

תרגיל בית 4 מבוא למערכות לומדות

דוח עבודה

מגישים:

אביב כספי – 311136691

יקיר יהודה - 205710528

שליב עבודה:

1. Pre-processing:

דבר ראשון שעשינו היה ביצוע עיבוד ראשוני לדאטה שלנו.
מימוש חלק זה נמצא בקובץ preprocessing.py.

- א. המרה ל one-hot של הפיצ'רים הנומינליים ללא חשיבות לסדר :
מתוך הפיצ'רים שנבחרו היה עלינו לשנות רק את 'Most_Important_Issue' ל-one-hot.
- ב. המרת כל הפיצ'רים הקטגוריאליים למספריים:
על מנת שנוכל לעבוד על הדאטה, עלינו לשנות את כל הפיצ'רים שיהיו מספריים.
- ג. תיקון ערכים שליליים בפיצ'רים:
בדומה לתרגיל קודם, גם כעת ישנו פיצ'ר בעל ערכים שליליים למרות שלא ייתכן כי הערכים של הפיצ'ר יהיו שליליים.
בצורה דומה לתרגיל הקודם, התמודדנו עם בעיה זו על ידי ביצוע abs על הפיצ'ר הנ"ל. (כמו בתרגיל קודם, זיהינו כי התפלגות הדאטה השלילי מתאימה בדיוק להתפלגות הפיצ'ר בערכים החיוביים, לכן החלטנו כי ייתכן כי בזמן הדגימה שונה הסימן בטעות, לכן החלטנו רק לבצע ערך מוחלט על הדאטה הנ"ל)
- ד. לאחר מכן ביצענו פיצול לדאטה שלנו:
בדומה לתרגיל הקודם בחרנו את הפיצול הבא – train(70%), val(15%), test(15%)
- ה. מחיקת outliers:
בשלב זה מצאנו את הדוגמאות בהן ערך הפיצ'ר במרחק גדול מ-4.5 סטיות תקן מהממוצע של הפיצ'ר, והצבנו NaN במקומות אלה.
פעולה זו ביצענו רק על סט האימון, והשתמשנו בלייבלים של הדאטה על מנת למצוא outliers יחסית ללייבל.
- ו. Imputation :
השלמת ערכים חסרים ביצענו בנפרד לסט האימון ולסט הולידציה והבדיקה - עבור סט האימון בחרנו בשיטה הנקראת Bootstrapping שלמדנו בכיתה.
בשיטה זו, עבור כל ערך חסר, השלמנו על ידי כך שדגמנו מתוך הדאטה שלנו עם חשיבות ללייבל, ערך חדש.
עבור סט הולידציה והבדיקה בחרנו לטפל בצורה שונה.
תחילה לקחנו את סט האימון, ועבור כל פיצ'ר חישבנו את ממוצע הפיצ'ר עבור ערכים נומריים ואת הערך הנפוץ ביותר עבור ערכים נומינליים (ללא חשיבות ללייבל).
לאחר מכן, השלמנו את סט הולידציה והבדיקה בעזרת הערכים שחושבו מסט האימון.
- ז. Scaling :
עבור פיצ'רים יוניפורמים ביצענו scaling לטווח [1 -1]
עבור פיצ'רים נורמלים ביצענו נורמליזציה עם ממוצע 0 וסטיות תקן 1.

2. מציאת קואליציות אפשריות על ידי מודל clustering:

על מנת לבצע חלק זה, בחרנו בשני מודלים של clustering אשר בעזרתם ננסה לבנות קואליציה טובה.

המודלים שבחרנו הינם: KMeans , GaussianMixture.

שלב 1: מציאת מספר הקלאסטרים האופטימלי עבור כל מודל.

שלב זה ביצענו על ידי ביצוע CV על מספר הקלאסטרים שכל מודל מחפש, ובנינו פונקציית ניקוד אשר מנקדת כל קלאסטר שהוצע על ידי המודל ולפי כך בחרנו את הערך הטוב ביותר. פונקציית הניקוד – בהינתן מודל מאומן, וסט ולידציה, תחילה ביצענו חיזוי עבור כל נקודה בדאטה, לאיזה מהקלאסטרים היא שייכת.

לאחר מכן עבור כל קלאסטר, עברנו על כל המפלגות שחזינו כי ישנם מצביעים הנמצאים באותו קלאסטר. עבור כל מפלגה בדקנו כמה אחוז מכלל המצביעים למפלגה זו נמצא בתוך הקלאסטר, אם אחוז זה גדול מ-60% (כלומר רוב המצביעים למפלגה נמצאים בקלאסטר זה) הוספנו אחוז זה לציון הקלאסטר. (אם אחוז המצביעים למפלגה מסוימת היה קטן מ-60% התעלמנו ממפלגה זו כי אינה שייכת לקלאסטר).

פונקציית ניקוד זו, בעצם מקשרת בין הקלאסטרים למפלגות הבעיה שלנו, כלומר ככל שהניקוד גדול יותר, כך יותר מפלגות שייכות ברובן לקלאסטר מסוים. לדוגמא אם בחנו אלגוריתם מסוים עם 3 קלאסטרים וקיבלנו כי מפלגה A נמצא בצורה שווה בכל קלאסטר (30% בכל קלאסטר), מודל זה לא יקבל ניקוד עבור המפלגה הנ"ל. אך אם מפלגה B שייכת בעיקר לקלאסטר מסוים (מעל 60%), המודל יקבל ציון לפי הגודל היחסי של המפלגה השייך לקלאסטר.

בסוף שלב זה קיבלנו כי שני המודלים מקבלים ניקוד אופטימלי עבור 2 קלאסטרים.

שלב 2: מציאת קואליציות אפשריות על ידי המודלים שמצאנו.

בשלב זה עבור כל מודל שמצאנו ביצענו:

עבור כל קלאסטר שהמודל מצא –

מצאנו את כל המפלגות אשר אחוז שייכות המצביעים לכל מפלגה בקלאסטר גדול מ-85% (בחרנו ערך זה מפני שהוא גבוה מספיק כדי להגיד שרוב גדול של מצביעים למפלגה שייכים לקלאסטר זה) לאחר מכן ספרנו את גודל סך המצביעים לכל מפלגות אלה, אם גודל זה גדול שווה ל-51%, הוספנו את המפלגות הנ"ל כקואליציה לרשימת הקואליציות האפשריות.

בסוף שלב זה קיבלנו רשימה של קואליציות אפשריות, כך שכל קואליציה מכילה מספר מפלגות שסך המצביעים הכולל למפלגות אלה גדול מ-51%.

שלב 3: מציאת קואליציה הומוגנית ביותר ושונה ביותר מהאופוזיציה.

