

אישה נעה נעה: מודלי עיבוד שפה טבעיות בעברית



ענבל יהב אביחי שריין

ד"ר ענבל יהב היא חברת סגל בפקולטה לניהול ע"ש קולר אוניברסיטת תל אביב. בעלת תואר ראשון במדעי המחשב ותואר שני במערכות מידע מהטכניון. קיבלה את הדוקטורט שלה באופטימיזציה וקריות נתונים מאוניברסיטת מרילנד בשנת 2010, 2011 והמשיכה לעבוד שם במשך שניםים כחברת סגל אורחת. עיקר עבודתה מתמקדת בפיותח והטאהה של מודלים סטטיסטיים לשימושם של חוקרים במערכות מידע. במחקריה היא משלבת אלגוריתמים לכריית נתונים, מודלי עיבוד שפה טבעיות ומודלי אופטימיזציה כדי ליצור מודלים של נiations נתונים עברו יישומים שונים, ובهم יישומי בריאות הציבור וניתוח רשותות חברותיות.

אביחי שריין הוא דוקטורנט בפקולטה לניהול ע"ש קולר אוניברסיטת תל אביב, מהת הנחייתה של ד"ר יהב. בעל תואר ראשון בכלכלת ולימודי מורה תיכון מאוניברסיטת חיפה ותואר שני בניהול מערכות מידע מאוניברסיטת תל אביב. תחום המחקר שלו הוא זההו שימושיות בטקסט, במיוחד לשפה העברית. הוא פיתח את HeBERT, מודל שפה עברי המבוסס על ברט ומודל זההו רגשות בעברית מtekst.

תקציר

עברית שפה קשה. למחשב, כמו לאדם, וקצת יותר. בארבע השנים האחרונות מודלי עיבוד שפה טבעיות נמצאים בשיא פריחתם עבור מגוון שפות ומגוון משימות מחשב, כגון תרגום, מענה על שאלות, ניתוח תחשויות וכיתבת מקציירים. העברית, לעומת זאת, נותרה קצת מאחור. זה לא מאד מפתיע מפני שהקל העד על עברית קטן שימושיות מזה של שפות אחרות, ומבנה השפה מורכב בהרבה. למעשה העברית נחשבת "שפה עשרה מורפולוגיות" – שפה שבה המידע המורפולוגי מכוון חלק מהמליה, ולא מופרד ממנה כמו במרבית השפות הלטיניות.

ב-2021 פותח על ידי כתבי מאמר זה מודל שפה מבוסס ברט ראשון לשפה העברית, שהיווה ירידת פתחה למחקרים רבים בתחום. במאמר זה נצין את האתגרים בפיתוח מודל השפה העברית, נסקרו את המודלים הקיימים והמאיצים המתmeshיכים לפיותח כלים ומודלים חדשים, ולאחר עוד אפשר ויכול לשאוף. בנוסף נצון כיצד קקרה כיצד ניתן, ללא ידע מוקדים עשיר, להשתמש במודל השפה בעברית ליהו תחשויות מותך שפה כתובה.

הקדמה: מה זה מודל שפה?

שקדמו לו (במודול דו-יכוני), או המונחים הקודמים למונה, ואלו שעקבם אחריו (במודול דו-יכוני). לבני עיתות ההקשר, שהיא משמעותית סבוכה יותר, נדרש הבנה מדויקת יותר של המונחים בשפה מעבר לשיחותם. את הבסיס הראשון לבניית ההקשר יצאו תומאס מיקולוב וחבריו (Mikolov et al., 2013, מוגול בשנת 2013, במודול הנקרא word2vec).

לפי מודל word2vec של מיקולוב וחבריו, מונח בשפה מייצג על ידי וקטור מסווג באורך קבוע, בთהליך שנקרא "שיכון מילים" (או אנגלית, word embedding). לדוגמה, המילה "ילד" יכולה להיות מוצגת במודול זה על ידי הווקטור [10.5, 9.9, ..., 3.2]. ערך הווקטור לנבדים בעוררת רשת נוירונים בעלת שתי שכבות, המאומנת על קורפוס גדול של טקסט, שptrתיה לשכנ את המונחים בקורסוס על המרחב הרציף באופן שבו מונחים דומים ימוקמו קרוב זה לזה במרחב. לפי מודל זה, משמעות המונח נמנעת מהקשה בשפה, כמו למשל מואוף המונחים שקדמים לו במשפטים שונים ומאלו שבאים לאחריו. בהתאם, רשת הנוירונים מקבלת קלט "הקשר": משפטים באורך שווה, שבהם מונח האמצע ממוקן "הקשר". משפטים באורך שווה, שבהם מונח האמצע ממוקן "משמעות": מונח בעל הסתברות גבוהה ביותר להשלמת המשפט (לדוגמה: "הילך הלך ___ הספר היסודי", וחוזה כפלו).

המשמעות (לדוגמה: "לבית").

אחרי הפיתוח של מיקולוב וחבריו, החלו לנצח (עדין ציטים) מודלי שפה שונים כפתרונות אחרי הגשם. אחד הנפוצים שבهم Bidirectional Encoder Representations (BERT) (from Transformers [BERT] 2018 Devlin et al., 2018) את המודל הזה, כפי שנפרט בהמשך המאמר, תרנמננו' לראשונה לעברית.

