

# מבוא לסטטיסטיקה

## קוזלובסקי

2	הרצאה מס' 1 - 30/10/00
2	מושגי יסוד
2	הרצאה 2 - 6/11/00
3	התפלגות
3	הרצאה 3 - 13/11/00
4	הרצאה 4 - 7/11/00
5	הרצאה 6 - 11/12/00
5	הסתברות
6	הרצאה 7 - 18/12/00
7	הרצאה 11 - 8/1/01
7	תיאורית בייס Baye's theory
8	הרצאה 12 - 15/1/1
8	מתאם קורלציה
9	Rank correlation / spearman
9	הרצאה 13 - 26/2/01
9	Standard error of estimate
9	השערות Hypotheses
10	הרצאה 14 - 5/3/1
11	הרצאה 15 - 19/3/1
11	רווח בר סמך Confidence interval
11	T test
11	הרצאה 16 - 26/3/1
12	מבחן זשל שני מדגמים בלתי תלויים
12	הרצאה 17 - 2/4/01
13	מחקר לפני אחרי match before after
13	הרצאה 18 - 16/4/01
15	מבחן r
15	הרצאה 19 - 23/4/01
15	הנחיות לביצוע מבחן t
15	מבחן $\chi^2$ - chi square
16	מבחן מנדל Mendel
16	הרצאה 20 - 30/4/01
16	מבחן $\chi^2$ לבדיקת תלות או אי תלות
17	ניתוח שונות Analysis of variance
17	הרצאה 21 - 7/5/01
18	הרצאה 22 - 14/5/01
18	הנוסחה של Sheffe
19	הרצאה 23 - 4/6/01
19	2 way anova ניתוח שונות דו כיווני

הרצאה מס' 1 - 30/10/00מושגי יסוד

Variable (משתנה) - ממד/תכונה בעלת אפשרויות התבטאות רבות. למשל גלאים. (מכיל טווח גדול של אפשרויות).

Constant (קבוע) - נתון קבוע שאינו משתנה. היות וכך, לרוב הוא פחות מעניין.

שני הנ"ל יכולים להתחלף במשמעויותיהם בהתאם למצבי הרקע הקיים. למשל ממדי חיים ומוות בכיתת לימוד מהווים מצב קבוע אך בבי"ח הם מהווים משתנה היות ומס' המתים/חיים שם הוא משתנה.

Continuous (רציף) - מספר אשר רמת הדיוק שלו היא אין סופית. למשל גובה של בן-אדם או המשקל שלו.

Discrete (בדיד) - נתון מוחלט. למשל מספר פריטים באוכלוסייה.

Infinite

Finite

Sample (מדגם) - על פי פרמטרים מסוימים.

Population (אוכלוסייה) - מספר פריטים בעלי תכונות זהות. מסיק מהמדגם על האוכלוסייה

Independent variable - משתנה בלתי תלוי.

Dependent variable - משתנה תלוי.  $Y=f(x)$  לעתים קיים קושי לזהות את המשתנה התלוי והבלתי תלוי.

אם יש פער בזמן בין המשתנה הבלתי תלוי ניתן להגדיר ביתר קלות מי הם המשתנים. באופן מקביל לסיבה ותוצאה, מהוא הפער בין השניים, מה הקשר ביניהם. במידה וקיים ספק בזיהוי אופי המשתנים, ניתן להתייחס אליהם כאל משתנים בלבד. (דוגמת חברת יצרניות הסיגריות - עישון לא מוביל לסרטן. קיים משתנה נוסף אשר הוא הגורם למחלה ויגרום לה בין אם האדם יעשן או לא).

