

# סטטיסטיקה ואקונומטריקה א'

פרק 6 - רגרסיה מרובה

תוכן העניינים

1. כללי .....

## רגרסיה מרובה:

### רקע:

ניבוי המשנה תלוי באמצעות יותר ממשנה ב"ית אחד.

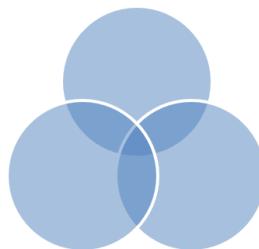
המודל אוכלוסייה:  $y = \alpha + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_n x_n$ .

מקדמי מודל הרגרסיה המרובה:

$\alpha$  = חותך אחד שמשמעותו: הציון המנווא כאשר כל המשתנים הב"ת = 0.

$\beta_1, \beta_2, \dots, \beta_n$  = מקדמי השיפוע. מס' הבוטות = מספר המשתנים הב"ת במודל.

משמעות מקדם השיפוע  $\beta_j$ : ההשפעה הייחודית של המשנה הב"ת מסוימת לניבוי המשנה תלוי, בנסיבות השפעתם של כל יתר המשתנים הב"ת האחרים המצויים במשוואת הרגרסיה.



### אמידת מודל הרגרסיה המרובה:

ברגרסיה מרובה, כמו ברגרסיה פשוטה, שיטת האמידה הטובה ביותר היא שיטת הריבועים הפחותים. כאמור, נרצה להביא את סכום הטעויות בניובי למינימום.

מפתרו פונקציית הריבועים הפחותים קיבל את אומדי הרגרסיה:  $\hat{\alpha}, \hat{\beta}_1, \dots, \hat{\beta}_n$ .

### מבחני מובהקות:

1. מבחן F למובהקות הרגרסיה:

בדיקה האם קיים קשר ליניארי בין המשנה תלוי  $Y$  לבין לפחות אחד מהמשתנים המסבירים.

$$H_0: \beta_1 = \beta_2 = \dots = \beta_k = 0$$

ההשערות הן:  $H_1: \text{Not } H_0: \text{at least one of the } \beta's \text{ is not } 0$

טבלת ניתוח שוונות (ANOVA) :

מקור	סכום ריבועים $SS$	דרגות חופש $d.f.$	ממוצע סכום ריבועים $MS = \frac{SS}{d.f.}$	$F_{st} \sim F_{k,n-k-1}$
מודל הרגression	$SSR$	$k$	$MSR = \frac{SSR}{K}$	$F_{st} = \frac{MSR}{MSE}$
שאריות	$SSE$	$n - k - 1$	$MSE = \frac{SSE}{(n - k - 1)}$	
סה"כ	$TSS$	$n - 1$		

$$\text{סטטיסטי המבחן : } F_{st} = \frac{MSR}{MSE}$$

כלל הכלראה : נדחה את  $H_0$  אם :  $F_{st} \geq F_{k,n-k-1}^{1-\alpha}$

חישוב סכומי הריבועים :

$$TSS = \sum_{i=1}^n y_i^2 - n\bar{y}^2$$

$$SSR = R^2 \cdot TSS$$

$$SSE = (1 - R^2) TSS$$

פרופורציצית השונות המוסברת :  $R^2$ 

$$R^2 = \frac{SSR}{SST} = 1 - \frac{SSE}{SST}$$

ברגression מרובה אומד זה לפרופורציצית השונות המוסברת הוא בעייתי שכן הוא מושפע ממספר המשתנים הב"ת במודל. אומד זה יכול רק לגודל בהוספה משתנים ב"ת למודל ולכן לא ניתן לנו אינדייקציה האם כדאי להוסיף אותם למודל או לא.

האומד המתוקן לפרופורציצית השונות המוסברת  $: AdjR^2$ 

$$\bar{R}^2 = 1 - \left[ \frac{(1 - R^2)(n - 1)}{n - k - 1} \right]$$

בניגוד ל-  $R^2$  לוקחים בחשבון את מספר המשתנים הב"ת במודל. יכול שלא לגודל ואף לקטונו בהוספה משתנה ב"ת שלא תורם תרומה משמעותית לניבוי.