מודל ברט

מודל ברט הוא מודל שפה מבוסס טרנספורמר דו-יכוני – ארכיטקטורה ייחודית המאפשרת למידת' יצוג של מונחים מתוך הקשר גלובלי (כגון מסמך) ולוקלי (כגון משפט). מודל זה מאפשר יצוג וקטורי שונה של מונחים דומים בעלי משמעות שונה, כפי שמשתמש מהקשרים במשפט, כגון המילה החזרת "געלה" במשפט "איישה נעלה נעלה נעללה".

1 ארכיטקטורת רשת נוירונים זו נקראת *Continuous Bag of Words* (CBOW) Skip-gram. אלטרנטיבית, בראכיטקטורת ה-words.

מודל שפה הוא בסיסו מודל הסתברותי המחשב את ההסתברות של מונח (token) להיות חלק "תקן" מהשפה שהוא למד. המונח הנלמד יכול לכלול מילה שלמה, חלק ממילה (כגון תחילית, סופית), או אפילו אות בודדת. המחשב לומד לשלים מונחים ממשפטים רבים, וכך לומד את מבנה השפה, הטויטה ומשמעות המילים בתוכה, וכך להוות בסיס למשימות שפה רבות.

ישנם סוגים שונים של מודלי השפה, והפשות שבהם הוא מודל האוניגרים (Unigrams) (Sebastiani, 2002) [WoBo], משפט מייצג על ידי "שיכון מילים" (Bag of Words),(Cloumer אוסף המונחים (או המילים, בישום הפשט של האלגוריתם) המרכיבים אותו. הסתברות משפט נמנעת מכפלת הסתברות המונחים (כלומר שכיחותם בשפה) בשך. לדוגמה, המשפט "הילדה הלכה לבית הספר" יוצן על ידי האוסף: {הילדה, הלכה, בית, הספר}, והסתברות המשפט תחושב באופן הבא:

$$(1) \quad = ("הילדה הלכה לבית הספר") \cap ("הספר") \cap ("לבית") \cap ("הליכה") \cap ("הילדה")$$

את הייצוג הזה ניתן כМОן לשפר על ידי עיבוד מקדים של המשפט; מילות ואותיות קישור ייסרו מהמשפט (שכן אין מושיפות מידע על הסתברות המילים המרכיבים אותו וויפויו באופן שכיח בכל משפט), המילים יהלטו בצוות השורשיות שלהן, ואוניגרים יחלפו בזמנים מילים אם לאלו שיש הסתברות גבוהה יותר. במקרה הזה הסתברות המשפט תוחלף במשווהה הבאה:

$$(2) \quad = ("הילדה הלכה בית הספר") \cap ("בית ספר") \cap ("הלך") \cap ("ילדה")$$

היתרונות של מודל האוניגרים הם הפשטות שבו והיכולת להבין את התוצר המתקבל מהאלגוריתם. החיסרון של המודל הוא החוסר הבהיר בסדר המונחים ובקשר שלהם. בהינתן אי מנבליה על סדר המונחים והקשרים, המשפט "הילדה הלכה בית הספר הלך בית הילדה". את המובלות הללו, כמו המשפט "הספר הלך בית הילדה". את המובלות הללו, פותרים מודלים מורכבים יותר, וביעית סדר המונחים נפתרת על ידי הסתברות מותנה: הסתברות מונח בהינתן המונחים

קדומים. המשימות נחלקות לשניים: משימות לא מופוקחות, שמטותן לראות כיצד המודל "مبין" את השפה, ומשימות סיווג שמטותן לראות כיצד הבנת השפה תורמת ליהו אלמנטים שונים בשפה. להלן פירוט המשימות הנפוצות:

1. משימות לא מופוקחות

א. משימת "מלא את החסר" (*Fill-in-the-blank*). משימה הבודקת את יכולת המודל להשלים מונחים חסרים בטקסט. ביצוע המשימה נבדקים על ידי מדדי perplexity, המכמתת את האיכות של חיזוי נוכנות משפט בשפה באמצעות המודל. מתמטי, איקות החיזוי של רצף (w), דוגמת משפט) בין N מונחים (גניך, מילים), מוגדרת על ידי הממוצע המעריכישן הנראות המקסימלית המותנית של המונחים ברכף:

$$PP(w) = \exp \left\{ -\frac{1}{N} \sum_{i=1}^N \log_{p_\theta}(w_i) \right\} \quad (3)$$

ב. הכללה למילים מחוץ למילון (*out-of-vocabulary*). משימה הבודקת את יכולת של המודל להזות מילים שאין מופיעות בקורסום שעליו הוא אומן. מדובר באחוז המילים שהמודל לא הצליח לשכנן (כלומר ליציר עברים וקטור מספורי).

2. משימות מופוקחות

ג. זיהוי ישיות (NER). Named-Entity Recognition (Named-Entity Recognition [NER]). משימה הבודקת את יכולת של המודל לסווג ישיותם שם בטקסט, כגון שמות, ארונים ומיקומים של אנשים. איקות הסיווג נמדדת על ידי מדדי F1:

$$2 \times \frac{\text{precision} \times \text{recall}}{\text{precision} + \text{recall}} \quad (4)$$

ד. זיהוי חלקו דיבור (POS). Part-of-Speech (POS). משימה הבודקת את יכולת של המודל לסווג את התפקיד הדקדוקי של מילה או ביטוי במשפט (לדוגמה: שם עצם, فعل, תואר השם). איקות הסיווג נמדדת אף היא על ידי מדדי F1. ניתוח תוצאות. משימה הבודקת את יכולת של המודל להזות את התפקידות המובאות בטקסט. משימה זו בדרכן.

