

# גרסיה ושיטות ניתוח ליניאריות



## תוכן העניינים

1. ניתוח שונות חד כיוונית ..... 1
2. ניתוח שונות דו כיוונית ..... 10
3. רגרסיה ליניארית חד משתנית ..... (ללא ספר)
4. משתנה דמי ..... 46
5. קו הרגרסיה במדגם ..... (ללא ספר)
6. מובהקות הרגרסיה באוכלוסיה ..... (ללא ספר)
7. מאפייני קו הרגרסיה המרובה במדגם ..... (ללא ספר)
8. מובהקות קו הרגרסיה המרובה ומקדמיו באוכלוסיה ..... (ללא ספר)
9. שיטות להרצת רגרסיה רבת משתנים ..... (ללא ספר)
10. רגרסיה לוגיסטית ..... 63
11. מבחן לדוגמה 1 ..... (ללא ספר)
12. מבחן לדוגמה 2 ..... (ללא ספר)

# רגרסיה ושיטות ניתוח ליניאריות

פרק 1 - ניתוח שונות חד כיוונית

תוכן העניינים

1. כללי..... 1

## ניתוח שונות חד כיוונית

### רקע תיאורטי

ניתוח שונות (חד כיוונית) הוא מבחן להשוואת תוחלות  $(\mu_1, \dots, \mu_k)$  של  $k$  אוכלוסיות שונות. לכן, בנייתוח שונות, השערות המחקר הן:

$$H_0: \mu_1 = \mu_2 = \dots = \mu_k \quad (\text{התוחלות של כל האוכלוסיות שוות})$$

$$H_1: \quad \text{אחרת} \quad (\text{לפחות שתיים מהתוחלות שונות})$$

### ההנחות הדרושות לביצוע התהליך:

(2) בכל אוכלוסייה מתוך  $k$  האוכלוסיות ההתפלגות נורמלית.

(3) כל האוכלוסיות הן עם אותה שונות  $\sigma^2$ .

(4) המדגמים בלתי תלויים זה בזה.

ישנו משתנה המבדיל בין הקבוצות השונות, הוא המשתנה הבלתי תלוי הנקרא גורם (factor). משתנה זה הוא קטגוריאלי עם  $k$  רמות (levels). כדי לבצע את התהליך יש לבצע מדגם מכל אוכלוסייה: נסמן ב- $n_i$  את גודל המדגם בקבוצה  $i$ .

$$n = \sum_{i=1}^k n_i \quad \text{- מספר התצפיות סך הכול (בכל המדגמים).}$$

$\bar{X}_1$  - ממוצע המדגם הראשון,  $\dots, \bar{X}_k$  - ממוצע המדגם ה- $k$ .  
 $\bar{X}$  - ממוצע כללי (של כל המדגמים).

$$SS_B = \sum_{i=1}^k n_i [\bar{X}_i - \bar{X}]^2 \quad \text{: סכום ריבועים בין הקבוצות}$$

$$SS_W = \sum_{i=1}^k n_i [n_i - 1] \cdot \hat{S}_i^2 \quad \text{: סכום ריבועים בתוך הקבוצות}$$

$$SS_T = \sum_{i=1}^k \sum_{j=1}^{n_i} [X_{ij} - \bar{X}]^2 \quad \text{: סכום ריבועים כללי}$$

$$SST = SSB + SSW$$

יש למלא את טבלת ניתוח השונות הבאה:

מקור השונות	סכום הריבועים SS	דרגות חופש $df$	ממוצע הריבועים MS	F
B - בין הקבוצות	SSB	$k - 1$	$\frac{SSB}{k - 1}$	$\frac{MSB}{MSW}$
W - בתוך הקבוצות	SSW	$n - k$	$\frac{SSW}{n - k}$	
T - סה"כ	SST	$n - 1$		

$$F = \frac{\frac{SSB}{k-1}}{\frac{SSW}{n-k}} \sim F(k-1, n-k)$$

אזור דחיית  $H_0$ :  $1 - \alpha : F > F_{(k-1, n-k)}$

**שאלות**

- (1) מחקר מעוניין להשוות בין שלוש תרופות לשיכוך כאבים במטרה לבדוק האם קיים הבדל בין התרופות מבחינת הזמן בדקות שלוקח עד שהתרופה משפיעה. לצורך הבדיקה נלקחו 15 אנשים שסובלים מכאבי ראש. אנשים אלה חולקו באקראי לשלוש: קבוצה 1 קיבלה "אקמול" קבוצה 2 קיבלה "אופטלגין" קבוצה 3 קיבלה "נורופן". כל אדם במחקר מסר את מספר הדקות עד שהתרופה השפיעה עליו.
- א. מהו המשתנה התלוי ומהו המשתנה הבלתי תלוי במחקר?
  - ב. מהו המבחן הסטטיסטי המתאים כאן? רשמו את ההשערות.
  - ג. מה הן ההנחות הדרושות כדי לבצע את המבחן הסטטיסטי שהצעת בסעיף הקודם?

- (2) בעיר מסוימת שלושה בתי ספר תיכון. ראש העיר התעניין לבדוק האם קיים הבדל בהצלחה של בתי הספר במקצוע מתמטיקה. לצורך כך הוא דגם מספר תלמידים שנבחנו במבחן הבגרות במתמטיקה ברמה של 3 יחידות בעירו ובדק עבור כל תלמיד מה ציון הבגרות שלו במתמטיקה. להלן הציונים שהתקבלו:

"הס"	"רבין"	"המתמיד"
85	98	78
83	62	65
74	55	70
85	80	90
75		56

- א. מהו המבחן הסטטיסטי המתאים? רשמו את ההשערות ואת ההנחות של המבחן.
- ב. מהו גודל המדגם? מהו המשתנה הבלתי תלוי (factor) כמה רמות יש לו?
- ג. חשבו את הממוצע ואת סטיית התקן של הציונים בכל אחד מהמדגמים.
- ד. מלאו את טבלת ANOVA.
- ה. רשמו את כלל ההכרעה למבחן שהוצע בסעיף א ברמת מובהקות של 5%.
- ו. האם קיים הבדל בין בתי הספר בעיר מבחינת רמת הצלחת התלמידים במקצוע המתמטיקה? ענה על סמך הסעיפים הקודמים.

- (3) מעוניינים לבדוק האם יש הבדל בהשפעה של שיטות טפול שונות על לחץ הדם הסיסטולי (SBP) באוכלוסייה של קשישים. נבדקו 4 שיטות שונות. בטבלה המצורפת מרוכזים ממצאי המחקר.

השיטה	A	B	C	D
גודל המדגם	12	14	8	12
הממוצע	178	172	180	182
סטיית התקן	4	8	5	3

- א. רשמו את השערות המחקר וההנחות הדרושות כדי לבצע את המבחן המתאים.
- ב. מה מסקנת המחקר ברמת מובהקות של 5%?
- ג. האם יש צורך לבצע השוואות מרובות?

- 4) שלושה אופים נתבקשו להכין עוגת שוקולד. לכל אופה בדקו את משך זמן ההכנה בדקות. כל אופה נדרש לאפות בכל יום 4 עוגות. האם קיים הבדל בין האופים מבחינת תוחלת זמני ההכנה של העוגות? בדקו ברמת מובהקות של 5%.

האופה	ניר	מוזס	שלום
סכום הזמנים	206	212	182
סכום ריבועי הזמנים	10644	11250	8982

- 5) להלן טבלת ניתוח שונות חד כיוונית. במחקר בחנו 4 סוגי סוללות. רצו לבדוק האם לסוג הסוללה השפעה על תוחלת אורך החיים שלה. הפעילו את כל הסוללות על אותו מכשיר ובדקו את אורך החיים של כל סוללה בשעות. מה המסקנה ברמת מובהקות של 10%? רשמו את ההשערות וההנחות הדרושות.

### ANOVA

	Sum of Squares	df	Mean Square	F	Sig.
Between Groups	10.317	3	3.439	1.361	.279
Within Groups	60.648	24	2.527		
Total	70.964	27			

- 6) להלן טבלת ANOVA בטבלה הושמטו חלקים. השלימו את החלקים בטבלה שהושמטו ומסומנים באותיות.

### ANOVA

	Sum of Squares	df	Mean Square	F	Sig.
Between Groups	357.450	ב	ג	ה	.000
Within Groups	א	17	ד		
Total	522.950	19			



7) חברת תרופות לקחה 15 אנשים ברמת בריאות דומה. החברה חילקה את האנשים ל שלוש קבוצות שוות בגודלן. לכל קבוצה ניתנה אותה תרופה במינון שונה (dosage). המינונים שניתנו הם: 10 מ"ג, 20 מ"ג ו-30 מ"ג. לאחר שעה מזמן לקיחת התרופה נבדק קצב פעימות הלב של כל אדם (pulse). הנתונים הוזנו לתוכנה סטטיסטית והתקבלו התוצאות הבאות:

ANOVA						pulse			
pulse						Tukey HSD <sup>a</sup>			
	Sum of Squares	df	Mean Square	F	Sig.	dosage	N	Subset for alpha = 0.05	
								1	2
Between Groups	414.400	2	207.200	19.733	.000	30.00	5	71.0000	
Within Groups	126.000	12	10.500			20.00	5		80.2000
Total	540.400	14				10.00	5		83.4000
						Sig.		1.000	.299

Means for groups in homogeneous subsets are displayed.  
a. Uses Harmonic Mean Sample Size = 5.000.

Post Hoc Tests

Multiple Comparisons

(I) dosage		(J) dosage	Mean Difference (I-J)	Std. Error	Sig.	95% Confidence Interval	
						Lower Bound	Upper Bound
10.00		20.00	3.20000	2.04939	.299	-2.2675	8.6675
		30.00	12.40000*	2.04939	.000	6.9325	17.8675
20.00		10.00	-3.20000	2.04939	.299	-8.6675	2.2675
		30.00	9.20000*	2.04939	.002	3.7325	14.6675
30.00		10.00	-12.40000*	2.04939	.000	-17.8675	-6.9325
		20.00	-9.20000*	2.04939	.002	-14.6675	-3.7325

\*. The mean difference is significant at the 0.05 level.

- א. בדקו ברמת מובהקות של 5% האם קיים הבדל בין המינונים השונים מבחינת תוחלת הדופק של האנשים? רשמו את ההשערות וההנחות הדרושות לצורך פתרון.
- ב. הסבירו ללא חישוב כיצד הייתה משתנה התשובה לסעיף הקודם אם הינו מעלים את הדופק של כל התצפיות במחקר ב-2.
- ג. האם יש צורך במחקר בהשוואת מרובות. נמקו!
- ד. לטבלת ANOVA צורפו טבלאות של השוואות מרובות בשיטה הנקראת "טוקי". ברמת בטחון של 95% מה הם הממצאים לפי שיטה זו?

- 8) בעיר מסוימת רצו לבדוק האם קיים הבדל ברמה של התלמידים בין בתי הספר השונים בעיר. ביצעו מדגם מכל בית ספר ונתנו מבחן זהה לכל הנדגמים. לאחר מכן ריכזו את הנתונים בתוכנה סטטיסטית והפעילו ניתוח שונות. מצורפים הפלטים שהתקבלו. ענו על הסעיפים הבאים:
- כמה בתי ספר יש בעיר?
  - כמה תלמידים השתתפו בסך הכול במחקר?
  - האם קיים הבדל בין בתי הספר בעיר מבחינה רמת הציונים? בדקו ברמת מובהקות של 1%
  - בביטחון של 95% אילו בתי ספר שונים זה מזה ברמת התלמידים? נמקו והסבירו.

### Oneway

#### ANOVA

grade

	Sum of Squares	df	Mean Square	F	Sig.
Between Groups	7799.600	4	1949.900	13.586	.000
Within Groups	2870.400	20	143.520		
Total	10670.000	24			

**Post Hoc Tests**

**Multiple Comparisons**

grade

Scheffe

(I) school	(J) school	Mean Difference (I-J)	Std. Error	Sig.	95% Confidence Interval	
					Lower Bound	Upper Bound
1.00	2.00	5.40000	7.57681	.971	-20.2543	31.0543
	3.00	36.80000*	7.57681	.003	11.1457	62.4543
	4.00	36.40000*	7.57681	.003	10.7457	62.0543
	5.00	-2.60000	7.57681	.998	-28.2543	23.0543
2.00	1.00	-5.40000	7.57681	.971	-31.0543	20.2543
	3.00	31.40000*	7.57681	.011	5.7457	57.0543
	4.00	31.00000*	7.57681	.013	5.3457	56.6543
	5.00	-8.00000	7.57681	.888	-33.6543	17.6543
3.00	1.00	-36.80000*	7.57681	.003	-62.4543	-11.1457
	2.00	-31.40000*	7.57681	.011	-57.0543	-5.7457
	4.00	-.40000	7.57681	1.000	-26.0543	25.2543
	5.00	-39.40000*	7.57681	.001	-65.0543	-13.7457
4.00	1.00	-36.40000*	7.57681	.003	-62.0543	-10.7457
	2.00	-31.00000*	7.57681	.013	-56.6543	-5.3457
	3.00	.40000	7.57681	1.000	-25.2543	26.0543
	5.00	-39.00000*	7.57681	.001	-64.6543	-13.3457
5.00	1.00	2.60000	7.57681	.998	-23.0543	28.2543
	2.00	8.00000	7.57681	.888	-17.6543	33.6543
	3.00	39.40000*	7.57681	.001	13.7457	65.0543
	4.00	39.00000*	7.57681	.001	13.3457	64.6543

\*. The mean difference is significant at the 0.05 level.

**Homogeneous Subsets**

grade

Scheffe<sup>a</sup>

school	N	Subset for alpha = 0.05	
		1	2
3.00	5	45.0000	
4.00	5	45.4000	
2.00	5		76.4000
1.00	5		81.8000
5.00	5		84.4000
Sig.		1.000	.888

Means for groups in homogeneous subsets are displayed.

a. Uses Harmonic Mean Sample Size = 5.000.

**תשובות סופיות**

1) א. משתנה בלתי תלוי : סוג התרופה. ב. ניתוח שונות חד כיוונית

$$H_0 : \mu_1 = \mu_2 = \mu_3$$

$$H_1 : otherwise$$

ג. 1. מדגמים בלתי תלויים.

2. שווין שונויות.

3. משתנים מתפלגים נורמלית.

2) א. המבחן לניתוח שונות חד כיוונית.

$$H_0 : \mu_1 = \mu_2 = \mu_3$$

$$H_1 : otherwise$$

הנחות:

1. מדגמים בלתי תלויים.

2. משתנים מתפלגים נורמלית.

3. שווין שונויות.

ב. גודל המדגם: 14. משתנה ב"ת: בית הספר, בעל 3 רמות.

$$g. \bar{X} = 71.8, \hat{S} = 12.93, \bar{X} = 73.75, \hat{S} = 19.29, \bar{X} = 80.4, \hat{S} = 5.46$$

ד. להלן טבלה:

F	MS	df	SS	מקור השונות
	100.3	2	200.6	B
	173.2	11	1904.75	W
0.58		13	2105.35	סה"כ

ה.  $F > 3.98$ .

ו. נקבל את  $H_0$ .

3) א.  $H_0 : \mu_1 = \mu_2 = \mu_3 = \mu_4$   
 ב. נדחה את  $H_0$ .  
 ג. כן.

$$H_1 : otherwise$$

הנחות:

1. מדגמים בלתי תלויים.

2. שווין שונויות.

3. משתנים מתפלגים נורמלית.

4) נקבל את  $H_0$  : נכריע שאין הבדל מובהק בין האופים מבחינת תוחלת זמן הכנה.

$$H_0 : \mu_1 = \mu_2 = \mu_3 = \mu_4 \quad (5)$$

$$H_1 : otherwise$$

הנחות :

1. מדגמים בלתי תלויים.

2. שוויון שונות.

3. משתנים מתפלגים נורמלית.

נקבל את  $H_0$  : לסוג סוללה אין השפעה של תוחלת החיים ברמת ביטחון של 10%.

6) להלן טבלה :

### ANOVA

	Sum of Squares	df	Mean Square	F	Sig.
Between Groups	357.450	ב 2	א 178.725	ה 18.36	.000
Within Groups	א 165.5	17	ד 9.735		
Total	522.950	19			

$$H_0 : \mu_1 = \mu_2 = \mu_3 \quad (7)$$

$$H_1 : otherwise$$

הנחות :

1. מדגמים בלתי תלויים.

2. משתנים מתפלגים נורמלית.

3. שוויון שונות.