2. מבוחן  $t$  לモביהקות משתנה ב'ית יחיד :

$$\begin{aligned} H_0: \beta_j &= 0 \\ \text{השערות:} \\ H_1: &\text{else} \end{aligned}$$

סטטיסטי המבחן וכלל הכרעת השערת האפס :

$$\cdot \left| t_{\hat{\beta}_j} \right| = \left| \frac{\hat{\beta}_j}{S_{\hat{\beta}_j}} \right| > t_{(T-k-1, 1-\frac{\alpha}{2})}$$

$$\cdot \hat{\beta}_j \pm t_{n-k-1, 1-\frac{\alpha}{2}} s.e.(\hat{\beta}_j) : \beta_j$$

3. מבוחן F חלקי (partial F) :  
 בודק את ההשערה שתוספת של משתנה אחד או קבוצה של משתנים מוסיפה תוספת מובהקת לניבוי המשתנה תלוי מעבר למשתנים אחרים שקייםים כבר במודול.

$$\begin{aligned} H_0: \beta_1 = \beta_2 = \dots = \beta_p &= 0 \\ \text{השערות:} \\ H_1: &\text{otherwise} \end{aligned}$$

ביצוע המבחן :

MRIIZIM SHETI RGRSIVOT :

1. UR : המודל המלא – רגרסיה עם כל המשתנים הב'ית ( $K$ ).

2. R : המודל החלקי – רגרסיה תחת  $H_0$  ( $K - P$ ).

$$\cdot F = \frac{\frac{R^2}{UR} - \frac{R^2}{P}}{\frac{1-R^2_{UR}}{n-k-1}} = \frac{\frac{SSR_{UR}}{P} - \frac{SSR_R}{P}}{\frac{1-SSR_{UR}}{n-k-1}} = \frac{\frac{SSE_R}{P} - \frac{SSE_{UR}}{P}}{\frac{1-SSE_{UR}}{n-k-1}}$$

סטטיסטי המבחן :

כלל הכרעה :

$$\cdot F > f_{1-\alpha}^{p, n-k-1}$$

**קשר בין מבוחן F חלקי למבוחן  $t$ :**

קיימים קשר בין מבוחן F חלקי לモביהקות תוספת משתנה ב'ית יחיד, למבוחן  $t$

$$\cdot F > f_{1-\alpha}^{1, n-k-1} = t_{\frac{1-\alpha}{2}, n-k-1}^2$$

לモביהקות אותו משתנה :

$$pvalue = pvalue$$

## מולטיקוליניאריות:

**בעיתת המולטיקוליניאריות:**

מצב שבו המשתנים המסבירים הם בעלי מתאם גובה ביןם ולבין עצם. מתאימים גובהם אלו יוצרים קושי ביכולת להבחין בהשפעה הבודדת של כל משתנה מסביר על המשתנה המוסבר.

**כיצד מזהים בעיתת מולטיקוליניאריות?**

- $R^2$  גדול וערכי קטנים, או שדוחים את השערת האפס ב מבחן F אך לא ב מבחני t.
- מקדמי מתאם זוגיים גבוהים בין המשתנים המסבירים (לפי טבלת קורלציות).

**איתור המולטיקוליניאריות:**

נבחר לרוב להוציא את אחד מהמשתנים לפיה הкрיטריונים הבאים:

1. המשתנה המתואם ביותר גם עם יתר המשתנים המסבירים.
2. המשתנה הכى פחות מובהק שהתקבל בהרצתה הראשונית (לפי t).

**אבחן קוליניאריות – מדד VIF (Variance Inflationary Factor) :**

מדד זה מבטא את הגידול (האינפלציה) בשינויו של מקדם הרגרסיה של משתנה ב'ית מסויים כתוצאה מהמתאים שלו עם משתנים ב'ית אחרים המצוים במשווה:

$$\cdot X_j = \frac{1}{(1-R_j^2)}$$

$R_j^2$  הוא מקדם ההסבר של  $X_j$ , כל יתר המשתנים ב'ית.

