

Bar Ilan University

אוניברסיטת בר אילן

Department of Mathematics

המחלקה למתמטיקה

848862401 סטטיסטיקה וניתוח נתונים

שם המרצה: ד"ר עינת מליק גדות

סמסטר אי, שנהייל תשפייד 2024

מגיש: דור גאון

R.f: XXXXXXXX

18/05/2024 : תאריך הגשה

1) מודל תיאורטי- א: שאלת מחקר

אילו משתנים משפיעים על השכר וכיצד?

בסט הנתונים המתקבל מקובץ "online foods" ישנם 12 משתנים, מבניהם ארצה לחקור את המשתנה "Monthly Income" המציר את דרגות השכר של המשתתפים בתצפית. המשתנה בעל 5 אופציות אפשריות במדגם המציגים טווח שכר של המשתתפים, לצורך ניתוח אגדיר את חמשת המשתנים הנמצאים כטווח למספר בעל ערך בודד-

הגדרת המשתנה לנתון מספרי בודד	שם המשתנה בקובץ הנתונים
0	No Income
10000	Below Rs.10000
17500	25000 to 10001
37500	50000 to 25001
50000	More than 50000

באמצעות המשתנה, ארצה לבחון אילו משתנים משפיעים על ההכנסה וכיצד זה בא לידי ביטוי-חיובי שלילי או נטרלי ביחס למשתני הבקרה, האם אוכל למצוא דפוס שכר לפי פרמטרים או שמדבור בנתון אקראי לחלוטין שלא נוכל להסביר באמצעות שאר המשתנים. אתייחס בעיקר ליחס הנישואים על השכר, כבחינת מחקרו של פרופי Robert I. Lerman

How Do Marital Status, Wage Rates, and Work Commitment Interact?

נכתב כי:

"We find that marriage increases men's earnings by about 20 percent ... These findings suggest that both marriage enhancing and earnings-enhancing policies can set off a virtuous circle, in which marriage and earnings reinforce each other over time"

ב: בכדי לענות על שאלת המחקר אשתמש ב3 משתני בקרה שתפקידם לבחון שוני בהכנסה, לצורך כד בחרתי 3 משתנים שלדעתי יהיו בעלי השפעה רבה על השכר בהתאם ל-

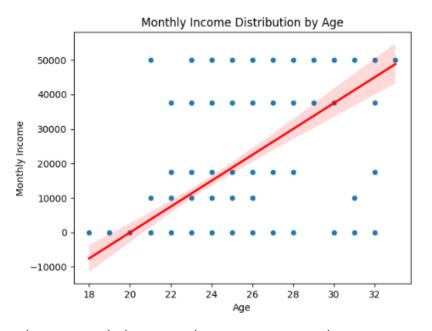
- Marital Status הנועד ליצור הבדלה בין הסטטוס הזוגי של המשתתפים, המשתנה בעל Marital Status (נשוי), Single (לא ידוע), Prefer not to say (נשוי), שלושה סוגי תוצאה: שלושה מוכל לדעת כיצד מצב משפחתי של נבדק משפיע על השכר, האם בעזרת המשתנה הזה נוכל לדעת כיצד מצב משפחתי של נבדק משפיע על השכר, האם נישואים מעלים את השכר או שלמעשה גורמים לירידה?
- Occupation משתנה המציג את הסטטוס התעסוקתי של המשתתפים, המשתנה בעל ארבעה סוגי תוצאה: Employee (שכיר), Housewife (עקרת בית), Student (עקרת בית), Housewife (עצמאי). בעזרת נתון עיסוק המשתתף נוכל לפלח את שכר המשתתפים Self Employed (עצמאי). בעזרת נתון עיסוק המשתתף ב2 קבוצות ספציפיות- עצמאים מול שכירים בהתאם למצב התעסוקתי שלהם, אתמקד ב2 קבוצות ספציפיות- עצמאים מול שכירים מאחר וידוע כי שכר סטודנט מוגבל, וסטטוס עקר בית מלווה לרוב בחוסר הכנסה.
 - Age הנועד לפלח את המשתתפים לפי הגילאים שלהם בתור משתנה רציף, בעזרת המשתנה גיל נוכל לבחון האם גיל גדול יותר גורם לשיפור בשכר- בהתאם לניסיון גבוה יותר. או שלמעשה אין קשר בין הגיל לשכר והוא לא בהכרח משפיע

