

## שאלה 1

א. קראו את הקובץ המצורף: "ex8\_data.csv". הפרידו את המשתנה  $Y$  משאר המשתנים המסבירים. בדקו והציגו את כל 4 האינדיקציות למולטיקוליאריות בין המשתנים המסבירים. האם ישנה עדות למולטיקוליאריות בנתונים? בין אילו משתנים? הסבירו בקצרה.

ב. אמדו את המודל כאשר  $Y$  הוא המשתנה המוסבר על כל המשתנים המסבירים ושמרו את התוצאות. אילו מהמשתנים יצאו מובהקים? האם התוצאה מפתיעה? בשלב הבא הוסיפו משתנה חדש:

$$Y_{new} = Y + rnorm(300,0,1)$$

אמדו את המודל המסביר את  $Y_{new}$  על כל שאר המשתנים המסבירים (ללא  $Y$ ).

ג. הציגו  $subplot$  המכיל שני גרפים: האחד  $barplot$  כפול, כאשר בצבע אחד מופיעים המקדמים שהתקבלו מהמודל הראשון, ובצבע אחר המקדמים שהתקבלו מאמידת המודל השני. השני, מאה נקודות שבחרתם מתוך וקטורי הערכים החזויים מכל אחת מהאמידות (כלומר 2 ערכים חזויים לכל תצפית -  $\hat{Y}_{new_i}, \hat{Y}_i$ ). בעבור כל אינדקס בין 1-100, הציגו בצבעים שונים את  $\hat{Y}_{new_i}, \hat{Y}_i$ . בכותרת של כל גרף, בהתאמה, הקפידו לציין את הגדלים:

$$\frac{||\hat{\beta} - \hat{\beta}_{new}||^2}{||\hat{\beta}||^2}, \frac{||\hat{Y} - \hat{Y}_{new}||^2}{||\hat{Y}||^2}$$

הסבירו את התוצאות ואת ההבדלים בין הגדלים.

## שאלה 2:

בשאלה הזאת נמצא פרשנות למדד VIF: אנחנו נראה ש-VIF עבור משתנה מסביר  $j$  זה למעשה היחס בין השונות של המקדם (הנאמד) ברגרסיה מרובה לבין השונות של המקדם ברגרסיה פשוטה.

כרגיל, אנחנו מסמנים ב- $X^{(j)}$  את העמודה במטריצה  $X$  שמתאימה למשתנה המסביר ה- $j$ , כך ש:  
 $X = [X^{(0)} \ X^{(1)} \ \dots \ X^{(p)}] \in \mathbb{R}^{n \times (p+1)}$ . נסמן  $X_{-j}$  את המטריצה  $X$  ללא העמודה  $X^{(j)}$ , עבור  $j = 1, \dots, p$ . כזכור, אם  $SSE(j)$ ,  $SST(j)$  ו- $R_j^2$  מסמנים, בהתאמה, את  $SSE$ ,  $SST$  ו- $R^2$  ברגרסיה מרובה של  $X^{(j)}$  על  $X_{-j}$ , אז:  $VIF = \frac{1}{1-R_j^2}$ .

א. רגרסיה פשוטה ללא חותך מתאימה לנתונים את המודל  $Y_i = \beta X_i + \epsilon_i$  כלומר זה אכן מודל רגרסיה עם משתנה מסביר בודד וללא האיבר הקבוע (שאותו אנחנו רגילים לסמן  $\beta_0$ ). הראו שהאומד של  $\beta$  ברגרסיה פשוטה ללא חותך נתון ע"י

$$\hat{\beta} = \frac{1}{\|x\|^2} x^T Y$$

כאשר  $x = (X_1, \dots, X_n)^T$ .

עובדה כללית (\*): ברגרסיה מרובה של  $Y$  על  $X$  אפשר לקבל את  $\hat{\beta}_j$  ע"י רגרסיה פשוטה ללא חותך של  $Y$  על הוקטור

$$z = (I - P_{-j})X^{(j)}$$

כאשר

$$P_{-j} = X_{-j}(X_{-j}^\top X_{-j})^{-1}X_{-j}^\top$$

היא מטריצת ההיטל על מרחב העמודות של  $X_{-j}$ .

ב. הראו/הסבירו שברגרסיה פשוטה (עם חותך) של  $Y$  על  $X^{(j)}$ , כלומר עבור  $j$  נתון ותחת המודל  $Y_i = \beta_0 + \beta_j X_{ij} + \epsilon_i$ , ניתן להציג

$$\hat{\beta}_j = \frac{1}{\|w\|^2} w^\top Y$$

כאשר

$$w = (I - P_0)X^{(j)}, \quad P_0 := \mathbf{1}_n(\mathbf{1}_n^\top \mathbf{1}_n)^{-1}\mathbf{1}_n^\top = \frac{1}{n}\mathbf{1}_n\mathbf{1}_n^\top$$

הדרכה: הסתכלו על הרגרסיה הפשוטה (עם חותך) בתור רגרסיה מרובה על המטריצה  $n \times 2$  שנתונה ע"י  $X = [\mathbf{1}_n \ X^{(j)}]$ , והשתמשו בעובדה הכללית (\*).

