני חשאי למחצה ששייך למשרד ראש הממשלה, רצוי לשייך את רגולטור הבינה המלאכותית למשרד המשפטים או למשרד הכלכלה או למשרד דיגיטל נפרד.

שנית, מבנה מסוג "רגולטור של רגולטורים", כלומר גוף רגולטורי שינחה או יפקח על הרגולטורים המגזריים, הוצע בתזכיר חוק הגנת הסייבר בשנת 2019, ושוורץ אלטשולר וארידור הרשקוביץ העלו לגביו חששות.¹¹⁴⁷ עיקרם: החשש מחוסר ההתאמה של מתווה זה לתרבות הפוליטית והרגולטורית בישראל,¹¹⁴⁸ שבאה לידי ביטוי בנכונות לשיתוף פעולה בין רגולטורים ובפערים בתרבות הארגונית. הסתייגויות דומות הועלו לאחרונה בתגובה להמלצה להקים רשות לרגולציה – גוף מרכזי לבקרה על רגולציה (regulatory oversight body).¹¹⁴⁹ לדברי המסתייגים, רגולטור-על עלול לסרב ולהאריך את עבודת הרגולציה¹¹⁵⁰ ולהחריף מאבקים בין רשויות על כוח וסמכויות.¹¹⁵¹

משפטיים על מערכות אלו. אין הכוונה שמערכות נבונות ביטחוניות יוחרגו בהכרח מסמכותו של רגולטור בינה מלאכותית, אלא שהן לא יהיו מרכז עיסוקו.

1146 על היעדר ההפרדה בין הביטחוני לאזרחי בתזכיר חוק הגנת הסייבר ראו למשל ארידור הרשקוביץ ושוורץ אלטשולר, לעיל ה"ש 454, בעמ' 165–169; מרכז המחקר להגנת הסייבר באוניברסיטה העברית, הערות לתזכיר חוק הגנת הסייבר ומערך הסייבר הלאומי התשע"ח–2018 (11.7.2018).

1147 ארידור הרשקוביץ ושוורץ אלטשולר, שם, בעמ' 163–165.

1148 למשל, בתחום התקשורת רגולטורים מתקשים לשחף פעולה ביניהם. ראו למשל נתי טוקר "מי הבוס של שוק התקשורת?" *TheMarker* (9.12.2019); תהילה שוורץ אלטשולר "קמצנות רגולטורית" *TheMarker* (3.12.2018). לביקורת על שיחוף הפעולה בין מערך הסייבר לרשות הגנת הפרטיות ראו מבקר המדינה, היבטים בהגנה על הפרטיות במאגרי מידע 6, 20–21 (דו"ח שנתי 69 2019).

1149 ראו הצוות הבין-משרדי לרגולציה חכמה, לעיל ה"ש 1084, בעמ' 105–106.

1150 עמירם ברקת "התוכנית שנחשפה בגלובס עולה מדרגה: מי צריך רגולטור שיילחם ברגולציה?" *גלובס* (6.7.2021).

1151 סמי פרץ "תנסו לא להירדם, זה נושא חשוב מדי" *הארץ* (7.7.2021).

11.3 המלצות

11.3.1 תפקיד הרגולטור א. קידום ותיאום אסדרה של מערכות נבונות בישראל, לרבות התנעת התהליך של גיבוש חקיקה רוחבית וריכוזה.

ב. התוויית מדיניות לפיתוח מערכות נבונות בישראל, להטמעתן ולשימוש בהן. בין השאר יגדיר הרגולטור סטנדרטים בעניין מקור הנתונים לאימון ולבדיקה של מערכות אלו ואיכותם; יקבע מבחנים לאומדן הטיות אלגוריתמיות; יבחן כיצד אפשר להבטיח מידה ראויה של שקיפות, הסברתיות או יכולות ביקורת (auditability)¹¹⁵² ונעקבות (traceability) של תהליכים במערכות אלו וכן יקבע סטנדרטים לניהול סיכונים.¹¹⁵³

ג. הנחיה מקצועית של רגולטורים מגזריים כדי להבטיח הרמוניזציה של הכללים החלים על פיתוח מערכות בינה מלאכותית, פרישה שלהן ושימוש בהן. הרגולטור ישמש גורם מנחה שיורי, כלומר יהיה ממונה על אסדרת מערכות נבונות פרטיות בתחומים שאין בהם מאסדר מגזרי.

אופי ההנחיה יכול להתבסס על מודל של ניהול סיכונים. לפי מודל זה תהיה ההנחיה של מפתחים, משתמשים או מפיצים של מערכות נבונות ברמת סיכון גבוהה נוקשה ומחייבת יותר מההנחיה החלה על מערכות נבונות ברמת סיכון נמוכה. נוסף על כך, הרגולטור יהיה הגוף המאשר את ההערכה של ארגוני החול הרגולטוריים שיוציאו אל הפועל רגולטורים מגזריים כדי למנוע חשש לניגוד עניינים, להגביר את השקיפות ולעורר מודעות להשלכות החברתיות הרחבות

1152 ראו גם ACM 2017, לעיל ה"ש 326, בעמ' 2 (עיקרון 6).

1153 AI HELG 2019, לעיל ה"ש 298, בעמ' 13. ל-traceability, ראו גם Defense

Innovation Board, לעיל ה"ש 342, בעמ' 8 (עיקרון מס' 3).

של המיזם בתום תקופת המבחן שלו.¹¹⁵⁴ על הרגולטור יהיה מוטל גם לוודא שהלקחים שהופקו ייושמו כהלכה באסדרה עתידית.¹¹⁵⁵

ד. הנחיה מקצועית של רשויות המדינה בתחום הבינה המלאכותית, ובכלל זה חוות דעת על מסמכי מכרזים ממשלתיים של מערכות בינה מלאכותית, במטרה להבטיח שהמדינה מצטיידת במערכות בינה מלאכותית העולות בקנה אחד עם הסטנדרטים הישראליים והעולמיים בתחום, וכן כדי לייצר בעבור השוק הפרטי תמריצים "רכים" לעבוד בהתאם לסטנדרטים.¹¹⁵⁶ זאת לנוכח העובדה שרכש תוכנה יכול להסוות קבלת החלטות מדיניות, למשל בחירת משתני המטרה, נתוני האימון, תכלית האופטימיזציה וסוג המודל.¹¹⁵⁷

ה. ייעוץ בנושאי משפט הבינה המלאכותית. בהיעדר גורם מרכזי המרכז את הייעוץ בנושאי משפט וטכנולוגיה בישראל (הוא היה מכלל תפקידיה של הרשות למשפט, טכנולוגיה ומידע במשרד המשפטים, אך נזנח כשזו הפכה לרשות להגנת הפרטיות), חוות דעתו של הרגולטור תחייב את הממשלה ורשויותיה בנושאים אלו.¹¹⁵⁸ אנו מציעים שהרגולטור יהיה גם בעל סמכות להציג עמדה עצמאית בתחומים אלו לכנסת ולבתי המשפט.

ו. מוקד ידע, הדרכה ושיתופי פעולה. הרגולטור יהיה ממונה על קידום האוריינות הדיגיטלית בתחומים אלו בקרב גורמי הממשלה וידון עם בעלי עניין מקומיים ובינלאומיים מתחומי התעשייה, הממשל והאקדמיה.

1154 ראו למשל עמדתם של Eric Brown and Dóra Piroška, *Governing Fintech and Fintech as Governance: The Regulatory Sandbox, Riskwashing, and Disruptive Social Classification*, NEW POLITICAL ECONOMY (2021). פלאטר-שנער (לעיל ה"ש 529) נדרשה גם לצורך בסקיפיות כלפי המשחמשים במהלך תקופת הניסוי.

1155 Katerina Yordanova, *The Shifting Sands of Regulatory Sandboxes For AI*, CITIP BLOG (18.7.201)

1156 לרגולציה באמצעות מכרזים ראו לעיל, ה"ש 1011.

1157 Deirdre K. Mulligan and Kenneth A. Bamberger, *Procurement as Policy: Administrative Process For Machine Learning*, 34 BERKELEY TECH. L.J. 781 (2019)

1158 השוו יהודה שופמן (יו"ר), הצוות לבחינת החקיקה בתחום מאגרי המידע: דין וחשבון 30 (2007); ס' 30 בהצעת חוק הגנת הפרטיות-2019, ודברי ההסבר לו. ארידור הרשקוביץ ושוורץ אלטשולר, לעיל ה"ש 403, בעמ' 116; ס' (5) לתזכיר חוק הגנת הסייבר.