נדחה את  $H_0$  : ברמת ביטחון של 5% קיים הבדל במינונים השונים מבחינת תוחלת הדופק.

$$\text{ב. ראה וידאו. ג. כן. ד. } \mu_{20} = \mu_{10} > \mu_{30} .$$

$$5 \text{ א. ב. } 25 \quad (8)$$

ג. נדחה את  $H_0$  : יש לפחות שני בתי ספר בעיר עם תוחלת רמת ציונים שונה.

$$\text{ד. } (\mu_3 = \mu_4) < (\mu_1 = \mu_2 = \mu_3) .$$

# רגרסיה ושיטות ניתוח ליניאריות

פרק 2 - ניתוח שונות דו כיווני

תוכן העניינים

10	1. הקדמה
20	2. אפקטים פשוטים, עיקריים ואינטראקציה
32	3. תהליך ניתוח שונות דו כיווני

## ניתוח שונות דו-כיווני - הקדמה

### רקע

ראשית, נחזור על עיקרי ניתוח השונות החד-כיווני (חד-גורמי).

בניתוח שונות חד-כיווני יש משתנה תלוי יחיד, שהוא כמותי, ומשתנה בלתי תלוי יחיד, שהוא משתנה קטגוריאלי (משתנה שהערכים שלו שייכים למספר סופי של קטגוריות). המשתנה הקטגוריאלי נקרא לעתים גם גורם (פקטור), והקטגוריות שלו נקראות רמות. המטרה בניתוח שונות חד-כיווני היא לבדוק האם לגורם יש השפעה מובהקת על המשתנה התלוי. השערת האפס של המחקר בניתוח שונות חד-כיווני היא שבכל הקטגוריות יש אותה התוחלת, והשערת המחקר טוענת שיש לפחות שתי קטגוריות שבהן התוחלות שונות.

### דוגמה: (פתרון בהקלטה)

נבדקו שלושה סוגי דיאטות על אנשים בעלי משקל עודף. נבחרו 30 מטופלים בעלי משקל עודף, והם חולקו באקראי לשלוש קבוצות שוות בגודלן, כך שכל קבוצה קיבלה דיאטה נחקרת אחרת. כעבור שלושה חודשים בדקו את מספר הקילוגרמים שהפחית כל מטופל ממשקלו בתקופה זו. מטרת המחקר הייתה לבדוק האם קיים הבדל בין הדיאטות מבחינת ההפחתה במשקל.

- מהו המשתנה התלוי במחקר?
- מהו המשתנה הבלתי תלוי במחקר? כמה רמות יש לו?
- מה הן השערות המחקר?
- מהו המבחן הסטטיסטי המתאים?

בניתוח שונות דו-כיווני אנו מוסיפים עוד משתנה בלתי תלוי למחקר, כלומר עוד גורם שאנו רוצים לבדוק איך הוא משפיע על המשתנה התלוי. לכן בניתוח שונות דו-כיווני יש משתנה תלוי כמותי יחיד ושני משתנים בלתי תלויים שכל אחד מהם קטגוריאלי. כזכור, למשתנים הבלתי תלויים אנו קוראים גם גורמים (פקטורים), ומספר הקטגוריות של כל גורם נקרא גם מספר הרמות שלו. ניתוח שונות רב-כיווני או רב-גורמי הוא ניתוח שונות שבו יש יותר מגורם אחד, כלומר יותר ממשתנה בלתי תלוי קטגוריאלי אחד. בניתוח שונות דו-כיווני יש שני גורמים, בניתוח שונות תלת-גורמי יש שלושה גורמים וכו'. ככל שנוסיף גורמים, הניתוח הסטטיסטי יהיה מורכב יותר ויידרשו יותר תצפיות למחקר אבל כיוון שהוא יקטין את שונות הטעויות (שונות מקרית) וייתן יותר הסבר לשונות הכללית, כך שהמבחן יהיה עוצמתי יותר.

המשך הדוגמה:

מבין 30 המטופלים שבמחקר 15 היו גברים ו-15 היו נשים. המטופלים חולקו כך שבכל דיאטה השתתפו 5 גברים ו-5 נשים.

מה הם המשתנים הבלתי תלויים? כמה רמות יש לכל משתנה?

בניתוח שונות דו-כיווני אנו בעצם רוצים לבדוק סימולטנית שלוש שאלות מחקר על אוכלוסיית כבדי המשקל:

- האם יש הבדלים משמעותיים בין שיעורי הפחתת המשקל של מטופלים כבדי משקל כתוצאה משימוש בדיאטות שונות?
- האם יש הבדלים משמעותיים בין שיעורי הפחתת המשקל של מטופלים כבדי משקל כתוצאה ממגדר שונה?
- האם יש השפעה משולבת (אינטראקציה) של שני הגורמים הנבדקים על הפחתת המשקל של מטופלים כבדי משקל, כלומר האם צירוף של דיאטה מסוימת ומגדר מסוים מביא להפחתת משקל גדולה יותר או קטנה יותר מצירופים אחרים?

נסמן גורם אחד ב- $a$  ואת מספר הרמות שלו ב- $A$ . באותו האופן הגורם האחר יסומן ב- $b$ , ואת מספר הרמות שלו נסמן ב- $B$ . מספר הקבוצות הכולל שאנו יוצרים הוא  $A \cdot B$ .

המשך הדוגמה:

- בחרו גורם אחד להיות  $a$  וגורם אחר להיות  $b$ . מהו  $A$  ומהו  $B$ ?
- כמה קבוצות שונות נוצרו במחקר?

נסמן ב- $m$  את מספר התצפיות בכל תא (בהנחה שהוא יהיה מספר קבוע). תא הוא שילוב של רמה מסוימת של גורם  $a$  עם רמה מסוימת של גורם  $b$ .

המשך הדוגמה:

- כמה תאים (קבוצות) יש במחקר?
- מה ערכו של  $m$ ?
- מהו הקשר המתמטי בין  $m$  לבין  $n$ , גודל המדגם?

נסמן ב- $a_1$  את הרמה הראשונה של  $a$ , ב- $a_2$  את הרמה השנייה שלו וכך הלאה.  
 נסמן ב- $b_1$  את הרמה הראשונה של  $b$ , ב- $b_2$  את הרמה השנייה שלו וכך הלאה.  
 נסמן ב- $\mu_i$  את התוחלת ברמה  $a_i$ . נסמן ב- $\mu_j$  את התוחלת ברמה  $b_j$ . נסמן ב- $\mu_{ij}$  את התוחלת של תא  $ij$ .

המשך הדוגמה :

- מה המשמעות של  $\mu_1$  ושל  $\mu_2$  ?
- מה המשמעות של  $\mu_{12}$  ושל  $\mu_{21}$  ?

השערות המחקר בניתוח שונות דו-כיווני

את השערות המחקר בניתוח שונות דו-כיווני אפשר לרשום בצורות רבות :

לגורם  $a$  אין השפעה על המשתנה התלוי :  $H_0$

אחרת:  $H_1$

לגורם  $b$  אין השפעה על המשתנה התלוי :  $H_0$

אחרת:  $H_1$

אין אינטראקציה בין שני הגורמים :  $H_0$

אחרת:  $H_1$

דרך אחרת היא שימוש בתוחלות:

$$H_0: \mu_1 = \mu_2 = \dots = \mu_A.$$

$H_1$ : אחרת

$$H_0: \mu_1 = \mu_2 = \dots = \mu_B.$$

$H_1$ : אחרת

$H_0$ : אין אינטראקציה בין שני הגורמים

$H_1$ : אחרת

המשך הדוגמה:

אם אנחנו מעוניינים לבצע ניתוח שונות דו-כיווני, מה הן ההשערות הנחקרות?

## שאלות

- (1) בחברת טקסטיל בחנו 4 סוגי בדים שונים מבחינת חוזקם. דגמו 5 חתיכות בד מכל סוג ובדקו את חוזק הקריעה של כל סוג בד.
- מהו המשתנה התלוי במחקר?
  - כמה משתנים בלתי תלויים יש במחקר? מה הם?
  - מהו המבחן הסטטיסטי המתאים במקרה זה?
- (2) במחקר בתחום הפסיכולוגיה נדגמו אנשים הסובלים מחרדות מסוגים שונים. כל מטופל סווג כסובל מאחד מסוגי החרדות הבאים: חרדה חברתית, חרדה כללית או אגורפוביה. במחקר השתתפו 6 מטופלים מכל סוג חרדה שצוין. המטופלים במחקר חולקו כך שכל מטופל היה צריך לעבור במשך שנה את אחד מהטיפולים הבאים: טיפול קוגניטיבי התנהגותי (CBT), טיפול קבוצתי או טיפול דיאלקטי התנהגותי (DBT). בכל סוג טיפול השתתפו 2 מטופלים מכל סוג חרדה. בסוף השנה נבדקו כל המטופלים וקיבלו ציון כמותי על השיפור במצבם הנפשי (משתנה כמותי). מטרת המחקר הייתה לבדוק האם סוג החרדה, סוג הטיפול והשילוב ביניהם משפיעים על המצב הנפשי של המטופלים.
- מהו גודל המדגם?
  - מהו המשתנה התלוי במחקר הזה ומה הם המשתנים הבלתי תלויים?
  - כמה קטגוריות יש לכל משתנה בלתי תלוי?
  - כמה קבוצות שונות יש במערך המחקרי?
  - מהו המבחן הסטטיסטי המתאים במערך מחקרי זה?

3) מחקר שיווקי בדק את השפעת גובה המדף בסופרמרקט והשפעת החומר שממנו עשוי הבקבוק (זכוכית או פלסטיק) על היקף המכירות של משקאות קלים. נבדקו שני סופרמרקטים. בכל סופרמרקט נבחן כל צירוף אפשרי של גובה המדף וחומר הבקבוק, ועבור כל צירוף כזה נבדק מספר בקבוקי המשקה הקל שנמכרו באותו סופרמרקט ביום מסוים. הנה התוצאות שהתקבלו:

פלסטיק	זכוכית	סוג בקבוק
		גובה המדף
59	23	נמוך
63	32	
88	47	בינוני
90	55	
51	40	גבוה
56	48	

- מהו המבחן הסטטיסטי המתאים? נמקו.
- מהו מספר הרמות של כל גורם מחקרי?
- מה יהיו השערות המחקר אם יתבצע ניתוח שונות דו-כיווני?
- מהו ערכו של  $m$  ומהו ערכו של  $n$ ?

4) יצרן של נוזל כביסה מעוניין לבחון שני נוזלי ניקוי מבחינת יעילותם בהסרת כתמים בשלוש רמות טמפרטורה. בכל אחד מששת הצירופים של סוג נוזל וטמפרטורה נבחנה יכולת הסרת הכתמים מבדים דומים, וניתן ציון בין 0 ל-15 (הציון הטוב ביותר).

מספר סידורי	סוג הנוזל	טמפרטורה במעלות צלזיוס	ציון הסרת כתמים
1	C	30	4
2	C	30	5
3	C	30	4
4	C	30	6
5	C	40	6
6	C	40	6
7	C	40	7
8	C	40	6
9	C	60	9
10	C	60	8
11	C	60	7
12	C	60	10
13	w	30	9
14	w	30	9
15	w	30	9
16	w	30	10
17	w	40	12
18	w	40	13
19	w	40	11
20	w	40	11
21	w	60	14
22	w	60	14
23	w	60	15
24	w	60	13

- א. כמה משתנים יש במחקר?  
 ב. לגבי כל משתנה קבעו האם הוא משתנה תלוי או בלתי תלוי.  
 ג. כמה רמות יש לכל גורם?  
 ד. אם נבצע ניתוח שונות דו-כיוונית, מה יהיו השערות המחקר?  
 ה. רכזו את נתוני המחקר בטבלה שבה בשורות גורם אחד, בעמודות גורם שני ובתאים התוצאות שהתקבלו למשתנה התלוי.
- 5) קבעו לגבי כל אחד מהבאים האם הוא משתנה קטגוריאלי:**
- א. מספר הניתוחים שעבר אדם בחייו.  
 ב. אחוז האבטלה בישראל בחודש זה.  
 ג. סוג הדם של חולה.  
 ד. שונות הציונים בבחינת הבגרות באנגלית במועד האחרון.  
 ה. משקל חבילה בדואר בגרמים.  
 ו. היבשת שאירחה את משחקי המונדיאל.
- בשאלות הבאות יש לבחור את התשובה הנכונה ביותר:**
- 6) משרד החינוך רוצה לבדוק עד כמה שיטת הוראה (יש 3 שיטות הוראה מקובלות) ומגדר משפיעים על ציוני הבגרות בהיסטוריה. מהו המבחן הסטטיסטי המתאים למחקר זה?**
- א. מבחן T להשוואת תוחלות.  
 ב. ניתוח שונות חד-כיוונית.  
 ג. ניתוח שונות דו-כיוונית.  
 ד. מבחן T לתוחלת אחת.
- 7) מחלקת שירות הלקוחות של חברת החשמל דגמה עובדים כדי לבחון האם ככל שמספר שנות הוותק של נותן השירות גדול יותר גדל גם מספר הלקוחות שבו הוא מטפל במהלך משמרת. מהו המבחן הסטטיסטי שיכול לבדוק זאת?**
- א. מבחן T להשוואת תוחלות.  
 ב. ניתוח שונות חד-כיוונית.  
 ג. ניתוח שונות דו-כיוונית.  
 ד. אף אחת מהאפשרויות שלעיל.

8) האיחוד האירופי המשותף דגם 10 עובדים מתחום ההוראה בכל אחת מהמדינות הבאות: הולנד, צרפת, בלגיה, גרמניה ואוסטריה. לכל עובד בדקו את גובה המשכורת החודשית שלו ביורו. אם נרצה להשוות בין המדינות הללו מבחינת תוחלת השכר של עובדי ההוראה במדינה, מה יהיה המבחן הסטטיסטי המתאים?

- א. מבחן T להשוואת תוחלות
- ב. ניתוח שונות חד-כיווני
- ג. ניתוח שונות דו-כיווני
- ד. אף אחת מהאפשרויות שלעיל

9) בקו ייצור 2 סוגים של מכונות ו-3 רמות ותק של מפעיל המכונה (עד שנתיים במפעל, בין שנתיים ל- חמש שנים במפעל, יותר מחמש שנים במפעל). מנהל הייצור רוצה לבדוק אם קיימת השפעה של סוג המכונה והוותק של המפעיל על מספר המוצרים הפגומים שיוצאים מהמכונה. מה יהיה המבחן הסטטיסטי המתאים במקרה זה?

- א. מבחן T להשוואת תוחלות.
- ב. ניתוח שונות חד-כיווני.
- ג. ניתוח שונות דו-כיווני.
- ד. ניתוח שונות תלת-כיווני.

10) במחקר נאספו הנתונים הבאים על קבוצת נחקרים:

1. כמה כוסות קפה הנחקר שותה ביום: לא שותה / 1-2 כוסות/ יותר מ-2 כוסות.
2. מין הנחקר: גבר/אישה.
3. דופק (מספר פעימות בדקה) שעתיים אחרי הקימה.

מטרת המחקר הייתה לבדוק האם מספר כוסות הקפה שאדם שותה ביום משפיע על הדופק אצל גברים אחרת מאשר אצל נשים. מה יהיה המבחן הסטטיסטי המתאים במקרה זה?

- א. מבחן T להשוואת תוחלות.
- ב. ניתוח שונות חד-כיווני.
- ג. ניתוח שונות דו-כיווני.
- ד. ניתוח שונות תלת-כיווני.

- 11) במחקר יש משתנה כמותי אחד ושני גורמים שלכל אחד מהם שתי רמות. אילו מהמשפטים הבאים אינו נכון?
- א. אפשר מבחינה טכנית לבדוק כיצד כל גורם בנפרד משפיע על המשתנה התלוי באמצעות ניתוח שונות חד-כיווני שייערך לכל גורם בנפרד.
- ב. אפשר מבחינה טכנית להשוות בין התוחלות של כל רמה בגורם הראשון על ידי מבחן T להשוואת תוחלות.
- ג. אפשר מבחינה טכנית לבצע ניתוח שונות דו-כיווני במערך מחקרי זה.
- ד. כיוון שבמחקר יש בסך הכול שלושה משתנים, אפשר מבחינה טכנית לבצע ניתוח שונות תלת-כיווני.

### תשובות סופיות

- 1) א. חוזק הקריעה.  
ג. ניתוח שונות חד גורמי.
- 2) א. 18  
ב. המשתנה התלוי: ציון במצב הנפש. המשתנים הב"ת: סוג חרדה, סוג הטיפול.  
ג. 3,3  
ד. 9
- 3) א. ניתוח שונות דו גורמי.  
ב. 3,2  
ג. ניתוח שונות דו גורמי.  
ד.  $H_0$ : אין אינטראקציה,  $H_1$ : יש אינטראקציה.  $m = 2, n = 12$
- 4) א. 3  
ב. משתנים ב"ת: סוג הנוזל, טמפרטורה. משתנה תלוי: ציון הסרת כתמים.  
ג. 3,2  
ד.  $H_0$ : אין אינטראקציה בין הגורמים,  $H_1$ : אחרת.  
ה. עיין בסרטון הוידאו.
- 5) א. כן.  
ב. לא.  
ג. כן.  
ד. לא.  
ה. תלוי.  
ו. כן.
- 6) ג.  
7) ד.  
8) ב.  
9) ג.  
10) ג.  
11) ד.

