

# The election challenge

1. בשלב ראשון טענו את המידע והורדנו את כל השורות בהן היו חסרים ערכים. הרצת describe הראתה לנו שיש שדות עם ערכים שליליים, שלא הגיוני שיהיו שם. לאחר שהורדנו גם את כל השורות עם הערכים השליליים נשארנו עם 8146 שורות.
  2. Feature selection - מאחר והמטרה היא לנקות פיצ'רים בצורה שמרנית, החלטנו לחפש לכל פיצ'ר סיבות שונות להשאיר אותו, ואת הפיצ'רים שלא מצאנו סיבה להשאיר, נוריד. התחלנו עם מבחן  $\chi^2$  לערכים הנומינליים, ומבחן ANOVA לערכים המספריים. החלטנו לשמור את הפיצ'רים שה-p-value שלהם קטן מ-0.05. לאחר מכן הדפסנו את פונקציות הצפיפות של המשתנים הרציפים, וההיסטוגרמות של המשתנים הדיסקרטיים, ולמדנו ש:
    - בפיצ'ר Last\_school\_grades כל הערכים מ-60 ומעלה מתנהגים אותו הדבר, לכן אפשר לצמצם אותם לערך אחד.
    - הערכים Avg\_monthly\_expense\_when\_under\_age\_21 ו-Avg\_Residency\_Altitude תלויים לינארית אחד בשני, לכן אפשר להוריד אחד מהם, אבל לפני זה כדאי להשתמש בו להשלמת ערכים עבור השני. במימוש הסופי הורדנו עוד שדות שתלויים לינארית.
    - יש התאמה בין ההיסטוגרמות של Last\_school\_grades ו-Most\_Important\_Issue
    - יש פיצ'רים שמתפלגים פואסונית (כמו Num\_of\_kids\_born\_last\_10\_years) ויש שמתפלגים (בקירוב כמובן) אקספוננציאלית (כמו AVG\_lottary\_expenses) אבל החלטנו לא להתייחס אליהם בצורה מיוחדת.
    - מסקנות נוספות נמצאות בnotebook לצד הגרפים הרלוונטיים.
  3. האלגוריתם:
    - a. מנקים את המידע.
    - b. מחפשים פיצ'רים דיסקרטיים עם התאמה בין הערכים, משלימים ערכים חסרים וזורקים אחד מהם (אם רק אחד מהם נומרי נעדיף לשמור אותו).
    - c. מחפשים פיצ'רים נומריים עם קורלציה לינארית גבוהה בכל הלייבלים, משלימים ערכים חסרים וזורקים אחד מהם.
    - d. בוחנים את הפיצ'רים הנומינליים בעזרת  $\chi^2$  ואת המספריים בעזרת ANOVA (פיצ'רים שמתפלגים לינארית הופכים קודם כל לנורמלים). ניסינו גם להריץ ANOVA על log values אבל לא קיבלנו פיצ'רים נוספים שכדאי לשמור.
    - e. מריצים SFS עם מספר מסווגים שונים, אם יש פיצ'ר שנבחר ע"י אחד מהם נוסף אותו לבחירה שלנו. (היחיד שנבחר הוא gender, אנחנו לא בטוחים לגביו אבל מכיוון שעדיף להשאיר פיצ'רים מיותרים מאשר לזרוק פיצ'רים חשובים אז נשאיר אותו)בפועל הרצנו כמה פעמים ובכל פעם הפיצ'רים שהתווספו בגלל SFS היו שונים (בין היתר קיבלנו את Avg\_education\_importance, Occupation, Occupation\_Satisfaction, %Time\_invested\_in\_work), ומהסתכלות בגרפים שלהם נראה שההבדלים בין הלייבלים די מקריים, לכן אנחנו חושבים שעדיף בלעדיהם.
- מכיוון ש Number\_of\_valued\_Kneset\_members לא חזר מאף ריצה של SFS, ומראש הוא נכנס עם p-value די גבוה (ב-0.0219284964384) ANOVA) אנחנו חושבים שעדיף לוותר עליו. בסוף האלגוריתם הדפסנו השוואה של התוצאות עם כל הפיצ'רים שבחרנו לעומת התוצאות כאשר מורידים כל אחד מהפיצ'רים בקבוצת הפיצ'רים שלא נבחרו באף הרצה של SFS. גם שם אנחנו רואים שהורדה שלו לא פוגעת בתוצאות. הפיצ'רים איתם נשארו:
- 'Political\_interest\_Total\_Score', 'Yearly\_ExpensesK', 'Avg\_monthly\_household\_cost', 'Number\_of\_valued\_Kneset\_members', 'Yearly\_IncomeK', 'Married', 'Last\_school\_grades',