11.3.2. תקומו הארגוני

של הרגולטור -

יחידה בתוך רשות

הרגולציה

אנו סבורים כי בעת הזאת יהיה נכון שרגולטור הבינה המלאכותית יהיה יחידה בתוך רשות האסדרה. אנחנו סבורים שהוא אינו צריך להשתייך למשרד המדע, שאיננו משרד רגולטורי באופיו ובפעולתו לקידום התעשייה מעדיף אפריורית קידום חדשנות מעקרונות כגון "האדם במרכז" ו"הדמוקרטיה במרכז".

רשות האסדרה הוקמה על פי חוק עקרונות האסדרה¹¹⁵⁹ והחלה את פעילותה בשנת 2022. זו יחידה במשרד ראש הממשלה ותפקידיה הם לספק לכל משרדי הממשלה שירותי ייעוץ והדרכה במהלך תהליכי גיבוש הרגולציה, לתת חוות דעת מקצועיות על תהליך המחקר המקדים שעשה המשרד הרלוונטי ולוודא שהמשרד מודע לכלל ההשלכות של החלטותיו. כל עוד אסדרת הבינה המלאכותית היא תחום מתפתח, תפקידים אלו של רשות האסדרה הולמים אותה מאוד.

לרשות תפקידים נוספים: לקדם למידה ועבודת מטה ומחקר, להרחיב את גבולות הדמיון הרגולטורי ולהציע פתרונות יצירתיים. לרשות אמורה אף להיות היכרות רוחב עם ארגו הכלים הרגולטורי ועם התכונות של כלים בו; למשל, היא אמורה לדעת מה היתרונות של רישיונות ומתי חובת דיווח עדיפה מחובת גילוי; מתי צריך רישיון קשיח ומתי יש להעדיף מנגנון של אישור בשתיקה. הרשות אמורה לקדם שיטתיות ואחידות בבחינת השפעות הרגולציה, כך שאפשר יהיה להשוות תהליכים בין מגזרים. היא אף אמורה לחתור לאמינות רבה יותר בשימוש בנתונים של הערכת תהליכים רגולטוריים.

יתרון נוסף הוא שגם הרשות היא סוג של "רגולטור של רגולטורים", אך היא רק מייעצת ולא מחייבת. במובן זה היא מתאימה יותר ממבנה מערך הסייבר, שמבוסס על הסדר כופה. הרשות ממילא עובדת מול כלל המשרדים והרשויות העוסקים באסדרה, ולכן הדיאגנטיקה הארגונית שלה מתאים. היא שייכת למשרד ראש הממשלה, מה שייכול לסייע לה לפתור בעיות של כפילויות בין רגולטורים ומשרדים בנוגע לחלוקת סמכויות בתחומים שונים של אסדרת בינה מלאכותית.

היא אמורה לאזן באופן מיטבי בין אינטרסים ציבוריים (כגון אקלים) לבין אינטרסים כלכליים.

החיסרון שבמיקום רגולטור הבינה המלאכותית בתוך הרשות לרגולציה הוא שאף על פי שאחד מתפקידיו הוא לתת ייעוץ משפטי רוחבי, הוא לא יהיה שייך למשרד המשפטים. חיסרון זה אינו משמעותי, שכן עבודת הרשות לרגולציה כוללת מטבע הדברים גם מתן ייעוץ משפטי רוחבי. כמו כן, השנים האחרונות מלמדות שמשרד המשפטים לא השכיל לקבל עליו אחריות של ממש בתחום הייעוץ הרוחבי במשפט ובטכנולוגיה, מה שגרם לעיכובים רבים בחקיקה ולחוסר ודאות ושקיפות.

ההתמודדות עם אתגרי הגמישות, עם הסביבה הרגולטורית העמוסה ועם הטכנולוגיה מצריכים הקצאת משאבים ראויה, שתאפשר ליצור תקנים הן למומחים טכניים הן למומחים במשפט ובמדיניות. אנו מציעים כי גם בהיעדר חוק בינה מלאכותית תתקצב רשות האסדרה בהחלטת ממשלה כדי שתוכל להקים את היחידה של רגולטור הבינה המלאכותית.

כפי שציינו לעיל, חקיקה לא נוצרת בן לילה, וגם אם טרם הבשילה העת לחוקק בישראל חוק בינה מלאכותית מן הראוי לעקוב אחר התפתחויות גלובליות רלוונטיות, לייצר ידע ארגוני ולגבש עמדות בעניין ההסדר הרצוי. בחלק זה נציע שלוש הצעות בנוגע לתקופת הביניים, בטרם תגובש חקיקה בישראל.

11.4 פיתוח אסדרת בינה מלאכותית בתקופת ביניים בסיס לאסדרה עתידית

11.4.1. יצירת "אקוסיסטם" אנו סבורים כי הקמת יחידה שתעסוק כבר כעת באסדרה של בינה מלאכותית תוכל לפרסם בשלב הביניים הנחיות וגילויי דעת שישמשו תשתית לאסדרה עתידית. כמתואר בתרשים שלהלן, אנו סבורים שראוי ליצור אקוסיסטם קדם-אסדרתי שבמסגרתו יפרסמו הרגולטור הייעודי, רגולטורים ענפיים ורגולטורים משלימים (כגון הרשות לפרטיות והרשות לתחרות) הנחיות

וגילויי דעת; והם בתורם ישמשו את התעשייה ושחקנים קונקרטיים בה, וגם את בתי המשפט. מסמכי הנחיות אלו יוכלו לשמש את התעשייה כדי לקבוע סטנדרטים של אסדרה עצמית, מתוך הנחה שרגולציה עתידית תהיה דומה להנחיות. הם יוכלו לשמש גם את חברות התקינה כדי ליצור את מסגרות התקינה.

סביב האקוסיסטם הקדם-אסדרתי יתפתחו תהליכים של חידוד דפוסי הבקרה (auditing) בהקשרים שונים, ייערכו תהליכי חשיבה בעניין סטנדרטים של זהירות ראויה וייעשה מאמץ סביר לקדם אחריות בפיתוח וביישום של מערכות לומדות; ייבנו ארגוני חול וייערכו ניסויים רגולטוריים שונים; ויקודם מאגר של מומחים שיוכלו לשמש את התעשייה, את הרגולטורים ואת בתי המשפט בהתמודדות עם סוגיות חדשות ומורכבות. במקביל, התעשייה והשחקנים יוכלו לתת משוב לגופים הרגולטוריים השונים וכך לטייב את ההנחיות. נוסף על כך, בהיעדר אסדרה, בתי משפט יוכלו להשתמש בהנחיות כהשראה לפרשנות בסכסוכים המובאים לפניהם. פסקי הדין יוכלו גם הם להעשיר את גוף הידע ולאפשר לשחקנים וגם לרגולטורים לטייב ולשייף את הנחיותיהם לקראת גיבוש אסדרה. בתרשים להלן מתואר האקוסיסטם הקדם-אסדרתי.

6 תרשים

האקוסיסטם הקדם-אסדרתי המוצע



11.4.2. קידום תהליכי בתעשיות רבות תקני בטיחות הם כלים מרכזיים להערכת סיכונים במוצרים חדשים – מזון, תרופות, מטוסים מסחריים, מכוניות וכיו"ב.

תקינה (סטנדרטיזציה) בתחום הבינה המלאכותית מאפשרת יכולת תאימות תפעולית (inter-operability) של מערכות ורכיבים, אימות וניידות של מודלים ונתונים, וממשק יעיל והעברה של טכנולוגיות בין פלטפורמות, סביבות ויישומים שונים. לכן יש לה פוטנציאל להפחית עלויות פיתוח ויישום של מערכות ולקדם ולטפח שיתוף פעולה בין בעלי עניין. נוסף על כך, תהליכי תקינה יכולים לסייע לארגונים לפתח אסטרטגיות פיקוח על פרישה אחראית של בינה מלאכותית וכך ליצור את מסגרת הערכים הסובבת כל הטמעה של טכנולוגיה חדשה, ביחוד במרחב שאין בו עדיין רגולציה. המכון הלאומי האמריקני לתקנים וטכנולוגיה (NIST), ארגון התקינה הבינלאומי ISO, הוועדה הבינלאומית לאלקטרו-טכניקה (IEC) והוועדה האירופית לתקינה אלקטרו-טכנית (CENELEC) נעשו בשנים האחרונות שחקני מפתח בפיתוח תקנים בתחום הבינה המלאכותית.