בשלב זה עבור כל קואליציה שמצאנו בשלב הקודם, חישבנו את שונות הפיצ'רים של מצביעה

(המצביעים של הקואליציה – מצביעים לכל אחת מהמפלגות בקואליציה)

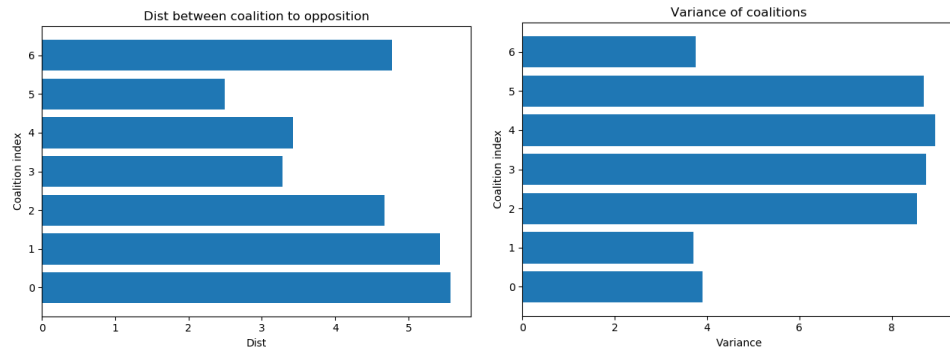
ובנוסף חישבנו את המרחק של ממוצע הפיצ'רים של הקואליציה לממוצע הפיצ'רים של האופוזיציה.

(כלומר כמה שונה האופוזיציה מהקואליציה)

ביצענו מיון על כל אחת מהרשימות הנ"ל (שונות ומרחק), והסתכלנו על התוצאות, כך שהרשימה של השונות ממוינת מהקטן לגדול (ערך קטן יותר – קואליציה הומוגנית יותר) והמרחק ממיון מהגדול לקטן (ערך גדול יותר – קואליציה שונה מהאופוזיציה יותר).

לאחר הסתכלות על התוצאות בחרנו בקואליציה אשר מיקומה היה בראשונים בכל רשימה, כך שערך השונות שלה קטן והמרחק מהאופוזיציה גדול.

הגרפים שקיבלנו עבור השונות והמרחק של הקואליציות:



מתוך רשימת הקואליציות הבאות:

[(0), (12, 11, 10, 9, 7, 6, 4, 2, 1, 0), (11, 9, 8, 6, 5, 4, 3, 2), (12, 11, 9, 8, 6, 5, 4, 3, 2), (9, 8, 6, 5, 4, 3, 2), (9, 7, 6, 4, 2, 1, 0), (10, 9, 7, 6, 4, 2, 1, 0), (11, 10, 9, 7, 6, 4, 2, 1, 0)]

כאשר:

'Violets' - 10
'Browns' - 1
'Greens' - 2
'Whites' - 11
'Reds' - 8
'Turquoises' - 9
'Purples' - 7
'Greys' - 3
'Blues' - 0
'Khakis' - 4
'Oranges' - 5
'Pinks' - 6
'Yellows' - 12

הקואליציה שנבחרה היא קואליציה 0

['Greens', 'Greys', 'Khakis', 'Oranges', 'Pinks', 'Reds', 'Turquoises', 'Whites', 'Yellows']

אשר מכילה 9 מפלגות מתוך 13.

בעלת מרחק מקסימלי ושונות כמעת מינימלית.

הערות:

בחיפוש הקואליציות האפשריות, ביצענו אימון של המודל על סט האימון, ולאחר מכן מצאנו את הקואליציות האפשריות מתוך סט הולידציה וגם מתוך סט האימון (ביצענו חיתוך בין התוצאות) זאת על מנת לזהות כי אכן הקואליציות שקיבלנו על ידי אימון על סט האימון, מתאימות גם לסט הולידציה. (סוג של בדיקת שפיות, על מנת לוודא שהמסווג שלנו עובד בצורה טוב מספיק).

בנוסף, לאחר בחירת הקואליציה האופטימלית, השתמשנו במודל חיזוי מהתרגיל הקודם (RandomForest) וביצענו חיזוי ללייבים של סט הבדיקה.

לאחר מכן ביצענו את אותו תהליך של בחירת קואליציה על סט הבדיקה, ווידאנו כי תוצאת הקואליציה מתאימה לקואליציה שמצאנו בחלקים הקודמים. (בדיקה נוספת כי המודל שלנו יודע להכליל גם על ידי חיזוי הלייבלים).

3. מציאת קואליציות אפשריות על ידי מודל generative:

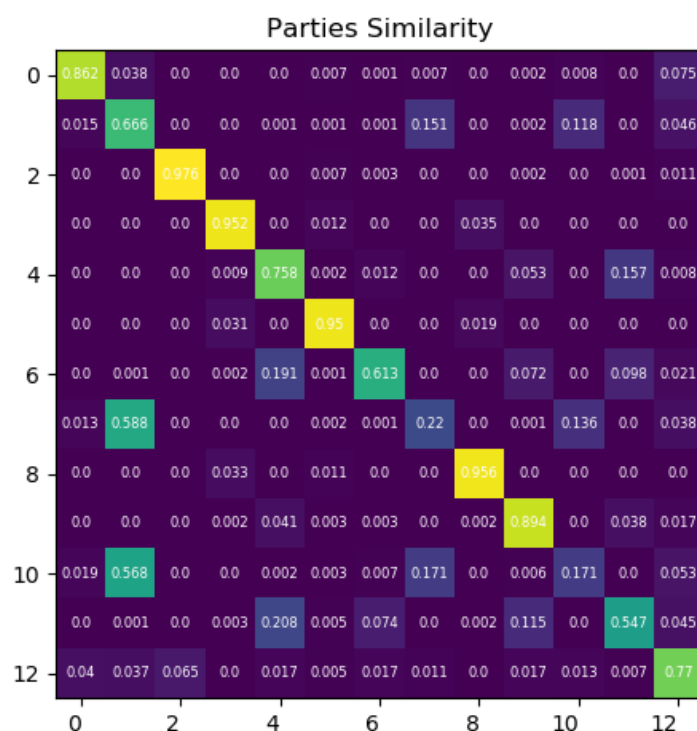
ראשית ביצענו k-folds cross validation על 2 מודלים גנרטיביים מסוג GaussianNB ו LinearDiscriminantAnalysis ובדקנו מי מקבל ציון גבוה ביותר בחיזוי הסיווג (לפי k-folds), קיבלנו GaussianNB מקבל דיוק גבוה יותר.

לאחר מכן בנינו מטריצת הסתברויות לפי המודל הטוב יותר (GaussianNB) על ידי שימוש בתכונה כי מודל גנרטיבי משערך את התפלגות הדאטה שלנו, השתמשנו בפונקציה predict_proba לבניית המטריצה.

מטריצה זו נבנתה בצורה שבה עבור כל מפלגה בדקנו את המצביעים אשר הצביעו למפלגה זו ומה ההסתברות לפי המודל שהם יצביעו למפלגה אחרת, זאת אומרת עבור כל מצביעי מפלגה A מה ההסתברות שהם הצביעו לכל אחד משאר המפלגות ובהתאם ממלאים את ערכים אלו במטריצה.