2) סטטיסטיקה תיאורית- א: סטטיסטיקה תיאורית של משתנה התוצאה ושל כלל המשתנים

נעבד את הנתונים הגולמיים המתקבלים מהקובץ "online foods" כדי לבחון סטטיסטיקה תיאורית, ראשית אנו רואים כי בקובץ ישנם 12 משתנים 388 תצפיות, הגיל הממוצע בקרב תיאורית, ראשית אנו רואים כי בקובץ ישנם 12 משתנים ו388 תצפיות, הגיל הממוצע בקרב המשתתפים הינו 24.6 (מקסימום), עם סטיית תקן של 3. בנוסף, ניתן לראות כי מבין הנבדקים, השכר הממוצע (לאחר התקנון בסעיף 1) עומד על \$17,322.47\$

כאשר אנו בוחנים את מצב המשפחתי של המשתתפים אני מגלים כי:

כאשר אנו בוחנים את השכר לפי קבוצות גיל:



מגרף זה ניתן להבחין בקו רגרסיה עולה בהתאם לגיל, מה שמעיד לנו על שכר גדל ככל שגיל המשתתף בתצפית גבוה, נתון זה מאמת את ההשערה כי גיל משפיע באופן חיובי על השכר. מכך ניתן לטעון כי ככל שאדם צובר יותר ניסיון הוא מתוגמל בשכר גבוה יותר, וכי השכר הממוצע מתקבל בסביבות גיל 24.

$$Y = \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_3 + \in$$

: כאשר

שכר חודשי = Y

 $4375.051 = מצב משפחתי = <math>\beta_1$

0ט מצב שואף = מצב מצב $= \beta_2$

1329.019 = גיל = β_3

שגיאות –∈

$$Y = 4375.051 * X_1 + 0 * X_2 + 1329.019 * X_3 + \in$$

בהתאם לתוצאות שקבלנו ניתן להסיק כי הערך המשמעותי ביותר לשכר הוא המצב המשפחתי של הנבדק נשוי ישנה עליה ניכרת בשכר המשפחתי של הנבדק נשוי ישנה עליה ניכרת בשכר לעומת רווק, אחריו גיל עם ערך נמוך יותר אך עדיין בעל השפעה חיובית שתואם את ההשערה כי גיל ואיתו ניסיון עוזרים לשפר את שכר העובד, יתכן ואם במדגם היו גילאים גדולים מ33 היינו רואים פערים גדולים יותר בהתאם. ניתן גם לראות כי השפעת המצב התעסוקתי באופן מפתיע כלל לא משפיעה על המודל, עם השפעה השואפת לאפס.

מהפלט שקיבלנו ניתן לראות R^2 גבוה מאוד העומד על 73.4%, דבר שמצביע לנו על קשר בין המשתנים לשכר אותו בדקנו, בנוסף ניתן לראות שכלל המשתנים מובהקים <0.05.

ב- המודל המצומצם

במודל המצומצם שערכנו ניתן לראות את השפעת הגיל על ההכנסה:

$$Y = 3756.212 * X_3$$

כאן אנו מוצאים השפעה גדולה יותר מאשר במודל המלא, מכך ניתן להסיק כי לגיל יש השפעה רבה על ההכנסה אך ביחס למצב המשפחתי ההשפעה פוחתת, נתון זה תואם למחקרם של-

"The Marriage Premium in the Michael Debowy, Gil S. Epstein, and Avi Weiss אשר מציג כיצד נישואים בעלי השפעה רבה על שכר העובדים, Israeli Labor Market" ושניתן לראות הפרש בין הנשואים לעמיתם הרווקים.

כאשר אנו משווים בין המודלים אנו רואים כי במודל המלא ה R^2 מקבל 73.4% המעיד על קשר חזק בין המשתנים לתוצאה, לעומת המודל המצומצם המציג R^2 השווה ל31.9%, פחות ממחצית ממה שקבלנו במודל המלא.