## אפקטים פשוטים, עיקריים ואינטראקציה

### רקע

בניתוח שונות דו-כיווני אנו דנים במשתנה כמותי תלוי יחיד ובשני משתנים בלתי תלויים (גורמים) המחולקים כל אחד למספר רמות. מטרת המחקר היא לבדוק שלוש השערות שונות:

לגורם  $a$  אין השפעה על המשתנה התלוי:  $H_0$

אחרת:  $H_1$

לגורם  $b$  אין השפעה על המשתנה התלוי:  $H_0$

אחרת:  $H_1$

אין אינטראקציה בין שני הגורמים:  $H_0$

אחרת:  $H_1$

נרצה להבין מה בדיוק כל השערה בודקת לגבי האוכלוסייה הנחקרת.

**אפקט עיקרי:** אם יש שתי קטגוריות (רמות) לפחות של גורם מסוים שהתוחלות שלהן שונות, נאמר שלגורם זה יש השפעה על המשתנה התלוי. השפעה זאת נקראת "אפקט עיקרי". למשל, אם יימצאו לפחות שתי תרופות נוגדות דיכאון שונות שמביאות לתוחלות שונות במצב הנפשי, נגיד שלסוג התרופה יש השפעה על המצב הנפשי, כלומר יש אפקט עיקרי. כמות האפקטים העיקריים שאפשר למצוא היא כמות הגורמים במחקר.

אפקט אינטראקציה: מצב שבו גורם אחד משפיע על המשתנה התלוי באופן שונה בקטגוריות שונות של הגורם השני. למשל, תרופה נוגדת דיכאון אחת מביאה את הגברים למצב רוח טוב יותר מאשר את הנשים לעומת תרופה אחרת שמביאה דווקא את הנשים למצב רוח טוב יותר מאשר את הגברים. אפקט האינטראקציה הוא יחיד, כלומר נאמר אם יש או אין אינטראקציה. כמו כן הוא אפקט סימטרי: אם קיימת אינטראקציה בין מגדר לסוג התרופה, יש גם אינטראקציה בין סוג התרופה למגדר.

אפקט פשוט: אפקט פשוט מתייחס להשפעת גורם אחד על המשתנה התלוי בתוך קטגוריה מסוימת של הגורם השני. למשל, נרצה לבדוק רק בקטגוריה של הגברים האם קיים הבדל בין התרופות נוגדות הדיכאון. אם נמצא הבדלים כאלה נאמר שיש

אפקט פשוט של סוג התרופה בקרב אוכלוסיית הגברים. כמות האפקטים הפשוטים שאפשר למצוא היא סכום מספר הקטגוריות (רמות) של כל גורם. למשל, אם יש שלושה סוגי תרופות ושתי אפשרויות למגדר, בסך הכול נוכל לבדוק 5 אפקטים פשוטים.

**דוגמה**

נבדקו שלושה סוגי דיאטות על אנשים בעלי משקל עודף. כעבור שלושה חודשים בדקו כמה קילוגרמים הפחית כל מטופל ממשקלו באותה התקופה. נניח שאנו יודעים את תוחלת הפחתת המשקל של כל דיאטה בחלוקה למגדרים.

נתאר כמה מצבים אפשריים לגבי האוכלוסייה הנחקרת וננתח כל מצב מבחינת ההשפעה של כל גורם על תוחלת המשתנה התלוי ומבחינת אפקט האינטראקציה.

שימו לב שהמצבים שנתאר להלן מתייחסים לתוחלות האמיתיות. בניתוח שונות אין לנו נתוני אמת, אלא רק נתוני מדגם, ונרצה לבדוק האם האפקטים שהתקבלו במדגם הם מובהקים, כנדרש בכל תהליך של הסקה סטטיסטית.

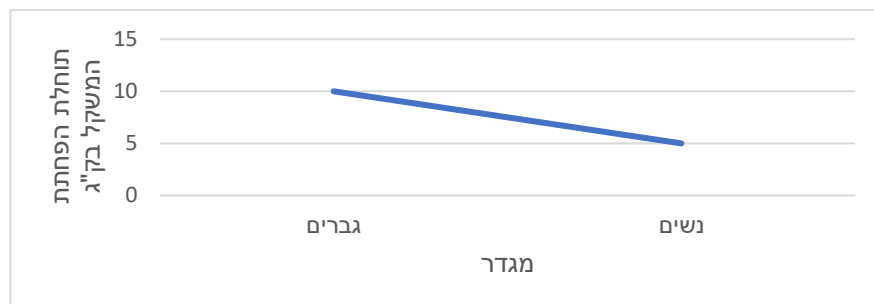
אם התוצאות שלנו יהיו ממוצעי מדגם ולא תוחלות, נוכל לבדוק אם קיימים אפקטים במדגם, אך אין זה אומר שקיימים אפקטים באוכלוסייה, כלומר לא נוכל לדעת אם האפקטים במדגם הם מובהקים. כדי לבדוק אם האפקטים הם מובהקים נצטרך לעשות את מבחן ניתוח השונות.

**מצב א:**

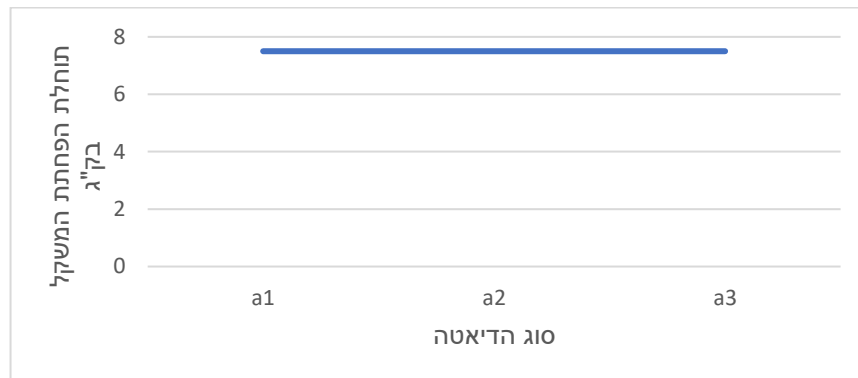
הטבלה הבאה מתארת את תוחלת הפחתת המשקל בק"ג לכל קבוצה:

נשים	גברים	
5	10	$a_1$
5	10	$a_2$
5	10	$a_3$

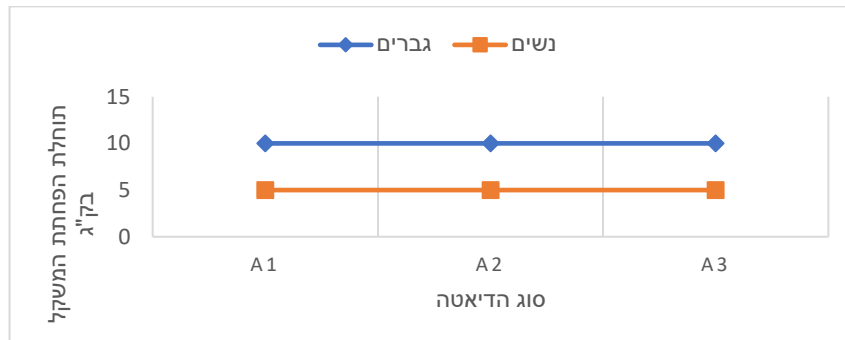
**תיאור גרפי לבדיקת אפקט למגדר**



### תיאור גרפי לבדיקת אפקט לסוג הדיאטה



### גרף אפקטים פשוטים



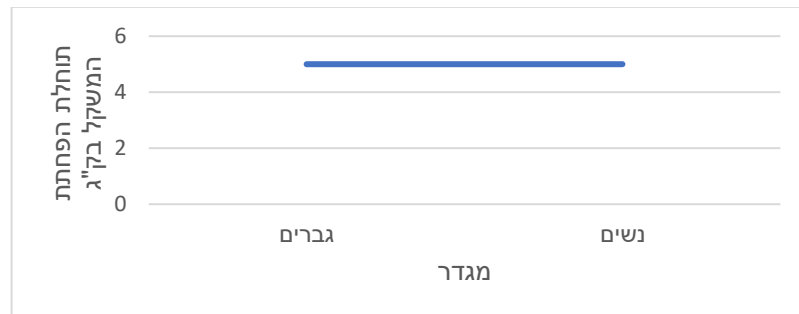
ניתוח המצב: למגדר יש אפקט, לסוג הדיאטה אין אפקט, אין אפקט אינטראקציה. הערה: אם הקווים הנוצרים בגרף האפקטים הפשוטים מקבילים או מתלכדים, אנו אומרים שאין אפקט אינטראקציה.

**מצב ב**

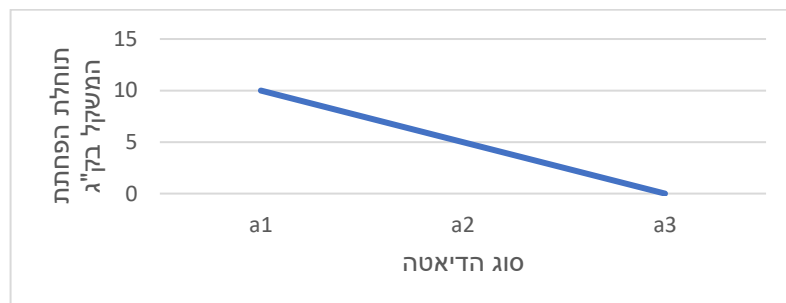
הטבלה הבאה מתארת את תוחלת הפחתת המשקל בק"ג לכל קבוצה:

נשים	גברים	
10	10	$a_1$
5	5	$a_2$
0	0	$a_3$

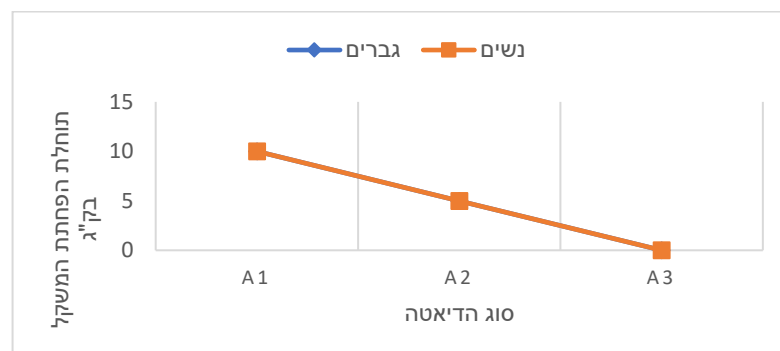
תיאור גרפי לבדיקת אפקט למגדר



תיאור גרפי לבדיקת אפקט לסוג הדיאטה



גרף אפקטים פשוטים



ניתוח המצב: למגדר אין אפקט, לסוג הדיאטה יש אפקט, אין אפקט אינטראקציה.

**מצב ג**

הטבלה הבאה מתארת את תוחלת הפחתת המשקל בק"ג לכל קבוצה:

נשים	גברים	
0	10	$a_1$
5	5	$a_2$
10	0	$a_3$

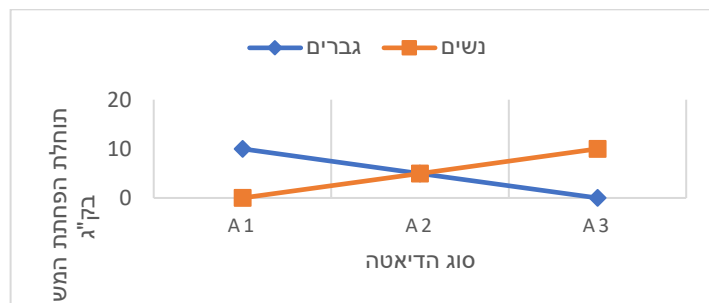
**תיאור גרפי לבדיקת אפקט למגדר**



**תיאור גרפי לבדיקת אפקט לסוג הדיאטה**



**גרף אפקטים פשוטים**



ניתוח המצב: למגדר אין אפקט, לסוג הדיאטה אין אפקט, יש אפקט אינטראקציה.

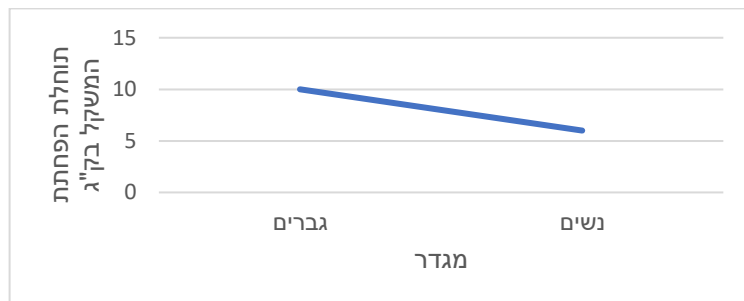
אינטראקציה דיסאורדינלית (נקראת גם "אינטראקציה מהותית"): אפשר לזהות מצב של אינטראקציה כזו באמצעות גרף של אפקטים פשוטים, כאשר נוצרים קווים נחתכים שאחד מהם עולה והאחר יורד. המשמעות היא שגורם אחד משפיע על המשתנה התלוי ברמה מסוימת של הגורם השני באופן הפוך משהוא משפיע על המשתנה התלוי ברמה אחרת של הגורם השני. במצב זה אין להתייחס לאפקטים עיקריים. יש להתייחס רק לאפקטים הפשוטים.

**מצבה**

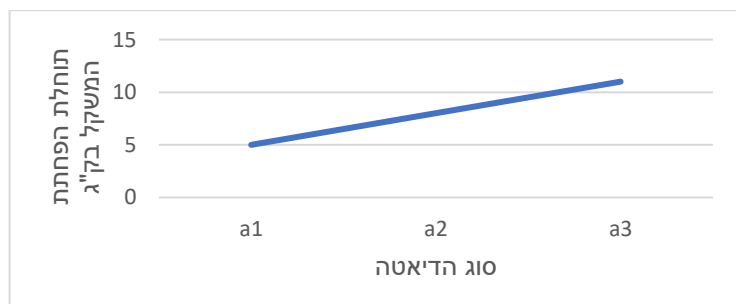
הטבלה הבאה מתארת את תוחלת הפחתת המשקל בק"ג לכל קבוצה:

נשים	גברים	
5	5	$a_1$
6	10	$a_2$
7	15	$a_3$

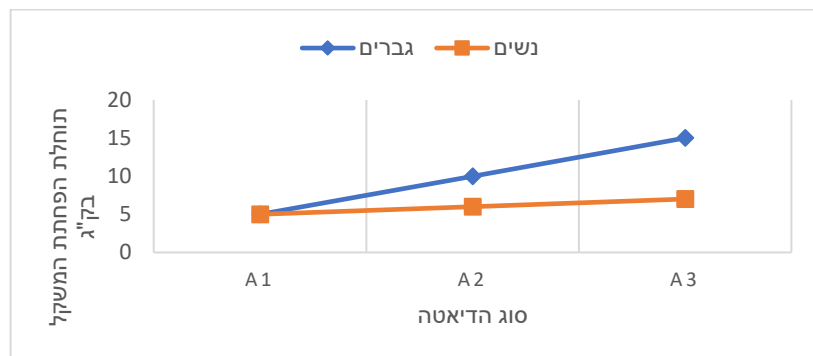
**תיאור גרפי לבדיקת אפקט למגדר**



**תיאור גרפי לבדיקת אפקט לסוג הדיאטה**



גרף אפקטים פשוטים



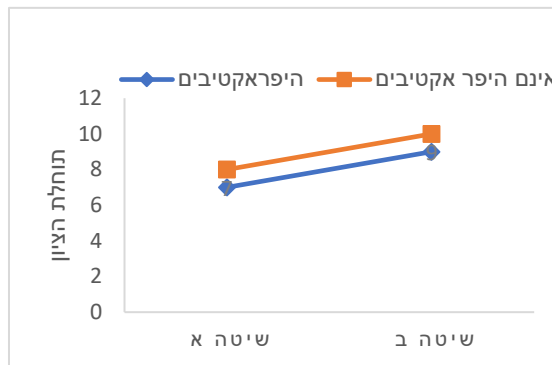
ניתוח המצב: למגדר יש אפקט, לסוג הדיאטה יש אפקט, יש אפקט אינטראקציה.

אינטראקציה אורדינלית (נקראת גם "אינטראקציה לא מהותית"): אפשר לזהות מצב של אינטראקציה כזו כאשר בגרף האפקטים הפשוטים נוצרים קווים נחתכים עם אותו הכיוון (כולם עולים או כולם יורדים אבל לא באותו השיפוע). המשמעות היא שבמעבר של גורם אחד מרמה אחת לרמה אחרת שלו הוא משפיע על המשתנה התלוי באותו אופן בכל רמה של המשתנה האחר אבל עם גודל אפקט שונה.

