

19 באפריל, 2021

לכבוד רונה קייזר מנהלת אגף בריאות דיגיטלית ומחשוב <u>משרד הבריאות</u>	לכבוד עו"ד טליה אגמון, המשנה ליועץ המשפטי <u>משרד הבריאות</u>	לכבוד פרופ' חזי לוי מנכ"ל משרד הבריאות <u>משרד הבריאות</u>
--	--	---

**בדואר אלקטרוני**

שלום רב,

**הנדון: תהליכי התממת מידע רפואי במשרד הבריאות**

אנו החתומים מטה, מומחים מתחום מדעי הנתונים, קריפטוגרפיה והגנת פרטיות, מתכבדים לפנות אליכם בנושא שבנדון, ובבקשה להטמיע שיטות התממה מודרניות של פרטיות דיפרנציאלית, לצורך הגברת ההגנה על הפרטיות במידע ולמען קידום המחקר והטכנולוגיה הרפואית.

**1. רקע ומסגרת נורמטיבית**

לאחרונה נוצר דיון ער על חשיבות המידע הרפואי שמצוי במערכת הבריאות הישראלית, ובפרט מידע רפואי שבידו לתרום תרומה משמעותית לפיתוח חדשנות רפואית, בישראל ובעולם כולו.

בעקבות הדיון בנושא, בחנו את המסגרת הנורמטיבית הפומבית לנושא התממת מידע רפואי בישראל, בכלל זה את [טיוטת תקנות זכויות החולה \(שימוש מחקרי במידע בריאות\), התש"ף – 2019](#) (להלן: "טיוטת התקנות"), [טיוטת הנחיות להתממת מידע בריאות לצרכי שימוש מחקרי מתאריך 6/10/19](#) (להלן: "טיוטת ההנחיות"), ו**מסקנות הועדה ליישום המלצות השימושים המשניים במידע בריאות.**

מעיון בחומרים המוזכרים לעיל, עולות מספר בעיות בסיסיות באשר לאופן בו מוצע לקיים את תהליך ההתממה בנוגע למידע מובנה (טבלאי). נראה שההנחיות הקיימות במשרד הבריאות, כפי שמפורסמות בציבור, מאמצות תהליכי התממה שעלולים להיות פגיעים למתקפות. בפרט, תהליכים אלו עלולים לאפשר שימוש במידע שהותמם על פי ההנחיות, כדי לחשוף זהות חולים ונתונים על מחלתם. יש לזכור כי חשיפת מידע רפואי פוגעת באופן קשה בפרטיותם של נושאי המידע ועשויה לגרום להם לנזק בלתי הפיך במישור האישי, הרפואי והכלכלי. בנוסף, פגיעה כזו תערער את אמון הציבור בהתממת מידע בכלל ובמשרד הבריאות בפרט, וכתוצאה יכולת המשרד לקבל נתונים ולתמוך במחקרים עתידיים תפגע גם היא.

כידוע, פרטיות היא זכות חוקתית בשיטה המשפטית בישראל. הזכות לפרטיות מעוגנת בסעיף 7 לחוק-יסוד: כבוד האדם וחירותו ולפיכך היא בעלת מעמד חוקתי על-חוקי. עוד בטרם הוכרה כזכות חוקתית, עוגנה הזכות לפרטיות בחקיקה ספציפית בחוק הגנת הפרטיות, התשמ"א-1981 ("חוק הגנת הפרטיות"). סעיף 92(9) לחוק הגנת הפרטיות קובע את עקרון צמידות המטרה, לפיו שימוש במידע של אדם שלא למטרה לשמה המידע נמסר, מהווה פגיעה בפרטיות.

הגדרת "מידע" בסעיף 7 לחוק הגנת הפרטיות הוכרה הן על-ידי בתי המשפט והן על-ידי הרשות להגנת הפרטיות כהגדרה רחבה שמכילה גם מידע מזהה וגם מידע שניתן לזיהוי (ע"א 1697/11 א. גוטסמן אדריכלות בע"מ ואח' נ' ורדי (23.1.2013)). עקרון יסוד בחוק הגנת הפרטיות הוא כי במקום שיש הסכמה לשימוש במידע, אין פגיעה בפרטיות.