בהסתמך על ניסיון העבר בתעשיות אחרות, יש להניח שכללי תקינה ימלאו תפקיד מכריע ביצירת סביבת הכללים. כללים אלו יוכלו להפוך אחר כך לעקרונות מנחים ביצירת רגולציה – הן בהיבט של יצירת טרמינולוגיה,¹¹⁶⁰ הן בהיבט של מתודולוגיות הערכה וניהול סיכונים,¹¹⁶¹ בקרות אבטחה והגנת פרטיות, הן בהיבט של הטמעת ערכים ועקרונות חברתיים. למשל, כאשר NIST פיתחה את מסגרת התקינה של ניהול סיכונים בינה מלאכותית היא הציעה, מלבד מסגרות לניהול סיכונים ומתודולוגיות הערכה, גם כללי שקיפות והסברות.¹¹⁶² יתר על כן, למוסדות תקינה יכולה להיות גם תרומה של ממש ליצירת תהליכים רגולטוריים ולהטמעתם. כך קרה בתהליך חקיקת חוק הבינה המלאכותית של האיחוד האירופי, כשהאיחוד נעזר ביכולותיו של מכון התקינה האירופי CENELEC,

1160 ראו למשל בתקן ISO לבינה מלאכותית: ISO/IEC JTC 1/SC 42 (2023)

1161 ראו למשל: *Artificial Intelligence Risk Management Framework (AI RMF 1.0)*, NIST (2023). נוסף על כך, בתקן AI NIST SP 800-53, המיועד לאבטח מידע, יש בקרו אבטחה ופרטיות שחוכננו במיוחד למערכות AI.

1162 ראו גם את התקינה של ISO/IEC 20547, AI ISO, המתמקדת בהגדרת הדרישות למערכות AI אמין ונדרשת להיבטים כגון יכולת הסבר, הוגנות וחוסן; וכן את תקן ISO/IEC 23894, שמשמש מסגרת פיקוח על בינה מלאכותית מתוך התייחסות לנושאים כגון אחריות, שקיפות ופיקוח אנושי.

שהוא עצמו עבד בשיתוף פעולה הדוק עם ISO ו-IEC כדי לנסח את החקיקה. מוסדות האיחוד קבעו מראש כי מכון התקינה האירופי ייצור תקינה שתכליתיה הטמעה הרמונית של העקרונות שנקבעו בחקיקה.¹¹⁶³

אנו מציעים לתרגם כבר כעת מסמכי תקינה מרכזיים בתחום הבינה המלאכותית; לבצע התאמה שלהם לתנאים הכלכליים, המשפטיים והרגולטוריים במדינת ישראל; ולהפיץ אותם בקרב השחקנים השונים. כל זה כדי לאפשר דיון ציבורי בעקרונות שהם מבטאים, דיון טכנולוגי בישימות שלהם ודיון רגולטורי בשאלה עד כמה יהיה אפשר לאכוף אותם אם יתורגמו לחקיקה. גם במדינות אחרות בעולם יש תקינה, ואנו סבורים שזהו השלב הראשון וההכרחי ביצירת מצע לאסדרה עתידית. גם אם אי-אפשר לאכוף אותה, קיומה מאפשר לתעשייה להבין את כיווני החשיבה של הרגולטור ולהתארגן בהתאם, לבעלי עניין לעורר סביבה דיון ציבורי ולרגולטור לבצע תהליכי למידה ואימון.

11.4.3. "ארגון חול" למושג "ארגון חול רגולטורי"¹¹⁶⁴ כמה משמעויות רגולטוריים
 המרכזית והמקורית שבהן היא יצירת סביבה מבוקרת שבה חברות יכולות לבדוק את המוצרים או את השירותים שלהן תחת פיקוח רגולטורי, אך רשאים לחרוג מהרגולציה או להגמיש אותה. לעיתים מדובר במשטר משפטי תחום בזמן ובהיקף, בניהול גורם מוסמך, שמשותפו זוכים לפטור מהוראות מסוימות בדין או להחרגות מסוימות והם יכולים ליישם טכנולוגיה או מודל עסקי בניגוד להוראות אלו. למשל, תקנות הבינה המלאכותית האירופיות פוטרות את החברות במסגרת

¹¹⁶³ ראו את תיאור התהליך באתר CENELEC; וכן את חוכנית העבודה של JTC21 מחודש מרץ 2022.

¹¹⁶⁴ ראו ההגדרה שהציעו Vladislav O. Makarov and Marina L. Davydova, *On the Concept of Regulatory Sandboxes, in "SMART TECHNOLOGIES" FOR SOCIETY, STATE AND ECONOMY 1014* (Elena G. Popkova and Bruno S. Sergi eds., 2021). ראו Dirk A. Zetsche et al., *Regulating a Revolution: From Regulatory Sandboxes to Smart Regulation*, 23 (1) FORDHAM JOURNAL OF CORPORATE & FINANCIAL Artificial Intelligence Act and Regulatory Law 31 (2017-2018); ראו גם Law 31 (2017-2018) Sandboxes (EU Parliament Briefing), THINK TANK (17.6.2022)

ארגז חול רגולטורי¹¹⁶⁵ מהחובה לציית לעקרונות צמידות המטרה שבתקנות הגנת המידע¹¹⁶⁶ ומאפשרות להן להשתמש בנתונים שנאספו למטרות אחרות כדי לפתח בינה מלאכותית חדשנית.¹¹⁶⁷ בפעמים אחרות מדובר בליוי רגולטורי שבו מוסכם כי אם המפוקח פועל בתום לב לא יאכפו לגביו הסדרים מסוימים.

מטרת ארגז החול היא לאפשר היערכות, איסוף נתונים ומדידה של ביצועים, בטיחות, חוסן ופגיעות של מודלים ומוצרים, בדיקת דיוק והטיות שלהם, וכן את עמידתם בעקרונות אתיים אחרים, לעיתים עקב קבלת משוב על המוצרים מצרכנים, והכול לפני יישום נרחב של מודלים ומוצרים וללא הסיכון של פגיעה בצרכנים או הפרת חוק ורגולציה. ארגזי חול נועדו לעודד חדשנות במובן זה שהם מספקים מרחב בטוח לניסויים – דרך זו מאפשרת לארגונים, לחברות ולעסקים לנווט את דרכם בין דרישות רגולטוריות ולזהות בעיות פוטנציאליות לפני יישום רחב יותר של מוצריהם. לכן הם נעשו נפוצים גם בהקשרים של מוצרי סייבר וממשל נתונים.¹¹⁶⁸ יתרונות נוספים של ארגזי חול: שילוב כוחות בין המאסדר ובין המפוקחים, במקום המצב הרגיל שהוא היררכי במהותו; האפשרות לשלב בתהליך הניסוי בעלי עניין שונים, כגון צרכנים ולקוחות, אזרחים מקבוצות שונות ושחקנים חדשים, וליצור מערכת המבוססת על אחריות ואמון ולא רק על פיקוח והגנה.

בישראל הוצעו ארגזי חול רגולטוריים בתחום הפיננסי¹¹⁶⁹ ובתחום האנרגיה הירוקה,¹¹⁷⁰ למשל, ונחקקו בפועל בתחום הטיס¹¹⁷¹ ובתחומים נוספים. בתחום

1165 ס' (6)-(5) להצעת תקנות הבינה המלאכותית האירופיות, לעיל ה"ש 53.

1166 ראו ס' (4) 6 של התקנות הכלליות בדבר הגנת מידע (GDPR).

1167 ס' (1)(a) 54 להצעת תקנות הבינה המלאכותית האירופיות, לעיל ה"ש 53.

1168 ראו למשל את ארגז החול הרגולטורי לניסויים בתחום הנתונים בסינגפור. *Data Regulatory Sandbox*, Infocomm Media Development Authority

1169 לאחרונה פורסמה הצעת חוק להקמת ארגז חול רגולטורי בתחום הפינטק (FinTech) בישראל. ראו הצעת חוק לעידוד פיתוח טכנולוגיה בתחום הפיננסי בישראל, החשפ"א-2021, ה"ח הממשלה 204; רות פלאטו-שנער "פינטק בישראל: תוכנית 'ארגז החול' הרגולטורי" *ICoN-S-IL BL06* (30.6.2021).

