



מטלה שלישית - חיזוי עתיד
כלכלה בעולם ה Big Data
ד"ר רועי ששון; אופיר בצר

להגיש לכל המאוחר עד:

2022-08-31 23:59

לא יתקבלו הגשות באיחור.

פרטים חשובים על המטלה:

- יש לבצע את המטלה **בזוגות**. לא יבדקו מטלות באופן יחיד, במידה ומתעוררת בעיה בנושא אנא הודיעו וימצא פתרון.
 - רק אחד מחברי כל קבוצה צריך להגיש את המטלה במודל.
 - נא לציין באופן ברור מי הם חברי הקבוצה בתוך הדוח ובשם של הקובץ של הדוח באופן הבא:
`EinBDW22_A3_{name #1}_{name #2}.pdf`
לדוגמא:
`EinBDW22_A3_roy_sasson_ophir_betser.pdf`
- משקל המטלה הוא **70%** מציון הקורס.
- את כל הניתוחים יש לבצע ב **R בלבד**.

יש להגיש:

- דו"ח מפורט ובו התשובות שלכם על השאלות, בצירוף גרפים וטבלאות רלוונטים. על הדו"ח להיות ממוקד - אין להגיש יותר מ 6 עמודים, וינתן משקל על אסתטיות של הגרפים ו/או הטבלאות המצורפים בו.
 - יש לצרף לכל אחד מהגרפים כותרות המסבירות את תוכנו.
 - יש לוודא שהגרפים שאתם מצרפים קריאים, וברורים, מובנים ואסתטיים.
 - אין צורך לכתוב הרבה מלל, אלא לענות באופן קצר וממוקד.
 - פורמט הקובץ צריך להיות **pdf**.
 - שפת הדוח צריכה להיות **עברית**, גופן דוד 12, רווח בין השורות 1.5.
- קובץ csv בפורמט הבא: **קובץ** () אשר בו יוגשו תחזיות המודל שלכם.
 - על הקובץ להיות בדיוק בפורמט המצורף - אותם שמות עמודות ואותם ערכים בעמודות ה `partition_datetime`.
 - על שם קובץ להיות בפורמט הבא: `EinBDW22_A3_predict_{name #1}_{name #2}.csv`.
- יש** לצרף את סקריפט ה R הכולל את הניתוחים שביצעתם.

על סט הנתונים:

במטלה זו תבצעו ניתוח וחיזוי של מחירים **הנתיב המהיר**. מחירי הנסיעה בנתיב מעודכנים בהתאם לביקוש הנהגים, כך בשעות העומס מחיר הכניסה לנתיב המהיר עולה. מחיר הכניסה לנתיב מעודכן באתר האינטרנט, וכך ניתן לעקוב על השינויים בעלויות הנסיעה בו.

את המחיר העדכני אנו אוספים בעזרת זחלן רשת. טבלאות עדכניות יפורסמו מידי שבוע **בדריב** (`price_data_yyyymmdd`). בכל טבלה יהיו 2 עמודות בלבד: זמן הדגימה (POSIXt) ומחיר הנסיעה הממוצע (`numeric`).

שימו לב שיתכנו רשומות חסרות בסט הנתונים, קחו זאת בחשבון בעת הניתוח וניקוי הנתונים. למתעניינים, **זהו סקריפט הזחלן** שאוסף את הנתונים (python).

שאלות מחקר: (חלקים 1 עד 4)

חלק 1 - אנליסט המוצר (45%)

חלק זה של המטלה כולל את ניתוח הנתונים המרכזי במטלה. עליכם "לספר את הסיפור שמאחורי הנתונים" באופן שלוקח בחשבון את משימת התחזית בחלק הרביעי של המטלה.

הכנת וניתוח הנתונים כולל את השלבים הבאים:

- טעינת הנתונים
- סידור וניקוי שלהם
- הוספת משתנים מסבירים (לדוגמא משתנים מבוססי תאריך, מזג אוויר, נתונים קורונה או כל מידע אחר שלדעתכם רלוונטי).
- ביצוע EDA

חלק 2 - הכלכלן (10%)

בחרו משתנה מסביר אחד מהרשימה, או אחר שאתם חושבים שיכול להיות מעניין: [מזג אוויר, סוג היום]. דונו ותארו כיצד ניתן לגלות מה הקשר הסיבתי בינו לבין מחירי הנתיב המהיר מבחינה כלכלית. בניתוח תארו כיצד הייתם ממליצים לבדוד את ההשפעה הסיבתיות. כחלק מהניתוח עשו שימוש במודל קלאסטריןג k-means ופרשו את תוצאותיו.

חלק 3 - האנליסט העסקי (5%)

האם הנסיעה בנתיב המהיר משתלמת?

ישראל ישראלי מתגורר בירושלים ועובד בתל אביב. הוא מרוויח שכר שעתי של 100 ש"ח, ויש לו את האפשרות לצאת מביתו בשעות 7, 8, 9, 10 או 11 (שעות עגולות בלבד) ברכבו הפרטי. נתחו את הכדאיות של הנסיעה בנתיב המהיר עבור מר ישראלי - כמה זמן של נסיעה על הנתיב לחסוך לעומת אלטרנטיבת הנסיעה בנתיב רגיל על מנת שהשימוש בו יהיה משתלם?

חלק 4 - מדען הנתונים (40%)

חלק זה הינו תחרותי בין הקבוצות. עליכם לתת תחזית בעלת mse נמוך ככל הניתן עבור נקודות הדגימה אשר נמצאות [הקובץ הזה](#) (בין ה 2022-09-01 עד ה 2022-09-13, בשעות 06:00 עד 12:00 מדי יום).

כחלק מהניתוח, תארו את אופן בניית המודל, אופן האימון שלו ודרך הבדיקה שלכם את תוצאותיו (למשל - האם יש overfitting).