

רגרסיה מרובה

ניבוי המשתנה התלוי באמצעות יותר ממשתנה ב"ת אחד.

$$y = \alpha + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_n x_n$$

המודל באוכלוסיה:

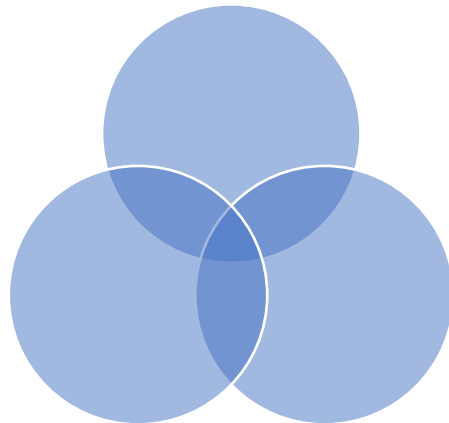
מקדמי מודל הרגרסיה המרובה:

α = חותך אחד שמשמעותו: הציון המנובא כאשר כל המשתנים הב"ת = 0.

$\beta_1 \dots \beta_j$ = מקדמי השיפוע. מס' הבטות = למספר המשתנים הב"ת במודל.

משמעות מקדם השיפוע β_j : ההשפעה הייחודית של המשתנה הב"ת המסוים לניבוי

המשתנה התלוי, בניכוי 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. כלומר, נרצה להביא את סכום הטעויות בניבוי למינימום.

מפיתרון פונקצית הריבועים הפחותים נקבל את אומדי הרגרסיה: $\hat{\alpha}, \hat{\beta}_1 \dots \hat{\beta}_j$

מבחני מובהקות:

(1) מבחן F למובהקות הרגרסיה

בדיקה האם קיים קשר ליניארי בין המשתנה התלוי Y לבין לפחות אחד מהמשתנים המסבירים.

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

השערות הן:

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

טבלת ניתוח שונות (ANOVA)

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

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

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

חישוב סכומי הריבועים:

$$TSS = \sum_{i=1}^n y_i^2 - n\bar{y}^2$$
$$SSR = R^2 \cdot TSS$$
$$SSE = (1 - R^2)TSS$$

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

פרופורציית השונות המוסברת R^2 :

ברגרסיה מרובה אומד זה לפרופורציית השונות המוסברת הוא בעייתי שכן הוא מושפע ממספר המשתנים ה"ב"ת במודל. אומד זה יכול רק לגדול בהוספת משתנים ב"ת למודל ולכן לא ייתן לנו אינדיקציה האם כדאי היה לי להוסיף אותם למודל או לא.

האומד המתוקן לפרופורציית השונות המוסברת $AdjR^2$:

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

בניגוד ל- R^2 לוקח בחשבון את מספר המשתנים ה"ב"ת במודל. יכול שלא לגדול ואף לקטון בהוספת משתנה ב"ת שלא תורם תרומה משמעותית לניבוי.

(2) מבחן t למובהקות משתנה ב"ת יחיד:

השערות:

$$H_0 : \beta_j = 0$$

$$H_1 : else$$

סטטיסטי המבחן וכלל הכרעת השערת האפס:

$$\left| T = \frac{\text{אמון מקדם}}{\text{סטיית תקן מקדם}} \right| > t_{1-\frac{\alpha}{2}}^{(n-k-1)} \quad for \quad n < 30$$

רווח בר סמך לאמידת ה- β_j :

$$\hat{\beta}_j \pm t_{n-k-1, 1-\alpha/2} \wedge s.e.(\hat{\beta}_j)$$

(3) מבחן F חלקי (partial F):

בודק את ההשערה שתוספת של משתנה אחד או קבוצה של משתנים מוסיפה תוספת מובהקת לניבוי המשתנה התלוי מעבר למשתנים אחרים שקיימים כבר במודל.

השערות:

$$p < k \quad H_0 : \beta_1 = \beta_2 = \dots = \beta_p$$
$$H_1 : \text{אחרת}$$

ביצוע המבחן:

מריצים שתי רגרסיות:

(1) UR - המודל המלא-רגרסיה עם כל המשתנים הב"ת (K)

(2) R - המודל החלקי-רגרסיה תחת H0 (K-P)

סטטיסטי המבחן:

$$F = \frac{R_{UR}^2 - R_R^2 / p}{1 - R_{UR}^2 / n - k - 1} = \frac{SSR_{UR} - SSR_R / p}{1 - SSR_{UR} / n - k - 1} = \frac{SSE_R - SSE_{UR} / p}{1 - SSE_{UR} / n - k - 1}$$

כלל הכרעה:

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

לדוגמא: נתונים 4 משתנים ב"ת לניבוי משתנה תלוי מסויים. רוצים לבדוק האם משתנה X1

ו-X2 מוסיפים תוספת משמעותית לניבוי Y מעבר למשתנים X3 ו-X4.

בהרצת רגרסיה עם כל המשתנים הב"ת התקבל $R^2 = 0.982$

בהרצת רגרסיה עם משתנים X3 ו-X4 בלבד התקבל $R^2 = 0.935$

השערות:

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

H1: אחרת

חישוב סטטיסטי המבחן:

$$F = \frac{R_{UR}^2 - R_R^2 / p}{1 - R_{UR}^2 / n - k - 1} = \frac{0.982 - 0.935 / 2}{1 - 0.982 / 12 - 4 - 1} = \frac{0.0235}{0.00257} = 9.144$$

כלל הכרעה:

$$F = 9.144 > F_{0.95}^{2,9} = 4.257$$

לכן יש סיבה מספקת לדחות את H_0 ברמת מובהקות של 0.05.