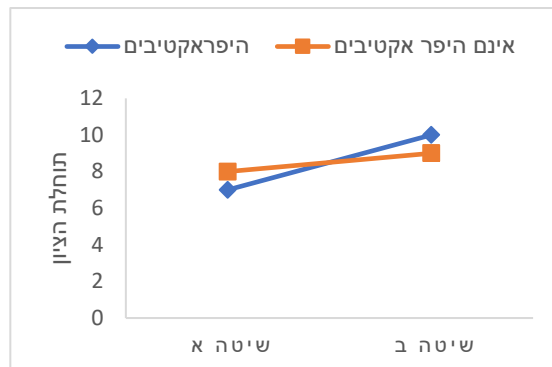
**שאלות**

1) בגני החובה יש שתי שיטות הוראה. שיטות אלו נוסו על ילדים היפראקטיביים וילדים שאינם היפראקטיביים. בתרשימים הבאים מיוצגים גרפים שמתארים את תוחלת הציון במבחן אוצר המילים שניתן לילדים בסוף השנה. בכל אחד מהמקרים יש לקבוע האם קיימת אינטראקציה בין שני הגורמים. אם קיימת אינטראקציה, יש לקבוע האם היא אינטראקציה אורדינלית או דיסאורדינלית.

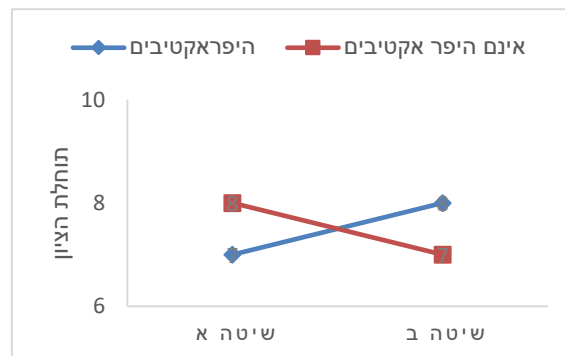
א.

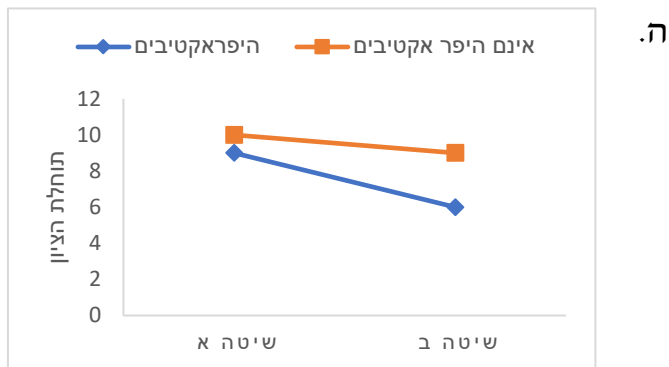
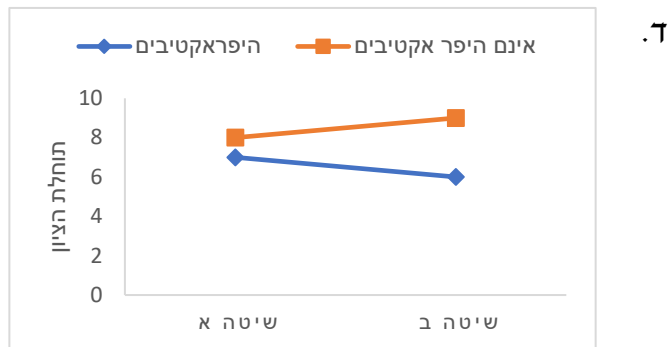


ב.



ג.





2) משרד האוצר פרסם נתונים על המחיר הממוצע של דירות גן ודירות גג של 4 חדרים ב-3 ערים בארץ. מחיר הדירות נמדד במיליוני שקלים. להלן התוצאות שהתקבלו:

דירות גג	דירות גן	
3	4	הרצליה
1	2	אשדוד
2	3	חולון

- א. מהו המשתנה התלוי ומה הם המשתנים הבלתי תלויים?
- ב. האם קיים אפקט לעיר? היעזרו בגרף מתאים.
- ג. האם קיים אפקט לסוג הדירה? היעזרו בגרף מתאים.
- ד. האם קיימת אינטראקציה בין הגורמים? אם כן, מהו סוג האינטראקציה? היעזרו בגרף מתאים.
- ה. האם יש אפקט פשוט לעיר עבור דירות גן?
- ו. האם יש אפקט פשוט לעיר עבור דירות גג?
- ז. האם יש אפקט פשוט לסוג הדירה בהרצליה?
- ח. האם יש אפקט פשוט לסוג הדירה באשדוד?
- ט. האם יש אפקט פשוט לסוג הדירה בחולון?

3) משרד החינוך פרסם נתונים על תוחלת הציונים בבחינת הבגרות באנגלית לפי עיר וסוג בית הספר (עיוני או מקצועי). להלן התוצאות שהתקבלו:

מקצועי	עיוני	
70	85	רעננה
75	75	תל אביב
85	70	פתח תקווה

- א. תארו את הנתונים באמצעות גרף אפקטים פשוטים.  
 ב. האם קיימת אינטראקציה בין הגורמים? אם כן, מה סוג האינטראקציה?  
 ג. באילו ערים קיים אפקט פשוט לסוג בית הספר?

4) משרד התחבורה פרסם נתונים על תוחלת מספר עבירות התנועה לבעלי רישיון נהיגה לפי עיר ולפי מגדר. להלן התוצאות שהתקבלו:

אישה	גבר	
1	2	חיפה
1	2	אשקלון
1	2	רמת גן

- א. האם קיים אפקט עיקרי לעיר?  
 ב. האם קיים אפקט עיקרי למגדר?  
 ג. האם יש אפקט פשוט לעיר אצל הגברים?  
 ד. האם קיימת אינטראקציה בין הגורמים? אם כן, מהו סוג האינטראקציה?

5) המשרד לאיכות הסביבה פרסם נתונים על תוחלת רמת זיהום האוויר בערים שונות בארץ בחורף ובקיץ. להלן התוצאות שהתקבלו:

חורף	קיץ	
20	20	חיפה
10	10	ירושלים
15	15	באר שבע

- א. האם קיים אפקט עיקרי לעיר?  
 ב. האם קיים אפקט עיקרי לעונה?  
 ג. האם קיימת עיר שבה יש אפקט פשוט לעונה?  
 ד. האם קיימת אינטראקציה בין הגורמים? אם כן, מה סוג האינטראקציה?

6) המשרד לאיכות הסביבה פרסם נתונים על תוחלת רמת זיהום האוויר בערים שונות בארץ בחורף ובקיץ. להלן התוצאות שהתקבלו:

חורף	קיץ	
10	10	רמת גן
10	10	גבעתיים
10	10	בת ים

האם קיים אפקט עיקרי לגורם כלשהו? האם קיימת אינטראקציה?

**בשאלות הבאות יש לבחור את התשובה הנכונה ביותר:**

7) במחקר נדגמו 5 אנשים מכל אחת מ-4 הקבוצות הבאות: 1. מתעמלים באופן קבוע ושומרים על תזונה בריאה; 2. מתעמלים באופן קבוע ולא שומרים על תזונה בריאה; 3. לא מתעמלים באופן קבוע ושומרים על תזונה בריאה; 4. לא מתעמלים באופן קבוע ולא שומרים על תזונה בריאה. להלן טבלה המסכמת את ממוצע הטריגליצרידים בדם (מ"ג לדציליטר) שנמצא בכל מדגם:

לא תזונה בריאה	תזונה בריאה	
100	90	מתעמלים
160	100	לא מתעמלים

- קיים אפקט עיקרי מובהק לגורם ההתעמלות.
- קיים אפקט עיקרי מובהק לגורם התזונה.
- קיים אפקט אינטראקציה מובהק בין שני הגורמים במחקר.
- אי אפשר לדעת אם קיים אפקט מובהק כלשהו על סמך תוצאות המדגם בלבד ללא ביצוע מבחן מתאים וללא קביעת רמת המובהקות של המחקר.

8) במחקר בדקו 3 טיפולים שונים לחולי פסוריאזיס. המחקר השווה גם בין גברים לנשים ובדק את זמן התגובה לטיפול. מסקנת המחקר הייתה שאצל גברים נמצאו הבדלים מובהקים בין הטיפולים השונים מבחינת תוחלת זמן התגובה. לאיזה סוג אפקט המסקנה מתייחסת?

- אפקט אינטראקציה.
- אפקט עיקרי של גורם המין.
- אפקט עיקרי של גורם סוג הטיפול.
- אפקט פשוט.

- 9) במחקר בדקו 3 טיפולים שונים לחולי פסוריאזיס. המחקר השווה גם בין גברים לנשים ובדק את זמן התגובה לטיפול. במדגם היה ממוצע זמן התגובה של הגברים שונה מממוצע זמן התגובה של הנשים.
- א. אפשר להגיד שבמדגם קיים אפקט עיקרי, אך אי אפשר לדעת אם האפקט העיקרי מובהק.
- ב. אפשר להגיד שבמדגם קיימת אינטראקציה, אך אי אפשר לדעת אם האינטראקציה מובהקת.
- ג. אפשר להגיד שקיים אפקט עיקרי מובהק.
- ד. אפשר להגיד שקיימת אינטראקציה מובהקת.
- 10) במחקר בדקו 3 טיפולים שונים לחולי פסוריאזיס. המחקר השווה גם בין גברים לנשים ובדק את זמן התגובה לטיפול. אחת המסקנות של המחקר הייתה שהטיפולים השונים משפיעים במידה משמעותית יותר על זמן התגובה של הגברים מאשר על זה של הנשים, אם כי באותו האופן.
- א. המסקנה היא שאין אינטראקציה בין הגורמים במחקר.
- ב. המסקנה היא שיש אינטראקציה אורדינלית בין הגורמים במחקר.
- ג. המסקנה היא שיש אינטראקציה דיסאורדינלית בין הגורמים במחקר.
- ד. המסקנה היא שיש אפקט עיקרי של המגדר.

### תשובות סופיות

- 1) א. אין אינטראקציה.  
 ג. אינטראקציה דיסאורדנלית.  
 ה. אינטראקציה אורדינלית.
- 2) א. המשתנים הבי"ת: העיר, סוג הדירה. המשתנה התלוי: מחיר.  
 ב. קיים.  
 ד. לא קיים.  
 ו. קיים.  
 ח. קיים.
- 3) א. עיין בסרטון הוידאו.  
 ג. רעננה ופתח תקווה.
- 4) א. לא.  
 ג. לא.
- 5) א. כן.  
 ג. לא.
- 6) לא, לא.
- 7) ד
- 8) ד
- 9) א
- 10) ב

## תהליך ניתוח שונות דו כיווני – הליך מבחן

### רקע

כפי שכבר ציינו, ניתוח שונות דו-כיווני נעשה כאשר יש שני גורמים מחקרניים ומשתנה כמותי תלוי אחד. מטרת המחקר היא לבדוק האם הגורמים משפיעים על המשתנה התלוי. מערך מחקר זה נקרא "מערך מחקר פקטוריאלי", כיוון שאנו בונים את המחקר לפי גורמים. מערך דו-גורמי יסומן כמעריך מסוג  $A \times B$ , כאשר  $A$  מייצג את מספר הרמות של גורם  $a$ , ו- $B$  מייצג את מספר הרמות של גורם  $b$ . במערך מחקרי תלת-גורמי נסמן את סוג המערך  $A \times B \times C$ , וכך הלאה.

### דוגמה

נבדקו שלושה סוגי דיאטות על אנשים בעלי משקל עודף. נבחרו 18 מטופלים בעלי משקל עודף, 9 מהם גברים ו-9 נשים. המטופלים חולקו כך שבכל דיאטה השתתפו 3 גברים ו-3 נשים. כעבור שלושה חודשים מתחילת הדיאטה נשקלו כלל המטופלים ונבדק המשקל בק"ג שהם הפחיתו. הטבלה הבאה מסכמת את המשקל שכל מטופל במדגם הפחית כעבור שלושה חודשים.

סוג הדיאטה \ מין	$b_1$	$b_2$	$b_3$	סה"כ
נשים	8	6	4	54
	4	8	6	
	0	10	8	
גברים	6	0	9	72
	10	2	12	
	14	4	15	
סה"כ	42	30	54	126

מטרת המחקר היא לבדוק האם יש השפעה של סוג הדיאטה, המין והשילוב ביניהם על ההפחתה במשקל.

- באיזה סוג מערך מחקרי מדובר?
- מהו המבחן הסטטיסטי המתאים לבדיקת ההשערות?
- מה הן השערות המחקר?

בדומה לניתוח שונות חד-כיווני גם התהליך של ניתוח שונות דו-כיווני דורש הנחות. ההנחות הן:

1.  $A \times B$  הקבוצות שנוצרות בלתי תלויות זו בזו.
2. בכל  $A \times B$  האוכלוסיות המשתנה התלוי מתפלג נורמלית.
3. בכל  $A \times B$  האוכלוסיות אותה שונות,  $\sigma^2$ .

הערה: ניתוח שונות הוא מבחן רובסטי, כלומר יש לו רגישות נמוכה להנחות. התיאוריה הסטטיסטית שפותחה התבססה על ההנחות האלה, אבל הלכה למעשה השיטה תעבוד טוב גם אם ההנחות הללו לא יתקיימו במדויק במלואן. זו הסיבה שהשיטה הזו נפוצה כל כך בעולם הסטטיסטיקה.

בהמשך לדוגמה

רשמו את כל ההנחות הדרושות לביצוע ניתוח השונות.

**הליך המבחן**

בניית טבלת ממוצעים

נבנה טבלת ממוצעים לכל רמה ולכל תא :

$\bar{X}_i$  – ממוצע המדגם ברמה  $i$  של גורם  $a$

$\bar{X}_j$  – ממוצע המדגם ברמה  $j$  של גורם  $b$

$\bar{X}_{ij}$  – ממוצע המדגם בתא  $ij$

**בהמשך לדוגמה**

- מלאו את טבלת הממוצעים הבאה :

סוג הדיאטה \ מין	$b_1$	$b_2$	$b_3$	$\bar{X}_i$
נשים				
גברים				
$\bar{X}_j$				

- שרטטו גרפים מתאימים לבדיקת אפקטים עיקריים ולבדיקת אינטראקציה במדגם. האם אפשר להגיד שיש אפקט מובהק?

**בניית טבלת ריבועי הפרשים מהממוצעים**

נמלא את הטבלה הבאה. בתוך תא  $ij$  נחשב:  $(\bar{X}_{ij} - \bar{X}_i - \bar{X}_j + \bar{X})^2$

**בהמשך לדוגמה**

- מלאו את טבלת הפרשי הממוצעים:

<div style="text-align: center;">סוג הדיאטה</div>	$b_1$	$b_2$	$b_3$	$(\bar{X}_i - \bar{X})^2$
מין				
נשים				
גברים				
$(\bar{X}_j - \bar{X})^2$				

**חישוב סכום ריבועי הסטיות מהממוצע**

מתוך טבלת ריבועי הסטיות מהממוצע נחשב את סכום ריבועי הסטיות מהממוצע הבאים :

הסימון  $SS$  הוא ראשי התיבות של "sum of squares" (סכום הריבועים).

סכום ריבועי הסטיות מהממוצע של גורם  $a$  :

$$SS_a = m \cdot B \sum_{i=1}^A (\bar{X}_{.i} - \bar{X})^2$$

סכום ריבועי הסטיות מהממוצע של גורם  $b$  :

$$SS_b = m \cdot A \sum_{j=1}^B (\bar{X}_{.j} - \bar{X})^2$$

סכום ריבועי הסטיות של האינטראקציה :

$$SS_{ab} = m \sum_{i=1}^A \sum_{j=1}^B (\bar{X}_{ij} - \bar{X}_{.i} - \bar{X}_{.j} + \bar{X})^2$$

סכום ריבועי השגיאות (סכום ריבועי הסטיות של התצפיות בתא מהממוצע בתא) :

$$SS_W = \sum_{i=1}^A \sum_{j=1}^B \sum_{k=1}^m (X_{ijk} - \bar{X}_{ij})^2 = (m-1) \sum_{i=1}^A \sum_{j=1}^B S_{ij}^2$$

סכום ריבועי הסטיות של כלל התצפיות מהממוצע הכללי :

$$SS_T = \sum_{i=1}^A \sum_{j=1}^B \sum_{k=1}^m (X_{ijk} - \bar{X})^2 = (n-1) \cdot S^2$$

הקשר המתמטי בין סכום הריבועים הללו הוא :

$$SS_T = SS_a + SS_b + SS_{ab} + SS_W$$

לכן אין אנו צריכים לחשב את כל חמשת המרכיבים הללו.

החלק הזה של הנוסחה מתייחס לשונות השיטתית :  $SS_a + SS_b + SS_{ab}$ . השונות השיטתית היא שונות שמקורה בגורמים עצמים.