1170 ראו אילנה קוריאלי "חדש: מועצה שתפתור חסמי רגולציה כדי לקדם סטארט-אפים ירוקים" *ynet* (10.11.2021).

1171 ראו ס' 16 לתקנות הטיס (נוהלי תיעוד כלי טיס וחלקיהם), תשל"ז-1977.

התחבורה נקבע הסדר מפורט בעניין הפעלה ניסיונית של רכב אוטונומי ובסעיף המטרה שלו יש הסבר ברור של תפיסת ארגו החוק הרגולטורי: "מטרתו של סימן זה קביעת הסדרים שיאפשרו הפעלה של רכב עצמאי, בדרך, למטרת ביצוע ניסוי, בלא נהג, בשמירה על הבטיחות ושימוש במגוון טכנולוגיות, כדי להביא לגיבוש תשתית ידע לגבי בטיחותו של הרכב העצמאי, יכולתו להשתלב בין עוברי הדרך ולתת שירות לנוסעים והשפעתו על התנועה בדרך, לאפשר הנגשה של הידע האמור לציבור ולבסס את אמון הציבור בו".¹¹⁷²

לשימוש בארגזי חול יש תכלית נוספת: לסייע בפיתוח אסדרה חדשה במקום שאין אסדרה קיימת, תהליך המכונה *regulatory prototyping*. בהיעדר מסגרות רגולטוריות ברורות בתחומי שיפוט מסוימים יכולים ארגונים, חברות ועסקים לזהות פערי אסדרה ולא לדעת אם על המוצרים שלהם חלות תקנות, אילו תקנות וכיצד לעמוד בהן. במקרה כזה הם יפנו לרגולטור ויציעו לו להקים ארגז חול שתכליתו לבדוק כיצד צריך לעצב רגולציה עתידית. ארגזי חול יכולים להיות מועילים גם במקום שבו תהליכים של הערכת השפעות וניהול סיכונים אינם מועילים. למשל, בתחום הבינה המלאכותית היוצרת המבוססת על מודלים בסיסיים שעליהם מורכבים מוצרים ויישומים, אין טעם לדרוש ממי שמשווק יישומים כאלה לבצע הערכת סיכונים של המודל הבסיסי שהוא לא יצר. במקרה כזה ייתכן שיהיה אפשר להפיק תועלת מארגז חול רגולטורי שידגים את החששות דווקא בהקשרי היישומים וייצור מסגרת לאסדרה עתידית. יש שמרחיבים את תחולת המושג אפילו יותר וכוללים בו שולחנות עבודה רכי-מגזריים שדנים בהם ברעיונות חדשים לרגולציה עתידית.

יש המפרשים את המושג "ארגז חול רגולטורי" ככלי שתכליתו להעצים את יכולות גופי האסדרה בתחומי ההנחיה, הפיקוח והאכיפה. פרשנות כזאת אמורה לאפשר לרגולטורים להבין את מידת התאימות של הרגולציה לטכנולוגיה, לאפשר להם למידה ולכן גם למזער את הסיכון לפיגור רגולטורי.

1172 ראו סימן 3ג (הפעלה ניסיונית של רכב עצמאי בלא נהג) פקודת התעבורה (תיקון 130) החשפ"ב 2022, בסעיף 14כד.

תקנות הבינה המלאכותית האירופיות מגדירות שתי מטרות מרכזיות להקמת ארגוני חול רגולטוריים: (1) לסייע לרשויות באמצעות מתן קווים מנחים בעניין הוראות התקנות לספקי מערכות בינה מלאכותית; (2) לסייע ללמידה רגולטורית בסביבה מבוקרת.¹¹⁷³ כלומר התפיסה האירופית רואה בארגוני החול כלי לסיוע בהטמעה של חוקי הבינה המלאכותית ובפרשנותה. לפיכך ארגוני החול הם החוליה המקשרת בין חקיקה חדשה לבין התקינה שתסייע ביישום שלה; תכליתם לסייע בקידומה של אכיפה אפקטיבית, הרמוניזציה וודאות משפטית.¹¹⁷⁴ לכן התקנות מורות על הקמה של ארגוני חול רגולטוריים אחד לכל הפחות בכל מדינה חברה; ועליו להיות פעיל לכל המאוחר ביום כניסת התקנות לתוקף.¹¹⁷⁵ לנוכח החשש שהובע במהלך החקיקה כי עמימות בעניין המתווה של ארגוני חול עלולה להביא לידי פרגמנטציה רגולטורית באיחוד,¹¹⁷⁶ נקבע בחקיקה הסופית כי הנציבות תאמץ חקיקה המגדירה את המתווים השונים להקמתם, פיתוחם והטמעתם של ארגוני חול רגולטוריים וכן את דרכי הפיקוח עליהם.¹¹⁷⁷

בחודש מרץ 2021, טרם פרסום התקנות המוצעות, החלה נורווגיה להפעיל ארגוני חול רגולטוריים כדי לסייע לחברות המפתחות בינה מלאכותית להתמודד עם הסביבה הרגולטורית המורכבת, ובפרט עם הקושי לציית להוראות התקנות הכלליות בדבר הגנת מידע (GDPR). שלא כמו המודל של התקנות האירופיות, ארגוני החול הנורווגיים אינו מגמיש את עקרון צמידות המטרה אלא נועד לאפשר לרגולטור לסייע ליזמים לפתח מערכות בינה מלאכותית המציינות ל-GDPR כדי להקנות מידה מסוימת של ודאות משפטית.¹¹⁷⁸

1173 ס' 53(1d) להצעה המחוקנת, לעיל ה"ש 57.

1174 פס' 71 למבוא להצעת תקנות הבינה המלאכותית האירופיות, לעיל ה"ש 53.

1175 ס' 53(1) להצעה המחוקנת, לעיל ה"ש 57.

1176 Sofia Ranchordas, *Experimental Regulations for AI: Sandboxes for Morals and Mores* (12.7.2021)

1177 ס' 53a להצעה המחוקנת, לעיל ה"ש 57.

1178 Dan McCarthy, *To Regulate AI, Try Playing in a Sandbox*, TECH BREW (26.5.2021); *Sandbox for Responsible Artificial Intelligence*, DATATILSYNET

השימוש בכלי של ארגז חול רגולטורי מעורר גם חששות. ראשית, הוא דורש הקצאה של משאבים לתשתיות נחוצות. הוא דורש גם כוח אדם מתאים, שכן נדרשות מומחיות טכנולוגית במנגנוני הערכת מודלים ובמתודולוגיות של ממשל נתונים והערכתם; וכן יכולת פרשנות ויכולת הסבר ומומחיות משפטית בשאלות של חבויות, קביעת אחריות משפטית והיבטי רגולציה נוספים. כל זאת, נוסף על מומחיות בענפים ספציפיים כגון בריאות, פיננסים ותחבורה. לכן בתקנות האירופיות נקבע כי על הרשויות המקימות ארגזי חול רגולטוריים להבטיח הקצאת נאותה של משאבים כדי שיוכלו לפעול בהתאם להוראות הדין.¹¹⁷⁹

שנית, יש חשש שארגזי חול רגולטוריים יובילו להתחמקות מרגולציה רוחבית, כמו הגנת פרטיות ואבטחת נתונים. לכן נקבע בתקנות האירופיות כי שימוש לצורכי ניסוי במידע שנאסף למטרות אחרות יכול להיעשות רק בכמה תנאים: מערכת הבינה המלאכותית שמבקשים לנסות נועדה להגן על אינטרסים ציבוריים חיוניים בתחום הביטחון הציבורי (כגון זיהוי עבירות פליליות ומניעתן), בתחום הבריאות (לרבות מניעה של מחלות, בקרה עליהן וטיפול בהן) או בתחום ההגנה על איכות הסביבה; הנתונים הגולמיים נחוצים ואי־אפשר להשתמש במידע סינתטי או במידע שעבר התממה; יש מנגנונים לניטור סכנות לזכויות יסוד של מושאי מידע שמתעוררות במהלך הניסוי; המידע לא מועבר או מוגש לצד שלישי כלשהו; עיבוד המידע במסגרת ארגז החול הרגולטורי אינו משפיע על החלטות שיתקבלו בעניין מושאי המידע; כל מידע אישי שמעובד במסגרת ארגז החול הרגולטורי נשמר בסביבה מאובטחת, נפרדת ומבוקרת שהגישה אליה מוקנית למורשים בלבד; כל פעולות עיבוד המידע האישי מתועדות; כל המידע האישי המשמש בארגז החול הרגולטורי מבוער עם סיום הפעילות; תוצאות הניסוי ותיעוד מלא של הבדיקות, האימון, התהליכים וההיגיון שהנחה אותם נשמרים.¹¹⁸⁰

שלישית, יש חששות מתחום המשפט המינהלי, משום שארגזי חול פועלים מטיבם על בסיס פרטני ויוצרים קשר ישיר וייחודי בין רגולטור למפוקח,

1179 ס' 53(1c) להצעה המתוקנת, לעיל ה"ש 57.