לכן חשוב להזכיר בהקשר זה, כי אין כיום בדין הישראלי דרישה לקבלת הסכמה לשימוש משני במידע שאינו מזהה או ניתן לזיהוי. על כן, קיימת חשיבות רבה לשימוש בשיטות התממה שאינן מאפשרות זיהוי מחדש של המידע – הן לצורך הגנה על נושאי המידע והן עבור משרד הבריאות בבואו לנהל סיכונים הגנת פרטיות ואבטחת מידע.

**על כן, אנו שמחים להציג בפניכם סקירה קצרה והמלצות בדבר מימוש תהליכי התממה, תוך שמירה על הזכות לפרטיות של נושאי המידע ותוך איזון יצירת הערך של הנתונים לקידום המחקר והטכנולוגיה הרפואית.**

התובנות והמסקנות המובאות במסמך זה מבוססות על מחקר עדכני בתחום טכנולוגיות הגנת פרטיות מידע ותקיפתה. תחום מחקר זה, שהינו תת-תחום של מדעי המחשב, ונושק לקריפטוגרפיה וסטטיסטיקה חישובית, התפתח רבות בשני העשורים האחרונים. **אנו בדעה כי בעידן המידע שבו אנו מצויים, שילוב מומחים מהתחומים המוזכרים לעיל בתהליכי פיתוח מדיניות, רגולציה וקבלת החלטות הוא חיוני למימוש אפקטיבי של טכנולוגיות הגנת הפרטיות.**

באופן טבעי כל שינוי של פרקטיקה מקצועית ורגולטורית אורך זמן ודורש היערכות ותכנון. לכן במסמך זה נציע גם קווים מנחים לתוכנית שמטרה הטמעת תהליכי התממה מיטביים המבוססים על המחקר המדעי העדכני בתחום טכנולוגיות הגנת פרטיות מידע ותקיפתה.

## 2. סיכונים במימוש תהליכי התממה המבוססים על הכללה ו/או K-Anonymity

### 2.1 שיטות "הכללה" אינן מספיק טובות וניתנות לתקיפה:

טיוטת ההנחיות מציגה דרישה רגולטורית לשימוש בשיטת התממה ספציפית מסוג "הכללה" עבור מידע מובנה (טבלאי), לאחר שלב ההסרה של שדות מידע מזהים (PII). בשיטות אלו הרזולוציה של המידע או רמת הדיוק שלו מופחתת בכוונה. למשל, שימוש במיקוד במקום כתובת מדויקת, חודש ושנה במקום תאריך מלא וכו'. אחד הקריטריונים הנפוצים יותר לבחינת אפקטיביות של שיטות התממה המשתמשות בהכללה הוא k-אנונימיות, ואכן זה הקריטריון המוצע בטיטת ההנחיות. קריטריון זה דורש שהמידע הלא קליני המופיע על כל אדם במאגר יהיה זהה למידע הלא קליני המופיע על לפחות k-1 אנשים אחרים במאגר.

שיטת ההכללה וקריטריון זה הוצעו על ידי החוקרת לטניה סוויני לפני יותר מ-20 שנה, והיוו פריצת דרך חשובה שהעלתה את נושא הפרטיות ובפרט את היכולת לזהות מחדש אנשים במאגר מידע שלא הותמם. אכן, שיטה זו נותנת הגנה מסוימת כנגד זיהוי מחדש, אך עם השנים הסתבר שגם שיטת הכללה זו אינה מקנה הגנת פרטיות מספקת. בפרט, כאשר מצליבים מידע משני ממאגרי נתונים שונים, ניתן לזהות מחדש אנשים שמידע עליהם מופיע בשני המאגרים - גם אם כל אחד מהמאגרים הותמם בשיטת ההכללה, וגם אם כל אחד מהמאגרים כשלעצמו מקיים k-אנונימיות. חוקרים הראו גם דרכים אחרות לזיהוי מחדש על ידי הצלבת מידע ממאגר נתונים אחד שהותמם ומקיים k-אנונימיות, עם מידע פומבי תמים על אנשים שמופיעים במאגר. פגיעות זו היא בעייתית במיוחד בימינו, מכיוון שכל אדם מופיע במאגרי מידע רבים בין אם פומביים או פרטיים, בין אם באותו התחום (למשל בריאות) או בתחום אחר לחלוטין (למשל תעסוקה, חינוך או פנקס הבחורים).