בנוסף נוכל לראות כי גם כאן המשתנים מובהקים 0.05>.

תוספת משתני הבקרה תרמו לנו לאמוד תוספת שכר של משתתף במדגם, על ידי כך שאימות מצבים אישיים יכולים להעיד לנו על תוספת בשכר- מצב משפחתי וגיל, אנחנו לא יכולים לומר זאת על המשתנה תעסוקה מאחר והתוספת שלו למודל אפסית ולכן היא לא יכולה לומר זאת על המשתנה תעסוקה מאחר והתוספת שלו למודל אפסית ולכן היינו יכולים לקבוע להעיד לנו על תוספת לשכר. יתכן ואם היינו מוסיפים עוד משתנים למודל היינו יכולים לקבוע בצורה יותר ברורה מהם המשתנים העיקריים המשפיעים על השכר, אך בהתאם למה שבדקנו נוכל לקבוע באופן חד משמעי שמצב משפחתי וגיל מקבלים ערך משמעותי לחיזוי השכר.

: המודל המלא

Dep. Variable:	Mon	thly Income	R-squared:			0.734
Model:		OLS	Adj. R-squ	ared:		0.732
Method:	Le	ast Squares	F-statisti	c:		353.9
Date:	Sat,	18 May 2024	Prob (F-st	atistic):	3.7	78e-110
Time:		14:51:40	Log-Likeli	hood:		4131.5
No. Observation	s:	388	AIC:			8271.
Df Residuals:		384	BIC:			8287.
Df Model:						
Covariance Type		nonrobust				
	coef	std err	t	P> t	[0.025	0.975]
const	-5.714e+04	4511.643	-12.666	0.000	-6.6e+04	-4.83e+04
Marital Status	4375.0510	1243.288	3.519	0.000	1930.547	6819.559
Occupation	2.555e+04	1087.134	23.505	0.000	2.34e+04	2.77e+04
Age	1329.0195	218.037	6.095	0.000	900.323	1757.716
omnibus:		47.419	 Durbin-Wat	======= son:		2.015
Prob(Omnibus):		0.000	Jarque-Ber	a (JB):	1	130.533
Skew:		0.571	Prob(JB):		4.	.52e-29
Kurtosis:		5.602	Cond. No.			217.

: סטטיסטיקה תיאורית

	Age	Monthly	Income	
count	388.000000	388	.000000	
mean	24.628866	17332	.474227	
std	2.975593	19802	.096490	
min	18.000000	0	.000000	
25%	23.000000	0	.000000	
50%	24.000000	10000	.000000	
75%	26.000000	37500	.000000	
max	33.000000	50000	.000000	
p-valu Degree Averag	uare statist e: 2.9270911 s of freedom e Monthly In of observat	09293357 : 8 come is:	e-09 17332.	

: המודל המצומצם

	OLS Regre	ssion Results	
Dep. Variable: Model: Method: Date: Time: No. Observations: Df Residuals: Df Model:	Monthly Income OLS Least Squares Sat, 18 May 2024 14:17:45 388 386	Adj. R-squared: F-statistic: Prob (F-statistic): Log-Likelihood: AIC:	0.319 0.317 180.5 5.05e-34 -4314.3 8633. 8641.
Covariance Type:	nonrobust		
C(pef stderr	t P> t [0.0	======= 25
const -7.518e-		10.838 0.000 -8.88e+ 13.434 0.000 3206.4	
Omnibus: Prob(Omnibus): Skew: Kurtosis:	4.693 0.096 0.267 3.019	Jarque-Bera (JB): Prob(JB):	2.052 4.614 0.0996 207.