1180 ס' (j)-(b)(1)54 להצעת תקנות הבינה המלאכותית האירופיות, לעיל ה"ש 53.

ששונה מן הפרוצדורה הרגילה של יישום ואכיפה של הרגולציה. מאפיין זה עלול אף לעורר חשש מאפליה ופגיעה בשוויון. לכן בתקנות האירופיות מוצעות הוראות מיוחדות לעידוד השתתפותם של ספקים קטנים וחברות הזנק בארגזי חול רגולטוריים.¹¹⁸¹ נוסף על כך, אם אין הקפדה על הליכים שקופים ומוגדרים בכירור, ועל מנגנונים לשיתוף בעלי עניין, הפגיעה של ארגזי החול באמון, באמינות ובלגיטימיות של הרגולטור עלולה לעלות על התועלת שבו. יש גם צורך לייצר הערכה ולמידה תקופתיות כדי לבחון את היעילות וההשפעה של ארגזי חול וכדי למנוע חששות מניצולם לרעה או משימוש בהם כמנגנון התחמקות מעול הרגולציה.¹¹⁸²

אנו סבורים שהשימוש בארגזי חול רגולטוריים שתכליתם לייצר החרגה מכללים או להציע ליווי שיש בו הסכמה להימנע מאכיפה של כללים יכול להיות אמצעי יעיל בתקופת הביניים בישראל. בשנים האחרונות נשמעו בישראל קולות בדבר הצורך לחוקק לארגזי חול רגולטוריים חוק מסגרת שיקבע עקרונות כלליים בעניין היקף ההחרגה מאסדרה קיימת (למשל אם אפשר לחרוג גם מחקיקה ראשית או רק מתקנות), מנגנוני אכיפה במקרים של הפרת הסכמי החרגה (למשל אם די בדיווח או שיש מקום לאכיפה מינהלית או פלילית), פרוצדורה (למשל דרכי הגשת בקשות, מגבלות זמן על שימוש בארגזי חול ונקודות יציאה). ואולם קבועי הזמן לחקיקה במדינת ישראל מלמדים כי הסיכוי לחקיקה כזאת בעתיד הקרוב נמוך ועל כן ספק אם אפשר לסמוך עליה, בייחוד בתקופת הביניים הקרובה. לכן אנו מציעים שאסדרת ארגזי החול תתבצע ברמה הענפית, אבל תיקבע בהנחיות חובת דיווח של הרגולטור המפעיל ארגז חול ליחידה העוסקת באסדרת בינה מלאכותית. המטרה – לאפשר למידת עמיתים רוחבית ולהעשיר את האקוסיסטם הקדם-אסדרתי. אנו סבורים שראוי שרגולטורים ענפיים המקימים ארגזי חול רגולטוריים יעשו זאת בקפידה ויעמדו בתנאי המשפט המינהלי. הטעם לכך הוא שמלכתחילה התרבות הרגולטורית בישראל ידועה

1181 שם, ס' 55.

1182 להרחבה ראו גם גיא מור וענת גרימלנד "רגולציה ניסיונית: כלים לשיפור איכות הרגולציה באמצעות תחולה מוגבלת" (נייר מדיניות, מאי 2019).

בגמישותה ובנטייתה לאפשר "חוזים רגולטוריים" עם מפוקחים, באופן שיכול ליצור העדפה בלתי ראויה של שחקנים מסוימים.

אם תהיה בעתיד בישראל חקיקה רוחבית בעניין בינה מלאכותית יהיה מקום לעסוק בשאלה אם יש צורך בארגזי חול כדי לקדם את תהליכי ההטמעה שלה, כפי שמוסדר בתקנות האירופיות. בשלב הביניים שבו אנו נמצאים דבר זה אינו נצרך.

נוסף על כך אנו מציעים ליזום "ארגזי חול מוסדיים" שתכליתם לנסות מוצרים מבוססי בינה מלאכותית בהקשרים ספציפיים, דוגמת בתי משפט או בתי חולים. מוסדות אלו יכולים לפרסם קולות קוראים להגשת בקשות להתנסות במוצרים או בשירותים מבוססי בינה מלאכותית הנוגעים לתחומי העיסוק שלהם. לשם כך עליהם ליצור מסגרות שיבחנו את מטרותיהן של המערכות, את תהליכי הפיקוח עליהן, את התועלת בהן ואת הפגמים המתגלים בהן, ולפרסם דוחות הערכה לעיון הציבור. במדינת יוטה מפעילה בשנים האחרונות יחידת החדשנות של בית המשפט העליון המדינתי ארגז חול מוסדי כזה. הצלחתו הייתה גדולה כל כך עד שבפברואר 2023 נאלצה היחידה להפסיק לקבל בקשות בשל עומס יתר.¹¹⁸³

לסיכום חלק זה, התרשים שלהלן מתאר את תהליך האסדרה המוצע בעת הזאת: שרטוט מפת דרכים וטיפולוגיה של הנושאים הטעונים תשומת לב; תרגום מסמכי תקינה ושילובם במסגרת הקדם-אסדרתית; העשרת המסגרת הזאת באמצעות ארגזי חול רגולטוריים; עבודה על עיצוב אסדרה וחקיקה רלוונטיות.

תרשים 7

תהליך האסדרה המוצע של בינה מלאכותית



פרק שנים עשר

סיכום

—

היא מלאכותית אבל היא "בינה", היא ממוחשבת אבל היא "ראייה", היא שייכת למכונה אבל היא "למידה", היא רשת אבל היא "נוירונית". אנחנו דורשים שהיא "תתאמן" היטב, שיהיה לנו "אמון" בה ושלא יהיו לה "הטיות". וכמובן – אנחנו מבקשים שהיא "תסביר" את עצמה כך שנבין.

בינה מלאכותית היא אגד של טכנולוגיות מסעירות ומרתקות שמתפתחות בשנים האחרונות בקצב מסחרר. מינוף יתרונותיהן והתמודדות עם האתגרים המשפטיים שהן מעמידות לפנינו מחייבים ערנות מצד ממשלות ורגולטורים, האקדמיה, החברה האזרחית וגם התעשייה. לא פחות מכך הן מחייבות ערנות מצד מי שיהיו מושאי ההחלטות האלגוריתמיות, המשתמשים במוצרים מבוססי בינה מלאכותית

וגם המושאים של שינויים אפשריים ברמה האישית והחברתית עקב שימושים אלו. לכולם יש מקום סביב השולחן, וחשיבות מקומם תיגזר מן הערנות שלהם, מנכונותם להרכיב את האוריינות הדיגיטלית שלהם ומהמאמץ שישקיעו בכל אלה בחשיבה ובמעשה.

מטרתו של ספר זה הייתה לספק לכלל השחקנים את האוריינות הדיגיטלית הדרושה להם כדי לקבל בשנים הקרובות החלטות בנוגע לאסדרה של בינה מלאכותית. הדעת נותנת שנודקק לעדכון מתמשך של הידע, ההבנה ויכולת הפרשנות שלנו בתחום זה. אבל לנוכח הניסיון המצטבר שרכשנו במהפכות הטכנולוגיות שהתרחשו בעשרים וחמש השנים האחרונות, ברור לנו כיום שהתפתחות ללא כל אסדרה איננה הגיונית.