מסקנה: המשתנים X_1 ו- X_2 מוסיפים תוספת מובהקת לניבוי של Y מעבר ליתר המשתנים ה"ב"ת במשוואה (X_3 ו- X_4).

קשר בין מבחן F חלקי למבחן t:

קיים קשר בין מבחן F למובהקות תוספת משתנה ב"ת יחיד למבחן t למובהקות אותו משתנה:

$$F_{1-\alpha}^{1, n-k-1} = t_{1-\frac{\alpha}{2}, n-k-1}^2$$

$$pvalue = pvalue$$

תירגול

לצורך בדיקת ההשערה שקיים קשר בין מספר המוניות בעיר באר שבע (y) לבין מספר התושבים בעיר באלפים (x_1) ומספר הרכבים הפרטיים באלפים (x_2). הוחלט לבנות מודל רגרסיה מהצורה: $y_i = \beta_0 + \beta_1 x_{1i} + \beta_2 x_{2i} + \varepsilon_i$, על סמך הנתונים הבאים:

$$MSE=119.789, \sum_i y_i^2 = 338657, \sum_i y_i = 1673$$

א. ע"ס הנתונים הנ"ל, השלימו את טבלת ניתוח השונות הבאה. איזו השערה ניתן לבדוק באמצעותה? כתוב את ההשערה ובחן אותה.

SOURCE	SS	DF	MS	F
Regression				
Error				
Total		8		

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

ג. נתונה טבלת המקדמים (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				

1. רשמו את האומדן למשוואת הרגרסיה ופרשו את מקדמיה.
2. בחנו את ההשערה כי קיים קשר בין מספר הרכבים הפרטיים לבין מספר המוניות ברמת מובהקות של 5%.
3. בנו רווח סמך למקדם של מספר התושבים בעיר.
4. ענה ללא חישוב (על סמך הסעיפים הקודמים) - האם קיים קשר בין מספר התושבים לבין מספר המוניות ברמת מובהקות 5%?
5. מהי תחזית מס' המוניות בבאר שבע עבור 100,000 תושבים ו-52,000 מכוניות פרטיות?
6. האם ניתן לסמוך על תחזית זאת?
7. חשב את סטטיסטי F חלקי של מס' הרכבים הפרטיים. האם מובהק (ענה ללא חישוב).

חישוב מובהקות התוספת (F חלקי) של משתנה ב"ת מסוים על פני האחרים:

במקרה של מולטיקוליניאריות במודל (מתאם חזק בין משתנים ב"ת), בכדי לדעת איזה משתנה ב"ת יש להוציא, ניתן לבחון את התוספת לניבוי של המשתנה ה"חשוד" על פני האחרים. אם היא איננה מובהקת, זוהי אינדיקציה שיש להוציא מהמודל.

במקרה שלנו נבחן את התוספת לניבוי של X4 על פני המשתנים הב"ת האחרים:

Model Summary

Model	R	R Square	Adjusted R Square	Std. Error of the Estimate
1	.991 ^a	.982	.974	2.44601

a. Predictors: (Constant), X4, X3, X1, X2

Model Summary

Model	R	R Square	Adjusted R Square	Std. Error of the Estimate
1	.991 ^a	.982	.976	2.31206

a. Predictors: (Constant), X3, X2, X1

$$F = \frac{(R_{UR}^2 - R_R^2) / P}{(1 - R_{UR}^2) / n - k - 1} = \frac{(0.982 - 0.982) / 1}{(1 - 0.982) / 13 - 4 - 1} = 0$$

מסקנה: X4 לא מוסיף תוספת מובהקת למודל.

תרגול מסכם:

מעוניינים למצוא קשר בין מחיר הדירה (ב-\$) לבין ארבעה משתנים מסבירים: (1) שטח הדירה ו- (2) גודל שטח האמבטיה (ב-Sqft) וכן (3) מרחק הדירה מהים ו- (4) מהאוניברסיטה (במיילים).

לשם כך נדגמו מספר דירות והריצו רגרסיה אשר בה המשתנה המוסבר הוא מחיר הדירה. להלן פלט הרגרסיה שהתקבל:

Model Summary

Model	R	R Square	Adjusted R Square	Std. Error of the Estimate
1	.952 ^a	○	○	○

a. Predictors: (Constant), Sea_Dist, Apartment, Bath

ANOVA^b

Model		Sum of Squares	df	Mean Square	F	Sig.
1	Regression					.000 ^a
	Residual					
	Total	1940484.615	25			

a. Predictors: (Constant), Univ_Dist, Bath, Sea_Dist, Apartment

b. Dependent Variable: Price

Coefficients^a

Model		Unstandardized Coefficients		Standardized Coefficients	t	Sig.
		B	Std. Error	Beta		
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 נמקו באופן מפורש מדוע לא ניתן).
- ב. כתבו את האומדן למשוואת מחיר הדירה בצורה מפורשת על סמך הפלט הנ"ל. פרשו את מקדמי הרגרסיה.
- ג. בדקו האם ארבעת הגורמים ביחד אכן מסבירים את מחיר הדירה. הסברו את המסקנה שהגעתם אליה. השתמשו ברמת מובהקות 5%.
- ד. הסברו מהו ערך ה-Pvalue ומה ניתן להסיק ממנו לגבי המשתנים המסבירים?
- ה. בנו רווח סמך למקדם גודל שטח האמבטיה. 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 מהים? השתמשו ברמת מובהקות 1%.
- ח. בדקו את השערה כי קיים קשר חיובי בין גודל הדירה למחירה ברמת מובהקות של 5%.

ט. נתונה מטריצת מקדמי המתאם 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 ממנה ומה משמעותה לגבי המודל?

י. האם משתנים X2 ו-X4 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 (ענה ללא חישוב)