לאחר מכן מקבלים מטריצה אשר בכל איבר i, j מקבלים מה ההסתברות שבוחרי i יצביעו למפלגה j , ועל ידי קשר זה אנו בונים את האפשרויות לקואליציה. מטריצה זו, נותנת לנו את הדמיון הקיים בין המפלגות (ככל שהערך במיקום מסוים (A,B) גדול יותר, כך יש הסתברות גבוהה יותר כי אדם אשר הצביע למפלגה A יצביע למפלגה B ולכן נוכל להסיק כי ישנו דמיון בין המפלגות)

המטריצה שקיבלנו הינה:



כעת אנו עוברים על כל האפשרויות לקואליציה שבה מתחילה ממפלגה A בודדת ואז מוסיפים לה את המפלגה B שהבחרים שלה הם עם הסתברות הכי גדולה לבחור במפלגה A. (דמיון גדול ביותר), נמשיך כך עד שמגיעים לפחות ל 51% מהמצביעים, כעת נוסיף קבוצת מפלגות זו לרשימת הקואליציות האפשריות.

לאחר שמגיעים ל 51% ממשיכים להוסיף מפלגות באותה צורה עד שלא נשאר מפלגות להוסיף וכל אפשרות כזו של קואליציה נכניס לרשימת הקואליציות האפשריות.

כעת, יש בידינו רשימה של קואליציות אפשריות, נרצה לבחור מתוכן את הקואליציה הטובה ביותר, כאשר נגדיר טובה ביותר בצורה הבאה:
קואליציה כמה שיותר גדולה, בעלת שונות קטנה, ומרחק כמה שיותר גדול מהאופוזיציה. לאחר

חישוב שלושת הערכים הנ"ל, ביצענו מיון כל רשימה (בצורה דומה לבחירת קואליציה בסעיף הקודם) ובחרנו בצורה ידנית את הקואליציה הטובה לדעתנו.

הקואליציה שבחרנו על ידי הנתונים מהמודל היא:

['Greys', 'Khakis', 'Oranges', 'Pinks', 'Reds', 'Turquoises', 'Whites', 'Yellows']

נשים לב כי זאת כמעט אותה קואליציה כמו שקיבלנו מהמודל של cluster למעט המפלגה Greens

4. זיהוי פיצ'רים מובילים עבור כל מפלגה:

על מנת לבצע חלק זה, השתמשנו בשיטת של feature selection אשר השתמשנו בתרגיל בית 2. שיטה זו הינה שיטת Wrapper אשר משתמשת במודל ExtraTreeClassifier המאפשר לקבל מידת חשיבות לכל פיצ'ר.

מימשנו חלק זה על ידי אימון מודל מסוג זה, על משימת סיווג בינארית בין מפלגה יחידה לכל השאר. כלומר, ביצענו עבור כל מפלגה:

חלוקת הדאטה לall vs one, כל שהמפלגה הנוכחית היא 1 וכל שאר המפלגות הם 0.

לאחר מכן אימנו את המודל על דאטה זה, והשתמשנו באופציה של feature_importances_ מנת לקבל ניקוד עבור כל פיצ'ר.

ציון זה בעצם מציין כמה פיצ'ר זה מפריד בין הדוגמאות המתאימות למפלגה שלנו לאלה שלא שייכות אליה.

לאחר הסתכלות על התוצאות הללו, בחרנו את מספר הפיצ'רים בעלי החשיבות הגבוה ביותר להיות הפיצ'רים המובילים של מפלגה זו.

להלן התוצאות עבור כל מפלגה:

מפלגה	פיצ'רים מובילים	Features importance
Greens	<ol style="list-style-type: none"> 1. Overall_happiness_score 2. Weighted_education_rank 3. Avg_size_per_room 4. Political_interest_Total_score 	
Greys	<ol style="list-style-type: none"> 1. Most_Important_Issue_Military 2. Most_Important_Issue_Foreign_Affairs 3. Overall_happiness_score 	

<p>Feature importance of Oranges vs all</p>	<ol style="list-style-type: none"> 1. Most_Important_Issue_Military 2. Political_interest_Total_score 3. Most_Important_Issue_Other 	<p>Oranges</p>
<p>Feature importance of Pinks vs all</p>	<ol style="list-style-type: none"> 1. Weighted_education_rank 2. Avg_size_per_room 3. Number_of_differnt_parties_voted_for 4. Political_interest_Total_score 5. Overall_happiness_score 	<p>Pinks</p>
<p>Feature importance of Purples vs all</p>	<ol style="list-style-type: none"> 1. Weighted_education_rank 2. Avg_size_per_room 3. Political_interest_Total_score 	<p>Purples</p>
<p>Feature importance of Reds vs all</p>	<ol style="list-style-type: none"> 1. Most_Important_Issue_Other 2. Most_Important_Issue_Foreign_Affairs 3. Overall_happiness_score 	<p>Reds</p>

