

Decision Trees

מצורף הקובץ votersdata.csv המכיל נתונים על בוחרים בארה"ב, כולל עמודת תוצאת הבחירה (רפובליקני/ דמוקרטי).

1. תחילה הגדירו random seed עם הערך 123 והשתמשו בו בפיצול הנתונים ובאתחול המודל.
2. הכרת הדאטה: נסו להכיר את הקשרים בין המשתנים השונים לבין עמודת המטרה vote.
2. א. שרטטו stacked bar plots עבור המשתנים הקטגוריים.
2. ב. שרטטו multivariate boxplot עבור הנומריים.
3. תקנו את הנתונים במידת הצורך. זה כולל טיפול בערכים חסרים, ערכים לא תקינים, נרמול והמרות נדרשות.
4. חלקו את הדאטה באופן רנדומי ל-70% training set ו-30% test set בשימוש ב seed שהגדרתם.
5. בנו מודל עץ בעזרת train set לחיזוי ערך המשתנה vote וענו על השאלות הבאות:
זכרו להציב Random seed

Model Evaluation

6. בנו Confusion matrix עבור חיזוי על ה- test set בעזרת המודל שבניתם בשאלה הקודמת.
הניחו כי "דמוקרטי" = Positive וחשבו את המדדים הבאים:
א. Accuracy
ב. Precision
ג. Recall
7. השוו את המדדים עבור ה- train set. האם מתקיימת תופעת ה-overfitting במודל החיזוי שבניתם? נמקו.
8. צרו מודל חדש משופר על סמך המסקנות מ-7. הגבילו את גובה העץ ל-5 ואת כמות הרשומות לחיתוך ל-40 וענו:

- א. מהו עומק העץ ?
 ב. כמה עלים יש בעץ ?
 ג. מהו פיצ'ר החלוקה הטוב ביותר בעץ ?
 ד. האם יש פיצ'רים שלא נכללו במודל ? מהם?
 ה. האם תצפית מס' 68 (בדאטה סט המקורי) סווגה נכונה במודל ? נמקו.
 9. בצעו שוב חיזוי על ה train set וה test sets ובנו מטריצות חדשות.

10. סטודנטים הריצו מודל משופר ויצאו להם המדדים הבאים:

Test set result:

Accuracy: 0.7946428571428571

Precision: 0.7142857142857143

recall: 0.9433962264150944

Train set results:

Accuracy: 0.7961538461538461

Precision: 0.7261904761904762

recall: 0.9457364341085271

מה אפשר להסיק מהתוצאות הנ"ל לגבי ביצועי המודל על הדאטה ? יש לפרט במילים.

Decision tree - Multiclass

- שנו את עמודת המטרה להיות "status" ובנו עץ החלטה לזיהוי הסטטוס.
 בנו confusion matrix והדפיסו את מדד ה **accuracy** לאחר בדיקה על ה test set.
 כתבו את התוצאה גם בקוד כהערה וענו:
 11. האם המודל יחזה טוב את הסטטוס המשפחתי על דאטה חדש ? מדוע?
 12. האם קיים חשד ל overfitting ?
 13. חשבו את מדד ה **precision** עבור הקטגוריה single.