החלק הזה של הנוסחה מתייחס לשונות המקרית :  $SS_W$ . השונות המקרית היא שונות שנקראת גם "שונות טעויות" או "שונות בתוך הקבוצות". זוהי שונות בין התצפיות שאינה נובעת מהגורמים הנחקרים. האות  $W$  מייצגת את המילה "Within", כלומר שונות בתוך התאים.



בהמשך לדוגמה

- חשבו את ריבועי הסטיות הבאים :

$$SS_a =$$

$$SS_b =$$

$$SS_{ab} =$$

$$SS_T =$$

$$SS_w =$$

**חישוב ממוצע ריבועי הסטיות וסטטיסטי המבחן**

MS הוא הסימון של ממוצע ריבועי הסטיות (Mean Square) שמהווה אומד לשונות של כל גורם. החישוב ייעשה על ידי חלוקת ה-SS המתאים בדרגות החופש המתאימות. לאחר מכן נחשב שלושה סטטיסטי מבחן, בהתאם לשלוש ההשערות הנבדקות.

נרכז את כלל החישובים הללו בטבלה הנקראת טבלת ניתוח שונות, ANOVA (Analysis of Variance).

מקור השונות Source of Variation	דרגות החופש Degrees of Freedom	סכום ריבועי הסטיות מהממוצע Sum of Squares	ממוצע ריבוע הסטייה Mean Square	F
a	A - 1	SS <sub>a</sub>	MS <sub>a</sub>	$F_a = MS_a / MS_w$
b	B - 1	SS <sub>b</sub>	MS <sub>b</sub>	$F_b = MS_b / MS_w$
ab	(A - 1)(B - 1)	SS <sub>ab</sub>	MS <sub>ab</sub>	$F_{ab} = MS_{ab} / MS_w$
Within	AB(m - 1)	SS <sub>w</sub>	MS <sub>w</sub>	
Total	n-1=ABm - 1	SS <sub>T</sub>		

בהמשך לדוגמה : מלאו את טבלת ניתוח השונות

מקור השונות Source of Variation	דרגות החופש Degrees of Freedom	סכום ריבועי הסטיות מהממוצע Sum of Squares	ממוצע ריבוע הסטייה Mean Square	F
a				
b				
ab				
Within				
Total				

**כללי ההכרעה לבדיקת ההשערות**

הסטטיסטי  $F_a$  מייצג את היחס בין השונות המדגמית של גורם  $a$  ובין השונות המקרית. לכן ככל שהערכים שלו גבוהים יותר, נרצה להגיד שלגורם  $a$  יש השפעה גדולה יותר על המשתנה התלוי.  $F_a$  יקבל ערכים גבוהים אם השונות המדגמית של גורם A תגדל או אם השונות המדגמית המקרית תקטן. הסטטיסטי מתפלג התפלגות F, ואזור הדחייה שלו יהיה בצד ימין.

- כלל ההכרעה לבדיקת המובהקות של גורם  $a$  :

דחה את השערת  $H_0$  ברמת מובהקות של  $\alpha$  אם

$$F_a > F_{1-\alpha}(df_a, df_w)$$

לפי אותו עיקרון שאר כללי ההכרעה יהיו :

- כלל ההכרעה לבדיקת המובהקות של גורם  $b$  :

דחה את השערת  $H_0$  ברמת מובהקות של  $\alpha$  אם

$$F_b > F_{1-\alpha}(df_b, df_w)$$

- כלל ההכרעה לבדיקת המובהקות של האינטראקציה :

דחה את השערת  $H_0$  ברמת מובהקות של  $\alpha$  אם

$$F_{ab} > F_{1-\alpha}(df_{ab}, df_w)$$

### בהמשך לדוגמה

רשמו את כל כללי ההכרעה המתאימים והסיקו מסקנות מתאימות ברמת מובהקות של 5%.

### הערות

1. אם מכריעים שקיימת אינטראקציה מובהקת, יש לבדוק האם היא אורדינלית או דיסאורדינלית. אם האינטראקציה דיסאורדינלית, יש לבדוק האם האפקטים העיקריים נמצאו מובהקים. אם לפחות אחד מהם נמצא מובהק נאמר שהוא אינו משמעותי כיוון שהוא נובע מהאינטראקציה בין הגורמים ולא מהגורם עצמו.
2. אם אחד מהאפקטים נמצא מובהק, אין זה אומר אילו רמות שונות זו מזו בתוחלת. למשל, אם נמצא הבדל מובהק בין סוגי הטיפולים, לא נוכל לדעת לפי זה איזה טיפול שונה מאחר באופן מובהק. לכן יש להמשיך בתהליך של השוואות מרובות כדי להסיק ממה נובע השוני.

### בהמשך לדוגמה

האם יש סיבה לבצע השוואות מרובות במחקר?

**שאלות**

1) מחקר שיווקי בדק את השפעת גובה המדף בסופרמרקט והשפעת החומר שממנו עשוי הבקבוק (זכוכית או פלסטיק) על היקף המכירות של משקאות קלים. נבדקו שני סופרמרקטים. בכל סופרמרקט נבחן כל צירוף אפשרי של גובה המדף וחומר הבקבוק, ועבור כל צירוף כזה נבדק מספר בקבוקי המשקה הקל שנמכרו באותו סופרמרקט ביום מסוים. הנה התוצאות שהתקבלו:

פלסטיק	זכוכית	סוג בקבוק
		גובה המדף
59	23	נמוך
63	32	
88	47	בינוני
90	55	
51	40	גבוה
56	48	

בצעו ניתוח שונות דו-כיווני על נתוני מחקר זה ברמת מובהקות של 5%. סכמו את המסקנות מתוך ניתוח השונות שביצעתם. מה הן ההנחות הדרושות לביצוע המבחן?

2) במחקר בתחום החקלאות נדגמו 8 חלקות אדמה : 4 חלקות בנגב ו-4 בעמק יזרעאל. בכל חלקה ההשקיה הייתה או באמצעות ממטרות או באמצעות טפטפות. בדקו את יבול העגבניות (בטונה לדונם) בכל חלקה. להלן התוצאות שהתקבלו :

מספר חלקה	מיקום החלקה	שיטת השקיה	יבול העגבניות
1	נגב	ממטרות	12
2	נגב	ממטרות	10
3	נגב	טפטפות	15
4	נגב	טפטפות	17
5	עמק יזרעאל	ממטרות	12
6	עמק יזרעאל	ממטרות	14
7	עמק יזרעאל	טפטפות	17
8	עמק יזרעאל	טפטפות	19

- א. רשמו את כלל המשתנים במחקר וציינו לגבי כל אחד מהם האם הוא משתנה תלוי או בלתי תלוי.
- ב. הציגו את נתוני המחקר באמצעות גרפים מתאימים. האם נראה שבמדגם יש אפקט עיקרי לכל גורם? האם יש אינטראקציה בין הגורמים במדגם? האם האפקטים מובהקים?
- ג. בדקו ברמת מובהקות של 5% האם האפקט העיקרי של כל גורם הוא מובהק והאם האינטראקציה היא מובהקת. מה הן ההנחות הדרושות?

3) חברה לייצור מוצרי שיער פיתחה נוסחה חדשנית לצבע לשיער שאינו דורש תוספת חמצן בעת תהליך הצביעה. החברה השוותה את צבע השיער החדש לצבע השיער הרגיל מבחינת כושר הכיסוי וזאת על שלושה סוגי שיער: בהיר, כהה ושיבה. ציון רמת הכיסוי הוא משתנה שמתפלג נורמלית עם שונות קבועה לכל סוג שיער ולכל סוג צבע. לכל קבוצה של סוג צבע וסוג שיער נדגמו 4 צביעות שנוסו על אנשים שונים, וניתן ציון מספרי על רמת הכיסוי. להלן סיכום תוצאות המדגם שהתקבלו:

שונות	ממוצע	הקבוצה
40	62	צבע רגיל על שיער בהיר
44	51	צבע רגיל על שיער כהה
42	45	צבע רגיל על שיער שיבה
46	60	צבע חדש על שיער בהיר
40	54	צבע חדש על שיער כהה
42	44	צבע חדש על שיער שיבה

בצעו ניתוח שונות דו-כיווני על הנתונים ברמת מובהקות של 5%. סכמו את כל המסקנות המתקבלות.

4) בוצע ניתוח שונות על נתונים. במערך המחקרי לגורם  $a$  יש 4 רמות ולגורם  $b$  יש 3 רמות. נערכו 3 תצפיות לכל אחת מ-12 הקבוצות שנוצרו. להלן טבלת ניתוח שונות דו-גורמי שבוצע:

מקור השונות	$df$	SS	MS	F
$a$	?	318	?	?
$b$	?	?	?	?
אינטראקציה	?	190	?	?
W	?	156	?	
T	?	674		

א. מלאו את כל התאים בטבלה המסומנים בסימני שאלה.

ב. בצעו את הבדיקות הבאות ברמת מובהקות של 5%:

- i. האם האינטראקציה מובהקת?
- ii. האם גורם  $a$  משפיע על המשתנה התלוי הנחקר?
- iii. האם לגורם  $b$  יש לפחות שתי רמות עם תוחלות שונות?

5) במחקר בדקו האם ארץ מוצא ומגדר של אדם משפיעים על שנות ההשכלה שלו. הנתונים סוכמו בטבלת ניתוח שונות:

מקור השונות	df	SS	MS	F
ארץ מוצא	4	34		
מגדר			2	
אינטראקציה		18	4.5	
W	10	12		
T				

- כמה ארצות מוצא נבדקו במחקר זה?
- מהו גודל המדגם הכולל במחקר זה?
- חשבו את ערכי F הסטטיסטי עבור ארץ המוצא, המגדר והאינטראקציה.
- מה הם האפקטים המובהקים במחקר זה ברמת מובהקות של 5%?

6) בטבלה הבאה מסוכמים הממוצעים של מערך מחקרי דו-גורמי עם משתנה כמותי תלוי:

	$b_1$	$b_2$	$b_3$
$a_1$	8	14	11
$a_2$	6	13	16

מספר התצפיות בכל תא הוא 5.  
הטבלה הבאה היא טבלה מסכמת של ניתוח השונות על סמך נתוני מחקר זה:

מקור השונות	df	SS	MS	F
$a$				
$b$		281.7		
$ab$		71.7		
W		190.1		
T				

- א. מלאו את טבלת ניתוח השונות.
- ב. הסיקו מסקנות ברמת מובהקות של 5%.
- ג. שרטטו גרף אינטראקציות והסבירו את משמעות הממצאים.

### תשובות סופיות

- 1) עיין בסרטון הוידאו.
- 2) א. משתנים ב"ת: מיקום החלקה, שיטת השקיה. משתנה תלוי: יבול בטונה לדונם.  
ב. עיין בסרטון הוידאו.  
ג. עיין בסרטון הוידאו.
- 3) עיין בסרטון הוידאו.
- 4) א. עיין בסרטון הוידאו. ב. i. כן. ii. כן. iii. לא.
- 5) א. 4. ב. 20. ג. עיין בסרטון הוידאו.
- 6) א. עיין בסרטון הוידאו. ב. עיין בסרטון הוידאו. ג. עיין בסרטון הוידאו.

# רגרסיה ושיטות ניתוח ליניאריות

פרק 3 - רגרסיה ליניארית חד משתנית

תוכן העניינים

1. כללי ..... (ללא ספר)

# רגרסיה ושיטות ניתוח ליניאריות

פרק 4 - משתנה דמי

תוכן העניינים

46 ..... 1. כללי

## משתנה דמי:

### רקע:

הכנסת משתנים ב"ת איכותיים למודל הרגרסיה.

למשל, נתונה משוואת הרגרסיה:  $W_t = \alpha + \beta \cdot S_t$ .

$W_t$  = השכר (התלוי).

$S_t$  = שנות לימוד (הבי"ת) שניהם כמותיים.

נניח שאנו סבורים שגם משתנה המגדר (משתנה איכותי) משפיע על השכר.

כדי להכניסו למשוואת הרגרסיה יש להגדיר משתני דמי (dummy variable):

נגדיר משתנה  $D$  שיקבל את הערך 0 אם מדובר ב"אישה" ואת הערך 1 אם מדובר ב"גבר". ניתן להכניס את משתנה הדמי למודל בשלושה אופנים שונים:

1. משתנה דמי לחותך – המגדר משפיע על השכר ההתחלתי בלבד.
2. משתנה דמי לשיפוע – המגדר משפיע על התוספת לשכר בגין שנות הלימוד.
3. משתנה דמי לכל הפונקציה – המגדר משפיע גם על החותך וגם על השיפוע.

### משתנה דמי לחותך:

המין משפיע על השכר ההתחלתי בלבד.

המודל:  $W_t = \alpha_0 + \alpha_1 D + \beta \cdot S_t + u_t$  החותך מייצג כאן את השכר ההתחלתי.

שכר ההתחלתי של אישה:  $\alpha_0$ .

שכר התחלתי של גבר:  $\alpha_0 + \alpha_1$ .

הבדל בשכר בין נשים וגברים:  $\alpha_1$  (הפרש בין החותכים).

בדיקת השערות על משתנה הדמי: מבחן  $t$  למובהקות הפרש החותכים:  $H_0: \alpha_1 = 0$ .

- השיפוע מייצג את התוספת בשכר כפונקציה של מס' שנות הלימוד והוא זהה עבור נשים וגברים.

### פונקציית רגרסיה המכילה משתנים איכותיים בלבד:

המגדר הוא המשתנה היחיד במשוואה:  $W_t = \alpha_0 + \alpha_1 D + u_t$ .

החותך מייצג כאן את השכר הממוצע עבור כל קטגוריה:

שכר הממוצע של אישה:  $\alpha_0$ .

שכר הממוצע של גבר:  $\alpha_0 + \alpha_1$ .

הבדל בשכר הממוצע בין נשים וגברים:  $\alpha_1$  (הפרש בין החותכים).

בדיקת השערות על משתנה הדמי: מבחן  $t$ :  $H_0: \alpha_1 = 0$  (מבחן זהה למבחן  $t$  להבדל בין ממוצעים).

**משתנה דמי לשיפוע:**

- המגדר משפיע על התוספת לשכר בגין שנות הלימוד:  $W_t = \alpha + \beta_0 S_t + \beta_1 DS_t + u_t$ .  
 השיפוע מייצג כאן את התוספת לשכר בגין שנות לימוד.  
 אצל אישה: התוספת לשכר בגין שנות לימוד:  $\beta_0$ .  
 אצל גבר: התוספת לשכר בגין שנות לימוד:  $\beta_0 + \beta_1$ .  
 הבדל בין גברים לנשים בתוספת לשכר בגין שנות הלימוד:  $\beta_1$  (הפרש השיפועים).  
 בדיקת השערות על משתנה הדמי: מבחן  $t$  למובהקות הפרש השיפועים:  $H_0: \beta_1 = 0$ .
- החותך, המייצג את השכר ההתחלתי, יהיה זהה עבור גברים ונשים.

**משתנה דמי לכל הפונקציה:**

- המין משפיע גם על החותך וגם על השיפוע – גם על השכר ההתחלתי וגם על התוספת לשכר ההתחלתי בגין שנות הלימוד.  
 המודל:  $W_t = \alpha_0 + \alpha_1 D + \beta_0 S_t + \beta_1 DS_t + u_t$ .  
 השכר ההתחלתי של אישה:  $\alpha_0$ .  
 השכר ההתחלתי של גבר:  $\alpha_0 + \alpha_1$ .  
 הבדל בשכר ההתחלתי בין המינים:  $\alpha_1$  (הבדל בחותכים).  
 אצל אישה: התוספת לשכר בגין שנות הלימוד:  $\beta_0$ .  
 אצל גבר: התוספת לשכר בגין שנות הלימוד:  $\beta_0 + \beta_1$ .  
 הבדל בין המינים בתוספת לשכר בגין שנות הלימוד:  $\beta_1$  (הבדל בשיפועים).

**2 דרכים לבדיקה האם יש השפעה למשתנה האיכותי:**

1. בדיקת השערות למשתני הדמי:  
 באמצעות מבחן WALT יש לבדוק:  $H_0: \alpha_1 = \beta_1 = 0$ .  
 לפחות אחד הפרמטרים שונה מ-0:  $H_1$ .  
 אם דוחים את השערת האפס, יש לבצע מבחני  $t$  עבור כל אחד מהפרמטרים  
 בנפרד:  $H_0: \alpha_1 = 0$  ו-  $H_0: \beta_1 = 0$ .
2. מבחן CHOW:  
 דרך נוספת לבדיקת ההבדל בין הקטגוריות בלא יצירת משתני דמי:  
 חלוקת המדגם לפי הקטגוריות של המשתנה האיכותי.  
 מדגם של גברים ( $T_m$ ) ושל נשים ( $T_f$ ).  
 עבור כל קבוצה לאמוד משוואות רגרסיה לניבוי שכר על ידי שנות לימוד:  
 נשים:  $W_t = \alpha_f + \beta_f X_t + u_t$ .  
 גברים:  $W_t = \alpha_m + \beta_m X_t + u_t$ .  
 השערות:  $H_0: \alpha_f = \alpha_m; \beta_f = \beta_m$ .