מודלי ברט נמצאים היום בשימוש נרחב ככלי להבנת שפות רבות, ובهن אנגלית, ערבית, רוסית ועוד שפות רבות אחרות. רוב המודלים מאומנים על קורפוס ענק, ובפרט קורפוס Open Super-large, וקיופידה אואסקר (Crawled Aggregated coRpus [OSCAR] Suarez et al., 2020), שעתוקם שלם קיימים בשפות השונות. בדומה למודל word2vec, מודל ברט מקבל קלט משפטים בשפה עם מונח מסוין (הקשר), ומטרתו לחזות את המונח הקשור (משמעות). בשונה מword2vec, מודל ברט כולל 12 שכבות של למידה عمוקה (הנקראות שכבות טרנספורמר), כאשר בכל שכבה יש למידה נוספת של הקשר המילויים במשפט. מידע נוסף על ארכיטקטורת ברט ניתן למצוא במאמר המקורי של דבלין וחבריו (Devlin et al., 2018).

לאימון מודל ברט על שפה חדשה נדרש החלטה אחת חשובה, והוא רמת הפירות של המונח שעליו יאות המודל (קלט המודל). נזכיר כי מונח הוא ייחידת טקסט היכלה לכלול מילה שלמה, חלקו מילה, או אות בודדת. רמת פירות של המילה "הילד" שהוצאה קודם. רמת פירות של אות בודדת, המונחים שייצגו מילה זו יהיו {"ה", "י", "ו", "ל", "ד", "ה"}. ברמת פירות של מילה שלמה, המונח יהיה שקול למילה השלמה: {"הילד"}*. חלקה לחלקי מילים יכול להתבצע בשתי דרכים: בעזרת ניתוח מורפולוגי {"יה", "ילד", "ה"}, או לפיה חלקו מילה שכוחים סטטיסטיים {"יה", "ויל", "דה"}.

לכל אחד מרמות הפירות יתרוננות וחסרונות משלו. ככל כל שרתת הפירות גבוהה יותר (כגון רמת אות), כך המודל מסוגל בקלות וחסית למדוד מילים חדשות העומדות בלונגה הכלילתית של השפה, אך מתקשה יותר למדוד את ההקשר (אנטואטיבית, המודל "ימבב" שכבות של למידה כדי להרכיב מילים), ולהיפך (Jawahar et al., 2019). בשפות שונות מצאו כי רמות פירות שונות מטיבותם עם המודל: באנגלית, המודל הנפוץ הוא מודל המבוסס מיליון כנ"ן 30 אלף חלקו מילה (Devlin et al., 2018); בעברית, לעומת זאת (הקרובה יותר לשפה העברית), הייצוג הנפוץ ביותר הוא מיליון מפורט בהרבה המכיל כ-60 אלף חלקו מילה (Antoun et al., 2020).

aicot מודל שפה

כדי ללמד איקות של מודל שפה, נהוג להריץ את המודל על מספר משימות שפה נפוצות ולהשוו את ביצועיו למודלים

לנורומים המרכיבים אותה וכן מוצלחות יותר בلمידת מבנה המילה (מורפולוגית). על פניו הציגו ש נם רמות ביןיהם – חלקן מילה, וחלקן מילה עם שימושות מורפולוגיות. האתגר בעברית הוא להבין הן את ההקשר והן את מבנה המילה.

כדי להבין את הפשרה בין החלופות, ערכנו ניסוי שבו אימנו מודל ברט לשפה העברית המאמן על משימת "מלא את החסר" עם רמות מונחים שונות (ספציפית, רמת מילה, מילים שונים של מילון חלקו מילים, חלקו מילה עם שימושות מורפולוגיות, ועוד). בחנונו את ביצועיו על שלוש המשימות המפוקחות שהוצעו בפרויקט הקודם (משימות לא מופוקחות אין בנות השוואה במקורה זה בשל גודלם השונה של המילונים Chriqui and Yahav, 2021). בשל המשאבים הרבים הנדרשים לאימון מודל לשפה, ביצענו את האימון על מסד נתונים קטן – מסד ויקיפדיה בעברית (מסד בגודל ~650靡נה). הממצאים מובאים בטבלה 1 לעיל. ניתן לראות בטבלה כי מודל "האמצע" המאמנים על חלקו מילה טוביים בהרבה ממודלי הלקצוט. בין חלופות האמצע, מודל המאמן על חלקו מילה עם שימושות מורפולוגיות מראה ביצועים טובים יותר במשימות הדורשות הבנת מילים בשפה – זהוי ישיות ויזיה חלקו דיבור, לעומת מודל המאמן על חלקו מילה, ללא שימושות מורפולוגיות, מצליח להבין את הרעיון המורכבי במשפט באופן טוב יותר, ובהתאם לחץ את התהוושה המתבטאת במשפט.

מודל השפה הסופי

כמו מודל שפה סופי, בחרנו לאמן מודל המבוסס על חלקו מילה. הסיבה לבחירה נובעת מהשימוש העיקרי בmorphological Rich Language (MRL) – שפה שבה המידע השופוט הלטניינית מהמילה ולא מופרד ממנה כמו במרבית השפות הלטנייניות הבאות: (i) ריבוי נטיות לכל מילה, על ידי הוספות מסווגות למילת בסיס (לדוגמה: יلد, ילדם); (ii) נטייה לא רציפה, שבה בסיס המילה משתנה ורק המוסףויות (לדוגמה: הילך, הולך); (iii) סדר המשפט לעתים חסר משמעות ולעתים רבמשמעות; (iv) מילים רבות יש שימושות כפולה המשתנה לפי ההקשר; (v) הניקוד בעברית, שאינו מופיע במרקビות הטקסטים הכתובים, מוסיף משמעות למילים בשפה.

