

## למידה חישובית – תרגיל 2

מיכל אברמוב, 301834297  
בוריס בורשבסקי, 311898746

### מה מצורף

1. מחברת המפרטת בשלבים את מה שעשינו
2. קובץ main.py המכיל את כל הקוד של התרגיל
3. קובץ html שמאפשר צפיה נוחה יותר במחברת
4. קבצי datan.

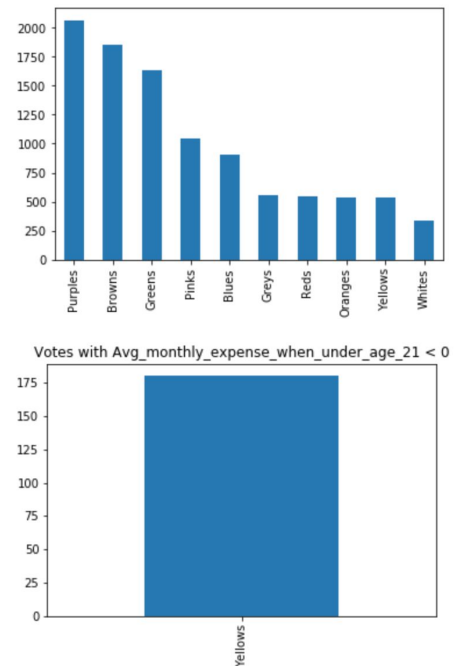
### Flown של התרגיל

#### 1. קריאת וחלוקת את המידע

- a. קראנו את המידע והוספנו עמודה בשם "Split" שתעזור לנו בחלוקה ל-Datasets.
- b. חילקנו את המידע ל60%, 20%, 20%, וייצאנו את raw data אחרי החלוקה (וללא עמודת split) ל3 קבצים:
  - i. Raw\_data\_train
  - ii. Raw\_data\_validation
  - iii. Raw\_data\_test
- c. בהמשך נשתמש בעמודה split כדי לוודא שחלוקת השורות על modified datan לקבצים מתאימה לחלוקת השורות בraw data.

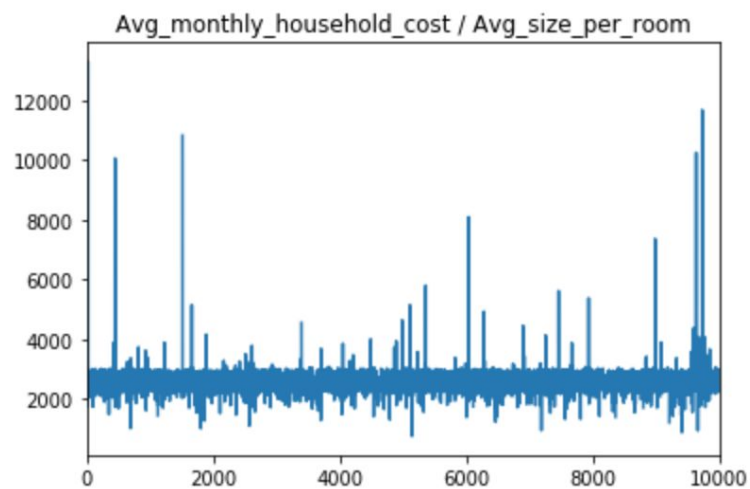
#### 2. מיפוי הפיצורים

- a. מיפינו את הפיצורים הקטגוריאליים והנומריים ע"מ לעבוד עליהם בהמשך
  - b. חיפשנו פיצורים שבהם ישנם ערכים לא הגיוניים, מצאנו שישנם מספר פיצורים עם ערכים שליליים שאינם מסתדרים, למשל:
    - i. Avg\_monthly\_expense\_when\_under\_age\_21
    - ii. AVG\_lottary\_expenses
- אבל ראינו שכל הערכים שקיימת קורלציה בינם לבין המטרה (כולם בחרו בצהוב) אז החלטנו להשאיר אותם במחשבה שהם ינורמלו לאחר מכן - ניתן לראות גרף.



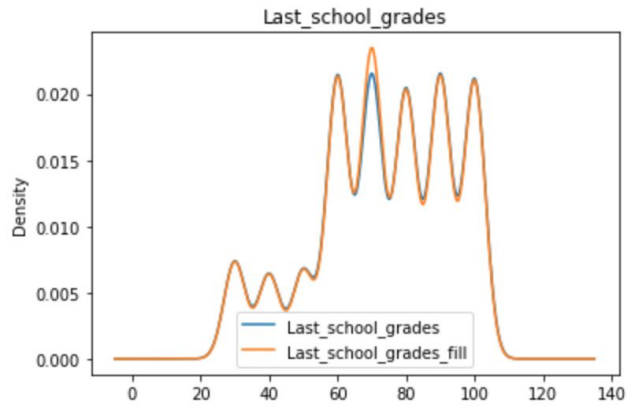
### 3. מילוי ערכים חסרים

a. בהתחלה בדקנו קורלציה בין הפיצ'רים וחיפשנו את כל הפיצ'רים הנומריים ביניהם ישנה קורלציה של לפחות 0.95 - במקרה זה השלמנו את כל הערכים החסרים בעזרת הקורלציה וסימנו את הפיצ'ר שהשלמנו ממנו כלא רלוונטי באיור: דוגמא לפיצ'רים עם קורלציה:



b. לאחר מכן עברנו על כל הפיצ'רים עם קורלציה של 70% ואותם השלמנו גם בעזרתה אך ללא סימון הפיצ'רים

c. לבסוף השלמנו את הערכים שנותרו חסרים ע"י הערך החציוני ובדקנו שלא פגענו בהתפלגות כמו בגרף המצורף:



d. את הערכים הקטגוריאליים השלמנו ע"י הערך הנפוץ ביותר תחת אותו label

4. טרנספורמציה של המידע

a. ערכים קטגוריאליים:

