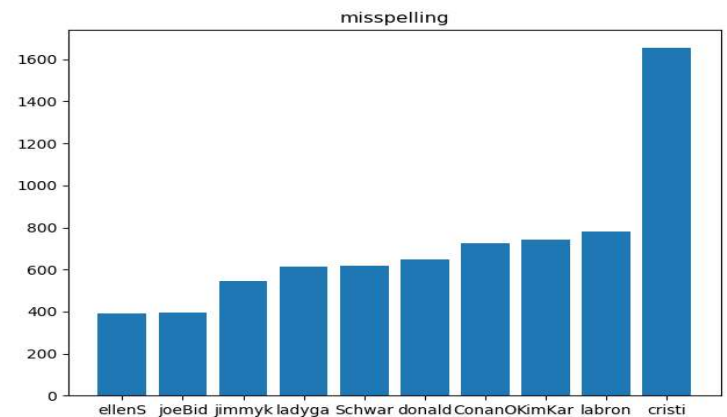
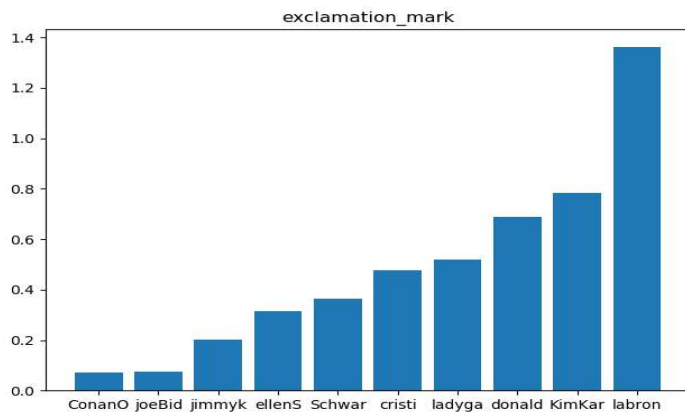
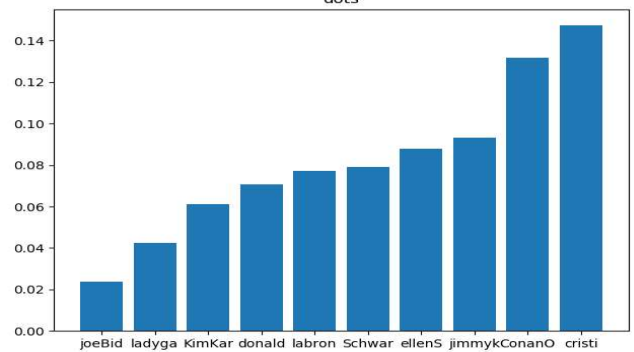
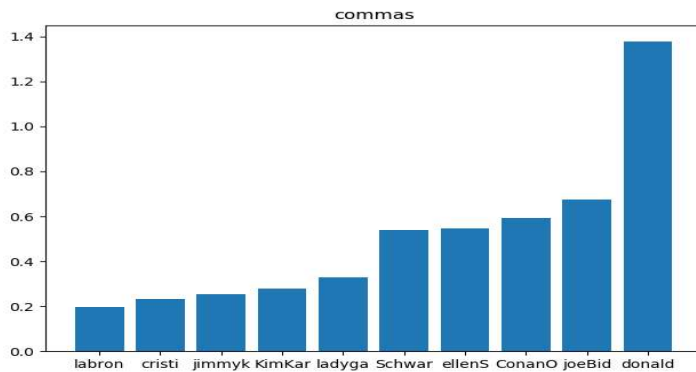
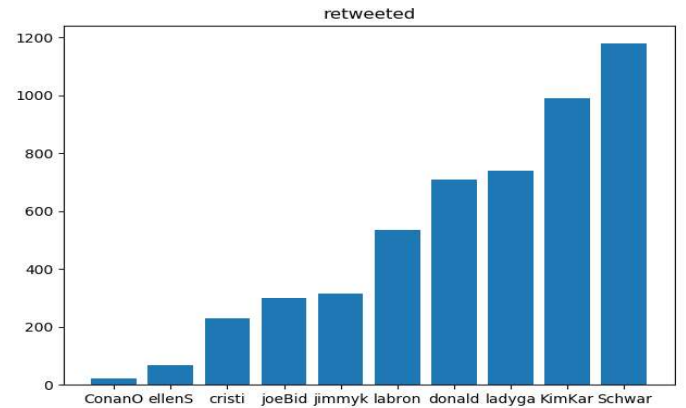
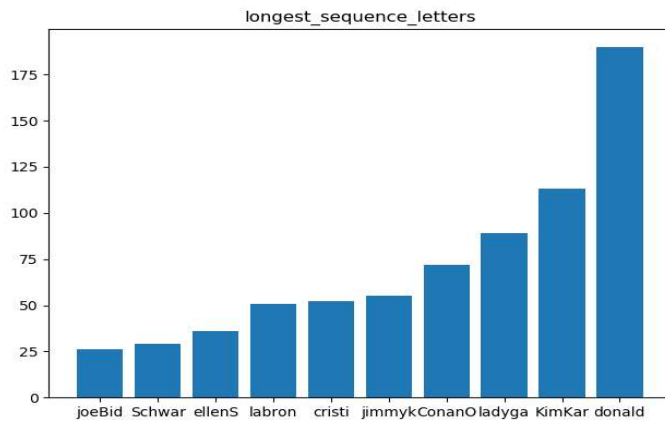
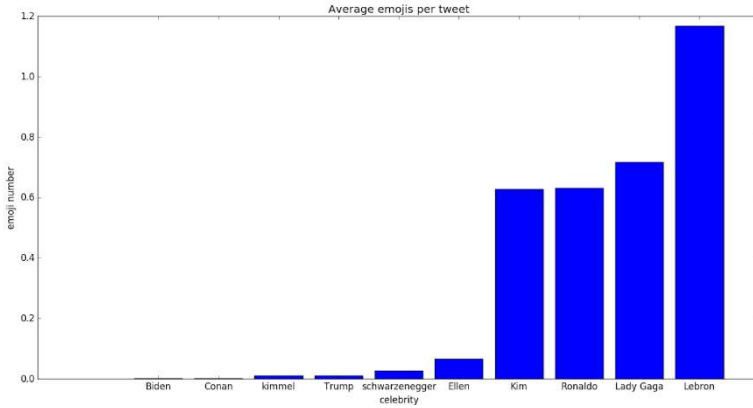