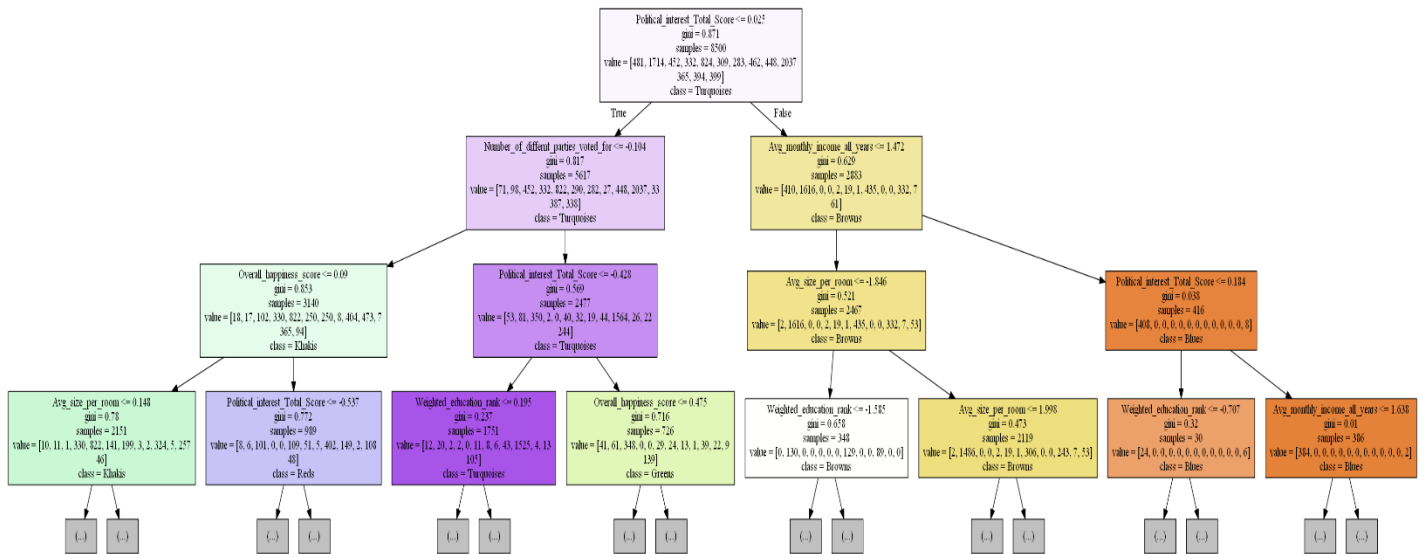
<p>Feature importance of Turquoises vs all</p> <table><thead><tr><th>Feature</th><th>Importance</th></tr></thead><tbody><tr><td>Most_Important_Issue_Social</td><td>0.01</td></tr><tr><td>Most_Important_Issue_Other</td><td>0.02</td></tr><tr><td>Most_Important_Issue_Military</td><td>0.01</td></tr><tr><td>Most_Important_Issue_Healthcare</td><td>0.01</td></tr><tr><td>Most_Important_Issue_Foreign_Affairs</td><td>0.02</td></tr><tr><td>Most_Important_Issue_Financial</td><td>0.01</td></tr><tr><td>Most_Important_Issue_Environment</td><td>0.01</td></tr><tr><td>Most_Important_Issue_Education</td><td>0.01</td></tr><tr><td>Weighted_education_rank</td><td>0.15</td></tr><tr><td>Avg_size_per_room</td><td>0.05</td></tr><tr><td>Overall_happiness_score</td><td>0.08</td></tr><tr><td>Avg_monthly_income_all_years</td><td>0.01</td></tr><tr><td>Avg_Satisfaction_with_previous_vote</td><td>0.11</td></tr><tr><td>Political_interest_Total_Score</td><td>0.21</td></tr><tr><td>Number_of_differnt_parties_voted_for</td><td>0.20</td></tr><tr><td>Yearly_IncomeK</td><td>0.02</td></tr></tbody></table>	Feature	Importance	Most_Important_Issue_Social	0.01	Most_Important_Issue_Other	0.02	Most_Important_Issue_Military	0.01	Most_Important_Issue_Healthcare	0.01	Most_Important_Issue_Foreign_Affairs	0.02	Most_Important_Issue_Financial	0.01	Most_Important_Issue_Environment	0.01	Most_Important_Issue_Education	0.01	Weighted_education_rank	0.15	Avg_size_per_room	0.05	Overall_happiness_score	0.08	Avg_monthly_income_all_years	0.01	Avg_Satisfaction_with_previous_vote	0.11	Political_interest_Total_Score	0.21	Number_of_differnt_parties_voted_for	0.20	Yearly_IncomeK	0.02	<p>1. Political_interest_Total_score 2. Number_of_different_parties_voted_f or 3. Weighted_education_rank</p>	<p>Turquoises</p>
Feature	Importance																																			
Most_Important_Issue_Social	0.01																																			
Most_Important_Issue_Other	0.02																																			
Most_Important_Issue_Military	0.01																																			
Most_Important_Issue_Healthcare	0.01																																			
Most_Important_Issue_Foreign_Affairs	0.02																																			
Most_Important_Issue_Financial	0.01																																			
Most_Important_Issue_Environment	0.01																																			
Most_Important_Issue_Education	0.01																																			
Weighted_education_rank	0.15																																			
Avg_size_per_room	0.05																																			
Overall_happiness_score	0.08																																			
Avg_monthly_income_all_years	0.01																																			
Avg_Satisfaction_with_previous_vote	0.11																																			
Political_interest_Total_Score	0.21																																			
Number_of_differnt_parties_voted_for	0.20																																			
Yearly_IncomeK	0.02																																			
<p>Feature importance of Khakis vs all</p> <table><thead><tr><th>Feature</th><th>Importance</th></tr></thead><tbody><tr><td>Most_Important_Issue_Social</td><td>0.01</td></tr><tr><td>Most_Important_Issue_Other</td><td>0.01</td></tr><tr><td>Most_Important_Issue_Military</td><td>0.01</td></tr><tr><td>Most_Important_Issue_Healthcare</td><td>0.01</td></tr><tr><td>Most_Important_Issue_Foreign_Affairs</td><td>0.01</td></tr><tr><td>Most_Important_Issue_Financial</td><td>0.01</td></tr><tr><td>Most_Important_Issue_Environment</td><td>0.01</td></tr><tr><td>Most_Important_Issue_Education</td><td>0.01</td></tr><tr><td>Weighted_education_rank</td><td>0.08</td></tr><tr><td>Avg_size_per_room</td><td>0.05</td></tr><tr><td>Overall_happiness_score</td><td>0.15</td></tr><tr><td>Avg_monthly_income_all_years</td><td>0.01</td></tr><tr><td>Avg_Satisfaction_with_previous_vote</td><td>0.05</td></tr><tr><td>Political_interest_Total_Score</td><td>0.15</td></tr><tr><td>Number_of_differnt_parties_voted_for</td><td>0.28</td></tr><tr><td>Yearly_IncomeK</td><td>0.02</td></tr></tbody></table>	Feature	Importance	Most_Important_Issue_Social	0.01	Most_Important_Issue_Other	0.01	Most_Important_Issue_Military	0.01	Most_Important_Issue_Healthcare	0.01	Most_Important_Issue_Foreign_Affairs	0.01	Most_Important_Issue_Financial	0.01	Most_Important_Issue_Environment	0.01	Most_Important_Issue_Education	0.01	Weighted_education_rank	0.08	Avg_size_per_room	0.05	Overall_happiness_score	0.15	Avg_monthly_income_all_years	0.01	Avg_Satisfaction_with_previous_vote	0.05	Political_interest_Total_Score	0.15	Number_of_differnt_parties_voted_for	0.28	Yearly_IncomeK	0.02	<p>1. Number_of_different_parties_voted_f or 2. Overall_happiness_score</p>	<p>Khakis</p>
Feature	Importance																																			
Most_Important_Issue_Social	0.01																																			
Most_Important_Issue_Other	0.01																																			
Most_Important_Issue_Military	0.01																																			
Most_Important_Issue_Healthcare	0.01																																			
Most_Important_Issue_Foreign_Affairs	0.01																																			
Most_Important_Issue_Financial	0.01																																			
Most_Important_Issue_Environment	0.01																																			
Most_Important_Issue_Education	0.01																																			
Weighted_education_rank	0.08																																			
Avg_size_per_room	0.05																																			
Overall_happiness_score	0.15																																			
Avg_monthly_income_all_years	0.01																																			
Avg_Satisfaction_with_previous_vote	0.05																																			
Political_interest_Total_Score	0.15																																			
Number_of_differnt_parties_voted_for	0.28																																			
Yearly_IncomeK	0.02																																			
