#### המשימה בה בחרנו היא המשימה הראשונה - משימת הסרטים.

#### :תיאור המידע

קיבלנו מידע עם עמודות רבות, אך עם מעט דגימות באופן יחסי. כמו כן, רוב מוחלט מהפיצ'רים היו טקסטואלים, רובם בפורמט JSON. בנוסף, ללא מעט מן הדגימות היה חסר מידע במאפיינים מסוימים. נציין שללא מעט דגימות היה חסר ה revenue עצמו מה שהופך את הדגימה לחסרת משמעות.

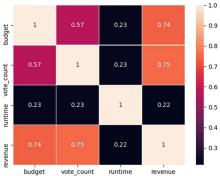
### אתגרים באפיון המידע:

- 1. הבנה כיצד להמיר JSON לתצורה נוחה לעבודה כ-DATA-FRAME. על אחת כמה וכמה כיצד לעבוד עם JSON מקונן. חלק גדול מהזמן היה מושקע בפרסור הטקסט.
  - 2. ניסיון להבין את תורת השימוש ב-free text לטובת בעיות למידה.
  - מילוי חוסרים במידעים של דגימות מסוימות (האם בכלל למלא את שורות אלו, אם כן,
    באיזה ערכים (ממוצע / מינימלי / חציון וכוי)).
- 4. מידע קטגוריאלי רב ברירתי (לדוגמה מספר רב של שחקנים) צריך להחליט אילו שחקנים הם המשמעותיים ביותר וכן לשאר הפיצ'רים בסגנון זה).
- חלק מהפיצ'רים היו תלוים זה בזה, והיה צורך להבין איך והאם הם מיתרים אחד את השני (לדוגמא, מיקום של חברת הפקה והמקום בו צולם הסרט).

# (preprocessing): אופן ניקוי המידע

- 1. הפיכת התאריך למשתנים אינדיקטיביים יותר: יום בשבוע, חודש ושנה.
  - 2. ניקוי שורות שהסתמנו כרעש.
  - companies, genres : הפיכת משתנים רבים לקטגוריאלים, לדוגמה
- בנוסף, את חלק מהמשתנים הקטגוריאלים איחדנו לאור תרומתם המעטה למודל. לדוגמה:
  רק חברות הפקה שמספר הסרטים שלהן בדאטה היה משמעותי, קיבלו קטגוריה משלהם
  (והקטנות אוחדו לקטגוריה אחרת) וכן ריכוז
  כל החודשים לפיצ׳ר של יחיד של חודשים

משמעותיים (חגים לדוגמה).



### <u>שיקולים שהנחו אותנו בעיצוב מודלי הלמידה:</u>

- 1. מספר ניסויים ראשוניים בניסיון למצוא קורלציות (ללא מידע מקדים) ניתן לראות גרף מצורף.
  - 2. בחינת מספר מודלים ליניאריים.

# השיטות שניסינו והתוצאות שהן הניבו:

- 1. ניסינו להשתמש ב-bag of words גם לפיציר overview גם לפיציר bag of words. ניסויים אלה העלו חרס.
  - 2. שימוש ב R.
- 2. כמו כן, לאור כמות המידע המועט, ניסינו להשתמש ב-ensemble. בסוף בחרנו להשתמש ב-random forest
  - לאור כמות הפיצירים הגדולה שנוצרה (לאור הפיכת משתנים רבים ל-dummies), ניסינו להשתמש ב-feature selection, כדי להוריד את כמות הפיצירים.
    בשני הניסיונות, lasso ו-selectKBest לא ראינו תוצאות חד משמעיות.
- Random forest, : על מנת למצוא את המודל הטוב ביותר, ניסינו להשתמש במודלים הבאים: Regression, Gradient boosting מאחר שהוא נתן .Regression, Gradient boosting לנו את התוצאות הטובות ביותר.

# שגיאת ההכללה של החיזוי של המודל שלנו:

במהלך האימון הגענו ל- $R^2$  של כ- 0.8. ול RMSE של כ- 62782878.193. (על הטסט שלא נגענו בו במהלך במהלך האימון הגענו ל- $R^2$  של כ- 0.8. ול שנות את כל האיימון כלל כמובן.) אנחנו חושדים שהיה לנו R סverfit כל האיימון כלל כמובן.) אנחנו חושדים שהיה לנו random-forest על הספקנו להגיע לזה.

לסיכום, רוב ההתעסקות שלנו היתה עם הדאטה, וחילוץ המידע ממנו היה מאוד מאתגר, עקב כך לדעתנו לא הספקנו להתעסק מספיק עם המודל עצמו ולשפר אותו.