נציין שהתקפות אלה אינן תיאורטיות גרידא. למשל, תקיפה המבוססת על עקרונות אלו בוצעה בהצלחה על מידע של פלטפורמת edX, מערכת שמנגישה קורסים מהמוסדות האקדמיים המובילים בעולם

בהובלת אוניברסיטאות הרווארד ו-MIT. לטובת מחקר, הפלטפורמה שחררה מידע אודות לומדים, שהותמם ע"י הכללה ועמד בדרישות k-אנונימיות עם  $k=5$ . עם זאת, חוקרים הצליחו לזהות מחדש מתוך המידע בעזרת הצלבה עם הרשת החברתית LinkedIn, ולחשוף מידע פרטי אודותיהם: ציונים בקורסים שלקחו במערכת.<sup>1</sup>

## 2.2 גם שיטות מתקדמות של "הכללה" אינן מספקות

עם השנים הוצעו שכלולים של שיטת ה"הכללה" וקריטריון ה-k-אנונימיות, ובפרט הוצעו קריטריונים שונים כמו l-diversity ו-t-closeness, שמנסים לפתור בעיות אלו ואחרות, אך סובלים מאותו עניין יסודי: צירוף של מספר מאגרי נתונים שהותמם כל אחד בנפרד, או לחלופין צירוף של מידע ממאגר אחד ומידע פומבי על אדם, יכולים לגרום לזיהוי מחדש. במילים אחרות, מבחינה טכנית כל מידע (ובמיוחד צירופים של ערכי מידע) הוא מזהה ועלול לחשוף מידע פרטי או לצמצם את האפשרויות למידע פרטי על אדם מסוים.

בהקשר זה מעניין לציין כי הוועדה ליישום המלצות השימושים המשניים במידע בריאותי התייחסה להיבט זה, אך עדיין בחרה להשתמש בשיטת k-אנונימיות. בפרט, הוועדה המליצה במפורש לא לבצע הכללות על השדות עם המידע הקליני, למרות הסכנה הברורה לזיהוי מחדש.

גישה אחרת היא לבצע "אגרזיה" למידע ולשחרר רק מידע סיכומי בפורמט של טבלאות, כמו כמויות וממוצעים. אולם, גם כאן חוקרי טכנולוגיות להגנת פרטיות הראו שלעיתים ניתן "לבנות מחדש" את חלק או רוב מהנתונים המקוריים ששימשו ליצירת המידע בעזרת הטבלאות בלבד ( **reconstruction attack**). תקיפה זו מתאפשרת כאשר כמויות או ממוצעים על משתנה מסויים (למשל מספר חולים במחלה מסוימת) מופיע במידע סיכומי בחיתוכים שונים (למשל, גיל, מגדר ואזור מגורים).

גם באשר ל-"אגרזיה" מסוג אחר, של שחרור מודל (כמו רגרסיה לינארית) שאומן על המידע, היא בעייתית. בין השאר, חוקרים פיתחו בהצלחה מספר "תקיפות שייכות" ( **membership attack**) שפענחו האם אדם מסוים הוא חלק ממאגר המידע ששימש לאימון המודל. למשל, אם מודל אומן רק על חולים במחלה כלשהי, "תקיפת שייכות" עלולה לחשוף האם אדם מסויים חולה באותה המחלה, למרות שמידע זה לא פורסם כלל.