חלק ב- ניתוח פלטים (1)

- א- הנחות המודל: מניחים כי המודל מתפלג נורמלית עם תוחלת אפס, בעל שונות קבוע ואי תלות.
 - ב- השערות המחקר:

אין הבדל התוחלת משקלי הבדל אין אין 10: $u_1=u_2=u_3$

אחרת ש הבדל בין משקלי התוחלת בקבוצה H1:

٦-

קבוצה	כמות
1 (בחורות)	13
2 (ילדות)	30
3 (מבוגרות)	18

- 7

קבוצה	ממוצע	סטיית תקן
1 (בחורות)	53.538	2.025
2 (ילדות)	13.8	2.441
3 (מבוגרות)	55.333	3.067

- ה- Variances of Homogeneity of Test הוא מבחן אשר בודק את ההומוגניות של שונות. באמצעותו ניתן לקבוע האם שונות של קובצות במדגם הם הומוגניות בהתאם לרמת המובהקות, ההומוגניות של שונות מוודאת האם ההתפלגות בכל קבוצה שוות או ניתנו להשוואה, האם יש קשר בין הקבוצות השונות.
 - ו- הבדל בין ממוצעי הקבוצות : ניתן לראות כי הסיגמה בטבלת ANOVA היא 0.000 ולכן ניתן לקבוע כי יש הבדל בין הקבוצות, כאשר הערך קטן מ0.05 אנו יודעים שקיים הבדל.
 - eta²

:ANOVA כדי לחשב נשתמש בטבלת

Total -מבין Between Groups לעומת סהייכ Sum of Squares נחשב את היחס של 25379.182/25761.213= **0.985**

מדובר בחישוב הבודק את פרופורציות שינוי המוסבר במשתנים במודל, ככל שמתקבלת תוצאה הקרוב יותר ל1 הקשר בין המשתנים חזק יותר (כי המשתנה בעל משקל גדול יותר), במקרה שלנו נמצא כי מעל ל98% מסך המדגם שייך לקבוצה אותה בדקנו.

ח- מבחן" Test Hoc Post" עוזר לנו להבין את ההבדלים וההשלכות בין הקבוצות שאותם אנו בודקים, באמצעות המבחן אנחנו יכולים להעריך האם ברמת מובהקות מוגדרת נמצא הבדלים בין המשתנים, במקרה שלנו 0.05 שעוזר לנו לזהות כי לפי ערכי הקבוצות, קבוצה 2 (ילדות) שונה מקבוצות 311 בהתאם לערך הסיגמה.

חלק ב- ניתוח פלטים (2)

- א- הנחות המודל: מניחים כי המודל מתפלג נורמלית עם תוחלת אפס, בעל שונות קבוע ואי תלות.
 - ב- השערות המחקר:

אין הבדל בין זמני המשלוח אין אוH0: $u_1=u_2=u_3=u_4=u_5$ ארת יש הבדל בין זמני המשלוח אחרת יש הבדל בין זמני המשלוח H1:

٦-

כמות	סניף
5	1
5	2
4	3
3	4
5	5

-7

סניף	ממוצע	סטיית תקן
1	606.400	206.139
2	240.400	32.292
3	385.000	104.527
4	850.333	8.144
5	632.800	231.598

ה- האם יש הבדל בין הקבוצות! כן.

ניתן לראות זאת לפי טבלת ANOVA כאשר הsig כאשר באר טבלת ANOVA לכן אנחנו יודעים שיש הבדל בין קבוצה אחת לפחות מבין כולם.

eta²

: ANOVA כדי לחשב נשתמש בטבלת

Total -מבין Between Groups נחשב את היחס של Sum of Squares מבין את היחס של 891975.997/1313581.864= **0.679**

מדובר בחישוב הבודק את פרופורציות שינוי המוסבר במשתנים במודל, ככל שמתקבלת תוצאה הקרוב יותר ל1 הקשר בין המשתנים חזק יותר (כי המשתנה בעל משקל גדול יותר), במקרה שלנו נמצא כי קרוב ל68% מסך המדגם שייך לקבוצה אותה בדקנו.

:י בעזרת מבחן Test Hoc Post יי בעזרת מבחן -ז Test Hoc Post

שונה מ	סניף
2	1
1,4,5	2
4,5	3
1,2,3	4
2,3	5

Tukey HSD בעזרת מבחן

שונה מ	סניף
2	1
1,4,5	2
4	3
2,3	4
2	5

ח- מבחן טוקי:

מדובר במבחן POST HOC לצורך השוואה בין הקבוצות שאנו בודקים באמצעות הממוצעים, מדובר בניתוח שונות הדומה למבחן כמו LSD אך עם תוצאה שונה. התוצאה השונה מתקבלת על ידי שינוי ההשוואות בין הקבוצות, בכדי שרמת המובהקות שלהם תהיה בערך מוגדר, כמו במקרה שלנו 0.05 שעזר לנו בסעיף הקודם. מבחן טוקי ביחס לLSD הוא בטוח יותר והסיכוי לטעות קטן.