לבדיקת ההשערה נשתמש במבחן CHOW (הזהה למבחן WALS) :  
 המודל המוגבל (R) לא לוקח בחשבון את השפעת המגדר ולכן יכול את  
 המדגם המאוחד :  $W_t = \alpha + \beta X_t + u_t$

המודל הלא מוגבל (U) כולל את שני חלקי המדגם :  
 $ESS_U = ESS_f + ESS_m$   
 $DF_U = DF_f + DF_m$

$$CHOW_{stat} = \frac{\frac{ESS_R - (ESS_f + ESS_m)}{DF_R - (DF_f + DF_m)}}{\frac{ESS_f + ESS_m}{DF_f + DF_m}} = WALS_{stat}$$

למרות התוצאות הזהות בשתי הדרכים, שיטת משתני הדמי עדיפה :

1. אם דחינו את  $H_0$  במבחן CHOW נתקשה לברר את מקור ההבדל שנמצא.
2. בהרצת שתי רגרסיות נפרדות אנו בודקים הבדל בכל הפונקציה ואילו שיטת משתני הדמי מאפשרת לבדוק הבדל רק בחותך או רק בשיפוע.

### סיכום ביניים :

משתנה דמי לכל הפונקציה	משתנה דמי לשיפוע	משתנה דמי לחותך	
$Y_t = \alpha_0 + \alpha_1 D + \beta_0 X_t + \beta_1 DX_t + u_t$	$Y_t = \alpha + \beta_0 X_t + \beta_1 DX_t + u_t$	$Y_t = \alpha_0 + \alpha_1 D + \beta \cdot X_t + u_t$	המודל
קיים הבדל בין הקטגוריות במשוואת הרגרסיה כולה (בחותך ובשיפוע).	קיים הבדל בין הקטגוריות בתוספת ל-Y בגין X (בשיפוע).	קיים הבדל בין הקטגוריות ב-Y ההתחלתי (בחותך).	ההשערה במילים
מבחן WALS להפרש בין הפונקציות (החותכים והשיפועים) : $H_0 : \alpha_1 = \beta_1 = 0$ **ניתן לבדוק את ההשערה בדבר הבדל בין הפונקציות גם במבחן CHOW. אם דוחים את $H_0$ יש לברר את מקור ההבדל באמצעות מבחני t (אפשרי רק ב- WALS) : $H_0 : \alpha_1 = 0$ $H_0 : \beta_1 = 0$	מבחן t להפרש השיפועים : $H_0 : \beta_1 = 0$	מבחן t להפרש החותכים : $H_0 : \alpha_1 = 0$	בדיקת ההשערה

**משתני דמי אם המשתנה האיכותי יכול לקבל יותר משני ערכים:**

כאשר המשתנה האיכותי כולל יותר משני ערכים/קטגוריות נגדיר מס' משתני דמי כמספר הקטגוריות פחות אחד.

למשל, את המשתנה האיכותי של עונות השנה הכולל 4 ערכים: אביב, קיץ, סתיו, חורף נייצג באמצעות 3 משתני דמי:

$D_1$  יקבל את הערך 1 אם מדובר באביב ו-0 אחרת.

$D_2$  יקבל את הערך 1 אם מדובר בקיץ ו-0 אחרת.

$D_3$  יקבל את הערך 1 אם מדובר בסתיו ו-0 אחרת.

אם מדובר בחורף אז כל משתני הדמי יקבלו את הערך 0 ולכן החורף היא קבוצת הייחוס. נניח שאנו רוצים לבדוק עונתיות במחירי הירקות:

$V_t =$  מדד מחירי הירקות.

$p_t =$  מדד המחירים לצרכן.

**1. משתני דמי לחותך:**

הטענה: יש הבדל בין עונות השנה במחיר ההתחלתי של הירקות.

המודל:  $V_t = \alpha_0 + \alpha_1 D_{1t} + \alpha_2 D_{2t} + \alpha_3 D_{3t} + \beta \cdot P_t + u_t$ .

כל עליה של יחידה אחת במדד המחירים לצרכן תעלה את מחירי הירקות ב- $\beta$ . למחיר זה יתווסף  $\alpha_0$  בחורף,  $\alpha_0 + \alpha_1$  באביב,  $\alpha_0 + \alpha_2$  בקיץ ו- $\alpha_0 + \alpha_3$  בסתיו.

ניתן לראות כי:  $\alpha_0$  - החותך בקטגוריה שהושמטה,  $\alpha_0 + \alpha_1$  - החותך בקטגוריה i.

בדיקת השערות:

$H_0: \alpha_1 = \alpha_2 = \alpha_3 = 0$

השערות:  $H_1: \text{OTHERWISE}$

המבחן הסטטיסטי – מבחן WALD:

(U)  $V_t = \alpha_0 + \alpha_1 D_{1t} + \alpha_2 D_{2t} + \alpha_3 D_{3t} + \beta \cdot P_t + u_t$

(R)  $V_t = \alpha + \beta \cdot P_t + u_t$

- שימו לב שהחותך במשוואה המוגבלת איננו  $\alpha_0$  שכן המשתנה המסביר של עונות השנה ירד.

אם נדחה את  $H_0$  במבחן הסטטיסטי של הסעיף הקודם, יש לבדוק מה מקור ההבדל בין החותכים על ידי מבחני  $t$ :

1. האם יש הבדל במחיר ההתחלתי של הירקות בין האביב לחורף:  
 $H_0: \alpha_1 = 0$

2. האם יש הבדל במחיר ההתחלתי של הירקות בין הקיץ לחורף:  
 $H_0: \alpha_2 = 0$

3. האם יש הבדל במחיר ההתחלתי של הירקות בין הסתיו לחורף:  
 $H_0: \alpha_3 = 0$

2. משתני דמי לשיפוע:

הטענה: יש הבדל בין עונות השנה בתוספת למחיר הירקות בגין המחיר לצרכן.

$$\text{המודל: } V_t = \alpha + \beta_0 P_t + \beta_1 (D_{1i} P_t) + \beta_2 (D_{2i} P_t) + \beta_3 (D_{3i} P_t) + u_t$$

המחיר ההתחלתי של הירקות שווה בין עונות השנה ( $\alpha$ ) אולם כל עליה

של יחידה אחת במדד המחירים לצרכן תעלה את מחירי הירקות

ב:  $\beta_0$  בחורף,  $\beta_0 + \beta_1$  באביב,  $\beta_0 + \beta_2$  בקיץ ו- $\beta_0 + \beta_3$  בסתיו.

ניתן לראות כי-  $\beta_0$ : השיפוע בקטגוריה שהושמטה  $\beta_0 + \beta_i$ :

השיפוע בקטגוריה i.

בדיקת השערות:

$$H_0: \beta_1 = \beta_2 = \beta_3 = 0$$

השערות:

$$H_1: \text{OTHERWISE}$$

המבחן הסטטיסטי – מבחן WALD:

$$(U) \quad V_t = \alpha + \beta_0 P_t + \beta_1 (D_{1i} P_t) + \beta_2 (D_{2i} P_t) + \beta_3 (D_{3i} P_t) + u_t$$

$$(R) \quad V_t = \alpha + \beta \cdot P_t + u_t$$

- שימו לב שהשיפוע במשוואה המוגבלת איננו  $\beta_0$  שכן המשתנה המסביר של עונות השנה ירד.

אם נדחה את  $H_0$  במבחן הסטטיסטי של הסעיף הקודם, יש לבדוק מה מקור ההבדל בין השיפועים על ידי מבחני  $t$ .

3. משתני דמי לכל הפונקציה :

הטענה : יש הבדל בין עונות השנה בפונקציית הרגרסיה לניבוי מחיר הירקות באמצעות המחיר לצרכן. המודל :

$$V_t = \alpha_0 + \alpha_1 D_{1t} + \alpha_2 D_{2t} + \alpha_3 D_{3t} + \beta_0 P_t + \beta_1 (D_{1t} P_t) + \beta_2 (D_{2t} P_t) + \beta_3 (D_{3t} P_t) + u_t$$

בדיקת השערות :

$$H_0 : \alpha_1 = \alpha_2 = \alpha_3 = \beta_1 = \beta_2 = \beta_3 = 0$$

המבחן הסטטיסטי - מבחן WALD :

(U)

$$V_t = \alpha_0 + \alpha_1 D_{1t} + \alpha_2 D_{2t} + \alpha_3 D_{3t} + \beta_0 P_t + \beta_1 (D_{1t} P_t) + \beta_2 (D_{2t} P_t) + \beta_3 (D_{3t} P_t) + u_t$$

$$V_t = \alpha + \beta \cdot P_t + u_t \quad (R)$$

אם דוחים את  $H_0$ , יש לבדוק במבחן WALD האם ההבדל הוא בין החותכים

$$H_0 : \beta_1 = \beta_2 = \beta_3 = 0, H_0 : \alpha_1 = \alpha_2 = \alpha_3 = 0$$

אם דוחים את  $H_0$  יש להמשיך לבדוק באמצעות מבחן  $t$  :

$$H_0 : \beta_j = 0, H_0 : \alpha_j = 0$$

### משתני דמי עבור שני משתנים איכותיים :

לדוגמא – שני משתנים איכותיים המשפיעים על פונקציית השכר : מגדר (אישה, גבר) וגזע (לבן, שחור).

נגדיר משתנה דמי  $G$  שיקבל 1 אם מדובר בגבר ו-0 אחרת (אישה).

נגדיר משתנה דמי  $R$  שיקבל 1 אם מדובר בלבן ו-0 אחרת (שחור).

נבדוק כיצד מגדר וגזע משפיעים על השכר ההתחלתי (החותך), כאשר השכר תלוי

גם בשנות לימוד  $(S_t)$ .

1. הבדל בחותך ללא אינטראקציה :

$$W_t = \alpha_0 + \alpha_1 G + \alpha_2 R + \beta \cdot S_t + u_t$$

במודל זה – אין השפעה משולבת של מגדר וגזע על השכר ההתחלתי.

ניתן לבדוק השערות על כל אחד מהמשתנים הבי"ת האיכותיים בנפרד :

$$1. H_0 : \alpha_1 = 0 \quad \text{הבדל בשכר ההתחלתי בין גברים לנשים}$$

$$2. H_0 : \alpha_2 = 0 \quad \text{הבדל בשכר ההתחלתי בין שחורים ללבנים}$$

2. הבדל בחותך עם אינטראקציה :

$$W_t = \alpha_0 + \alpha_1 G + \alpha_2 R + \alpha_3 G \cdot R + \beta \cdot S_t + u_t$$

המודל :  $W_t = \alpha_0 + \alpha_1 G + \alpha_2 R + \alpha_3 G \cdot R + \beta \cdot S_t + u_t$ .  
במודל זה הטענה היא כי קיימת, בנוסף להשפעה של מגדר וגזע בנפרד על השכר, גם השפעה משולבת (אינטראקציה) של מגדר וגזע על השכר ההתחלתי.

במודל זה, לעומת הקודם, נוספת ההשערה לבדיקת השפעת האינטראקציה בין מגדר לגזע על השכר ההתחלתי :

$$H_0 : \alpha_3 = 0$$

3. דרך נוספת ליצירת מודל עם אינטראקציה :

הגדרת משתני דמי המייצגים שילוב בין המשתנים האיכותיים גזע ומגדר באופן הבא :

$D_1$  יקבל 1 אם מדובר בגבר לבן ו-0 אחרת.

$D_2$  יקבל 1 אם מדובר בגבר שחור ו-0 אחרת.

$D_3$  יקבל 1 אם מדובר באשה לבנה ו-0 אחרת.

הנשים השחורות מהוות כאן את קבוצת הייחוס.

$$W_t = \gamma_0 + \gamma_1 D_1 + \gamma_2 D_2 + \gamma_3 D_3 + \delta \cdot S_t + u_t$$

נעזר בטבלה בכדי לנסח את ההשערות לבדיקת האינטראקציה :

הפרש	אישה	גבר	
$\gamma_1 - \gamma_3$	$\gamma_0 + \gamma_3$	$\gamma_0 + \gamma_1$	לבן
$\gamma_2$	$\gamma_0$	$\gamma_0 + \gamma_2$	שחור
	$\gamma_3$	$\gamma_1 - \gamma_2$	הפרש

ההשערות לבדיקת קיום האינטראקציה :  $H_0 : \gamma_1 - \gamma_3 = \gamma_2$  או  $H_0 : \gamma_1 - \gamma_2 = \gamma_3$   
התוצאות שיתקבלו כאן יהיו כמובן זהות לחלוטין לתוצאות שהתקבלו בדרך

$$WALD = t^2$$

$$PF = Pt$$

## שאלות:

## משתנה דמי לחותך:

- (1) על בסיס מדגם של 50 איש העובדים בחברה מסוימת התקבלו התוצאות הבאות:  

$$W_t = 5500 + 1043 \cdot D + 119 \cdot S_t$$
 (S.E) (134) (56) (24)  
 המספרים בסוגריים הם טעויות התקן של מבחני המובהקות לפרמטרים.  
 א. מהו השכר ההתחלתי של גבר בעל 12 שנות לימוד?  
 ב. מה ההבדל בשכר ההתחלתי בין גברים לנשים?  
 ג. האם הבדל זה מובהק באוכלוסייה?  
 ד. בדקו את הטענה כי השכר ההתחלתי של גברים גבוה ביותר מ-500 ₪ מזה של נשים.  
 ה. בדקו את הטענה שהשכר ההתחלתי של נשים נמוך ב-600 ₪ מזה של גברים.

## פונקציית רגרסיה המכילה משתנים איכותיים בלבד:

- (2) על אותו המדגם של 50 איש העובדים בחברה מסוימת ביקש החוקר לבדוק האם יש הבדל בשכר הממוצע בין גברים לנשים.  
 תוצאות האמידה:  $W_t = 5200 + 1120 \cdot D$   
 נתון:  $S_{\hat{\alpha}_1} = 63$   
 בדקו האם קיים הבדל מובהק בשכר הממוצע בין נשים וגברים?

## משתנה דמי לשיפוע:

- (3) על בסיס אותו מדגם, ביקש החוקר לדעת האם קיים הבדל מובהק בין גברים לנשים בתוספת לשכר בגין שנות הלימוד.  
 תוצאות האמידה נתונות להלן:  

$$W_t = 5000 + 110 \cdot S_t + 120 \cdot D \cdot S_t + u_t$$
 (68) (23) (25)  
 בדוק את ההשערה.

## משתנה דמי לכל פונקציה:

(4) חוקר רצה לבדוק את הטענה שסוג הכביש משפיע על מס' תאונות הדרכים בקטעי כביש בינעירוניים, בהינתן נפח התנועה. החוקר בדק האם הפונקציה של מס' התאונות בהינתן נפח התנועה, שונה בין כבישים מהירים לבין כבישים שאינם מהירים. לשם כך אמד החוקר את ארבע המשוואות הבאות:

$$1. \quad NUM_t = \gamma_1 + \delta_1 \cdot AVGD_t + \varepsilon_{1t} \quad \text{כבישים מהירים בלבד.}$$

$$2. \quad NUM_t = \gamma_2 + \delta_2 \cdot AVGD_t + \varepsilon_{2t} \quad \text{כבישים לא מהירים בלבד.}$$

$$3. \quad NUM_t = \gamma_3 + \delta_3 \cdot AVGD_t + \varepsilon_{3t} \quad \text{שני סוגי הכביש (כל המדגם).}$$

$$4. \quad NUM_t = \alpha + \beta_1 \cdot TYPE_t + \beta_2 \cdot AVGD_t + \beta_3 \cdot (AVGD \cdot TYPE)_t + U_t$$

כאשר:

$NUM_t$  - מס' תאונות הדרכים הקטלניות בקטע כביש  $t$  בשנה.

$AVGD_t$  - נפח התנועה בקטע כביש  $t$  ליום באלפים.

$TYPE_t$  - משתנה דמי המקבל את הערך 1 כאשר הכביש מהיר, ו-0 כאשר הכביש לא מהיר.

תוצאות אמידת המשוואות מופיעות בהמשך השאלה.

א. בדקו את טענת החוקר בשתי דרכים שונות. ציינו איזה מן המשוואות רלוונטיות עבור כל דרך.

ב. חשבו את הערכים המספריים עבור אומדני משוואה (4).

ג. מהו האומדן הנקודתי למס' התאונות בכביש מהיר כאשר נפח התנועה עומד על ארבעת מכוניות ליום בקטע הכביש האמור?

הועלתה הטענה כי המקדם להשפעה של נפח התנועה בדרכים מהירות הינו כפול מזה שבדרכים לא-מהירות.

ד. מהי השערת האפס לבדיקת הטענה (במונחי משוואה (4))?

ה. מהי הרגרסיה "תחת"  $H_0$  למבחן WALS?