כאשר:  $VIF_j > 5$  המשמעות היא ש-  $X_j$  מתואם בצורה חזקה עם יתר המשתנים הב'ית במשווה.

**חישוב מובהקות התוספת (F חלק)** של משתנה ב'ית מסוים על פני האחרים:

במקרה של מולטיקוליניאריות במודל (מתאים חזק בין משתנים ב'ית) בכדי לדעת איזה משתנה ב'ית יש להוציא, ניתן לבדוק את התוספת לניבוי של המשתנה ה"חשוד" על פני האחרים. אם היא אינה מובהקת, זהה אינדיקציה שיש להוציאו מהמודל.

**שאלות:****המודל ומבחן המובוקות:**

1) לצורך בדיקת ההשערה שקיים קשר בין מספר המוניות בעיר באර שבע ( $y$ ) לבין מספר התושבים בעיר באלפיים ( $x_1$ ) ומספר הרכבים הפרטיים באלפיים ( $x_2$ ).

$$\text{הוחלט לבנות מודל רגרסיה מהצורה: } y_i = \beta_0 + \beta_1 x_{1i} + \beta_2 x_{2i} + \varepsilon_i, \text{ על סמך}$$

$$\text{הנתונים הבאים: } .MSE = 119.789, \sum_i y_i^2 = 338657, \sum_i 1673$$

א. ע"י הנתונים הניל, השלימו את טבלת ניתוח השונות הבאה.

אייזו השערה ניתן לבדוק באמצעותה? כתוב את ההשערה ובן אותה.

SOURCE	SS	DF	MS	F
Regression				
Error				
Total		8		

ב. חשבו את מדד טיב ההתאמנה. הסבר את משמעותו.

ג. נתונה טבלת המקדים (החלקית) הבאה:

	Coefficients	Standard Error	t Stat	P-value	Lower 95%	Upper 95%
Intercept	-511.727	114.9476				
X 1	9.208785		3.732167			
X 2	-8.79921	4.420456				

i. רשמו את האומדן למשוואת הרגרסיה ופרשו את מקדמיה.

ii. בחנו את ההשערה כי קיים קשר בין מספר הרכבים הפרטיים לבין מספר המוניות ברמת מובהקות של 5%.

iii. בנו רוח סמך למקדם של מספר התושבים בעיר.

iv. ענה ללא חישוב (על סמך הסעיפים הקודמים).

האם קיים קשר בין מספר התושבים לבין מספר המוניות ברמת מובהקות של 5%?

v. מהי תחזית מס' המוניות באאר שבע עבור 100,000 תושבים ו-52,000 מכוניות פרטיות?

vi. האם ניתן לסמן על תחזית זאת?

vii. חשב את סטטיסטי F חלק של מס' הרכבים הפרטיים. האם מובהק? (ענה ללא חישוב).

**תרגול מסכם:**

2) מעוניינים למצוא קשר בין מחיר הדירה (ב-\$) לבין ארבעה משתנים מסבירים:

1. שטח הדירה.

2. גודל שטח האמבטיה (Sqft). (b-).

3. מרחק הדירה מהים.

4. מהאוניברסיטה (במיילים).

לשם כך נדגומו מספר דירות והריצו רגرسיה אשר בה המשתנה המסביר הוא מחיר הדירה.

להלן פلت הרגסיה שהתקבל:

**Model Summary**

Mode	R	R Square	Adjusted R Square	Std. Error of the Estimate
1	.952 <sup>a</sup>			

a. Predictors: (Constant), Sea\_Dist, Apartment, Bath

**ANOVA<sup>b</sup>**

Model	Sum of Squares	df	Mean Square	F	Sig.
1 Regression					
Residual					
Total	1940484.615	25			

a. Predictors: (Constant), Univ\_Dist, Bath, Sea\_Dist, Apartment

b. Dependent Variable: Price

**Coefficients<sup>a</sup>**

Model	Unstandardized Coefficients		Beta	t	Sig.
	B	Std. Error			
1 (Constant)	-265.514	146.673		-1.810	.085
Apartment		.449	.722	6.572	
Bath	4.256		.297	2.687	.014
Sea_Dist	-32.114	11.090	-.223		.009
Univ_Dist	11.746	9.439	.095	1.244	.227

a. Dependent Variable: Price

**ענה על הסעיפים הבאים :**

א. מלאו את התאים החסרים בטבלה (אם לא ניתן למלא את כל התאים החסרים באופן מלא ונמקו באופן מפורש מדוע לא ניתן).

