

NLP Assignment 3 – Report

Nir Rahav – 316275437

Yael Berkovich – 324879501

המשימה והמודלים

נתונים

המטלה שלנו היא לבצע קלסיפיקציה על ציוצים ב-Twitter. נרצה לתייג האם הציוץ נכתב על ידי דונלד טראמפ או לא. באמצעות עיבוד מקדים של הטקסט ובחינת מודלים שונים נרצה לבצע את תיוג הדאטה. לצורך כך ישנן הנחה מקדימה כי חלק גדול מהציוצים שנכתבו על ידי דונלד טראמפ לא נכתבו על ידו אלא על ידי הצוות שלו. הדאטה שלנו מורכב מהעמודות הבאות:

1. Tweet id
2. User handle
3. Tweet text
4. Time stamp
5. Device

את המחלקות נגדיר בצורה הבאה: ציוצים שנכתבו לפני אפריל 2017, מהיוזר "realDonaldTrump", דרך מכשיר אנדרואיד יוגדרו כדונלד טראמפ. נגדיר עמוד חדשה בשם "author" כאשר לפי היוריסטיקה המוזכרת במידה וציוץ נכתב על ידי טרמפ מוגדר 0 ובמידה והצוות שלו כותב את הציוץ יוגדר 1 כאשר ניתן לראות את החלוקה שלהם בצורה הבאה:

0	2255
1	1250

נרצה לבחון האם מודלים של למידת מכונה מסוגלים לבצע את ההבחנה הזו בין התיוגים השונים בצורה טובה. לצורך כך ישנן שתי הנחות שאנו מסתמכים עליהן

1. לאחר אפריל 2017 לא ניתן לקבוע את כותב הציוץ מכיוון שטראמפ עבר לאייפון
2. אין לנו חוסרים משמעותיים בדאטה שלנו

המודלים

לטובת המשימה נבדקו 5 מודלים:

1. מודל LogisticRegression - חלק מדרישת הבסיס - מודל פשוט, קל לפירוש הסתברותי, ומהיר לחישוב.
2. מודל SVM - חלק מדרישת הבסיס - מודל יותר מורכב חישובית אך בעל ביצועים טובים באופן אמפירי לבעיות קלסיפיקציה פשוטות וגם מורכבות.
3. מודל XGBoost - מודל רובסטי בעל ביצועים טובים.
4. מודל FFNN - חלק מדרישת הבסיס - נבחר בזכות גמישותו ויכולת הלמידה שלו דפוסים לא לינאריים ומורכבים.
5. מודל BERT

Feature Engineering

למודלים השונים הוזנו פיצורים בצורה של וקטור אשר הורכבו ממחלקה עצמאית (FeatureVectorizer). מטרת המחלקה הייתה ליצור וקטור פיצורים שישלב לפי קלט המשתמש בין TF-IDF, Word2Vec ווקטור שנבנה בצורה עצמאית עם הפיצורים הללו:

מספר המילים, מספר האותיות הגדולות, מספר סימני השאלה ("?"), מספר סימני הקריאה ("!"), אורך הממוצע של מילה, האם כל המילה באותיות גדולות, האם נכתב בחג או לא, שעת הציון, היום בשבוע של הציון, החודש שבו נכתב הציון, השנה שבה נכתב הציון

הערכת המודלים

לטובת הבחינה הדאטה סט חולק ל-80% עבור train set ו-20% עבור test set, לכלל המודלים בוצע שימוש באופטימיזציה על מנת לבצע אופטימיזציה להיפר פרמטרים הטובים ביותר, בנוסף במודל SVC בוצע תנאי עצירה מקדים על פונקציית הloss במידה ולא יבוצע שיפור עד כדי אפסילון.

