

תרגיל בית 4

הנחיות להגשת התרגיל

- יש להגיש שני קבצים נפרדים, אחד עבור הקוד (קובץ ipynb) ואחד לפתרון החלק היבש. בתוך קובץ פתרון החלק היבש הוסיפו גם את הפלט של הקוד.

שאלה 1

$$OR = \frac{odds(E=1|D=1)}{odds(E=1|D=0)}$$

א. הוכיחו כי OR במדגם Case Control

שאלה 2

הראו והסבירו מה ההשפעה של עליה ביחידה אחת ב- X_1 במודלים הבאים.

$$1. \mathbb{E}[Y|X] = \pi(\beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_{1,2} X_1 X_2)$$

$$2. \mathbb{E}[Y|X] = \pi(\beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_1^2)$$

שאלה 3

הנתונים איתם נעבוד בתרגיל זה מתארים ספירה של מקרי פציעות ושברים שהתרחשו במכרות פחם באזור האפלצ'ים במערב וירג'יניה. בסך הכול נאספו 44 תצפיות על מכרות באזור זה. נמדדו ארבעה משתנים הקשורים למאפייני המכרות:

- X_1 – Inner burden thickness in feet (INB)
- X_2 – Percent extraction of the lower previously mined seam (EXTRP)
- X_3 – lower seam height in feet (HEIGHT)
- X_4 – Time that the mine has been opened in years (TIME)

הנתונים מצורפים לתרגיל הבית בקובץ mine_fracture.csv

- כתבו משוואת מודל רגרסיה פואסונית עבור הקשר בין מספר הפציעות לבין ארבעת המשתנים המסבירים האלה.
- השתמשו ב statsmodels.GLM או בכל פונקציית ספריה אחרת כדי לאמוד את מקדמי הרגרסיה הפואסונית.
- פי כמה משתנה תוחלת מספר הפציעות כאשר גובה השכבה התחתונה (X_3) גדל ביחידה אחת? מצאו אומד נקודתי לגודל זה ורווח סמך.
- מצאו תת מודל בגודל 2 (מודל עם שני משתנים מסבירים מבין X_1, X_2, X_3, X_4) עבורו מתקבל מדד ה-AIC הטוב ביותר. אם אתם נעזרים בפונקציית ספריה, שימו לב כיצד ה-AIC מוגדר בספריה זו ואם צריך למזער או למקסם אותו. לפי מדד AIC, האם עדיף לבחור בתת המודל שמצאתם בסעיף זה או במודל המלא מסעיף א'?
- בצעו מבחן שבודק אם המודל המלא מסעיף א' עדיף על המודל המצומצם שמצאתם בסעיף ד'.