ב. כתבו את האומדן למשווהת מחיר הדירה בצורה מפורשת על סמך הפלט הניל. פרשו את מקדמי הרגסיה.

ג. בדקו האם ארבעת הגורמים ביחד אכן מסבירים את מחיר הדירה. הסבירו את המסקנה שהגעתם אליה. השתמשו ברמת מובהקות 5%.

- ד. הסבירו מהו ערך ה- Pvalue ומה ניתן להסיק ממנו לגבי המשתנים המסבירים?
- ה. בנו רוח סמך למקדם גודל שטח האטבטיה. השתמשו ברמת מובהקות של 2%.
- ו. ברמת מובהקות של 5% יש לבדוק האם המרחק מהאוניברסיטה אכן משפיע על מחיר הדירה.
- ז. האם במודל הרגרסיה הנוכחי ניתן לוותר על גורם המרחק מהיים? השתמשו ברמת מובהקות 1%.
- ח. בדקו את ההשערה כי קיים קשר חיובי בין גודל הדירה למחירה ברמת מובהקות של 5%.
- ט. נתונה מטריצת מקדמי המתאים הבאה:

	$X1$	$X2$	$X3$	$X4$
$X1$	1			
$X2$	0.228579	1		
$X3$	-0.22413	-0.13924	1	
$X4$	-0.24545	-0.97295	0.029537	1

- מה ניתן ללמוד ממנו ומה משמעותו לגבי המודל?
- י. הניחו כי השאריות המתקבלות מניתוח הרגרסיה הן בעלות הערכים הבאים (סדר הקראיה הוא משמאל לימין) הנח כי אלו כל השאריות הקיימות במודל: -1, 7, 9, -3, -7, -3, 12, 18, 6, -5, 10, 5, -4, 6, 3, -12, 7, 9, -3, -7, -3, 12, 18, 6, -5, 10, 5, -4, 6, 3, -12. האם משתנים  $X_2$  ו-  $X_4$  מוסיפים תוספת משמעותית לניבוי? אם לא ניתן לענות על השאלה, ציין מדוע.
- יא. מה יהיו תוצאות מבחן F לבדיקת התוספת לניבוי של המרחק מהאוניברסיטה על פני המשתנים האחרים (עננה ללא חישוב).

## תשובות סופיות:

. א. (1)

SOURCE	SS	DF	MS	F
Regression	26945.784	2	13472.892	
Error	718.733	6	119.789	112.414
Total	27664.517	8		

$$H_0: \beta_1 = \beta_2 = 0$$

. נדחה את השערת האפס,  
 $H_1: \text{at least one of the } \beta's \text{ is not } 0$

.  $\hat{y}_i = -511.727 + 9.208x_{1i} - 8.799x_{2i}$  ג.י. .97.4%

.  $p(3.17 \leq \beta_1 \leq 15.24) = 0.95$  iii. ii. אין עדות לכך.

.  $F = 3.96$  vi. v. 321 מוניות. iv. מובהק.

.1760019.5 .4. א. .92 .3. א. .0.89 .2. א. .0.907 .1. א. (2)

.440004.875 .8. א. .21 .7. א. .4 .6. א. .180465 .5. א.

. sig < 0.001 .12. א. .2.95 .11. א. .51.2 .10. א. .8593.57 .9. א.  
 .-2.896 .14. א. .1.58 .13. א.

.  $\hat{y}_i = -256.514 + 2.95x_{1i} + 4.256x_{2i} - 32.114x_{3i} + 11.746x_{4i}$  ב.

.  $p(1.016 \leq \beta_2 \leq 7.496) = 0.98$  ה. ד. ראו סרטון.

. א. אין עדות לכך. ג. לא.

. יא. ראו סרטון. ט. בעיית קולינニアריות. י. ראו סרטון.