ג. בעובדה כללית (\*) הגדרנו  $z = (I - P_{-j})X^{(j)}$ . הראו שבביטוי הזה ניתן להחליף את  $X^{(j)}$  ב- $w$ , כלומר שבאופן כללי מתקיים:

$$z = (I - P_{-j})w$$

עבור הוקטור  $w$  שהוגדר בסעיף ב', והסיקו בעזרת עובדה כללית (\*) שברגרסיה מרובה של  $Y$  על  $X$ , אפשר להציג:

$$\hat{\beta}_j = \frac{1}{\|(I - P_{-j})w\|^2} [(I - P_{-j})w]^\top Y$$

ד. את הסעיפים הקודמים אפשר לסכם באופן הבא:

$$\hat{\beta}_j = \frac{1}{\|w\|^2} w^\top Y, \quad \text{ברגרסיה פשוטה (עם חותך):}$$

$$\hat{\beta}_j = \frac{1}{\|(I - P_{-j})w\|^2} [(I - P_{-j})w]^\top Y, \quad \text{ברגרסיה מרובה:}$$

כאשר  $w$  זה הוקטור שהוגדר בסעיף ב'. חשבו את  $\text{Var}(\hat{\beta}_j)$  בכל אחד משני המקרים. בנוסף, הראו שהיחס בין  $\text{Var}(\hat{\beta}_j)$  ברגסיה מרובה לבין  $\text{Var}(\hat{\beta}_j)$  ברגסיה פשוטה נתון ע"י:

$$\frac{\|w\|^2}{\|(I - P_{-j})w\|^2}.$$

ה. הסבירו שמתקיים:

$$\|(I - P_0)X^{(j)}\|^2 = \|w\|^2 = SST(j)$$

$$\|(I - P_{-j})X^{(j)}\|^2 = \|(I - P_{-j})w\|^2 = SSE(j)$$

ו. הסיקו:

$$VIF = \frac{1}{1 - R_j^2} = \frac{\|w\|^2}{\|(I - P_{-j})w\|^2},$$

כלומר זה בדיוק היחס בין השונות שחושב בסעיף ד'.

### שאלה 3:

בשאלה זו נניח כרגיל שהנתונים הם  $X \in \mathbb{R}^{n \times (p+1)}$ , מטריצה דטרמיניסטית (לא מקרית), ו-  $Y \in \mathbb{R}^n$ , וקטור מקרי שהתפלגותו תלויה ב-  $X$ , אבל המודל עבור  $Y$  יהיה

$$Y = \mu + \varepsilon, \quad \varepsilon \sim (0, \sigma^2 I_n), \quad (4)$$

כאשר  $\mu \in \mathbb{R}^n$  דטרמיניסטי ויכול להיות פונקציה של  $X$ .  
יהי  $\hat{Y} \in \mathbb{R}^n$  וקטור ערכים חזויים שהתקבלו מהנתונים  $(X, Y)$  בשיטה כלשהי (לא בהכרח מריבועים פחותים). נגדיר את תוחלת השגיאה הריבועית בחיזוי (MSPE) המתאימה בתור  $MSPE := \mathbb{E} \|\hat{Y} - \mu\|^2$ .  
כאן נתרכז בשיטות חיזוי שעבורן  $\hat{Y} = PY$ , כאשר  $P \in \mathbb{R}^{n \times n}$  היא מטריצת ההטלה על איזשהו תת-מרחב של  $\mathbb{R}^n$  מממד  $0 \leq r \leq n$ .  $P$  דטרמיניסטית, ויכולה - אך לא חייבת - להיות תלויה ב-  $X$ .

א. הוכיחו את התוצאה הכללית הבאה. יהיו  $Z, W \in \mathbb{R}^n$  זוג וקטורים מקריים בעלי משותף כלשהו. הראו כי מתקיים:

$$\text{tr}[\mathbb{E}[ZW^T]] = \text{tr}[\text{Cov}(Z, W)] + \text{tr}[\mathbb{E}(Z) \cdot [\mathbb{E}(W)]^T]$$

ב. נגדיר  $SSE(P) = \|Y - PY\|^2$ . הוכיחו:

$$MSPE = \mathbb{E}[SSE(P) + 2\sigma^2 r - n\sigma^2]$$

הדרכה - בדומה למה שראינו בכיתה, הציגו  $\|\hat{Y} - \mu\|^2 = \mathbb{E} \|(\hat{Y} - Y) + (Y - \mu)\|^2$ , הרחיבו את הביטוי בעזרת כללים אלמנטריים  $\|a + b\|^2 = \|a\|^2 + \|b\|^2 + 2a^T b$ , והשתמשו בסעיף א' כדי לסיים.

ג. תהי  $Y^*$  התממשות חדשה ממודל (4) ( $Y, Y^*$  הם iid). הראו:

$$\mathbb{E} \|Y^* - \hat{Y}\|^2 = MSPE + n\sigma^2 \quad (5)$$

$$\mathbb{E} \|Y - \hat{Y}\|^2 = \mathbb{E}[SSE(P)] \quad (6)$$

והסיקו

$$\mathbb{E} \|Y^* - \hat{Y}\|^2 - \mathbb{E} \|Y - \hat{Y}\|^2 = MSPE - \mathbb{E}[SSE(P)] \geq 0 \quad (7)$$

הערה: את הביטוי ב-(7) אפשר להבין בתור תוחלת ההפרש בין שגיאת החיזוי "מחוץ למדגם" ("Out - Of - Sample") לבין שגיאת החיזוי "בתוך המדגם" ("In - Sample"). אם בגישה נאיבית היינו מצפים שתוחלת ההפרש הזה היא אפס,

אז (7) מראה שבאופן כללי, תוחלת ההפרש הזה חיובית (כלומר,  $\|Y - \hat{Y}\|^2$  זה אומד "אופטימי" (optimistic) עבור  $\mathbb{E} \|Y^* - \hat{Y}\|^2$ ).