האסדרה המשפטית של מערכות בינה מלאכותית ברחבי העולם עורונה בחיתוליה. במחקר זה נסקרו הצעות חוק שונות, אך רובן יבשילו לכלל חקיקה רק בשנים הבאות ועדיין לא ידוע כמובן מה יהיו השלכותיהן בפועל. אך אין פירוש הרבר שיש לשבת בחיבוק ידיים – בזמן ההמתנה להתפתחויות הרגולטוריות מעבר לים יכולים מעצבי המדיניות בישראל להכשיר את הקרקע לחוק בינה מלאכותית ישראלי. אפשר להתחיל לעגן בחקיקה זכויות של מושאי החלטות אוטומטיות ולגבות אותן בסעדים ובהתאמות פרוצדורליות. אפשר לנסות לגבש עמדות בעניין איסורים קטגוריים על פיתוח מערכות בינה מלאכותית ברמת סיכון גבוהה ועל שימוש בהן. איסורים זמניים כאלו יקנו למחוקק די זמן לעצב את מודל האסדרה המתאים לישראל. אפשר להקים ארגון חול רגולטוריים, שבמסגרתם ילווה רגולטור מגזרי פיתוח של מערכות בינה מלאכותית בתחום שהוא מאסדר במטרה לגבש מדיניות אסדרה מתאימה.

בעיקר חשוב לגבש צוות עבודה בין-משרדי שיהיה מופקד על ניסוח הצעת החוק ותיקוני חקיקה אחרים בהסתמך על לקחים שהופקו ממקרי הבוחן בארגון החול הרגולטוריים, על עקרונות אתיים מוסכמים מראש ועל דוגמאות מהדין הזר (ואף על יישומם בפועל). תהליך זה יכול להניב מודל רגולטורי בשל בפרק זמן סביר.

בשנים הקרובות ילך ויגבר השימוש במערכות אלגוריתמיות, ובהן במערכות לומדות ובמערכות בינה מלאכותית, ויחדור לכל תחומי חיינו. שמיכת הטלאים

המשפטית שמאסדרת את התחום הדיגיטלי בישראל קצרה מכדי לכסות את כל הנדרש כדי להבטיח שתושבי המדינה ואזרחיה יהיו מוגנים מפני הסיכונים הגלומים במערכות אלו. קידומה של תשתית משפטית ומוסדית לבניה מלאכותית בטוחה היא צו השעה.

בתרשים שלהלן מסוכמות המלצותינו בנוגע לאסדרה של בינה מלאכותית בישראל. במרכזו, שני העקרונות המרכזיים האמורים להדריך את כלל פעולות האסדרה בתחום: האדם במרכז והדמוקרטיה במרכז. סביבם עקרונות כלליים נוספים: חשיבותה של האוריינות הדיגיטלית של מקבלי ההחלטות; הצורך להתאים את הרגולציה הישראלית למתרחש בעולם – לא להקדים אותה, אך גם לא ליצור פערים רגולטוריים שיהפכו את מדינת ישראל לחצר אחורית במה שנוגע לטיפול בנתונים ולשימושי בינה מלאכותית; הצורך לתת פרשנות חדשה לזכויות אדם קיימות, כגון הזכות לפרטיות וחופש הביטוי, והחשיבות של פיתוח זכויות אדם חדשות, כגון הזכות לקבל שירות מבן אנוש; וכן האתגר הטמון בפיתוח רגולציה חסינת עתיד, כלומר כזאת שאין בה התייחסות לטכנולוגיה ספציפית אלא לעקרונות, והחשיבות של שילוב עקרונות אתיים ועקרונות משפטיים. סביב העקרונות האלה שתי רצועות קונקרטיות: האחת, הנוגעת למסגרת המוסדית של הרגולציה של בינה מלאכותית, מציעה לשלב בין אסדרה ענפית לאסדרה רוחבית ומצביעה על הצורך לראות ברגולציה של בינה מלאכותית חלק אחד בתצרף אסדרתי רחב יותר הכולל חובות לחקיקה עדכנית ולעדכון חקיקה קיימת. השנייה, החיצונית, עניינה הצעות קונקרטיות למיני התערבות לצורך פיקוח על בינה מלאכותית, ובהן ההצעה לפתח מתודות לניהול סיכונים, להחריג (החרגה מוגבלת) איסורים עקרוניים על שימושים מסוימים, הצעה לארגזי חול רגולטוריים, משילות נתונים, הנדסת הוגנות, שקיפות ופיקוח על הקשר בין משימה לתוצאה במערכות של למידה לא מפוקחת.

8 תרשים

סיכום מודל האסדרה המוצע לבינה מלאכותית



תקופות של מהפכות טכנולוגיות גורמות לאנשי רוח לומר, כפי שכתב צ'רלס דיקנס בפתחה הנודעת לרומן "בין שתי ערים", שעלילתו מתרחשת בתקופת המהפכה הצרפתית, "היה זה הטוב בזמנים, היה זה הרע בזמנים". המשורר והמחזאי הגרמני ברטולט ברכט כתב בתחילת המאה העשרים את השורות האלה:

הזמנים החדשים / ברטולט ברכט

זמנים חדשים אינם מתחילים בכת-אחת.

סבי כבר חי בזמן החדש.

נכדי בודאי עוד יחיה בישן.

בשר חדש אוכלים במזלגות ישנים.

לא כלי-הרכב הנוסעים-לבד עשו זאת

לא הטנקים

לא המטוסים מעל לגגותינו עשו זאת

ולא המפציצים.

מן האנטנות החדשות בא הטמטום הישן.

החקמה נמסרת הלאה מפה אל פה.¹¹⁸⁴

ברכט, שכמו מרפרר לדיקנס, ניסה ללמדנו שלא הטכנולוגיה היא שאחראית לתוצאותיה. לא "כלי-הרכב הנוסעים-לבד עשו זאת", לא הטנקים ולא המפציצים. אלה היו הכוחות האנושיים. וכיוון שזמנים חדשים אינם מתחילים בבת אחת חשוב כל כך לחקור אותם, לעמוד על מקורותיהם ולהגות באחריות שצומחת מהם. ולשם כך עלינו להתבונן גם בטכנולוגיה. האם נצעד היישר לגן עדן או שמא היישר לצד השני?

1184 ברטולט ברכט גלות המשוררים 201-202 (מגרמנית: בנימין הרשב, ספרי סימן קריאה/הקיבוץ המאוחד, 1978).



Human, Machine, State

Toward the Regulation of
Artificial Intelligence

Amir Cahane | Tehilla Shwartz Altshuler

Abstract

In recent years, artificial intelligence (AI) systems have increasingly become part of the fabric of daily life. They recommend travel routes and the next song to be played, support medical diagnoses, and, lately, even take an active part in doing homework. Public entities around the world are assimilating algorithmic systems that make, or support, administrative decisions on resource allocation, planning, crime prediction, and protection of the public space—from personal digital assistants to autonomous cars, from robots that carry out simple tasks to monitoring, detection, and forecasting systems.

However, despite their inherent advantages, algorithmic systems may menace human rights and basic freedoms unless there is oversight of their use, development, and deployment. These hazards may emerge at various stages of their development and use—from defining the purpose of the system, via reliance on incomplete, erroneous, corrupted, or biased data, to non-application of post-deployment oversight of the systems' outputs. Furthermore, the more AI develops, the more its systems tend to present new capabilities that were neither intended nor predicted by their developers. Some of these capabilities are inspiring, and therefore called “sparks,” but others have the potential to cause harm, such as carrying out offensive cyber actions, manipulating people by discursive means, and disseminating erroneous and misleading artificial

information. Therefore, the ability to identify these capabilities and limit their associated hazards has become a supremely important challenge.

What is artificial intelligence? What are its advantages? And what fears does it evoke, particularly when used by the authorities? These questions are the focus of this book.

Decision-makers, industry, academia, and civil-society organizations in Israel and around the world have all identified AI as a disruptive technology for which national strategy and regulatory policy must be prepared in advance. The end of the previous decade saw the publication of dozens of ethics documents regarding AI which sought to lay down principles for the development, use, and application of algorithmic systems. The ethical values proposed in these documents are founded on several principles: transparency; fairness; damage prevention and safety; responsibility and accountability; privacy; promoting the common good and prioritizing people; and freedom and autonomy.