<p>Feature importance of Violets vs all</p> <table><thead><tr><th>Feature</th><th>Importance</th></tr></thead><tbody><tr><td>Most_Important_Issue_Social</td><td>0.01</td></tr><tr><td>Most_Important_Issue_Other</td><td>0.01</td></tr><tr><td>Most_Important_Issue_Military</td><td>0.01</td></tr><tr><td>Most_Important_Issue_Healthcare</td><td>0.01</td></tr><tr><td>Most_Important_Issue_Foreign_Affairs</td><td>0.01</td></tr><tr><td>Most_Important_Issue_Financial</td><td>0.01</td></tr><tr><td>Most_Important_Issue_Environment</td><td>0.01</td></tr><tr><td>Most_Important_Issue_Education</td><td>0.01</td></tr><tr><td>Weighted_education_rank</td><td>0.11</td></tr><tr><td>Avg_size_per_room</td><td>0.16</td></tr><tr><td>Overall_happiness_score</td><td>0.17</td></tr><tr><td>Avg_monthly_income_all_years</td><td>0.09</td></tr><tr><td>Avg_Satisfaction_with_previous_vote</td><td>0.11</td></tr><tr><td>Political_interest_Total_Score</td><td>0.14</td></tr><tr><td>Number_of_differnt_parties_voted_for</td><td>0.08</td></tr><tr><td>Yearly_IncomeK</td><td>0.09</td></tr></tbody></table>	Feature	Importance	Most_Important_Issue_Social	0.01	Most_Important_Issue_Other	0.01	Most_Important_Issue_Military	0.01	Most_Important_Issue_Healthcare	0.01	Most_Important_Issue_Foreign_Affairs	0.01	Most_Important_Issue_Financial	0.01	Most_Important_Issue_Environment	0.01	Most_Important_Issue_Education	0.01	Weighted_education_rank	0.11	Avg_size_per_room	0.16	Overall_happiness_score	0.17	Avg_monthly_income_all_years	0.09	Avg_Satisfaction_with_previous_vote	0.11	Political_interest_Total_Score	0.14	Number_of_differnt_parties_voted_for	0.08	Yearly_IncomeK	0.09	<p>1. Overall_happiness_score 2. Avg_size_per_room 3. Political_interest_Total_score</p>	<p>Violets</p>
Feature	Importance																																			
Most_Important_Issue_Social	0.01																																			
Most_Important_Issue_Other	0.01																																			
Most_Important_Issue_Military	0.01																																			
Most_Important_Issue_Healthcare	0.01																																			
Most_Important_Issue_Foreign_Affairs	0.01																																			
Most_Important_Issue_Financial	0.01																																			
Most_Important_Issue_Environment	0.01																																			
Most_Important_Issue_Education	0.01																																			
Weighted_education_rank	0.11																																			
Avg_size_per_room	0.16																																			
Overall_happiness_score	0.17																																			
Avg_monthly_income_all_years	0.09																																			
Avg_Satisfaction_with_previous_vote	0.11																																			
Political_interest_Total_Score	0.14																																			
Number_of_differnt_parties_voted_for	0.08																																			
Yearly_IncomeK	0.09																																			
<p>Feature importance of Whites vs all</p> <table><thead><tr><th>Feature</th><th>Importance</th></tr></thead><tbody><tr><td>Most_Important_Issue_Social</td><td>0.02</td></tr><tr><td>Most_Important_Issue_Other</td><td>0.01</td></tr><tr><td>Most_Important_Issue_Military</td><td>0.01</td></tr><tr><td>Most_Important_Issue_Healthcare</td><td>0.01</td></tr><tr><td>Most_Important_Issue_Foreign_Affairs</td><td>0.01</td></tr><tr><td>Most_Important_Issue_Financial</td><td>0.01</td></tr><tr><td>Most_Important_Issue_Environment</td><td>0.02</td></tr><tr><td>Most_Important_Issue_Education</td><td>0.01</td></tr><tr><td>Weighted_education_rank</td><td>0.12</td></tr><tr><td>Avg_size_per_room</td><td>0.13</td></tr><tr><td>Overall_happiness_score</td><td>0.13</td></tr><tr><td>Avg_monthly_income_all_years</td><td>0.05</td></tr><tr><td>Avg_Satisfaction_with_previous_vote</td><td>0.06</td></tr><tr><td>Political_interest_Total_Score</td><td>0.13</td></tr><tr><td>Number_of_differnt_parties_voted_for</td><td>0.19</td></tr><tr><td>Yearly_IncomeK</td><td>0.05</td></tr></tbody></table>	Feature	Importance	Most_Important_Issue_Social	0.02	Most_Important_Issue_Other	0.01	Most_Important_Issue_Military	0.01	Most_Important_Issue_Healthcare	0.01	Most_Important_Issue_Foreign_Affairs	0.01	Most_Important_Issue_Financial	0.01	Most_Important_Issue_Environment	0.02	Most_Important_Issue_Education	0.01	Weighted_education_rank	0.12	Avg_size_per_room	0.13	Overall_happiness_score	0.13	Avg_monthly_income_all_years	0.05	Avg_Satisfaction_with_previous_vote	0.06	Political_interest_Total_Score	0.13	Number_of_differnt_parties_voted_for	0.19	Yearly_IncomeK	0.05	<p>1. Number_of_different_parties_voted_f or 2. Avg_size_per_room 3. Overall_happiness_score 4. Political_interest_Total_score 5. Weighted_education_rank</p>	<p>Whites</p>
Feature	Importance																																			
Most_Important_Issue_Social	0.02																																			
Most_Important_Issue_Other	0.01																																			
Most_Important_Issue_Military	0.01																																			
Most_Important_Issue_Healthcare	0.01																																			
Most_Important_Issue_Foreign_Affairs	0.01																																			
Most_Important_Issue_Financial	0.01																																			
Most_Important_Issue_Environment	0.02																																			
Most_Important_Issue_Education	0.01																																			
Weighted_education_rank	0.12																																			
Avg_size_per_room	0.13																																			
Overall_happiness_score	0.13																																			
Avg_monthly_income_all_years	0.05																																			
Avg_Satisfaction_with_previous_vote	0.06																																			
Political_interest_Total_Score	0.13																																			
Number_of_differnt_parties_voted_for	0.19																																			
Yearly_IncomeK	0.05																																			