כל נבדקה על זיהוי קווטביות בתיחסות (תחווה חיובית, ניטרלית, או שלילית), ונאמדת אף היא בעורף מדריך F1.

מודל לשפה בעברית: HeBERT

על מנת לפתח מודל לשפה נדרשות שלוש החלטות. הראשתו היא רמת המונח שהמודל לומד (מילה, חלקו מילה, או אות). השניה היא ארכיטקטורת הלמידה העמוקה – אללו סוני שכבות למידה לכלול, כמה שכבות למידה, מה הקשרים ביןם, וכדומה. ההחלטה השלישייה שלישית נוגעת למשימה שעליה מאמן המודל, כלומרஇזוי משימה המודל לומד לבצע כדי להבין את השפה.

בפיתוח מודל לשפה בעברית, מודל HeBERT (& Chriqui & Yahav, 2021), בחרנו להשתמש בארכיטקטורת הבסיס של מודל ברט, שהוכחה אלטרנטיבית מובילה לביעות לשפה שונות (Radford et al., 2018). משימת האימון שנבחרה הייתה משימת "מלא את החסר", המוגדרת כמשימת הבסיס של מודל ברט. שתי החלטות אלו הן תלוות לשפה, וניתן להחילפן בארכיטקטורות אחרות (דגםת RoBERTa (Liu et al., 2019) ומשימות שונות (Dongmei- pointwise-mutual-information (PMI) masking (Levine et al., 2020).

לבחירה רמת המונח המתאימה לשפה העברית, נדרשת הבנה של מאפייני השפה הייחודיים. עברית נחשבת "שפה עשירה מורפולוגית" (באנגלית: Morphologically Rich Language, MRL) – שפה שבה המידע המורפולוגי מקיים חלקן מהמילה ולא מופרד ממנה כמו במרבית השפות הלטנייניות (Tsarfaty et al., 2010). מאפייני השפה כוללים את התרונות הבאות: (i) ריבוי נטיות לכל מילה, על ידי הוספות מסווגות למילת בסיס (לדוגמה: יلد, ילדם); (ii) נטייה לא רציפה, שבה בסיס המילה משתנה ורק המוסףויות (לדוגמה: הילך, הולך); (iii) סדר המשפט לעתים חסר משמעות ולעתים רבמשמעות; (iv) מילים רבות יש שימושות כפולה המשתנה לפי ההקשר; (v) הניקוד בעברית, שאינו מופיע במרקビות הטקסטים הכתובים, מוסיף משמעות למילים בשפה.

הה תלבות בין רמות המונח השונות נעה על הציג שבין רמת מונח נבואה (מילה), השמרת על מבנה המילה המורכבת תוך המשפט שבו היא מוגנת ולן אפשרות למידה טובה יותר של הקשר, לבין רמת מונח נמוכה (אות), המחלקת מילה

טבלה 1: השוואת בין ביצועי מודל ברט בעברית, תחת רמות מונחים שונות

רמת מונה	זיהוי ישיות ¹	זיהוי חלקיק דיבור ²	ניתוח תחושים ³
אות	0.74	0.92	0.69
חלקי מילה עם משמעות מורפולוגית (More et al., 2019)	0.92	0.95	0.65
חלקי מילה	0.79	0.90	0.79
מילה	0.86	(לא ניתן לחישוב בשל ביצועים נמוכים במשימת מלאי-אתה-הolson)	0.43

1. מתוך המסל שפורסם על ידי Mordecai and Elhadad (2005)

2. מתוך המסל שפורסם על ידי Simaan et al. (2002)

3. מתוך המסל שפורסם על ידי Amram et al. (2018).

מודל זיהוי רגשות: HebEMO

משימת זיהוי רגשות היא אחת המשימות הנפוצות בעיבוד שפה טבעי. מטרת משימה זו היא להוות קשת רחבה של רגשות כפי שהוא בא לידי ביטוי בתוכן כתוב, ובכללها אושר, כעס, פחד ועוז. זיהוי רגשות אלו יכול לשפוך אוור על אמוןנות, התנהלות צפויות, ומצבים נפשיים שבהם שרויים אנשים. בספרות הפסיכולוגית ישנן מספר תיאוריות המגדירות את קשת התחושים של אדם. הנפוצה שבנה היא זו שפורסמה על ידי פלוטצ'יק (Plutchik, 1980) ומגדירה מעגל של ארבעה תכונות מנוגדות: עצב-שמהה, כעס-פחד, אמון-גועל, והפתעה-ציפייה.