### משוואה (1) - כבישים מהירים בלבד:

The REG Procedure

Model: MODEL1

Dependent Variable: num num

Number of Observations Read 344

Number of Observations Used 344

#### Analysis of Variance

Source	DF	Sum of Squares	Mean Square	F Value	Pr > F
Model	1	4700.81174	4700.81174	89.12	<.0001
Error	342	18039	52.74684		
Corrected Total	343	22740			

Root MSE	7.26270	R-Square	0.2067
Dependent Mean	5.10465	Adj R-Sq	0.2044
Coeff Var	142.27617		

#### Parameter Estimates

Variable	DF	Parameter Estimate	Standard Error	t Value	Pr >  t
Intercept	1	1.55289	0.54303	2.86	0.0045
avgd	1	0.02098	0.00222	9.44	<.0001

## משוואה (2) - כבישים לא מהירים בלבד:

The REG Procedure

Model: MODEL1

Dependent Variable: num num

Number of Observations Read 410

Number of Observations Used 410

### Analysis of Variance

Source	DF	Sum of Squares	Mean Square	F Value	Pr > F
Model	1	971.99073	971.99073	145.83	<.0001
Error	408	2719.34830	6.66507		
Corrected Total	409	3691.33902			

Root MSE	2.58168	R-Square	0.2633
Dependent Mean	1.38780	Adj R-Sq	0.2615
Coeff Var	186.02612		

### Parameter Estimates

Variable	DF	Parameter Estimate	Standard Error	t Value	Pr >  t
Intercept	1	0.14978	0.16360	0.92	0.3605
avgd	1	0.02877	0.00238	12.08	<.0001

### משוואה (3) - שני סוגי הכביש (כל המדגם):

The REG Procedure

Model: MODEL1

Dependent Variable: num num

Number of Observations Read 754

Number of Observations Used 754

#### Analysis of Variance

Source	DF	Sum of Squares	Mean Square	F Value	Pr > F
Model	1	8052.00804	8052.00804	288.84	<.0001
Error	752	20964	27.87730		
Corrected Total	753	29016			

Root MSE 5.27990 R-Square 0.2775

Dependent Mean 3.08355 Adj R-Sq 0.2765

Coeff Var 171.22758

#### Parameter Estimates

Variable	DF	Parameter Estimate	Standard Error	t Value	Pr >  t
Intercept	1	0.73903	0.23665	3.12	0.0019
avgd	1	0.02330	0.00137	17.00	<.0001

**משוואה (4):**

The REG Procedure

Model: MODEL1

Dependent Variable: num num

Number of Observations Read 754

Number of Observations Used 754

## Analysis of Variance

Source	DF	Sum of Squares	Mean Square	F Value	Pr > F
Model	3	8256.966	2752.322	99.44	<.0001
Error	750	20759	27.678		
Corrected Total	753	29016			

Root MSE	5.26102	R-Square	0.2846
Dependent Mean	3.08355	Adj R-Sq	0.2817
Coeff Var	170.61553		

## Parameter Estimates

Variable	DF	Parameter Estimate	Standard Error	t Value	Pr >  t
Intercept	1	0.14978	0.33340	0.45	0.6534
type	1				0.0067
avgd	1				<.0001
avgdtype	1				0.1283

## משתנה איכותי עם יותר משתי קטגוריות:

(5) ענה על הסעיפים הבאים:

- א. הועלתה הטענה כי יש הבדל במחיר ההתחלתי בין האביב לקיץ.  
 i. מהי השערת האפס לבדיקת הטענה?  
 ii. פרטו שני מבחנים סטטיסטיים בעזרתם ניתן לבדוק את הטענה.
- ב. הועלתה הטענה כי יש רק שתי עונות המשפיעות על מחיר הירקות ההתחלתי: קיץ + אביב, חורף + סתיו.  
 i. מהי השערת האפס לבדיקת הטענה?  
 ii. מהו המבחן הסטטיסטי המתאים? פרטו.

## משתנה דמי עבור שני משתנים איכותיים:

(6) חוקר בדק השפעות של השכלה, גזע (שחור, לבן) וניסיון (EXP) על לוג השכר ( $\ln(Y)$ ) במדגם בן 306 תצפיות:

$$\ln(Y)_t = \alpha_0 + \alpha_1 D_1 + \alpha_2 D_2 + \alpha_3 D_3 + \beta_1 EXP_t + \beta_2 EXP_t^2 + u_t$$

$\ln(Y)$  - לוג השכר.

EXP - שנות ניסיון.

$D_1$  - מקבל את הערך 1 עבור שחורים בעלי השכלה גבוהה (ו-0 אחרת).

$D_2$  - מקבל את הערך 1 עבור שחורים בעלי השכלה נמוכה (ו-0 אחרת).

$D_3$  - מקבל את הערך 1 עבור לבנים בעלי השכלה גבוהה (ו-0 אחרת).

תוצאות אמידת משוואת הרגרסיה מוצגות בבלט להלן:

Analysis of Variance					
Source	DF	Sum of Squares	Mean Square	F Value	Pr > F
Model	5	-----	-----	-----	-----
Error	300	140	-----		
Corrected Total	305	210			
Root MSE			-----	R-Square	-----
Dependent Mean			-----	Adj R-Sq	-----
Coeff Var			-----		

Parameter Estimates

Variable	DF	Parameter		t Value	Pr >  t
		Estimate	Standard Error		
Intercept	1	-----	-----	60.84	0.00
D1	1	-----	-----	-3.20	0.00
D2	1	-----	-----	-5.56	0.00
D3	1	-----	-----	7.23	0.00
EXP	1	-----	-----	8.11	0.00
EXP <sup>2</sup>	1	-----	-----	-7.45	0.00

- א. לפי המשוואה הניסיון זהה עבור שחורים ולבנים :  
 נכון/לא נכון/ לא ניתן לדעת
- ב. בדוק את הטענה כי בקרב אנשים בעלי השכלה נמוכה אין השפעה לגזע.
- ג. בדוק את הטענה כי אין השפעות השכלה בקרב לבנים.
- ד. מהי השערת האפס לבדיקת הטענה כי אין אינטראקציה בין גזע להשכלה?
- ה. לבדיקת ההשערה של הסעיף הקודם בוצע מבחן W.L.D.  
 הרגרסיה המוגבלת תחת השערת האפס הינה :  

$$Z_0 = \gamma_0 + \gamma_1 Z_1 + \gamma_2 Z_2 + \gamma_3 Z_3 + \gamma_4 Z_4 + v$$
 מהם ה-Zים?
- ו. בדוק את ההשערה אם ידוע שבמודל המוגבל  $R^2 = 0.33$ .
- ז. החוקר החליט לאמוד במקום את המשוואה המקורית את המשוואה :  

$$\ln(Y)_t = \lambda_0 + \lambda_1 S + \lambda_2 E + \lambda_3 (S \cdot E) + \delta_1 EXP + \delta_2 EXP^2 + \omega_t$$
 כאשר :  
 S מקבל את הערך 1 עבור שחורים ו-0 אחרת (לבנים).  
 E מקבל את הערך 1 עבור השכלה גבוהה ו-0 אחרת (השכלה נמוכה).  
 מה הקשר בין המקדמים של שני המודלים?
- ח. אם יאמוד החוקר את המשוואה :  

$$\ln(Y)_t = \lambda_0 + \lambda_1 S + \lambda_2 E + \delta_1 EXP + \delta_2 EXP^2 + \omega_t$$
 האם תהיה טעות ספציפיקציה של השמטת משתנה רלוונטי (היעזר בסעיפים ד', ו' ו-ז').

(7) חוקרת בדקה השפעות השכלה, מגדר וניסיון על הכנסה מעבודה לפי המשוואה הבאה :

$$\ln(MWAGE) = \alpha_0 + \alpha_1 \cdot S + \alpha_2 \cdot E + \alpha_3 \cdot (S \cdot E) + \beta_0 \cdot EXP + \beta_1 \cdot (EXP \cdot S) + \beta_2 \cdot (EXP \cdot S) + \beta_3 \cdot (EXP \cdot S \cdot E) + U$$

כאשר :

$S$  משתנה דמי : 1 = עבור נשים, 0 = גברים.

$E$  משתנה דמי : 1 = עבור השכלה גבוהה ( $scl > 12$ ), 0 = השכלה נמוכה.

א. רשמו את הפונקציה לחישוב :

- i. תחזית לוג השכר עבור גבר בעל השכלה נמוכה ו-10 שנות ניסיון.
- ii. תחזית לוג השכר ההתחלתי עבור נשים משכילות.
- iii. לאחר כמה שנות ניסיון ישתווה השכר של נשים משכילות לזה של גברים משכילים?

ב. רשמו את השערות האפס המתאימות לבדיקת הטענות הבאות :

- i. אין השפעה של מגדר והשכלה על השכר.
- ii. השפעת ההשכלה אינה תלויה במגדר.
- iii. אין השפעות השכלה אצל גברים.
- iv. אין הבדל בשיעורי התשואה לניסיון, בקרב הנשים.

## תשובות סופיות:

- (1) א.  $W_t = 7971$  . ב. 1,043 נח. ג. כן. ד. יש עדות לכך.
- ה. יש עדות לכך. (2)  
יש עדות לכך. (3)  
יש עדות לכך. (3)
- (4) א. יש עדות לכך, מבחן CHOW : 1, 2 ו-3, משתנה דמי : 3 ו-4.  
ב.  $\hat{\alpha} = 0.14978$ ,  $\hat{\beta}_1 = 1.40311$ ,  $\hat{\beta}_2 = 0.002877$ ,  $\hat{\beta}_3 = -0.008$ .  
ג.  $NUM_t = 1.532398$ . ד.  $H_0 : \beta_2 + \beta_3 = 2 \cdot \beta_2$   
 $H_0 : \beta_3 = \beta_2$   
ה.  $NUM_t = \alpha + \beta_1 \cdot TYPE_t + \beta_3 \cdot (AVGD_t + AVGD \cdot TYPE_t) + U_t$ .
- (5) א.i.  $H_0 : \alpha_1 = \alpha_2$  . ii. WALD t-1 . ב.i.  $H_0 : \alpha_1 = \alpha_2, \alpha_3 = 0$  . ii. WALD .
- (6) א. נכון. ב. יש עדות לכך. ג. יש עדות לכך. ד.  $H_0 : \alpha_3 = \alpha_1 - \alpha_2$  או  $H_0 : \alpha_2 = \alpha_1 - \alpha_3$ .  
ה.  $Z_0 = \ln(Y)_t$ ,  $Z_1 = D_1 + D_3$ ,  $Z_2 = D_2 - D_3$ ,  $Z_3 = EXP_t$ ,  $Z_4 = EXP_t^2$ .  
ו. אין עדות לכך. ז.  $\lambda_0 = \alpha_0$ ,  $\lambda_1 = \alpha_2$ ,  $\lambda_2 = \alpha_3$ ,  $\lambda_3 = \alpha_1 - \alpha_2 - \alpha_3$ . ח. לא.
- (7) א.i.  $\hat{\ln}(MWAGE) = \hat{\alpha}_0 + \hat{\beta}_0 \cdot 10$  . ii.  $\hat{\ln}(MWAGE) = \hat{\alpha}_0 + \hat{\alpha}_1 + \hat{\alpha}_2 + \hat{\alpha}_3$  . iii.  $EXP_t = \frac{-(\alpha_1 + \alpha_3)}{\beta_1 + \beta_3}$  . ב.i.  $H_0 : \alpha_1 = \alpha_2 = \alpha_3 = \beta_1 = \beta_2 = \beta_3 = 0$  . ii.  $H_0 : \alpha_3 = \beta_3 = 0$  . iii.  $H_0 : \alpha_2 = \beta_2 = 0$  . iv.  $H_0 : \beta_2 + \beta_3 = 0$  .

# רגרסיה ושיטות ניתוח ליניאריות

פרק 5 - קו הרגרסיה במדגם

תוכן העניינים

1. כללי ..... (ללא ספר)

# רגרסיה ושיטות ניתוח ליניאריות

פרק 6 - מובהקות הרגרסיה באוכלוסיה

תוכן העניינים

1. כללי ..... (ללא ספר)

# רגרסיה ושיטות ניתוח ליניאריות

פרק 7 - מאפייני קו הרגרסיה המרובה במדגם

תוכן העניינים

1. כללי ..... (ללא ספר)

# רגרסיה ושיטות ניתוח ליניאריות

פרק 8 - מובהקות קו הרגרסיה המרובה ומקדמיו באוכלוסיה

תוכן העניינים

1. כללי ..... (ללא ספר)

# רגרסיה ושיטות ניתוח ליניאריות

פרק 9 - שיטות להרצת רגרסיה רבת משתנים

תוכן העניינים

1. כללי ..... (ללא ספר)

# רגרסיה ושיטות ניתוח ליניאריות

פרק 10 - רגרסיה לוגיסטית

תוכן העניינים

1. רגרסיה לוגיסטית.....63

## הגרסה לוגיסטית:

רקע:

מתי נבצע רגרסיה לוגיסטית?

כאשר המשתנה המנובא הוא דיכוטומי (Binary Logistic):  
 יכול לקבל ערכים של 0 או 1.  
 הפונקציה הלוגיסטית מתארת את הסיכויים לקבל "1" במשתנה התלוי כתלות במשתנים ה"ב"ת.

הלוגיקה בניתוח רגרסיה לוגיסטית:

השוואת ניבוי Y ללא המשתנים המנבאים במודל לניבוי Y במודל הכולל את המשתנים המנבאים (סטטיסטי  $\chi^2$ ).

טיב מודל הרגרסיה ("Goodness of fit"):

1. מובהקות המודל:

Omnibus Tests of Model Coefficients

	Chi-square	df	Sig.
Step 1 Step	12.225	4	.016
Block	12.225	4	.016
Model	12.225	4	.016

מבחן  $\chi^2$  - תחת שורת ה-model נמצא את חי בריבוע ואת מובהקות המודל.

2. אחוז שונות מוסברת:

Model Summary

Step	-2 Log likelihood	Cox & Snell R Square	Nagelkerke R Square
1	96.524	.139	.189

$Nagelkerke R^2$  – מקביל ל- $R^2$  כללי ברגרסיה. אחוז שונות Y המוסברת ע"י כל המנבאים יחד (בטווח מוכר של 0-1).

## 3. דיוק בניבוי :

Classification Table<sup>a</sup>

Observed				Predicted		Percentage Correct
				whether mom believes course will help		
				no	yes	
Step 1	whether mom believes course will help	no	yes	46	5	90.2
				17	14	45.2
	Overall Percentage					73.2

a. The cut value is .500

- סגוליות (true negative) – ביחס ל-  $Y=0$  במדגם, כמה המודל דייק בניבוי (90.2%).
- רגישות (true positive) – ביחס ל-  $Y=1$  במדגם, כמה המודל דייק בניבוי (45.2%).
- אחוז הניבוי הכללי – בכמה בסה"כ המודל מדייק בניבוי (73.2%).

## מושגים חשובים להבנת טבלת המקדמים :

: ODDS

"הסיכוי להתרחשות אירוע מסוים" – ההסתברות שהאירוע יקרה לעומת ההסתברות

$$ODDS = \frac{p}{1-p} \text{ : יקרה לא אירוע}$$

ODDS=1 – הסיכוי שהאירוע יתרחש שווה לסיכוי שהוא לא יתרחש  $(\frac{0.5}{0.5})$ .

ODDS>1 – הסיכוי שהאירוע יתרחש גבוה מהסיכוי שלא יתרחש (למשל-  $\frac{0.75}{0.25}$ ).

ODDS<1 – הסיכוי שהאירוע יתרחש נמוך מהסיכוי שלא יתרחש (למשל-  $\frac{0.25}{0.75}$ ).

: ODDS RATIO (OR)

$$OR = \frac{ODDS(A)}{ODDS(B)} \text{ - יחס בין סיכויים}$$

כיצד משתנה ההסתברות במעבר מקבוצה A לקבוצה B.

OR=1 – הסיכוי להתרחשות האירוע שווה בין שתי הקבוצות- אין קשר בין המב"ת למ"ת.

OR>1 – הסיכוי להתרחשות האירוע בקבוצה A גבוה מאשר בקבוצה B – קשר חיובי.

OR<1 – הסיכוי להתרחשות האירוע בקבוצה A נמוך מאשר בקבוצה B – קשר שלילי.

## טבלת המקדמים – תרומות ייחודיות של כל מנבא:

(מקביל לטבלת Coefficients בגרסיות לינאריות)

		Variables in the Equation					
		B	S.E.	Wald	df	Sig.	Exp(B)
Step 1 <sup>a</sup>	EDU_YRS	-.107	.138	.603	1	.438	.898
	AGE	-.029	.020	2.078	1	.149	.971
	SATISFAC	.118	.175	.457	1	.499	1.126
	BIRTH#	.882	.321	7.530	1	.006	2.415
	Constant	.001	1.796	.000	1	.999	1.001

a. Variable(s) entered on step 1: EDU\_YRS, AGE, SATISFAC, BIRTH#.