חלק ג- מבחנים אפרמטריים

: יעילותו של מלטונין כטיפול לבעיית שינה (1

בשביל לחשב את יעילות מלטונין כטיפול לבעיות שנה נעזר במבחן וילקוקסון- מבחן שמטרתו השוואה בין שתי קבוצות על סמך מדגמים בלתי תלויים, המשתנה התלוי הינו כמותי שלא מתפלג נורמלית וכאשר מדובר במדגם קטן הפחות מ30 תצפיות.

: נגדיר

יתרופת הפלסיבו שווה לתרופת המלטונין, לא נמצא יעילות למלטונין OH

: 1H תרופת הפלסיבו שונה מהמלטונין, נמצאה יעילות למלטונין

נבדק	Placebo	Melatonin	הפרש (בערך מוחלט)	דירוג
1	2	4	2	5
2	3	3	0	-
3	5	7	2	5
4	2	4	2	5
5	4	5	1	2
6	6	9	3	7
7	7	6	1	2
8	3	9	6	8.5
9	8	7	1	2
10	1	7	6	8.5

סהייכ דירוג: 45

n=9 : מספר תצפיות בדירוג

סכום הדירוגים החיוביים: 41

סכום הדירוגים השלילים: Ts=4 (הקטן מבין השניים)

הערך הקריטי המתקבל בטבלת וילקוקסון בהתאם לנתונים הוא 8, הערך שקבלנו 4 קטן יותר

. ולכן נדחה את השערה האפס, נמצאה יעילות בשימוש מלטונין וולכן נדחה את השערה האפס, וולכן נדחה את וולכן נדחה את השערה האפס,

: האם יש הבדל בין גברים ונשים בכמות הטסטים עד קבלת רישיון

בשביל לבדוק האם קיים הבדל בין כמות הטסטים בין גברים ונשים עד קבלת הרישיון נבצע בדיקה באמצעות מבחן מו-ויטני מאחר וגדלי המדגם בין הקבוצות אינם שוות, רמת מובהקות 0.05 ונבדוק מבחן חד צדדי כדי למצוא הבדל בין ההתפלגויות

: נגדיר

אין הבדל בין גברים ונשים =0H

שים ונשים =1H

7.5 3 4.5 2 9.5 4 1.5 1 7.5 3 12 5 14.5 6 4.5 2 2 2 2 2 2 2 3 2 2 2 3 2 3 2 4.5 2 4.5 5 14.5 6 4.5 2 12 5 12 5 15 1	דרוג	גברים
9.5 4 1.5 1 7.5 3 12 5 14.5 6 4.5 2 current 2 9.5 4 12 5 14.5 6 4.5 2 12 5 14.5 6 4.5 2 12 5 14.5 6 4.5 2 12 5	7.5	3
1.5 1 7.5 3 12 5 14.5 6 4.5 2 - נשים 4.5 2 9.5 4 12 5 14.5 6 4.5 5 14.5 6 4.5 5	4.5	2
7.5 3 12 5 14.5 6 4.5 2 - נשים 4.5 2 9.5 4 12 5 14.5 6 4.5 2 12 5 5 2 12 5 14.5 6 4.5 2 12 5	9.5	4
12 5 14.5 6 4.5 2 - נשים 4.5 2 9.5 4 12 5 14.5 6 14.5 6 12.5 5	1.5	
14.5 6 4.5 2 נשים 4.5 9.5 4 12 5 14.5 6 4.5 2 12 5	7.5	
לשים - 2 נשים - 3 4.5 2 9.5 4 12 5 14.5 6 4.5 2	12	
רשים		
4.5 2 9.5 4 12 5 14.5 6 4.5 2 12 5	4.5	2
9.5 4 12 5 14.5 6 4.5 2 12 5	-	נשים
12 5 14.5 6 4.5 2 12 5	4.5	
14.5 6 4.5 2 12 5	9.5	4
4.5 2 12 5	12	
12 5	14.5	6
	4.5	
1.5	12	5
	1.5	1

