

פרויקט מסכם

בפרויקט זה עלייכם לחתולק לצוותים של שלושה סטודנטים לכל היוטר. כל צוות יבחר באחד הפרויקטים המופיעים בטבלה וישתבז באמצעות רישום עצמי (לא כפיליות- צוות לכל פרויקט), בכתובת הבאה :

https://docs.google.com/spreadsheets/d/1vwlnCG2zuCpmVRUsGqYno2CooBBxnRTVm1u_gBViokg/edit?usp=sharing

במודול תוכלו למצוא תקינה שכוללת את כל סטי הנתונים הרלוונטיים, מחברת לדוגמא לפרויקט -Titanic וכן מחברת "load data", שוראת את סטי הנתונים - pandas data frame.

עליכם להציג את הפרויקט במסגרת יום הצגת הפרויקטים. שעוט ההצגה תקבעה בהמשך, גם בהתאם לכמות הסטודנטים שיירשוו לכל אחד מהימים. אני שריםנו את היום שבחרתם לצורך העניין .

דרישות המציג- במסגרת הצגת הפרויקט, עליכם להתייחס לאלמנטים הבאים :

1. הגדרת הבעיה- מה מטרת הפרויקט? עם איזו בעיה הוא מתמודד? מהם הנתונים הזמינים לצורך הפרויקט (נקודות) ? מה המוטיבציה לפרויקט ואפליקציות רלוונטיות.

2. הערכת ביצועים- מהי המטריקה שבה בחרתם להשתמש לצורך הערכת ביצועי המודל ומדוע ; כיצד חילקתם את סט הנתונים?

3. נתונים- תאור סט הנתונים :

(1) כמה דוגמאות יש בסט?

(2) כמה מאפיינים יש בסט?

(3) מהי התפלגות ה- labels?

(4) האם ישנים ערכים חסריים?

(5) הצג 2-3 גרפים המתארים את האспектים השונים של סט הנתונים.

4. הנדסת נתונים (-engineering). data engineering

(1) האם הסרתם מאפיינים כלשחים?

(2) האם הוספתם מאפיינים כלשחים?

(3) כיצד התמודדתם עם ערכים חסריים?

5. ביצועי המודל- נתחו והשו את ביצועי של המודלים הבאים-

(1) מודל "baseline"- בעיית רגרסיה מודל זה יכול להתבסס, לדוגמא, על הערך המומוצע ; בעיית סיוג- מודל זה יכול להתבסס על אחד המטוגנים הבסיסיים וסט המאפיינים הראשוני.

(2) מודלים מוכרים-

- a. K nearest neighbors, try 3 different k values.
 - b. K nearest neighbors with scaled values, try 3 different k values.
 - c. Decision tree with 3 different max depth values.
 - d. Random forest with 3 different max depth values and 100 estimators.
 - e. Ada Boost with 3 different max depth values and 100 estimators.
 - f. Lasso regression with scaled values, try 3 different alpha values (only for regression).
6. בוחן את המודלים שקיבלה. לדוגמה, במקרה של random forest כדאי לבדוק את ה- feature importance.
7. היפר-פרמטרים- השתמש במודל שהניב את התוצאות הטובות ביותר. חפש את ההיפר-פרמטרים של המודל אשר משפרים בצורה הגדולה ביותר את ביצועי המודל. הצג השוואת ביצועים עבור ערכי היפר-פרמטרים אלו.
8. ניתוח נוסף (בחירת אחד מהבאים)-
- a. ביצועים כפונקציה של גודל סט הנתונים- השתמש במודל המצליח ביותר. הצג גרף של ביצועי המודל כאשר משתמשים ב- 10%, 30%, 50%, 70% או 100% מסט האימון. האם הייתה ממליצה להגדיל את סט הנתונים?
 - b. מיצוע מודלים- בחר שלושה מודלים שונים שבחנת קודם לבן. צור מודל "חדש" אשר מבוסס על המוצע של מוצאי שלושת המודלים שבחرت. השווה את הביצועים של מודל זה ביחס לכל יתר המודלים שבחנת בפרויקט.

אין צורך להגיש דו"ח ייעודי, אלא להזכיר מחברת גייפטר מסודרת, עם תיעוד מתאים של הקוד. את הקישור למחברת יש לעדכן בכתבות שצוינה מעלה (אין להגיש את המחברת).

בהצלחה!

צוות הקורס