- i. בערכים בהם יש חשיבות לסדר מיפינו את הערכים לפי הסדר
- ii. בערכים בוליאנים הגדרנו ש-NO מייצג 1, yes מייצג 0, maybe 1
- iii. בערכים ללא חשיבות לסדר יצרנו פיצ'רים בוליאני עבור כל אחד מהאפשרויות הקיימות כך שנוצר בעצם עבור כל פיצ'ר קבוצה נוספת של פיצ'רים למשל: עבור Occupation יצרנו ls\_occupation\_hightech וכו'

b. ערכים נומריים:

- i. בחרנו לעשות scale לינארי לכל הערכים כך שכולם יהיו בטווח בין 0 ל-1. (בהנחה שכבר בדקנו קורלציה בין הפיצ'רים)
- c. שינינו את voten לערך נומרי לפי כל צבע

5. ניקוי רעשים

- a. ביצענו ניקוי רעשים בעזרת outlier detection ע"מ לנקות את כל המופעים שבהם ערך הפיצ'ר רחוק בלפחות 3 סטיות תקן מהערך הממוצע של אותו פיצ'ר.
- b. הסתכלנו על כל השורות שבהם הורדנו פיצ'ר אך לא מצאנו משהו חשוד.

6. בחירת פיצ'רים

- a. בשלב זה הורדנו את הפיצ'רים שסימנו מיותרים בשלב הקורלציה והורדנו מופעים בהם יש עדיין ערכים ריקים
- b. השתמשנו בVarianceThreshold על מנת להעיף את כל הפיצ'רים שאין להם שונות (לא הצלחנו להעיף כלום)
- c. שלב ה Filter Method:
  - i. הרצנו select Percentile עבור הפיצ'רים בעזרת f-classif ובעזרת mutual-information בחרנו את הפיצ'רים ב-25 אחוז העליונים בשני המקרים
- d. שלב ה Wrapper Method:
  - i. השתמשנו בRFE בשני ריצות שונות, אחת עם 3 cross validation ובשניה עם cross validation מסוג StratifiedKFold הוספנו את כל הפיצ'רים שנוספו בכל אחת מהריצות.
- e. השתמשנו ב tree based feature selection, גם שם בחרנו את 25% הפיצ'רים עם tree weight הגבוה ביותר והוספנו לרשימת הפיצ'רים הנבחרים.

## 7. שמירת המידע

a. שמרנו את המידע המעובד לקבצים לפי עמודת split שהגדרנו בהתחלה לקבצים

- i. processed\_test.csv
- ii. processed\_train.csv
- iii. processed\_validation.csv

## 8. הפיצורים שנבחרו לבסוף:

a. פיצורים חדשים:

- i. Married'
- ii. 'Avg\_monthly\_income\_all\_years'
- iii. 'Is\_Most\_Important\_Issue\_Financial'
- iv. 'Yearly\_IncomeK'
- v. 'Number\_of\_valued\_Kneset\_members'
- vi. 'Is\_Most\_Important\_Issue\_Environment'
- vii. 'Garden\_sqr\_meter\_per\_person\_in\_residency\_area'
- viii. 'Will\_vote\_only\_large\_party'
- ix. 'Avg\_monthly\_expense\_when\_under\_age\_21'
- x. 'Is\_Most\_Important\_Issue\_Military'
- xi. 'Is\_Most\_Important\_Issue\_Education'
- xii. 'Is\_Most\_Important\_Issue\_Foreign\_Affairs'
- xiii. 'AVG\_lottary\_expanses'
- xiv. 'Is\_Most\_Important\_Issue\_Other'
- xv. 'Last\_school\_grades'
- xvi. 'Looking\_at\_poles\_results'
- xvii. 'Weighted\_education\_rank'
- xviii. 'Overall\_happiness\_score'
- xix. 'Is\_Most\_Important\_Issue\_Social'

b. פיצורים מקוריים (שעליהם מבוססים הפיצורים החדשים):

- i. Number\_of\_valued\_Kneset\_members'
- ii. 'Avg\_monthly\_expense\_when\_under\_age\_21'
- iii. 'Avg\_monthly\_income\_all\_years'
- iv. 'Yearly\_IncomeK'
- v. 'Married'
- vi. 'Last\_school\_grades'
- vii. 'Weighted\_education\_rank'
- viii. 'Garden\_sqr\_meter\_per\_person\_in\_residency\_area'
- ix. 'Will\_vote\_only\_large\_party'
- x. 'Looking\_at\_poles\_results'
- xi. 'Overall\_happiness\_score'
- xii. 'Most\_Important\_Issue'
- xiii. 'AVG\_lottary\_expanses'