מודל	וקטור	היפר פרמטרים אופטימליים	Accuracy	F1	Recall	Precision
Logistic Regression	tfidf	'C': 9.864948922299297, 'penalty': 'l1'	0.86	0.79	0.73	0.87
SVM	tfidf + additional info	'kernel': 'rbf', 'C': 4.31289165277942, 'gamma': 0.030396135757067042	0.82	0.73	0.69	0.79
XGBoost	tfidf + additional info	n_estimators': 477, 'max_depth': 9, 'learning_rate': 0.016548545657454578, 'subsample': 0.8444703426347393, 'colsample_bytree': 0.6238505376857482, 'reg_lambda': 0.4156047528735853	0.89	0.83	0.78	0.88
FFNN	tfidf	'hidden_dim = 256, n_epoch = 15, dropout = 0.4, batch_size = 512 lr = 1e-3	0.80	0.72	0.69	0.75
BERT	Plain text	epoch = 3, batch_size = 8	0.88	0.81	0.71	0.94

ניתוח

תהליך preprocess

הטקסט והתאריך עברו תהליך מקדים אשר כלל את הצעדים הבאים:

- הסרת תווים מיוחדים והחלפתם בתווים תקינים.
- הרחבת קיצורים באנגלית (לדוגמה: "don't" → "do not").
- החלפת HTML, URL, תאריכים, מספרים וערכים עשרוניים בתוויות מתאימות, וניקוי תווים פגומים.
- מחלק את הטקסט ל-tokens וכולל אפשרות ל:
 - המרת אותיות לקטנות עבור tokens שאינם placeholders.
 - הסרת stopwords (למעט placeholders).
 - ביצוע lemmatization על tokens (למעט placeholders).

ההחלטה להמרה של האותיות לקטנות עבור tokens שאינם placeholders וביצוע של הסרת stopwords נתונה להחלטת היוזר ואינה מחייבת. המרה של טקסט לאותיות קטנות זוהי שיטה מקובלת אך ההחלטה על כך היא לשמר את ה"טון הדיבור" בציוץ שכן זהו יכול להיות מאפיין חשוב לכותב הציוץ.

הגדרות, פיצ'רים ופרמטרים

הפרמטרים הטובים ביותר נקבעו באופטימיזציה בשיטת Cross Validation באמצעות ספריית Optuna, שבדקה טווחים רבים של הייפר פרמטרים ובחירה את המוצלחים ביותר בזמן ריצה קצר, האופטימיזציה כללה פרמטרים כמו penalties, גודל וכמות ה-hidden layers, כמות epochs, עומק מקסימלי ומספר מעריכים. חשוב לציין שבהינתן יותר זמן ומשאבים היה ניתן להגדיל את מספר הניסויים. המודל של BERT מיובא מ Hugging Face בצורה של pretraining לטובת משימת הסיווג הבינארי של מחבר הציוץ, הוא מקודד את הטקסטים באמצעות טוקניזציה, מחלק את הנתונים לסט אימון וסט ולידציה, ומאמן את המודל תוך שימוש ב-Hugging Face Trainer שכולל מעקב אחרי ביצועים, עצירה מוקדמת ושמירה של המודל הטוב ביותר. לאחר האימון, המודל מסוגל לחזות את מחבר הציוץ החדש על סמך טקסט בלבד.

מסקנות ותובנות

המודל הטוב ביותר - מודל XGBoost ומודל BERT, שני המודלים הוציאו תוצאות גבוהות ברוב המדדים אם כי מודל BERT היה עם זמן חישוב גבוה משמעותי, ועל כן אם צריך לבחור את המודל הטוב ביותר מודל XGBoost עדיף בגלל זמני הריצה.

וקטוריזציה - TF-IDF היה יעיל ביותר במודלים שהשתמשו בו, הוא הצליח לשפר ביצועים כאשר שולב עם פיצ'רים מהונדסים, מכיוון שגודל הדאטה קטן יחסית Word2Vec לא היה טוב במיוחד.

היפר-פרמטרים: שימוש ב-Optuna אפשר למצוא היפר-פרמטרים אופטימליים שהובילו לשיפור משמעותי בביצועים.

תובנות מהדאטה: הסרת מילות עצירה ולמטיזציה תרמו לשיפור הביצועים על ידי הפחתת רעשים והכנסת משמעות לשפה הטבעית. עיבוד הנתונים והניקוי היו חיוניים להשגת תוצאות מדויקות ואמינות. הפיצ'רים הנוספים הוסיפו שיפור מסוים, אך לא הצליחו לספק תוצאה חיזוי טובה בפני עצמם. אנו מאמינים שבהינתן יותר דאטה, מודלים גדולים כמו BERT (ואולי גם FFNN) היו מבצעים טוב יותר.