1R=61.5 : סכום הדירוגים גברים

2R=58.5 : סכום הדירוגים נשים

$$Ui = Ri - \frac{n1(n1+1)}{2}$$
 : נחשב את

$$U1 = 61.5 - \frac{8(8+1)}{2} = 25.5$$
 $U2 = 58.5 - \frac{7(7+1)}{2} = 30.5$

25.5= Us = min{U1, U2} : ערך סטטיסטי

:לפי טבלת מן-וטיני

n ₂	α									11	1								
		3	4	5	6	7	-8	9	10	11	12	13	14	15	16	17	18	19	20
3	.05		0	0	1	1	2	2	3	3	4	4	5	5	6	6	7	7	8
	.01		0	0	0	0	0	0	0	0	1	1	1	2	2	2	2	3	3
4	.05		0	1	2	3	4	4	5	6	7	8	9	10	11	11	12	13	14
	.01		-	0	0	0	1	1	2	2	3	3	4	5	5	6	6	7	8
5	.05	0	1	2	3	5	6	7	8	9	11	12	13	14	15	17	18	19	20
	.01	-	-	0	1	1	2	3	4	5	6	7	7	8	9	10	11	12	13
6	.05	1	2	3	5	6	8	10	11	13	14	16	17	19	21	22	24	25	27
	.01	-	0	1	2	3	4	5	6	7	9	10	11	12	13	15	16	17	18
7	.05	1	3	5	6	8	10	12	14	16	18	20	22	24	26	28	30	32	34
	.01	-	0	1	3	4	6	7	9	10	12	13	15	16	18	19	21	22	24
8	.05	2	4	6	8	10	13	15	17	19	22	24	26	29	31	34	36	38	41
	.01		1	2	4	6	7	9	11	13	15	17	18	20	22	24	26	28	30

מתקבל הערך הקריטי 10.

Uקריטיש, $\alpha c < Us$

ולכן לא נוכל לדחות את השערת האפס, נגיע למסקנה שאין הבדל בכמות הטסטים בין המינים.

3) האם קיים הבדל בין שני המועמדים בעימות הרדיופוני, במספר הקולות שסחפו לטובתם? מאחר והמשתנים שלנו הם משתנים נומינליים המקבלים 2 ערכים בלבד $2X^2$ נשתמש במבחן מקנמר:

: נגדיר

שינה את דעת המצביעים =0H

1H העימות שינה את דעת המצביעים

1 אם דרגת עם איב איב ובהתפלגות אמה במבחן איב במבחן יש 20<C+B מאחר מאחר יש 20

$$\chi_s^2 = \sum_{i=1}^n \frac{(o_i - e_i)^2}{e_i} = \frac{\left(b - \frac{b+c}{2}\right)^2}{\frac{b+c}{2}} + \frac{\left(c - \frac{b+c}{2}\right)^2}{\frac{b+c}{2}} = \dots = \frac{(b-c)^2}{b+c}$$

 $\chi_s^2 = \frac{(18-23)^2}{18+23} = 0.609$: במקרה שלנו הערך הסטטיסטי

 $\chi^2_{(1,5\%)} = 3.84$: נחשב את הערך הקריטי

 $\chi^2_{(1,5\%)} > \chi^2_s$: קריטי

לכן לא נוכל לדחות את השערת האפס, העימות לא שינה את דעת המצביעים ברמת מובהקות 5%.

1) How Do Marital Status, Wage Rates, and Work Commitment Interact?

Robert I. Lerman / Urban Institute

https://docs.iza.org/dp1688.pdf

- 2) Michael Debowy, Gil S. Epstein, and Avi Weiss "The Marriage Premium in the Israeli Labor Market" https://www.taubcenter.org.il/wp-content/uploads/2022/12/Marriage-Premium-ENG-2022.pdf
- 3) Dr. Aviel Einat

https://lemida.biu.ac.il/

4) Wikipedia

5) Alexandar osipov

https://u.math.biu.ac.il/~osipova/index.html