1. מבחן WALD למובהקות המשתנים:  
מבטא את מובהקות המשתנה מבחינת תרומתו הייחודית לניבוי Y.
2. B – מקדמי המשתנים ב-log odds:  
בטא חיובית – עלייה ב-log odds של Y כפונקציה של עליה ביחידה אחת של X.  
בטא שלילית – ירידה ב-log odds של Y כפונקציה של עליה ביחידה אחת של X.
3. משוואת הרגרסיה:

$$\log\left(\frac{p}{1-p}\right) = \hat{\alpha} + \hat{\beta}_1 x_{1i} + \hat{\beta}_2 x_{2i} + \hat{\beta}_3 x_{3i} + \hat{\beta}_4 x_{4i}$$

$$p = \frac{1}{1+e^{-\log odds}} : (p) \text{ חישוב הניבוי במונחי הסתברות}$$

$$ODDS = e^{\log odds} : (ODDS) \text{ חישוב הניבוי במונחי סיכויים}$$

4. Exp(B) - יחס הסיכויים (Odds Ratio):  
מבטא את העלייה (אם גדול מ-1) או את הירידה (אם קטן מ-1) בסיכויים להיות בעלי ערך '1' ב-Y כאשר הערך במשתנה המנבא גדל ביחידה אחת.

$$\log \text{Exp}(B) = B \quad ; \quad e^B = \text{Exp}(B) : \text{Exp}(B) \text{ ל-} B \text{ היחס בין } B$$

## שאלות:

- 1) חוקרת בחוג למגדר ביקשה לבדוק האם מגדר משפיע על תעסוקה. היא התבססה על סקר של הלמ"ס שדגם 826 מבוגרים בגילאי העבודה המרכזיים (25-55). היא הגדירה את המשתנים באופן הבא:  
 WOMEN - "1" = אישה ; "0" = גבר.  
 WORKING - "1" = כן ; "0" = לא.  
 מהצלבה של שני המשתנים התקבלה הטבלה הבאה:

		women		Total
		.00	1.00	
working	.00	13	130	143
	1.00	338	345	683
Total		351	475	826

- על סמך הטבלה חשבו:
- מה ההסתברות של אישה לעבוד?
  - מה הסיכוי של אישה לעבוד?
  - מה ההסתברות של גבר לעבוד?
  - מה הסיכוי של גבר לעבוד?
  - מה יחס הסיכויים (OR) של נשים לעבוד לעומת גברים?
  - מה הלוגריתם של יחס הסיכויים?
  - מה יהיה ערך מקדם השיפוע B בגרסיה הלוגיסטית לניבוי תעסוקה על פי מגדר ומה משמעותו?
  - מה יהיה ערך  $Exp(B)$  בגרסיה הלוגיסטית ומה משמעותו?

2) במחקר ביקשו לבדוק כיצד מצב משפחתי וגובה המשכורת משפיעים על בעלות על דירה.

משתני המחקר:

apartm - בעלות על דירה: "1" - כן; "0" - לא.

status - מצב משפחתי: status (0) - רווק; status (1) - בזוגיות;

status (2) - בזוגיות עם ילדים; status(3) - פרוד או גרוש.

incom - הכנסה (בעשרות אלפי שקלים).

התקבלו הממצאים הבאים:

Observed	Predicted	apartm		Percentage Correct
		.00	1.00	
		Step 1 apartm .00	22	11
1.00	10	22	68.8	
Overall Percentage			67.7	

a. The cut value is .500

Omnibus Tests of Model Coefficients

	Chi-square	df	Sig.
Step 1 Step	10.218	4	.037
Block	10.218	4	.037
Model	10.218	4	.037

Model Summary

Step	-2 Log likelihood	Cox & Snell R Square	Nagelkerke R Square
1	79.876 <sup>a</sup>	.145	.194

a. Estimation terminated at iteration number 4 because parameter estimates changed by less than .001.

Variables in the Equation

	B	S.E.	Wald	df	Sig.	Exp(B)
Step 1 <sup>a</sup> Status			.682	3	.877	
Status(1)	-.498	.713	.487	1	.485	.608
Status(2)	-.520	.784	.441	1	.507	.594
Status(3)	-.180	.748	.058	1	.810	.835
income	.000	.000	8.580	1	.003	2.536
Constant	-2.734	1.079	6.417	1	.011	.065

a. Variable(s) entered on step 1: Status, income.

- א. האם ניתן לדחות את השערת האפס הטוענת כי אין קשר בין בעלות על דירה להכנסה ולסטטוס משפחתי?
- ב. כמה אחוזים מצליחים המשתנים הבי"ת להסביר מהשונות של המשתנה "בעלות על דירה"?
- ג. באיזה אחוז מצליח המודל לנבא באופן מדויק בעלות על דירה מתוך כלל המקרים?
- ד. באיזה מידה מצליח המודל לנבא בהצלחה בעלות על דירה מתוך בעלי הדירה במדגם? כיצד נקרא המדד המתאים?
- ה. באיזה מידה מצליח המודל לנבא בהצלחה אי-בעלות על דירה מתוך אלו שאינם בעלי דירה במדגם? כיצד נקרא המדד המתאים?
- ו. מהי המשוואה לניבוי בעלות על דירה על סמך המשתנים הבי"ת?
- ז. לאיזה מהמשתנים הבי"ת יש תרומה ייחודית מובהקת לניבוי בעלות על דירה? מהי משמעות מקדם B ו- $\text{Exp}(B)$  של משתנה זה?
- ח. על כל עליה ב-10,000 ₪ בהכנסה, בכמה אחוזים יעלה הסיכוי לבעלות על דירה?
- i. 53.6%
- ii. 253.6%
- iii. 153.6%
- iv. 93%
- ט. על כל עליה של 20,000 ₪ בהכנסה, בכמה אחוזים יעלה הסיכוי לבעלות על דירה?
- i. 307%
- ii. 423%
- iii. 542%
- iv. 642%
- י. מה ההסתברות של רווק המשתכר 20,000 ₪ להיות בעלים של דירה?
- יא. האם ההסתברות של אותו רווק להיות בעל דירה גבוהה / שווה / קטנה מההסתברות שלו לא להיות בעל דירה?
- יב. מהם הסיכויים (ODDS) שלו להיות בעל דירה?
- יג. עבור איזה משכורת הסיכוי (הסתברות) של רווק להיות בעל דירה עולה על הסיכוי שלו לא להיות בעל דירה?
- יד. במידה ומשתנה ההכנסה היה נמדד באלפי שקלים (ולא בעשרות אלפי שקלים), כיצד הדבר היה משפיע על ההשפעה השולית של מקדם ההכנסה, אם בכלל?

- 3) חוקרים בחנו את המאפיינים שעשויים לנבא את הביצוע של חניכים במבחן הסיום של קורס פקחי טיסה. הביצוע במבחן נמדד על סולם של הצלחה/כשלון והמשתנים הבלתי תלויים כללו מין (1-זכר 0-נקבה), השכלה קודמת (0-ריאלית, 1-לא ריאלית) וביצוע במהלך הקורס (1-7).  
להלן תוצאות ניתוח הרגרסיה:

	Chi-square	Df	Sig.
Step 1 Step	20.982	3	.000
Block	20.982	3	.000
Model	20.982	3	.000

## Model Summary

Step	-2 Log likelihood	Cox & Snell R Square	Nagelkerke R Square
1	17.209 <sup>a</sup>	.503	.699

a. Estimation terminated at iteration number 7 because parameter estimates changed by less than .001.

## Variables in the Equation

	B	S.E.	Wald	df	Sig.	Exp(B)
Step 1 <sup>a</sup> מין	4.445	2.611	2.897	1	.089	85.161
השכלה קודמת	-.146	2.054	.005	1	.943	.864
ביצוע במהלך הקורס	2.283	.944	5.846	1	.016	9.810
Constant	-19.284	8.056	5.731	1	.017	.000

a. Variable(s) entered on step 1: מין, השכלה, הקורס, הביצוע במהלך הקורס.

- האם למודל הכולל את שלושת המנבאים יכולת הסבר משמעותית?
- כמה אחוזים מתוך השונות של Y מצליח המודל להסביר?
- מהי משוואת הניבוי?
- לאיזה מן המשתנים הב"ת תרומה מובהקת לניבוי?
- הסבירו את משמעות המקדמים (b) שהתקבלו עבור המשתנים הב"ת: מגדר, השכלה קודמת והביצוע במהלך הקורס.
- בטאו את המקדמים במונחי הסיכויים להצלחה בקורס (odds) והסבירו אותם.
- הועלתה הטענה כי ההסתברות ההצלחה של נשים בקורס היא נמוכה ביותר, גם אם הן בעלות השכלה ריאלית ושביצוען במהלך הקורס מקסימאלי. אנא בדקו את הטענה.
- עבור זכר, בעל השכלה ריאלית, מהי ההשפעה השולית של עליה ביחידה אחת בדירוג הביצוע במהלך הקורס על הסיכוי להצליח בקורס?

(4) לפי מדגם של 20 זוגות נשואים, נאספו נתונים על המשתנה  $Y$  השווה ל-1 אם הזוג נוהג לצאת למסעדה לפחות פעם בשבוע ו-0 אחרת.

$$\text{נאמד המודל: } p = \frac{1}{1+e^{-z}} \text{ כאשר } p = P(Y=1).$$

התקבלו התוצאות הבאות:  $z = -9.456 + 0.368INCOM - 1.207BABY$ .  
 $INCOM$  - ההכנסה של שני בני הזוג (באלפים). ההכנסה במדגם נעה בין 17 אלף ל-44 אלף.

$BABY$  - משתנה דמי המקבל את הערך '1' אם הזוג צריך להיעזר בשמרטפית ו-'0' אחרת.

ענה נכון/לא נכון:

- זוג הנעזר בשמרטפית ומשתכר 30.5 אלף, יוצא למסעדה לפחות פעם בשבוע בהסתברות גבוהה מ-0.5.
- עבור זוג שאינו נעזר בשמרטפית, עליה של אלף שח בהכנסה, מעלה את ההסתברות לצאת למסעדה ב-0.368.
- כל אחד מערכי  $P$  הנאמדים כאן איננו גבוה יותר מ-0.99.
- הסיכוי של זוג, שהכנסתו עלתה ב-3000 שח, לצאת למסעדה יעלה ב-200% בערך.
- המשכורת שצריך להרוויח זוג, אשר אינו נעזר בשמרטפית, כדי שהסיכוי שלו לצאת למסעדה יהיה שווה לסיכוי שלא לצאת למסעדה הוא 27,000.
- זוג, שלא נעזר בשמרטפית, צריך להרוויח יותר מ-28,000 שח כדי שהסיכוי שלו לצאת למסעדה יהיה גבוה פי 3 מהסיכוי שלו לא לצאת למסעדה.
- עבור odds ratio של משתנה "שמרטפות" התקבל רווח בר סמך הבא:  
 $[0.123 ; 1.01]$  ברמת ביטחון של 95%.
- לפיכך ניתן לומר כי למשתנה "שמרטפות" תרומה מובהקת לניבוי הסיכוי לצאת למסעדה.

(5) בשנה מסוימת הוגשו 750 בקשות לקבלת משכנתא ורק חלק מהן אושר. המשתנה התלוי  $Y=1$  אם הבקשה למשכנתא אושרה ול-0 אם נדחתה. המנבאים:

$S$  משתנה דמי השווה ל-1 אם מבקש המשכנתא הוא רווק ול-0 אחרת.  
 $AGE =$  גיל בשנים.

$$\text{המודל הנאמד הינו: } p = \frac{1}{1+e^{-z}} \text{ כאשר } p = P(Y=1).$$

$$z = \alpha + \beta_1 age + \beta_2 age^2 + \beta_3 S$$

תוצאות אמידת המודל:  $z = -9.3 + 0.52age - 0.006age^2 - 0.314S$

א. הסבירו את השפעת הגיל והמצב המשפחתי על ההסתברות לאישור המשכנתא.

ב. מה ההסתברות שתאושר משכנתא לרווק בן 30?

ג. עבור איזה גיל ההסתברות של אדם נשוי לקבל משכנתא היא מקסימאלית?

6) משרד הקבלה של האוניברסיטה רצה לבדוק באיזה מידה ניתן לחזות את ההצלחה של הסטודנט בקורס בסטטיסטיקה על סמך נתונים של מבחן פסיכומטרי, ציון ממוצע של תעודת בגרות וסוג תעודת הבגרות: ריאלית או לא ריאלית.

במדגם של 50 סטודנטים נאספו נתונים על המשתנה Y השווה ל-1 אם הסטודנט הצליח במבחן בסטטיסטיקה ו-0 אם נכשל.

כמו כן נרשמו עבור כל סטודנט ציון הפסיכומטרי, ממוצע הבגרות וסוג הבגרות (1 - בגרות ריאלית, 0 - לא ריאלית).

להלן התוצאות שהתקבלו:

	B	S.E.	Wald	df	Sig.
פסיכומטרי	.090	.046	3.723	1	.054
ציון בגרות		2.070	1.089	1	.297
<u>בגרות ריאלית</u>	4.535	2.519	3.241	1	.072
Constant	-84.892	42.858	3.923	1	.048

- א. באיזה שיטת ניתוח הייתם ממליצים להשתמש ומדוע?
- ב. נתון כי ההסתברות להצלחה בקורס בסטטיסטיקה עבור סטודנט שעשה בגרות הומנית, קיבל 690 בפסיכומטרי וציון 9 בבגרות הינה: 0.034.
- ההסתברות של סטודנט שקיבל אותו ציון בפסיכומטרי, עם בגרות הומנית אבל ציונו בבגרות הוא 10 הינה: 0.233.
- על סמך הנתונים הללו השלם את הערך החסר בפלט המקדמים.
- ג. לאיזה משתנים השפעה מובהקת על הסיכוי להצלח במבחן לסטטיסטיקה? (אלפא 10%)
- ד. מה ההסתברות של סטודנט להצליח במבחן אם קיבל 680 בפסיכומטרי, ציון 10 בבגרות ולמד במגמה ריאלית?
- ה. מהו השינוי בסיכויים (odds) להצליח במבחן בסטטיסטיקה כפונקציה של שינוי ביחידה אחת בפסיכומטרי?
- ו. מהי ההשפעה השולית של נקודה נוספת בציון הבגרות על הסיכוי להצליח במבחן בסטטיסטיקה עבור סטודנט שקיבל 640 בפסיכומטרי ולמד במגמה ריאלית?
- ז. רותי שיפרה את הפסיכומטרי שלה ב-20 נקודות.
- בכמה יעלה הסיכוי שלה להצליח בקורס בסטטיסטיקה?
- ח. אם החוקר היה מחליט לקודד בגרות שאינה ריאלית כ-1 ובגרות ריאלית כ-0, האם הדבר היה משפיע על ערכו של  $Exp(b)$  של סוג בגרות ועל המשמעות שלו?

## תשובות סופיות:

- (1) א. 0.73    ב. 2.7    ג. 0.96    ד. 0.24    ה. 0.11  
 ו. -2.207    ז.  $B = -2.207$     ח.  $\text{Exp}(B) = 0.11$
- (2) א. כן.    ב. 19.4%    ג. 67.7%    ד. רגישות = 68.8%  
 ה. סגוליות = 66.7%
- ו.  $\ln(odds) = -2.734 - 0.498status(1) - 0.52status(2) - 0.18status(3) + 0.93 \cdot incom$   
 ז. משתנה "הכנסה"  
 יא. קטנה.    יב. 0.42    יג. 29,400    יד. 0.093    ט. 3    י. 0.3
- (3) א. כן.    ב. 69.9%    ג.  $\log(odds) = -19.284 + 4.445x_{1i} - 0.146x_{2i} + 2.283x_{3i}$   
 ד. המשתנה – "ביצוע במהלך הקורס".    ה. ראו סרטון.  
 ו. מגדר -  $\text{Exp}(b) = 85.19$ , השכלה קודמת -  $\text{Exp}(b) = 0.864$   
 ז. הטענה נכונה ( $p = 0.035$ ).  
 ח.  $\text{Exp}(b) = 9.81$
- (4) א. נכון.    ב. לא נכון.    ג. לא נכון.    ד. נכון.    ה. לא נכון.  
 ו. נכון.    ז. לא נכון.
- (5) א. ראו סרטון.    ב. 0.533    ג. 40
- (6) א. רגרסיה לוגיסטית.    ב.  $B = 2.16$     ג. "פסיכומטרי" ו-"בגרות ריאלית".  
 ד. 0.914    ה. 1.09    ו. 8.67    ז. 504%    ח. ראו סרטון.

# רגרסיה ושיטות ניתוח ליניאריות

פרק 11 - מבחן לדוגמה 1

תוכן העניינים

1. כללי ..... (ללא ספר)

# רגרסיה ושיטות ניתוח ליניאריות

פרק 12 - מבחן לדוגמה 2

תוכן העניינים

1. כללי ..... (ללא ספר)