Ethical values, however, are not enough. To assure the upholding of human rights and basic freedoms, these principles must be moored in legislation. Indeed, we are now seeing around the world initial signs of legislation that seeks to regulate the use and development of AI technologies. Some follow an across-the-board legislative model applied to AI systems at large, as in the European AI Act; others are regulatory patchworks that target specific sectors and uses of AI, such as laws that aim specifically to cope with algorithmic discrimination in systems that hire and promote workers.

This book offers guiding principles for the creation of rights-oriented AI policy and of a toolbox for AI regulation in Israel.

Main Recommendations

Guidelines for the creation of rights-oriented AI policy

Prioritizing people. The main purpose of developing artificial intelligence and learning systems should be serving the human race—individually and collectively—in a way that enhances its welfare. This principle of human freedom and autonomy connects with the principles of promoting the common good and preventing harm, and underscores all the more the importance of the principle of human-centering, freedom, and autonomy. What this anthropocentric outlook means in practice is that the development, deployment, and use of intelligent systems should be based on an approach that considers the defense of basic rights and civil liberties a paramount principle, rather than simply paying lip service to them while actually giving preference to principles such as “promoting innovation,” pursuing economic interests such as the advancement of high-tech industry, “making Israel a global technological leader,” or even making public-sector processes more efficient.

Prioritizing democracy. AI-based systems have a great deal of potential to infringe on democracy in its broad sense by affecting public discourse and circulating ideas; serving as an instrument of control, surveillance, and policing; and sowing doubt and subverting the very ability to determine reality and distinguish between original and counterfeit and between truth and falsehood. Therefore, the principle of “democracy at the center” should be given much weight, even if this sometimes comes at the expense of technological progress and innovation.

Digital literacy among decision-makers. Digital literacy means the ability to analyze the market and understand the direction in which technology is developing, at least in the near term. For example: Where do the tech giants keep their research and development funds? What patents have they registered in order to secure new technological

developments? It also includes understanding commercial and regulatory possibilities for guiding technological development and, by extension, understanding policymakers' responsibility to influence technological development and not merely observe it from the sidelines. A liminal stratum is needed between understanding technology and making technology policy—a framework for understanding the implications of technological systems, being able to imagine the new possibilities that they offer, and gauging their implications for social ethics and the contours of the judicial method. The frequent lack of framework often results in lacunae in understanding, particularly in matters that have broad implications such as AI.

Reconceptualizing systems and capabilities in the AI field. The concept of artificial intelligence is a highly powerful politico-technological metaphor. Comparing AI systems to the human brain creates proximity and similarity, leading to the social assimilation of the idea that machines operate like the human brain, perform human actions the same way people perform them, and, in fact, compete with people. We recommend that machines' actions and the traits attributed to them be conceptualized in a manner that is not contingent on this comparison.

Developing rights for the objects of AI decisions. Basic rights have to be rethought in two senses. First, the constitutional theory of existing human rights, such as freedom of speech and the right to privacy, needs an injection of new meanings. And second, new digital rights, unneeded in the past, should be created: foremost the rights of individuals who come into contact with algorithm-based machine systems.

Israel must not become a “digital backyard.” In recent years, proposals for AI regulation have been advanced in leading countries and the European Union, in what is presumably the onset of a global trend. Even if there are values-based differences among regulatory mechanisms in different places—whether in the choice of systems defined as dangerous or in the

declared goals of the legislation itself—they also have much in common. Therefore, a situation should not be allowed to develop in which Israel lacks legislation that is well aligned with accepted legislation abroad. Admittedly, a lower regulatory standard than the convention abroad may stimulate innovation, but not necessarily of the desired kind; it may turn Israel into a technological backyard, a place where systems are created whose development, dissemination, or use are banned in other western countries.

Future-proof regulation. Technology is advancing rapidly; in countries such as Israel, where legislative processes are very slow, an especially wide gap is being created. Therefore, regulation should not target any specific technology, as that is a sure prescription for regulatory obsolescence. Instead, an effort should be made to establish guiding principles and general definitions that would invest future enforcement with flexibility.

A hybrid regulatory framework: principles, rights, and legislation. Principle-based frameworks present a matrix of ethical core principles; rights-based frameworks focus on protecting the human rights and liberties of those affected by AI-based technological applications; and legislation-based frameworks make it possible not to rely solely on voluntary regulation predicated on the goodwill of economic actors. The three frameworks do not clash with each other; they should be integrated into a hybrid structure that combines soft regulation (ethical principles) with a risk-management approach manifested in rigid legislative provisions and regulatory rules.

Flexibility in the timing of regulatory intervention. In certain cases, there are clear advantages to early regulation, introduced before products based on a certain technology enter the market. If a technology is perceived as especially dangerous—physically (as in autonomous cars) or ethically (such as artificial creation of content that incites to terrorism)—

then advance regulation makes sense. Even in less extreme cases, prior intervention may be useful in shaping the directions of research and in planning resource investment in development. Furthermore, since investment is relatively small and sunk costs are smaller at this stage, regulatory intervention may meet more limited resistance from interested parties. Contrastingly, in certain cases it may be better to wait and cope with problems when they arise instead of trying to anticipate them.

Sectorial vs. across-the-board regulation. Across-the-board regulation attains goals of governance, creates regulatory harmony, and, accordingly, may enhance public trust and increase regulatory certainty for industry. Sectorial regulation, conversely, allows the use of existing regulators and their enforcement powers, does not require the establishment of new institutional frameworks, facilitates more accurate tailoring of enforcement arrangements and methods to a given industry, enhances regulatory clarity and certainty, and also allows stakeholders in each sector to participate in formulating these arrangements. The problem with sectorial regulation, however, is that it may result in discrepancies among industries, create inconsistent standards, exacerbate gaps among regulators, and leave behind unregulated spheres that fall into the cracks. Therefore, we recommend a combination of across-the-board and sectorial regulation, such as having a general regulator with instructional and advisory powers vis-à-vis sectorial regulators.

The AI regulation toolbox

Regulation of learning systems should be based on an understanding of their “lifecycle.” To create effective regulation of learning systems, all components of their lifecycle should be taken into account. Given that the principles of regulation—such as fairness, privacy, transparency, accountability, and risk management—are manifested in different

contexts in each component of the lifecycle, inattention to these components may result in excessive regulation of certain elements and disregard of others, compromising regulatory effectiveness.

An integrated approach, by contrast, gives consideration to the full set of lifecycle components and their interrelations. The purpose of a learning system and the framing of the problem it addresses, for example, should influence the choice of model used (whether or not to allow the choice of a more opaque model in terms of the ways it makes decisions). Evaluation outcomes at the stage at which the model is built energize risk-assessment processes—which, in turn, require decisions on how the model should be trained, deployed, and protected. The purpose of the model affects the choice of user interfaces: Should a model that dispenses medical advice notify users that it may be wrong? Should users be told that they are connecting with an artificial system and not with a human one?

An important part of understanding the lifecycle of learning systems relates to the need to monitor them after they are deployed in the real world (e.g., when systems are integrated into a product or an interface). This is because a learning system—unlike other products, such as pharmaceuticals—can, by its very nature, change even after it is applied due to the feedback-loop that it receives from its users.

Development of risk-management methodologies. The application of risk-management methodologies to algorithmic systems, despite being vitally needed, is still in its infancy. We propose a risk-management model that requires a double observation in order to assess the dangerousness of a system, first to assess the potential dangerousness of the system as designed and then to assess the strength of the alignment of its task with its outcome, namely, the potential of a system to manifest its dangerousness outside the role its designers intended for it.

This double layer assessment should be the basis for decision-making. To make it practicable, rules of governance and safety need to be formulated, relating *inter alia* to responsible training (whether to train a new model that exhibits early indications of danger—and in what way) and responsible application (whether, when, and how to implement models that may be dangerous); requisite levels of transparency and documentation in cases of models that may pose extreme danger; and the auditing and cyber-security systems that should be applied to them.

Data documentation, data governance, and post-deployment auditing procedures. The complexity of AI systems generally, and of learning systems particularly, imposes special difficulties on policymakers and regulators when it comes to formulating rules of liability and identifying the actual chain of causation that precipitates an infringement of rights, particularly in consideration of variances among sectors and among applications.

What all these have in common is that they are impossible without proper documentation. The basis for every factual examination of a concrete failure in AI systems is data governability and painstaking documentation of working procedures, information sources, labeling, models, coding processes, risk assessment and databases, and detection of discrepancies in each. Good documentation design also serves the interests of entrepreneurs and developers because it allows them to investigate failures and unexpected phenomena after their occurrence, and also to satisfy regulatory documentation obligations that originate elsewhere.