<p>Feature importance of Yellows vs all</p> <table><thead><tr><th>Feature</th><th>Importance</th></tr></thead><tbody><tr><td>Most_Important_Issue_Social</td><td>0.010</td></tr><tr><td>Most_Important_Issue_Other</td><td>0.010</td></tr><tr><td>Most_Important_Issue_Military</td><td>0.010</td></tr><tr><td>Most_Important_Issue_Healthcare</td><td>0.010</td></tr><tr><td>Most_Important_Issue_Foreign_Affairs</td><td>0.010</td></tr><tr><td>Most_Important_Issue_Financial</td><td>0.010</td></tr><tr><td>Most_Important_Issue_Environment</td><td>0.010</td></tr><tr><td>Most_Important_Issue_Education</td><td>0.010</td></tr><tr><td>Weighted_education_rank</td><td>0.120</td></tr><tr><td>Avg_size_per_room</td><td>0.110</td></tr><tr><td>Overall_happiness_score</td><td>0.110</td></tr><tr><td>Avg_monthly_income_all_years</td><td>0.140</td></tr><tr><td>Avg_Satisfaction_with_previous_vote</td><td>0.140</td></tr><tr><td>Political_interest_Total_Score</td><td>0.110</td></tr><tr><td>Number_of_differnt_parties_voted_for</td><td>0.060</td></tr><tr><td>Yearly_IncomeK</td><td>0.110</td></tr></tbody></table>	Feature	Importance	Most_Important_Issue_Social	0.010	Most_Important_Issue_Other	0.010	Most_Important_Issue_Military	0.010	Most_Important_Issue_Healthcare	0.010	Most_Important_Issue_Foreign_Affairs	0.010	Most_Important_Issue_Financial	0.010	Most_Important_Issue_Environment	0.010	Most_Important_Issue_Education	0.010	Weighted_education_rank	0.120	Avg_size_per_room	0.110	Overall_happiness_score	0.110	Avg_monthly_income_all_years	0.140	Avg_Satisfaction_with_previous_vote	0.140	Political_interest_Total_Score	0.110	Number_of_differnt_parties_voted_for	0.060	Yearly_IncomeK	0.110	<p>1. Avg_Satisfaction_with_previous_vote 2. Weighted_education_rank 3. Political_interest_Total_score</p>	Yellows
Feature	Importance																																			
Most_Important_Issue_Social	0.010																																			
Most_Important_Issue_Other	0.010																																			
Most_Important_Issue_Military	0.010																																			
Most_Important_Issue_Healthcare	0.010																																			
Most_Important_Issue_Foreign_Affairs	0.010																																			
Most_Important_Issue_Financial	0.010																																			
Most_Important_Issue_Environment	0.010																																			
Most_Important_Issue_Education	0.010																																			
Weighted_education_rank	0.120																																			
Avg_size_per_room	0.110																																			
Overall_happiness_score	0.110																																			
Avg_monthly_income_all_years	0.140																																			
Avg_Satisfaction_with_previous_vote	0.140																																			
Political_interest_Total_Score	0.110																																			
Number_of_differnt_parties_voted_for	0.060																																			
Yearly_IncomeK	0.110																																			
<p>Feature importance of Blues vs all</p> <table><thead><tr><th>Feature</th><th>Importance</th></tr></thead><tbody><tr><td>Most_Important_Issue_Social</td><td>0.005</td></tr><tr><td>Most_Important_Issue_Other</td><td>0.005</td></tr><tr><td>Most_Important_Issue_Military</td><td>0.005</td></tr><tr><td>Most_Important_Issue_Healthcare</td><td>0.005</td></tr><tr><td>Most_Important_Issue_Foreign_Affairs</td><td>0.005</td></tr><tr><td>Most_Important_Issue_Financial</td><td>0.005</td></tr><tr><td>Most_Important_Issue_Environment</td><td>0.005</td></tr><tr><td>Most_Important_Issue_Education</td><td>0.005</td></tr><tr><td>Weighted_education_rank</td><td>0.080</td></tr><tr><td>Avg_size_per_room</td><td>0.080</td></tr><tr><td>Overall_happiness_score</td><td>0.080</td></tr><tr><td>Avg_monthly_income_all_years</td><td>0.480</td></tr><tr><td>Avg_Satisfaction_with_previous_vote</td><td>0.080</td></tr><tr><td>Political_interest_Total_Score</td><td>0.080</td></tr><tr><td>Number_of_differnt_parties_voted_for</td><td>0.020</td></tr><tr><td>Yearly_IncomeK</td><td>0.080</td></tr></tbody></table>	Feature	Importance	Most_Important_Issue_Social	0.005	Most_Important_Issue_Other	0.005	Most_Important_Issue_Military	0.005	Most_Important_Issue_Healthcare	0.005	Most_Important_Issue_Foreign_Affairs	0.005	Most_Important_Issue_Financial	0.005	Most_Important_Issue_Environment	0.005	Most_Important_Issue_Education	0.005	Weighted_education_rank	0.080	Avg_size_per_room	0.080	Overall_happiness_score	0.080	Avg_monthly_income_all_years	0.480	Avg_Satisfaction_with_previous_vote	0.080	Political_interest_Total_Score	0.080	Number_of_differnt_parties_voted_for	0.020	Yearly_IncomeK	0.080	<p>1. Avg_monthly_income_all_years</p>	Blues
Feature	Importance																																			
Most_Important_Issue_Social	0.005																																			
Most_Important_Issue_Other	0.005																																			
Most_Important_Issue_Military	0.005																																			
Most_Important_Issue_Healthcare	0.005																																			
Most_Important_Issue_Foreign_Affairs	0.005																																			
Most_Important_Issue_Financial	0.005																																			
Most_Important_Issue_Environment	0.005																																			
Most_Important_Issue_Education	0.005																																			
Weighted_education_rank	0.080																																			
Avg_size_per_room	0.080																																			
Overall_happiness_score	0.080																																			
Avg_monthly_income_all_years	0.480																																			
Avg_Satisfaction_with_previous_vote	0.080																																			
Political_interest_Total_Score	0.080																																			
Number_of_differnt_parties_voted_for	0.020																																			
Yearly_IncomeK	0.080																																			
<p>Feature importance of Browns vs all</p> <table><thead><tr><th>Feature</th><th>Importance</th></tr></thead><tbody><tr><td>Most_Important_Issue_Social</td><td>0.010</td></tr><tr><td>Most_Important_Issue_Other</td><td>0.010</td></tr><tr><td>Most_Important_Issue_Military</td><td>0.010</td></tr><tr><td>Most_Important_Issue_Healthcare</td><td>0.010</td></tr><tr><td>Most_Important_Issue_Foreign_Affairs</td><td>0.010</td></tr><tr><td>Most_Important_Issue_Financial</td><td>0.010</td></tr><tr><td>Most_Important_Issue_Environment</td><td>0.010</td></tr><tr><td>Most_Important_Issue_Education</td><td>0.010</td></tr><tr><td>Weighted_education_rank</td><td>0.140</td></tr><tr><td>Avg_size_per_room</td><td>0.130</td></tr><tr><td>Overall_happiness_score</td><td>0.160</td></tr><tr><td>Avg_monthly_income_all_years</td><td>0.080</td></tr><tr><td>Avg_Satisfaction_with_previous_vote</td><td>0.160</td></tr><tr><td>Political_interest_Total_Score</td><td>0.190</td></tr><tr><td>Number_of_differnt_parties_voted_for</td><td>0.040</td></tr><tr><td>Yearly_IncomeK</td><td>0.050</td></tr></tbody></table>	Feature	Importance	Most_Important_Issue_Social	0.010	Most_Important_Issue_Other	0.010	Most_Important_Issue_Military	0.010	Most_Important_Issue_Healthcare	0.010	Most_Important_Issue_Foreign_Affairs	0.010	Most_Important_Issue_Financial	0.010	Most_Important_Issue_Environment	0.010	Most_Important_Issue_Education	0.010	Weighted_education_rank	0.140	Avg_size_per_room	0.130	Overall_happiness_score	0.160	Avg_monthly_income_all_years	0.080	Avg_Satisfaction_with_previous_vote	0.160	Political_interest_Total_Score	0.190	Number_of_differnt_parties_voted_for	0.040	Yearly_IncomeK	0.050	<p>1. Political_interest_Total_score 2. Overall_happiness_score 3. Avg_Satisfaction_with_previous_vote</p>	Browns
Feature	Importance																																			
Most_Important_Issue_Social	0.010																																			
Most_Important_Issue_Other	0.010																																			
Most_Important_Issue_Military	0.010																																			
Most_Important_Issue_Healthcare	0.010																																			
Most_Important_Issue_Foreign_Affairs	0.010																																			
Most_Important_Issue_Financial	0.010																																			
Most_Important_Issue_Environment	0.010																																			
Most_Important_Issue_Education	0.010																																			
Weighted_education_rank	0.140																																			
Avg_size_per_room	0.130																																			
Overall_happiness_score	0.160																																			
Avg_monthly_income_all_years	0.080																																			
Avg_Satisfaction_with_previous_vote	0.160																																			
Political_interest_Total_Score	0.190																																			
Number_of_differnt_parties_voted_for	0.040																																			
Yearly_IncomeK	0.050																																			