'Avg\_monthly\_expense\_when\_under\_age\_21', 'Looking\_at\_poles\_results',  
'Financial\_agenda\_matters', 'Overall\_happiness\_score', 'AVG\_lottary\_expanses'

המשמעות של הפיצ'רים שבחרנו ביחס ללייבל (כפי שניתן לראות מהגרפים בנוטבוק):

- Political\_interest\_Total\_Score, Avg\_monthly\_household\_cost, Yearly\_IncomeK: בכל המפלגות התוחלת נראית דומה, אבל במפלגות מסוימות ההתפלגות מאוד צפופה סביב התוחלת (הכי קיצוני במפלגה 3, ) ובמפלגות אחרות, כמו 1,5,8, השונות הרבה יותר גדולה. 3 הפיצ'רים האלה מאוד דומים וכנראה שאפשר לוותר על חלקם, גם Yearly\_ExpensesK ו-Overall\_happiness\_score נראים די דומה, אבל פחות. יתכן ש Yearly\_ExpensesK, Overall\_happiness\_score, ואולי פיצ'רים נוספים, נגזרים מ Yearly\_IncomeK עם רעש מסוים. עם המסווגים שהרצנו בסוף, הורדה של Avg\_monthly\_household\_cost לא פגעה בציונים, לכן אנחנו חושבים שאפשר להוריד אותו. גם כשמורידים את Political\_interest\_Total\_Score הפגיעה היא מינורית מאוד, אבל נהיה זהירים ונשאיר אותו.
  - Married – רווקים נוטים מאוד להצביע למפלגות 2 ו-5, הרבה יותר מנשואים. נשואים לעומת זאת מצביעים ל-4 ו-9 הרבה יותר מרווקים.
  - Last\_school\_grades – אנשים עם ציון 30 מצביעים רק ל-5,8. אנשים עם ציון 40 מצביעים רק ל-2,8. אנשים עם ציון 50 מצביעים רק ל-2,5. ההצבעה בשאר הציונים מתחלקת די אחיד בין שאר המפלגות (ביחס לגודלן כמובן, הגרפים מנורמלים). כמו כן מצאנו התאמה מלאה בין כל ציון לאחת הקטגוריות ב-Most\_Important\_Issue. העדפנו את הפיצ'ר המספרי, אבל אולי למספרים אין באמת משמעות.
  - Looking\_at\_poles\_results – הפוך בדיוק מ-Will\_vote\_only\_large\_party. האלגוריתם שלנו הוריד במקרה את Will\_vote\_only\_large\_party אבל זו בדיוק המשמעות, מצביע מסתכל בסקרים אמ"מ הוא לא מצביע למפלגה גדולה, אמ"מ הוא יצביע בפועל למפלגה קטנה (1 או 6). בתוך כל קטגוריה ההצבעות מתחלקות אחיד בין כל המפלגות (שוב, ביחס לגודלן).
  - Financial\_agenda\_matters – מפריד בדיוק בין מפלגות 2,5,6,8 לשאר.
  - Avg\_monthly\_expense\_when\_under\_age\_21 – יש התאמה לינארית מלאה בינו לבין Avg\_Residency\_Altitude, לכן האלגוריתם זרק את Avg\_Residency\_Altitude. מהגרף רואים שהפיצ'ר מפריד טוב ל-3 קבוצות: 1,6, וכל השאר.
- מיפוי המספרים לשמות:

0=Blues  
1=Browns  
2=Greens  
3=Greys  
4=Oranges  
5=Pinks  
6=Purples  
7=Reds  
8=Whites  
9=Yellows

רשימת הפיצ'רים הסופית שבחרנו:

'Political\_interest\_Total\_Score', 'Yearly\_ExpensesK', 'Married', 'Yearly\_IncomeK',  
'Last\_school\_grades', 'Looking\_at\_poles\_results',  
'Avg\_monthly\_expense\_when\_under\_age\_21', 'Financial\_agenda\_matters',  
'Overall\_happiness\_score', 'AVG\_lottary\_expanses'

אנחנו מאמינים שגם ברשימה הזאת יש יתירות, אבל לשם הזהירות נשאיר אותה כך.

בנוסף, מצורף קובץ "looking at the data" ובו הגרפים שהשתמשנו בהם (קיימת גם גרסאת html לקריאה נוחה ללא צורך בהרצה).