Development of tools to contend with biases and “fairness engineering.”

Although the algorithmic-bias problem defies prevention, particularly when the bias originates in reality itself as embodied in data, it may be identified and mitigated. Several strategies to mitigate algorithmic biases should be pursued, including having in place statistically fair

procedures, diversifying human capital among developers of AI systems, and applying after-the-fact auditing procedures.

Coping with the challenges of algorithmic transparency: We suggest a model that is based on the classic conceptualization of transparency, but includes an alternative tailored to the technological limitations of algorithmic systems that cannot provide explainability for specific outputs. The model is meant for cases in which an output cannot be provided but society needs to see the obligation of transparency upheld so as not to thwart the development and use of certain technologies.

Recommendations for an AI regulatory institution in Israel

We recommend establishing an AI regulation authority in Israel for the purpose of advancing and systematizing the regulation of intelligent systems in the country, including the formulation of across-the-board legislation with an eye on corresponding developments abroad. This authority should make policy on the development, implementation, and use of AI-based products in Israel; provide sectorial regulators with professional guidance in order to assure the consistency of the rules that apply to the development, deployment, and use of these systems; and serve as a residual guiding entity, that is, be responsible for regulating those products in fields where there is no sectorial regulator.

The regulator should also provide state authorities with professional AI guidance. Among other things, it should express its professional opinion about government tender documents pertaining to artificial intelligence in order to make sure that the state procures AI systems that meet Israeli and foreign standards in the field, and should give the private market “soft” incentives to align its work with standards.

The regulator should also be tasked with consultation in matters of AI law and should be empowered to present the Knesset and the courts with its own position in these fields. In addition to the aforementioned roles, the proposed regulator should be a source of knowledge, instruction, and cooperation and should promote digital literacy among government players in these regards. It should also discuss the possible social effects of these technologies with local and international stakeholders in industry, government, and academia.

At the present time—at least for the next few years, until the field stabilizes—the AI regulator should be established in the form of a unit within the Regulatory Authority, because the latter's roles are highly suited to the evolving world of AI regulation. To cope with the challenges of flexibility, the burdened regulatory environment, and the technology, enough resources should be allocated to allow positions to be created for experts in technical fields and also in law and policy. Even in the absence of a comprehensive artificial intelligence law, the Authority should be budgeted by government resolution so that it can set up the AI regulator's unit.

Recommendations for the interim period until AI regulation is introduced

Supplementary legislation and legislative amendments. Right now, even in the absence of a broad and dedicated artificial intelligence law, designated decision-makers and regulators should be obligated to update existing legislation and pass supplementary legislation—mainly to statutes such as the Competition and Consumer Protection Law, the Copyright Law, the Protection of Privacy Law, the Evidence Ordinance, and the Government Procurement Law.

Create a “pre-regulatory ecosystem.” A “pre-regulatory ecosystem” should be created, in which the designated regulator, sectorial regulators, and supplemental regulators (such as the Protection of Privacy Authority and the Competition Authority) would issue guidelines and professional opinions while serving industry, its players, and the courts. These guideline documents would help industry set self-regulation standards on the assumption that future regulation would resemble the guidelines. They would also serve standards-setting companies in constructing their standards frameworks.

Under the influence of this pre-regulatory ecosystem, auditing patterns in various contexts would undergo fine-tuning, thinking about due caution standards would take place, an effort to encourage responsibility in developing and implementing learning systems would be made, sandboxes would be built, various regulatory trials would be undertaken, and a pool of experts who could serve industry, the regulators, and the courts in coping with new and complex issues would take shape. Concurrently, industry and its actors would be able to provide the various regulatory bodies with feedback, thus improving the guidelines. In addition, in the absence of regulation, the courts could use the guidelines for inspiration as they interpret conflicts presented to them. Their rulings would also enrich the body of knowledge and enable players and regulators alike to enhance and polish their guidelines until regulation can be formulated.

Series & Cover Design: Studio Alfabees
Typesetting: Nadav Shtechman Polischuk
Chart Graphics: Dana Berger
Printed by Maor Wallach Print

ISBN: 978-965-519-433-3

No portion of this book may be reproduced, copied, photographed, recorded, translated, stored in a database, broadcast, or transmitted in any form or by any means, electronic, optical, mechanical, or otherwise. Commercial use in any form of the material contained in this book without the express written permission of the publisher is strictly forbidden.

Copyright © 2023 by the Israel Democracy Institute (RA)

Printed in Israel

The Israel Democracy Institute

4 Pinsker St., P.O.B. 4702, Jerusalem 9104602

Tel: (972)-2-5300-800

Website: <http://en.idi.org.il>

Online Book Store: en.idi.org.il/publications

E-mail: orders@idi.org.il

All IDI publications may be downloaded for free, in full or in part, from our website.

The views expressed in this book do not necessarily reflect those of the Israel Democracy Institute.

מערכות בינה מלאכותית ממליצות על נתיבי נסיעה או בוחרות את השיר הבא שיושמע, תומכות באבחון רפואי ומשתתפות באופן פעיל בהכנת שיעורי הבית ובכתיבת מאמרים מדעיים. גופים ציבוריים ברחבי העולם מטמיעים מערכות אלגוריתמיות שמקבלות החלטות מינהליות הנוגעות להקצאת משאבים, לתכנון, לחיזוי פשיעה או להגנה על המרחב הציבורי. מעוזרים אישיים דיגיטליים ועד מכוניות אוטונומיות, מרובוטים המבצעים משימות פשוטות ועד מערכות מעקב, זיהוי וחיזוי – מהפכת הבינה המלאכותית בעיצומה.

ואולם על אף היתרונות הגלומים במערכות אלגוריתמיות, ללא בקרה על השימוש בהן, על פיתוחן ועל פרישתן הן עלולות לסכן זכויות אדם וחירויות יסוד, לגרום נזק ולבצע פעולות סייבר התקפיות, לתמרן אנשים או להפיץ מידע מלאכותי מוטעה ומטעה. האפשרות לזהות יכולות אלו ולהגביל את הסיכונים שהן מביאות איתן נעשתה אפוא אתגר חשוב מאין כמוהו. ספר זה, הראשון מסוגו בעברית, בוחן מהי בינה מלאכותית, מהם יתרונותיה ואילו חששות היא מעוררת בארץ ובעולם. מחבריו קוראים למקבלי ההחלטות, התעשייה, האקדמיה וארגוני החברה האזרחית בישראל להיערך לקראתה באמצעות אסטרטגיה לאומית ומדיניות רגולטורית. כדי

להבטיח שאנשים הם שישלטו במכונות ולא להפך, נדרש פיתוח של עקרונות מנחים ליצירת מדיניות בינה מלאכותית מכוונת זכויות וארגון כלים לאסדרת בינה מלאכותית.

ד"ר תהילה שוורץ אלטשולר היא עמיתה בכירה וראשת התוכנית "דמוקרטיה בעידן המידע" במכון הישראלי לדמוקרטיה. מומחית לאתיקה, מדיניות ומשפט של תקשורת וטכנולוגיה. תחומי המחקר שלה הם פרטיות, אסדרת סייבר, רשתות חברתיות, ריכוזיות בשוק התקשורת ומדיניות תקשורת. חברה במועצת הארכיונים העליונה ובארגוני חברה אזרחית העוסקים באתיקה עיתונאית ובזכויות דיגיטליות ובעלת טור בנושאי טכנולוגיה ורגולציה במגזין "דה מרקר".

עו"ד עמיר כהנא הוא חוקר בתוכנית "דמוקרטיה בעידן המידע" במכון הישראלי לדמוקרטיה; סטודנט לתואר שלישי בפקולטה למשפטים באוניברסיטה העברית בירושלים, בוגר תואר ראשון במשפטים מהמרכז הבינתחומי הרצליה ותואר שני במשפטים מאוניברסיטת קיימברידג'. עמית מחקר בתוכנית לסייבר ומשפט במרכז פדרמן לחקר הסייבר באוניברסיטה העברית בירושלים.

מחיר מומלץ: 82 ש"ח

יולי 2023

מסת"ב: 3-433-519-965-978

www.idi.org.il



0 450001257 8

דאנאקוד 450-1257