5. מציאת פיצ'רים שעל ידי שינויים נקבל תוצאות בחירות שונה:

חלק זה נמצא בקובץ fourth_prediction.

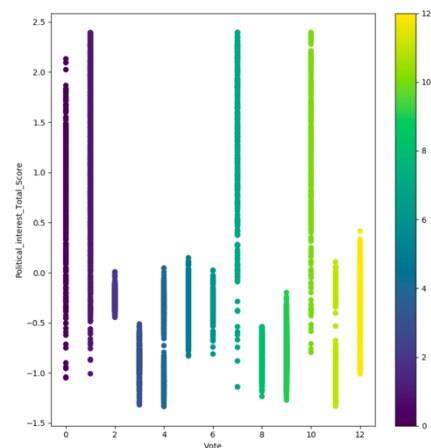
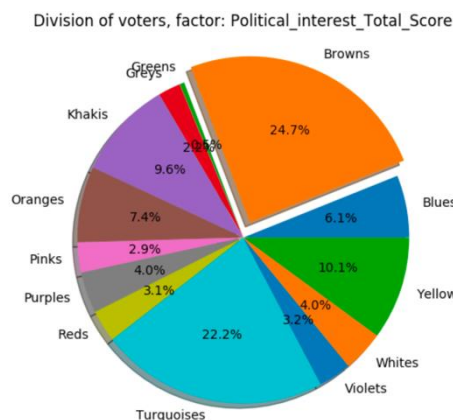
ראשית הסתכלנו על השכבות הראשונות של עץ החלטה וראינו מי מהתכונות בעלי משקל גדול על תוצאות המודל והסתכלנו על התכונות איך הן משפיעות על ההצבעה ואיך הן מפולגות. לכל תכונה בנפרד שינינו את קבוצת ה test עבור אותה תכונה כדי להשפיע על התוצאות של המנצח על ידי שהמסווג יטעה בגלל תכונה זו בסיווג. שינינו כל תכונה ולאחר שינוי קבוצת ה test בדקנו את המודל (שמאומן על train תקין) ואת המנצח.



בנוסף ראינו מתרגיל קודם שהמנצח הוא Turquoises (9) ומקום שני הוא Brown (1) וקיבלנו את התכונות הבאות כקריטריות עבור תוצאות הסיווג

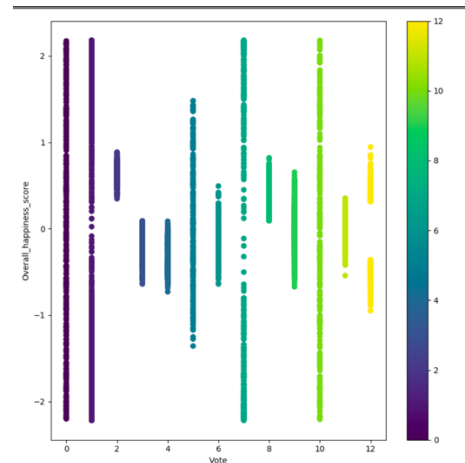
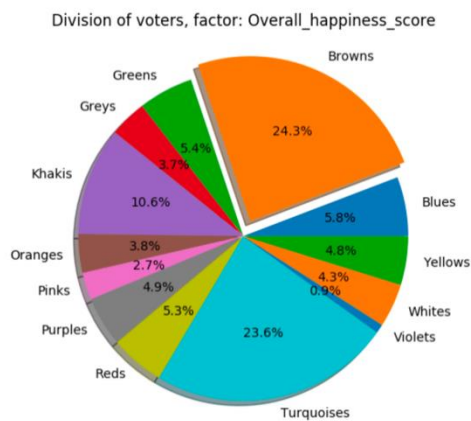
Political_interest_Total_Score : ניתן לראות כי המצביעים של 1 ו-9 מפולגים בצורה כזאת שעבור 1 יש בוחרים בטווח גדול יותר ובאזור הגבוהה של הערכים והטווח עבור 9 קטן יותר ולמטה.

פילוג של train data על תכונה זו ביחס להצבעות: brown מנצח לאחר שינוי תכונה זו test:



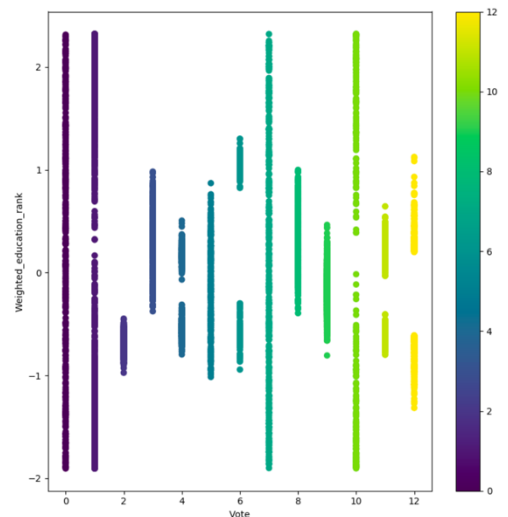
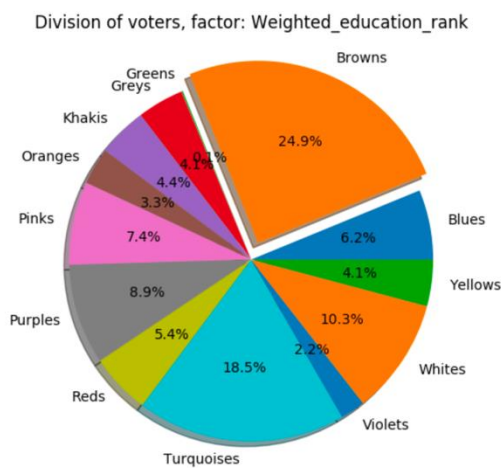
Overall_happiness_score : ניתן לראות כי המצביעים של 1 ו-9 מפולגים בצורה כזאת שעבור 1 יש בוחרים בטווח גדול יותר בכל הטווח מלבד האמצע ב והטווח עבור 9 קטן יותר ובאמצע.