במחקר שלנו (Chriqui & Yahav, 2021), התבוסנו על הנדרת התחושים לפיו פלוטצ'יק כדי ליצור מודל זיהוי רגשות אוטומטי. לשם כך אספנו מסד נרחב הכלול כ-4,000 העותות שנכתבו בתגובה לכתבות חדשותיות (אפיו האיסוף מפורט במאמר). בחרנו לאוסף תגניות לכתבות שנכתבו בנושאים הקורונה בשנת 2020, תקווה המוגדרת כתקופה עמוסה רגשית (Pedrosa et al., 2020), ופורסמו באחד משלושת אתרי החדשנות הבאים: *ynet*, *ישראל היום*, ובוחורי חרדים. בחירת האתרים נבעה מהאופן שבו הם מיצגים את המנגנונים השונים בארץ: שני הראשונים פונים לעיר לקהל היהודי החילוני משני קצוות הקשת הפוליטית, ואילו השלישי פונה למגזר החradi. את התגניות שהלחנו לתיוג אנושי באתר *Prolific*,⁶ כך שכל תגינה תיוגה על ידי כעשרה מתינים דוברי השפה.

את מודל השפה המאמון HeBERT ניתן למצוא בניט,² באתר מודלי השפה *huggingface*³ ובשירותי הענן של AMAZON⁴ (AWS, Amazon Web Services).

התפתחות כל עיבוד שפה טבעית נוספים בשפה עברית

במשך השנים האחרונות התפתח תחום עיבוד השפה הטבעית בעברית לאין שיעור. חלק לא מבוטל מההתפתחות זו מיחסות לתפקידו של שביבון לפרויקטם בתחום על ידי משרד החדשנות, המדע והטכנולוגיה, שמהם צמח נס האינוד*ישראלי* לטכנולוגיות שפת אנושי, המאנד נציגי אקדמיה ותעשייה במטרה לפתח את תחום עיבוד השפה בעברית.

כדי להדנים את התפתחות השפה, נפרט במאמר זה אודות שני מאמצים מחקרים שביצעו כתבי מאמר זה: פיתוח מודל זיהוי רגשות בעברית המשמש במודל השפה הבסיסי בעברית לצורכי משימת סיווג, ופיתוח מודל שפה משפטית ה"מלמד" את מודל השפה "להבין" את המונחים המשפטיים בעברית.

- | | |
|---|---|
| https://github.com/avichaychriqui/HeBERT | 2 |
| https://huggingface.co/avichr/heBERT_sentiment_analysis | 3 |
| https://github.com/aws-samples/aws-lambda-docker-serverless-inference/tree/main/hebert-sentiment-analysis-inference-docker-lambda | 4 |
| https://www.iahlt.org | 5 |

משפטית), כמות המשאבים שהוא דורש היא נדולה – בהיבט מקורות המידע, בזמן העבודה, ובמשאבי החומרה הנדרשים. לעומת זאת, מודל שעובר התאמת מידע פחות משאים, ארך עלול להיות פחות מדויק בהבנת הקשרים משפטיים פורמליים.

לצורך פיתוח מודל השפה העברית המשפטי ובחינת האופן המתאים לאימון המודל, במחקר אחרון („Chiriqui et al., 2022“) אספנו מקורות מידע משפטיים רבים בגודל כולל של ~3.7 ג'ינה, המכילים את סוף החוקים בישראל, מאגרי פסקין דין, מאגרי החלטות בית דין ועוד. על מנת זה אימנו שני מודלים – אחד המאמן מהתחלת על המסלול החדש, ואחד המבוסס על מודל HeBERT ומתאים אותו לתחום המשפט. שני המודלים המאמנים הועלו לניט לשימוש הציבור.⁷

את ביצועי המודלים בחנו על שתי מישיות סיוג. למשימת הסיוג הראשונה בחרנו את מסד חוק ההסדרים („Kosti, 2021“). מטרת המשימה, כפי שהונדרה במאמר המקורי, הייתה לסייע משפטיים ככלו המכילים סמכיות לרשותם ולאלו שלא. יכולת היהוי האוטומטי בסיס דזה נאמדה על 0.87, כאשר המודל שאומן מהתחלת הראה את הביצועים הטובים ביותר. במשימת הסיוג השנייה בדקנו את יכולת המודלים להזיהות תכונות מסוימות אצלינו חקיקה בכנסת ישראל.⁸ במקרה זהה המודל טוב יותר היה דווקא המודל המותאם, וביצועיו נאמדו על $F_1=0.74$. להערכתנו, הסיבה לביצועים הטובים יותר של המודל השני בקרה זה טמונה באופי השפה בדיוני חקיקה – שפה המשלבת שפה טבעית ושפה משפטיית – וישורה עם אופי מודל HeBERT המותאם לשפה המשפטית.

בימים אלו אנו עוסקים על פיתוח נרחב של השפה העברית המשפטית, הכולל איסוף והנגשת מקורות מידע, תינוק ישיות משפטיות, ובנויות מודדי שפה (ראו מימון מחקר ושותפים בפרק התווות).

⁷ <https://github.com/avichaychiriqui/Legal-HeBERT?fbclid=IwAR3sFizNjEfPIXm0Agg5HpELUm49v11kfksjes72-Q-9CxMwv8hdR815ahg>
⁸ בסיס שנאפס וונגן על ידי פרופ' איתן בר-סימן-טוב מאוניברסיטת בר-אילן וורם פורסם.

לבסוף, ביצענו ציון (Fine-tuning) של מודל השפה HeBERT למשימת הסיוג של זיהוי התחששות (כלומר, עדכנו את משקלות ארבע השכבות האחרונות במודל לזיוי אופטימלי של התחששות). איקות המודל לפי מדד F_1 נעה בין 0.78 ל-0.97 לתחששות השונות, מלבד הרגש "הפטעה" שאוות המודל לא הצליח לזהות ($F_1 = 0.41$), כמוポート בטבלה 2.