## 3. האתגר: מתן הגנה עמידה כנגד הצלבות מידע, הפתרון: פרטיות דיפרנציאלית

חוקרי הגנת פרטיות פיתחו מדד אלטרנטיבי לרמת הפרטיות המובטחת על ידי מנגנוני התממת מידע - מדד שמבטיח שרמת הפרטיות המובטחת תשמר גם כאשר מאגר הנתונים המותמם מוצלב עם מידע נוסף כלשהו, בין אם מידע שכזה זמין בהווה או בכל נקודת זמן בעתיד. המדד, שנקרא  $\epsilon$ -פרטיות דיפרנציאלית<sup>2</sup> פותח על ידי החוקרים דוורק, מקשרי, ניסים וסמית ב 2006 ומבטיח, בעקרון, שאף תוקף לא יצליח לנחש בהסתברות הצלחה יותר טובה מערך התלוי ב- $\epsilon$  (אפסילון) אם נתונים על אדם מסוים מופיעים במאגר או לא.

מאז הצעת המדד החדש, פותחו עשרות שיטות התממה שעונות על הדרישה, כשכל שיטה מותאמת לסיטואציה וסוג נתונים אחרים. חוקרים פיתחו שיטות שעונות לדרישת הפרטיות הדיפרנציאלית

<sup>1</sup> Cohen, A. (2019). [New guarantees for cryptographic circuits and data anonymization](#). MIT PhD Thesis

<sup>2</sup> Wood, A., Altman, M., Bembeneq, A., Bun, M., Gaboardi, M., Honaker, J., Nissim, K., O'Brien, D.R., Steinke, T. and Vadhan, S. (2018). [Differential privacy: A primer for a non-technical audience](#). Vand. J. Ent. & Tech. L., 21, 209.

ומתאימות למנעד רחב של תרחישים: החל מפרסום נתונים סטטיסטיים פשוטים, דרך אימון גרסיות לינאריות ועד לאימון מערכות מורכבות של בינה מלאכותית.

נציין שלכל השיטות האלה יש רעיון משותף: הזרקה של כמות מדודה ומכויילת של רעש אקראי לתוך הנתונים. הרעש הוא אקראי ולא ניתן לחיזוי, כך שאי-אפשר לגלות את ערכם המדויק של הנתונים המקוריים. מצד אחד הרעש צריך להיות גדול דיו כדי למנוע זיהוי פרטני והצלבות מידע, ומצד שני קטן דיו כדי לא לפגוע בניתוחים סטטיסטיים של הנתונים (כמו קורלציה בין מצב רפואי לנתונים ביוגרפיים). לדוגמה, אם אחד השדות במאגר נתונים הוא "הכנסה", במקום לבצע "הכללה" ולמשל לייצג הכנסה כטווח של 100 (כמו 1100-1100, 1200-1100 וכיו"ב), בשיטות פרטיות דיפרנציאלית נוסף רעש אקראי לכל אחד מהערכים - עבור אדם אחד נוסף 50 ולאחר נחסיר 100. בגלל שהרעש מתווסף באופן מדוד, עדיין נוכל להסיק מסקנות סטטיסטיות תקפות על הנתונים, כמו ממוצע הכנסה, תוך שמירה על הפרטיות של כל אחד מהאנשים במאגר מול כל תוקף. באופן כללי, ככל שבמאגר יש יותר נתונים על יותר אנשים (גודל המדגם), מידת הרעש היא קטנה יותר. פרטיות דיפרנציאלית מאפשרת למקבלי החלטות לכייל באופן מדויק את המתח הבלתי-נמנע בין הגנה על פרטיות לבין דיוק הנתונים.

נציין שדו"ח הועדה ליישום המלצות השימושים המשניים במידע בריאות התייחס לפרטיות דיפרנציאלית והסביר היטב את עקרונות השיטה. אולם, הקביעה בדו"ח, שאין דוגמאות בולטות ליישום בפועל של השיטה להגנה על פרטיות, אינה עדכנית. זאת משום שבשנים האחרונות התרחשה פריחה מחקרית בנושא וגישה זו יושמה במספר מקרים משמעותיים. דוגמה דומיננטית היא ההחלטה של לשכת מפקד האוכלוסין בארה"ב להתמים את טבלאות המידע של מפקד 2020 רק בעזרת שיטות שמקיימות פרטיות דיפרנציאלית. ההחלטה נלקחה לאחר שעורכי המפקד גילו התקפת זיהוי מחדש על תוצאות הסקר של 2010 - למרות העובדה שתוצאות אלה עברו התממה על ידי שילוב של מספר שיטות. לטענתם, המשך שימוש באותן שיטות הכללה יהווה עבירה על החוק האמריקאי שמורה על עורכי הסקר לשמור על פרטיות המידע של הנסקרים.