פילוג של train data על תכונה זו ביחס להצבעות: brown מנצח לאחר שינוי תכונה זו test:



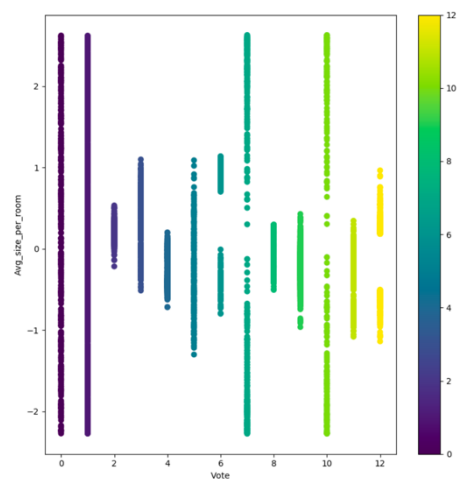
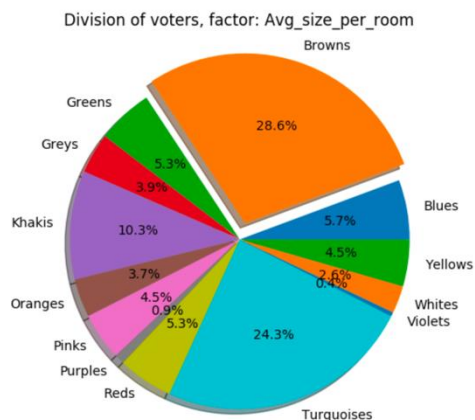
Weighted_education_rank : ניתן לראות כי המצביעים של 1 ו-9 מפולגים בצורה כזאת שעבור 1 יש בוחרים בטווח גדול יותר בכל הטווח מלבד האמצע והטווח עבור 9 קטן יותר ובאמצע.

פילוג של train data על תכונה זו ביחס להצבעות: brown מנצח לאחר שינוי תכונה זו test:



Size per room : ניתן לראות כי המצביעים של 1 ו-9 מפולגים בצורה כזאת שעבור 1 יש בוחרים בטווח גדול יותר בכל הטווח והטווח עבור 9 קטן יותר ובאמצע.

פילוג של train data על תכונה זו ביחס להצבעות: brown מנצח לאחר שינוי תכונה זו בtest:

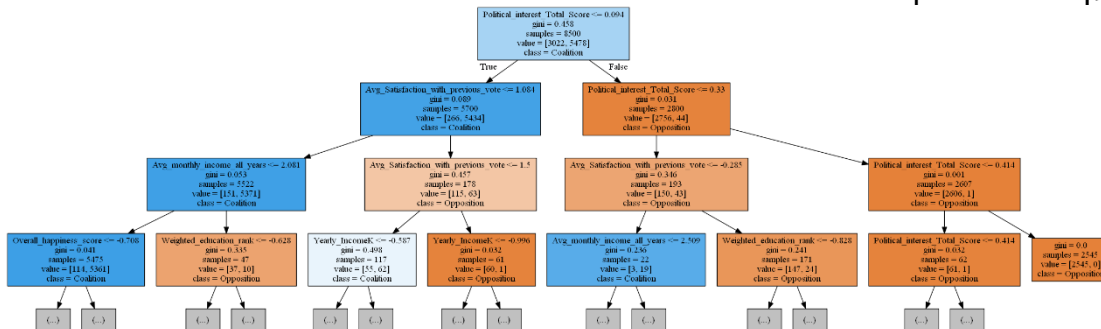


6. זיהוי הפיצ'רים שעל ידי שינויים נוכל לחזק את הקואליציה שלנו:

שלב זה ביצענו בצורה דומה לסעיף קודם. תחילה בחרנו את הקואליציה שהחזיר לנו אלגוריתם ה cluster, כעת שינינו את הדאטה שלנו כך שני לייבלים, קואליציה ואופוזיציה. כלומר עברנו על סט האימון, ועבור כל לייבל השייך לקואליציה הצבנו 1 ובכל שאר הלייבלים 0-1. כעת ביצענו אימון של עץ החלטה על סט אימון זה. כעת עלינו לחזק את הקואליציה, הגדרנו חיזוק זה בצורה הבאה: חיזוק הקואליציה – הגדלת מספר המצביעים לקואליציה, הקטנת שונות הקואליציה (יותר הומוגנית), הגדלת מרחק הקואליציה מהאופוזיציה. כעת על מנת למצוא את הפיצ'רים שעלינו לשנות השתמשנו ב 3 דרכים:

- הסתכלנו על עץ ההחלטה שקיבלנו, כאשר הסתכלנו על הרמות העליונות ביותר. רמות אלה מבצעות את ההפרדה הגדולה ביותר ולכן על ידי שינוי של אחד מהפיצ'רים המופיעים שם נוכל להשפיע בצורה גדולה יותר.
- חישוב ה Variance של כל הפיצ'רים של מצביעי הקואליציה. לאחר שחישבנו מידע זה, הסתכלנו איזה מהפיצ'רים בעל שונות גדולה ביותר וכך נכון להשפיע עליו על מנת להקטין את השונות של הקואליציה.
- חישוב מרחק ממוצע כל פיצ'ר ממצביעי הקואליציה לממוצע כל פיצ'ר ממצביעי האופוזיציה. לאחר חישוב זה הסתכלנו איזה פיצ'רים מבדילים כמה שיותר בין הקבוצות ונסה לשנות אותם עוד יותר כדי להגדיל את המרחק בין הקבוצות. תוצאות כל אחד מהדרכים:

א. עץ ההחלטה שקיבלנו:



- ניתן לראות כי הפיצ'ר Political_interest_total_score הינו המשפיע ביותר על פיצול הדוגמאות, לכן על ידי הקטנת פיצ'ר זה בכל הדוגמאות, נהפוך יותר מצביעים לקואליציה.
- פיצ'ר Number_of_differnt_parties_voted_for בעל השונות הגדולה ביותר, לכן נקרב את ערכיו לממוצע על מנת להקטין את השונות. (נשים לב כי הכפלה זו תגרום גם להקטנת המרחק בין הקואליציה והאופוזיציה אבל בצורה קטנה יותר מאשר הקטנת השונות)
- פיצ'ר Political_interest_Total_Score בעל ההבדל הגדול ביותר בין הקואליציה לאופוזיציה, כאשר ממוצע הקואליציה גדול מ-0 וממוצע האופוזיציה קטן מ-0. ננסה לפצל עוד יותר את הדאטה על ידי הכפלה של הערכים. (נשים לב כי הכפלה זו תגרום גם להגדלת השונות של פיצ'ר זה, אז בצורה פחות משמעותית)

תוצאות:

לפני השינוי –

גודל הקואליציה: 64.13%

שונות הקואליציה: 3.843

מרחק מהאופוזיציה: 5.843

אחרי השינוי –

גודל הקואליציה: 68.2%

שונות הקואליציה: 3.213

מרחק מהאופוזיציה: 6.375

נשים לב כי הצלחנו להגדיל את גודל הקואליציה, להקטין את שונות הקואליציה (יותר הומוגנית), ולהגדיל את המרחק של הקואליציה מהאופוזיציה. כלומר קיבלנו קואליציה יותר יציבה וחזקה.