טבלה 2: ביצועי מודל ניתוח התחששות בעברית – HebEMO

תחשושה	Recall	Precision	F1
עצב	0.84	0.83	0.84
שמחה	0.87	0.89	0.88
כעס	0.97	0.97	0.97
פחד	0.77	0.84	0.80
אמון	0.70	0.88	0.78
גועל	0.95	0.97	0.96
הפתעה	0.37	0.47	0.41
ציפייה	0.87	0.83	0.85

מודל שפה משפטי Legal-HeBERT בעברית

הצורך בכלל עיבוד שפה טבעית למקרים בתחום המשפט והחקיקה, בעולם בכלל ובעברית בפרט, הלך והתעצם בשנים האחרונות (לדוגמה, „Katz & Nay, 2019; Sarne et al., 2021“). לצד צורצוזה, ההתפתחו מודלי שפה המותאמים לענה המשפטית בשפות רבות, ובראשן צפוי בשפה האנגלית (Chalkidis et al., 2020).

התאמת מודל שפה לעולם תוקן חדש יכולה להתבצע בשתי דרכים. הראשונה היא אימון מודל שפה חדש (נניח ברט), שמקורות המידע שעליהם הוא מאמין הם מקורות משפטיים בלבד. הדרך השנייה, הנזכרת "התאמת תחום" (domain adaptation, ראו Devlin et al., 2018), נסמכת על מודלים ותוקן חדש. במקרה זהה נקודת ההתחלת תהיה מודל שפה מאומן על קורפוס כלל (דגם HeBERT), שחלק מהשכבות שלו עוביות אימון מחודש בעזרת מקורות המידע המשפטיים. הותרונות והחסרונות של כל שיטה בוררים: בעוד אימון מחדש מייצר מודל שמיוכן בצורה טובה יותר את השפה המשפטית ואת הקשי המילים בעולם תוקן זה (ambil להיות "מובלבל" מהקשרים של אותן מילים בשפה שאינה

מה הלהא?

עתה, תיון אלף ישיות משפטיות במסדים אלו, וויתוח מספר מודלי שפה המותאמים לישומים שונים בתחום המשפט. הפרוייקט השני יתמקד בפיתוח מודלים בעברית פיננסית לצורך ניתוח אוטומטי של דיווחי אנידים ציבוריים. פרוייקט זה יעשה בשיתוף פעולה עם הרשות לנויות ערך ועתיד להתחילה השנה הקרובה.

תודות

מחקר זה התאפשר בזכות מספר מקורות מימון. הראשון הוא תקציב ממשרד החדשות, המדע והטכנולוגיה (גראנט מס' #17991-3) בהובלת ד"ר ענבל יהב ופרופ' איתן בר-סימן טוב (אוניברסיטת בר אילן). השני, תקציב רשות החדשנות (גראנט מס' 01103654#) בהובלת פרופ' לב מוצניך (האוניברסיטה העברית) וד"ר ענבל יהב.

בנוספ', החוקרים מבקשים להודות לתמיכת האדיבה של קרן גירמי קולר ומכוון הנרי קראון למחרק עסקו בישראל.

inbalyahav@tauex.tau.ac.il

ד"ר ענבל יהב

מעבדות המחקר באוניברסיטה השונות ממשיכות ומשיכו לפתח כלים עבור השפה העברית, ובכללם מודלי שפה לתחומים שונים, כגון מודל השפה המשפטי (Chriqui et al., 2022), איסוף מסדים ייעודיים לאיתון, דוגמת מסד השאלות (Shmidman et al., 2022) ומודל לשפה רבנית (Keren and Levy, 2021) Parashoot ותשבות פרוייקט תיון היישיות ותיון מורפולוגי שמוביל האינוד הישראלי לטכנולוגיות שפת אנוש.

במקביל, עוד ועוד תחומי דעת עושים וישו שימוש במודלי שפה בעברית כחלק אינטגרלי מהמחקר שהם מבצעים. דוגמאות ראשונות לכך ניתן לראות בפסיכולוגיה (Hachohen-Kerner et al., 2022), בפוליטיקה (Litvak et al., 2022) ובמדעי החברה (Bialer et al., 2022).

אצלנו, מעבדה לבינה מלאכותית וניתוח עסקי בפקולטה לנויה ע"ש קולר באוניברסיטת תל אביב, אנו עתדים להוביל בשנים הקרובות שני מחקרים מרכזים בתחום עיבוד השפה הטבעית בעברית. האחד, המשך פיתוח השפה המשפטי, פרוייקט זה נמצא כתה שלבי הראויים, כולל איסוף מסדי