כדוגמה נוספת, מאז אפריל 2020 גוגל מפרסמת באופן תדיר 'דו"חות מוביליות בקהילה' (Community Mobility Reports) המאפשרים לחוקרים ולציבור הרחב לנתח מידע שינוי בהתנהגות הניידות של אנשים בעקבות נגיף הקורונה, בתגובה למדיניות אפידמיולוגיות. המידע מבוסס על נתוני-המיקום שגוגל אוספת, והוא מותמם בעזרת הוספת רעש אקראי מדוד באופן שמבטיח פרטיות דיפרנציאלית.<sup>3,4</sup>

לסיכום, מדד ה-k-אנונימיות לרמת ההתממה הינו מדד ישן שהוכח כלא אפקטיבי בעולם של היום, ונדחה על ידי ארגונים מובילים בעולם כלא מקיים את דרישות הדין להגנת הפרטיות ו/או להתממת מידע. לעומת זאת, מדד הפרטיות הדיפרנציאלית הוא מדד מודרני, שמאפשר לצמצם את סיכוני הפרטיות ואבטחת המידע, באופן שימנע זיהוי מחדש של יחידים מתוך מידע מותמם. כל שיטת התממה, בפרט k-anonymity, פוגעת ברמת הדיוק של הנתונים, אך שיטות המבוססות על מדד הפרטיות הדיפרנציאלית מאפשרות למקבלי החלטות לשלוט בדיוק על המתח בין דיוק להגנת פרטיות. בנוסף, המדד ישים: קיימת ספרות עשירה על שיטות התממה שעומדות בדרישות המדד, והמימוש הוא ישיר וקל כאשר מדובר בשחרור נתונים סטטיסטיים פשוטים כמו כמויות וממוצעים. חלק משיטות אלה כבר מיושמות בעולם במגוון של אפליקציות וסוגי מידע.

<sup>3</sup> Google (2019). [Enabling developers and organizations to use differential privacy](#). Also: [Video](#).

<sup>4</sup> Google (2020). [COVID-19 Community Mobility Reports](#).

**4. הצעת קווים מנחים אפשריים לתוכנית שמטרתה הטמעת פרטיות דיפרנציאלית עבור התממת מידע**

**רפואי במשרד הבריאות**

4.1 אנו מציעים לעבוד במודל של פיילוט: לזהות פרויקטים רלוונטיים של שיתוף מידע רפואי במערכת הבריאות וללמוד באופן ממוקד כיצד יש לעשות שימוש בשיטות פרטיות דיפרנציאליות. גישה שכזו תאפשר למשרד הבריאות ללמוד כיצד להנגיש את הטכנולוגיה של פרטיות דיפרנציאליות בפרויקטים מבוססי נתונים. אחת התועלות החשובות של הפיילוט הראשון, שהיא רלוונטית מאוד לכל היישומים העתידיים של פרטיות דיפרנציאלית, היא קיום דיון עקרוני על רמת ההגנה על הפרטיות הנדרשת, וקביעה בפועל של מדד הפרטיות ( $\epsilon$ -אפסילון).

4.2 לפרויקטים של שיתוף מידע רפואי יש מגוון של לקוחות אפשריים, כמו למשל חברות מסחריות, חוקרים מהאקדמיה, גורמים ממשלתיים והציבור. אנו מציעים למקד את הפיילוט הראשונים בשיתוף מידע רפואי עם חברות מסחריות, כיוון שיש להם מומחיות טכנולוגית, משאבים ותמריצים רבים יותר בהשוואה ללקוחות האחרים כדי להשתמש במידע שעבר התממה באמצעות שיטות של פרטיות דיפרנציאלית.