## תהליך העבודה:

התחלנו משקילת הפיצורים שבהם נרצה להשתמש, עשינו זאת בשלבים:

1. בנינו רשימה של כל הפיצורים שחשבנו שיכולים להיות רלוונטיים בתהליך הסיווג.
  2. פירסרנו את הציוצים לפי הפיצורים הללו ובדקנו את ההיסטוגרמה של כל אחד מהם.
- a. להלן ההיסטוגרמות שבלטו לעינינו:



b. ניתן אם כן לראות כי בהיסטוגרמות הללו ההפרשים בין אישויות שונות נראה לעין ולכן תהליך הסיווג עם מאפיינים אלו יהיה יעיל יותר.

לאחר מכן, עברנו ליצירת מאפיינים על בסיס מילים, בדרך הבאה:

1. לכל ציוץ קודדנו וקטור בשיטת BOW. וקטורים אלה מכילים את שכיחות המילה בתוך הטוויט (כלומר, יכולים בהחלט להיות ערכים גדולים מ-1. בדאטא שלנו השכיחות הכי גבוהה בתוך טוויט אחד היתה 10). סך כמות המילים השונות בהם השתמשו המשתמשים בכל הטוויטים (כלומר 'המילון') הוא גדול מאוד.

2. אולם, הוספת למעלה מ-40 אלף מאפיינים לתהליך הסיווג לא היה יעיל במיוחד ופינו לבחון דרכים לסינון המילים. כדי להקטין את גודל הוקטורים (כלומר את גודל הפיצורים) השתמשנו בSTEMMING. שיטה זו מחליפה כל מילה ליגזע' אליו היא שייכת. כך למשל מילים כמו "programming", "programmer" שייכות כולן לגזע "program". בכך, הכנסנו כפי שרצינו עוד BIAS למודל, והקטנו את VARIANCE.

3. בשלב הבא בדקנו קורלציה בין מילים שונות, שמופיעות בכלל הציוצים, לבין אישויות מסויימת. בדרך זו יכלנו לקבוע אילו מילים משפיעות יותר על הניבוי ועל אילו נוכל לוותר כאשר נשתמש בשק המילים בתור מאפיין. יחד עם זאת, קורלציה של מאפיינים יחידים עם הלייבלים התבררה כלא יעילה בשלב בחירת המאפיינים. דבר זה נובע למשל מהעובדה שמילה בודדת אולי לא תשפיע במיוחד על הסיווג, אבל שילוב של כמה מילים, יכול מאוד להתקשר לאדם ספציפי.

• בדקנו אם אפשר להשיג עוד דאטא מטוויטר - לא מצאנו דרך לעקוף את הטוויטר API בלי לשלם וזאת נראתה השקעה קצת מוגזמת:)

### קלסיפיירים שחשבנו לנסות:

Random forest\SVM\Logistic regression

### קלסיפייר נבחר:

בחרנו להשתמש בGLM. ראשית, בהרצות ראשונית על הדאטא הביצועים שלו היו טובים יותר מקלסיפיירים אחרים. שנית, מספר הפיצורים בBOW הוא גדול (כ-40000 לאחר STEMMING), הוא עלול לעבור את מספר נק' הטריינינג דאטא ( $n < d$ ) וגם במקרים שאינם כאלה ליצור בעיה של אוברפיטינג. השימוש ברגרסיה לוגיסטית מאפשר לנו להשתמש גם ברגוריצית LASSO שמורידה באופן משמעותי את הפיצורים. בנוסף לכך, בחרנו לקודד את הדאטא בעזרת BOW בתוספת הפיצורים הידניים שיצרנו עליהם חשבנו מראש (ר' לעיל). כיון שחלק מהפיצורים מתייחסים להופעתה של מילה בודדת (כמו למשל הופעתם של סימני קריאה בטוויט), אנחנו חשופים לבעיית קולינאריות. בעיה זו נפתרת בעקבות איפוס רוב הפיצורים (ובאופן ספציפי אלו הגורמים לקולינאריות). סה"כ הורדנו בעזרת לאסו את מס' הפיצורים לכ-3000 (הדבר תלוי בפרמטר למבדא, בתחילה ניסינו אחד ולאחמ"כ ירדנו לחצי).

בחיפוש אחר classifier מתאים נתקלנו באינטרנט במאמר הבא:

Document author classification using Generalized Discriminant Analysis (moon, howland, gunther)  
שהעלה את האפשרות ששימוש ב-bag of words שמכיל את כל המילים שמופיעות באחד המסמכים ובשיטה של discriminant analysis, יאפשר להגיע לתוצאה אופטימלית. ניסינו שיטה זו והגענו לתוצאות סבירות, אבל פחות טובות משמעותית מאשר עם ה-GLM.

**הערה:** רצינו להשתמש ב-validation k-cross אבל כשניסינו לאמץ את הלומד על כל הדטה הוא רץ הרבה מאוד זמן ולא הספקנו לעלות על הבעיה.