רשימת מקורות

- Amram, A., David, A. B., & Tsarfaty, R. (2018, August). Representations And Architectures in Neural Sentiment Analysis For Morphologically Rich Languages: A Case Study From Modern Hebrew. In *Proceedings of the 27th International Conference on Computational Linguistics* (pp. 2242-2252).
- Antoun, W., Baly, F., & Hajj, H. (2020). AraBERT: Transformer-based Model for Arabic Language Understanding. *arXiv preprint arXiv:2003.00104*.
- Bialer, A., Izmaylov, D., Segal, A., Tsur, O., Levi-Belz, Y., & Gal, K. (2022). Detecting Suicide Risk in Online Counseling Services: A Study in a Low-Resource Language. *arXiv preprint arXiv:2209.04830*.
- Chalkidis, I., Fergadiotis, M., Malakasiotis, P., Aletras, N., & Androutsopoulos, I. (2020). LEGAL-BERT: The Muppets Straight Out of Law School. *arXiv preprint arXiv:2010.02559*.
- Chriqui, A., & Yahav, I. (2021). HeBERT & HebEmo: A Hebrew BERT Model and A Tool For Polarity Analysis and Emotion Recognition. *INFORMS Journal on Data Science* 1(1):81-95.
- Chriqui, A., Yahav, I., & Bar-Siman-Tov, I. (2022). Legal HeBERT: A BERT-based NLP Model for Hebrew Legal, Judicial and Legislative Texts. *Judicial and Legislative Texts* (June 27, 2022).
- Devlin J, Chang MW, Lee K, Toutanova K (2018) Bert: Pre-Training of Deep Bidirectional Transformers For Language Understanding. *arXiv preprint arXiv:1810.04805*.
- Hacohen-Kerner, Y., Manor, N., Goldmeier, M., & Bachar, E. (2022). Detection of Anorexic Girls-In Blog Posts Written in Hebrew Using a Combined Heuristic AI and NLP Method. *IEEE Access*, 10, 34800-34814.
- Jawahar, G., Sagot, B., & Seddah, D. (2019, July). What does BERT learn about the structure of language?. In *ACL 2019-57th Annual Meeting of the Association for Computational Linguistics*.
- Katz, D. M., & Nay, J. J. (2021). Machine Learning and Law. *Legal Informatics*.
- Keren, O., & Levy, O. (2021). ParaShoot: A Hebrew Question Answering Dataset. *arXiv preprint arXiv:2109.11314*.
- Kosti, N. (2021). Centralization Via Delegation: The Long-Term Implications of The Israeli Arrangements Laws. In *Comparative Multidisciplinary Perspectives on Omnibus Legislation* (pp. 73-94). Springer, Cham.
- Levine, Y., Lenz, B., Lieber, O., Abend, O., Leyton-Brown, K., Tennenholz, M., & Shoham, Y. (2020). PMI-Masking: Principled Masking of Correlated spans. *arXiv preprint arXiv:2010.01825*.
- Litvak, M., Vanetik, N., Talker, S., & Machlouf, O. (2022, October). Detection of Negative Campaign in Israeli Municipal Elections. In *Proceedings of the Third Workshop on Threat, Aggression and Cyberbullying (TRAC 2022)* (pp. 68-74).
- Liu, Y., Ott, M., Goyal, N., Du, J., Joshi, M., Chen, D., Levy, O., Lewis, M., Zettlemoyer, L. & Stoyanov, V. (2019). RoBERTa: A Robustly Optimized BERT Pretraining Approach. *arXiv preprint arXiv:1907.11692*.

- Mikolov, T., Chen, K., Corrado, G., & Dean, J. (2013). Efficient Estimation of Word Representations In Vector Space. *arXiv preprint arXiv:1301.3781*.
- Mordecai, N. B., & Elhadad, M. (2005). Hebrew Named Entity Recognition. *MONEY*, 81(83.93), 82-49.
- More, A., Seker, A., Basmova, V., & Tsarfaty, R. (2019). Joint Transition-Based Models for Morpho-Syntactic Parsing: Parsing Strategies for MRLs And a Case Study From Modern Hebrew. *Transactions of the Association for Computational Linguistics*, 7, 33-48.
- Pedrosa, A. L., Bitencourt, L., Fróes, A. C. F., Cazumbá, M. L. B., Campos, R. G. B., de Brito, S. B. C. S., & Simões e Silva, A. C. (2020). Emotional, Behavioral, And Psychological Impact of The COVID-19 Pandemic. *Frontiers in psychology*, 11, 566212.
- Plutchik, R. (1980). A General Psychoevolutionary Theory of Emotion. In *Theories of emotion* (pp. 3-33). Academic press.
- Radford, A., Narasimhan, K., Salimans, T., & Sutskever, I. (2018). Improving Language Understanding by Generative Pre-Training.
- Sarne, D., Schler, J., Singer, A., Sela, A., & Bar Siman Tov, I. (2019, May). Unsupervised Topic Extraction from Privacy Policies. In *Companion Proceedings of The 2019 World Wide Web Conference* (pp. 563-568).
- Sebastiani, F. (2002). Machine Learning in Automated Text Categorization. *ACM computing surveys (CSUR)*, 34(1), 1-47.
- Seker, A., Bandel, E., Bareket, D., Brusilovsky, I., Greenfeld, R., & Tsarfaty, R. (2022, May). AlephBERT: Language Model Pre-training and Evaluation from Sub-Word to Sentence Level. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)* (pp. 46-56).
- Shmidman, A., Guedalia, J., Shmidman, S., Shmidman, C. S., Handel, E., & Koppel, M. (2022). Introducing BEREL: BERT Embeddings for Rabbinic-Encoded Language. *arXiv preprint arXiv:2208.01875*.
- Sima'an, K., Itai, A., Winter, Y., Altman, A., & Nativ, N. (2001). Building A Tree-Bank of Modern Hebrew Text. *Traitement Automatique des Langues*, 42(2), 247-380.
- Suárez, P. J. O., Romary, L., & Sagot, B. (2020). A Monolingual Approach to Contextualized Word Embeddings for Mid-Resource Languages. *arXiv preprint arXiv:2006.06202*.
- Tsarfaty, R., Seddah, D., Goldberg, Y., Kübler, S., Versley, Y., Candito, M., Foster, J., Rehbein, I., & Tounsi, L. (2010, June). Statistical Parsing of Morphologically Rich Languages (SPMRL) What, How And Whither. In *Proceedings of the NAACL HLT 2010 First Workshop on Statistical Parsing of Morphologically-Rich Languages* (pp. 1-12).