4.3 מבחינה טכנולוגית, הפרויקטים המתאימים ביותר לפיילוט הראשונים הם כאלה בהם המידע הוא טבלאי וכולל שיתוף של (1) מידע סיכומי רב (למשל, הרבה חיתוכים שונים), ו/או (2) מידע פרטני על אנשים, ו/או (3) מודלים סטטיסטיים או של למידת מכונה (ובפרט גרסיה לינארית ורגרסיה לוגיסטית, אך לא רק).

**5. המלצות**

5.1 לאור סעיף 13(ו) בטיטת התקנות, הקובע כי "התממה תבוצע בשיטות המקצועיות המיטביות הזמינות בתחום זה באותה עת", אנו ממליצים להחיל דרישה ששיטות התממת מידע רפואי של מידע מובנה (טבלאי)<sup>5</sup> תעמודנה כלל במדד הפרטיות הדיפרנציאלית עם רמת הגנה ( $\epsilon$  - אפסילון) שתקבע על ידי משרד הבריאות בשיתוף הרשות להגנת הפרטיות. חריגה מהכלל תאושר באופן פרטני, לאחר שהוכח שעמידה בדרישה הכללית אינה אפשרית במקרה הנתון ולאחר שהרשות להגנת הפרטיות השתכנעה שהפגיעה בפרטיות המוצעת היא מינימלית, וסבירה ביחס לתועלת הציבורית שתופק משיתוף הנתונים במקרה זה. במקרה זה, יש לפעול בהתאם לחוק הגנת הפרטיות והתקנות שהותקנו מכוחו.

5.2 פיתוח וניתוח אפקטיביות של שיטות התממה הינו תחום שדורש מומחיות מקצועית גבוהה, גם בשיטות הגנה וגם בשיטות התקפה. לכן אנו ממליצים שילוב של צוות מייעץ רב-תחומי שמורכב, בין היתר, ממומחים מתחומי הגנה פרטיות (ובפרט בתקיפת פרטיות), מדעי הנתונים ומשפט, לצורך מעורבות בהחלטות בנושא התממת מידע בגופים ציבוריים ורגולטוריים. נציין שמספר חוקרים ישראלים בארץ ובחו"ל הם בין המובילים העולמיים בתחום הזה, גם בתיאוריה וגם במעשה, כך שהמלצה זו ישימה באופן מיידי, ונשמח לעמוד לרשותכם לצורך כך.

<sup>5</sup> מידע לא-מובנה (כמו טקסט חופשי, דימות וסאונד) הוא לרוב מורכב לאין ערוך ממידע מובנה (טבלאי), ומציב קשיים משמעותיים לפרטיות. מכתב זה אינו עוסק במידע לא-מובנה, אך המלצה 5.2 מאפשרת לעמוד באתגר תוך שימוש בשיטות היישומיות העדכניות ביותר העומדות לרשותנו.

מדינת ישראל מובילה כיום בכמות, איכות וריכוז המידע הרפואי הנגיש למחקר. אולם היא עדיין בפער ניכר בתחום הגנת הפרטיות. ישראל יכולה להפוך למובילה בינלאומית גם בתחום זה, על-ידי מימוש שיטות התממה מודרניות קיימות ופיתוח שיטות חדשות.

נודה אפוא לקיום פגישה בנושא עם האנשים הרלוונטיים מטעמכם,

בכבוד רב,

רן קנטי, פרופסור למדעי המחשב, אוניברסיטת בוסטון

אלוני כהן, חוקר בתר-דוקטורט למדעי המחשב ומשפטים, אוניברסיטת בוסטון

שלומי הוד, דוקטורנט למדעי המחשב, אוניברסיטת בוסטון

עו"ד נעמה מטרסו, מנכ"לית פרטיות ישראל

#### העתקים:

עו"ד רז נורי, המשנה ליועץ המשפטי לממשלה (ציבורי-חוקתי), משרד המשפטים

ד"ר שלומית וגמן, עו"ד, מ"מ ראש הרשות להגנת הפרטיות, משרד המשפטים

עו"ד איל זנדברג, ראש תחום, מחלקת ייעוץ וחקיקה (ציבורי חוקתי), משרד המשפטים

עוז שנהב, מנהל מחלקת חדשנות ופיתוח מדיניות, הרשות להגנת הפרטיות, משרד המשפטים