נספח א: שימוש במודל שפה בעברית ליזוי תחושים ורנשות בשפה כתובה

בחלק זה של המאמר נסביר כיצד ניתן לבצע ניתוח תחושים ורנשות על מסמך המכיל משפטים בעברית. בהדרכה זו נניח כי קובץ הקלט נקרא "data.csv" וכי העמודה שנרצה לנתח נקראת "text", כמפורט בדיון באור 1. שמו לב שכך לשמר את העברית בקובץ csv, יש לשמור את הקובץ בפורמט "CSV UTF-8". קוד הנិធות בפרק זה כתוב בשפת פ'יתון, וניתן להריצתו בעזרת מחברת Jupyter⁹ או בכל תוכנת פ'יתון אחרת. בהדרכה נניח כי קוד הנិធות וקובץ הקלט נמצאים באותה תיקייה.

אייר 1: דוגמה לקלט לנិזהות

	A	B	C	
1	rowNum	CommentName	text	
2	1	Gilad	"זיננות מפוזרת ליזבוק בקווינה"	
3	2	Mika	"הסכנות במנת חישון אחת"	
4	3	Sarit	"המצב הבריאותי השתרף בצורה יצאת דופן?"	

על מנת לנתח את הקובץ באור 1, יש לטען מהילה את הספריות הנדרשות ולהרוא את הקובץ לזכרון. בקטע קוד 1 נמצאות שורות הקוד הרלוונטיות.

קוד 1: טיענת ספריות וקובץ לזכרון

```
## Required installations (run only once)
!pip install transformers

## Upload libraries
import pandas as pd
from transformers import pipeline

## Read file
input_path='data.csv'
df=pd.read_csv(input_path)
```

בשלב הבא נוכל להריץ את ניתוח התיחסות (קוד 2) וניתוח הרנשות (קוד 3). שורות הקוד של הנិזהות כוללות שלושה שלבים: (1) טיענת המודל לזכרון (תחושים או רנשות), (2) חישוב התיחסה או הרנש לכל שורה במסד, -(3) שימירת התוצר לקובץ פ'ט. קובצי הפ'ט מופיעים באורוים 2-3 בהתאם, ומיכלים מידע לבני התיחסה (חויבי, שלולי, ניטרלי) ורמת הביטחון בתיחסה זו, וכן לבני הרנש (קיים [LABEL_1] או לא [LABEL_2]), ורמת הביטחון ברנש זה. לצורך ההדגמה, קובץ התיחסות מכיל מידע מעודע על תיחסות שמחה וכעס בלבד. ניתן לראות שה毛病ט הראשון מכיל אלמנט של כעס, ואילו האחרון מכיל אלמנט של שמחה.

קיד 2: ניתוח תחושים

```
# Load sentiment pipeline
sentiment_cls=pipeline(
    'sentiment-analysis',
    model='avichr/heBERT_sentiment_analysis',
    tokenizer='avichr/heBERT_sentiment_analysis',
    device=0) #run on GPU (if there is no GPU
                #installed on your machine, change to 'device = -1')

# Add sentiment score
df = df.join(
    pd.DataFrame([sentiment_cls(df['text'][i])[0] for i in range(len(df))]).rename(columns={'label':'label','score':'confidence'}))

df.to_csv('data_polarity_analysis.csv', encoding='utf-8-sig')
```

אייר 2: פלט ניתוח תחושים

	A	B	C	D	E	
1	rowNum	CommentName	text	label	confidence	
2	1	Gilad	"ווננות מפחדות להידבק בקורונה"	negative	0.9999	
3	2	Mika	"הסכנות במנת חיסון אחת"	negative	0.9373	
4	3	Sarit	"המצב הבריאותי השתפר בצורה ויצאת דופן?"	positive	0.9996	

קיד 3: ניתוח רגשות

```
# Compute each emotion separately
emotions = ['anticipation', 'joy', 'trust', 'fear',
            'surprise', 'anger', 'sadness', 'disgust']

for emotion in emotions:
    # Load emotions pipeline
    sentiment_cls=pipeline(
        'sentiment-analysis',
        model='avichr/hebEMO_' + emotion,
        tokenizer='avichr/heBERT',
        device=0) #run on GPU (if there is no GPU
                    #installed on your machine, change to 'device = -1')

    # Add emotions score
    df = df.join(
        pd.DataFrame([sentiment_cls(df['text'][i])[0] for i in range(len(df))]).rename(columns = {'label':emotion, 'score':emotion+'_confidence'}))

    # Save to file
df.to_csv('data_emotions_analysis.csv', encoding='utf-8-sig')
```

אייר 3: פלט חלקו של ניתוח הרגשות

	A	B	C	D	E	F	G	
1	rowNum	CommentName	text	joy	joy-confidence	anger	anger_confidence	
2	1	Gilad	"ווננות מפחדות להידבק בקורונה"	LABEL_0	0.9999	LABEL_1	0.9997	
3	2	Mika	"הסכנות במנת חיסון אחת"	LABEL_0	0.9999	LABEL_0	0.9991	
4	3	Sarit	"המצב הבריאותי השתפר בצורה ויצאת דופן?"	LABEL_1	0.9374	LABEL_0	0.9999	