הרצאה 2 - 6/11/00סולמות - Scales

Nominal (נומינאלי) - משתנה. למשל מין (gender) אמנם אין יותר מדי אפשרויות אבל ציון של משתנה זה המדגם לא מציין או מהווה העדפה של הנתון לכיוון זה או אחר. הוא אפיון קבוצתי של קטגוריה תיאורית בלבד. אין היררכיה כזו או אחרת. מספר ת.ז. הוא שמי. אין בו שימוש יום יומי אבל יחד עם זאת ניתן לקטרג אותו. המיון הוא לא מדור, אין משמעות לסדר.

ordinal (אורדינאלי/דירוגי) - יש סדר בהגדרות. מאפיין בסדר אורדינאלי הוא שיש מדרג אך ההפרש בין דרג אחד לשני אינו ידוע/לא קיים מבחינת המדע. (לדוגמה - דירוג מלכות יופי או העדפה למשקה). אין דרך להגיע למשמעות של ההפרש.

Interval (אינטרוואלי) - בין אורדינלי ליחסי. קיים מדרג אבל לא בייחוס לאפס מוחלט משום שהוא בעייתי להגדרה. קיימת משמעות להפרשים בין המדדים ולהשוואה בין פערים על אותה סקלה. נקודת היחס היא שרירותית על הסקלה. למשל - ציונים במבחן מול רמת הידע של הנבחן. מה רמת הידע שמגלים הציון? האם מי שקיבל 90 יודע פי 2 ממי שקיבל 45? האם קיימת משמעות להפרש שבין שני הציונים?

ratio (יחס/רציונלי) - כמו בפיסיקה, היחס בין מדדים הוא בעל משמעות כמותית. הוא בהשוואה לאפס המוחלט (absolute zero) המהווה את נקודת היחס, בניגוד למספרים הסידוריים בסולם האורדינלי או למספרים הסמליים בסולם הנומינלי.

## התפלגות

### IQ

מאגר גדול של נתונים ניתן להגדיר ע"פ סולם של מדרג, רמת השכיחות ( $f$ ) המופיעה במדרג היא ההתפלגות. שימוש במדרג / אינטרבלי הופך את תוצאות ההתפלגות לנכונות יותר. מה גם שעבוד נתונים להתפלגות הוא נוח יותר לעבודה. השאיפה היא לכל היותר 20 אינטרבליים ובעדיפות אפילו לא יותר מ-10. טווח האינטרבליים תלוי באופי נתונים ובטווח הבסיסי שלהם. טווחים מקובלים הם 5, 10, 50 או 100. התחלת הסדרה תהיה במספר שמתחלק בגודל האינטרבלי ובספרה הגדולה יותר:

Class interval	f
125-129	0
120-124	2
115-119	2
110-114	1
105-109	1

אינטרבלי אמיתי / Real Exact Interval - אינו משאיר פערים בין שלב ושלב באנטרבליים. מתמודד בעיקר עם מדדים רציפים. במקרה כזה האינטרבליים יהיו:

124.5-129.5
124.5-120.5

מאחר ואנו מוגבלים למכשירי מדידה אנו משתמשים ב class interval. סביר להניח שמדדים שאנו מתעסקים איתם הם רציפים אבל אין לנו דרך אמיתית לבדוק זאת.

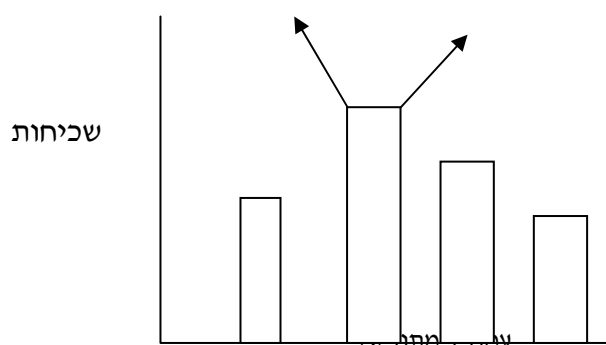
### Cumulative Frequency

שכיחות מצטברת מציגה את סה"כ השכיחות:

f	Cf
0	6
2	6
2	4
1	2
1	1

גבול אמיתי תחתון

גבול אמיתי עליון



## הרצאה 3 - 13/11/00

צורת גרף התפלגות (המקרה זה של משתנה איכותי בדיד): העמודה בנויה על mid point סכום הגבול האמיתי העליון עם הסכום האמיתי התחתון

חלקי 2. ע"פ הספר יש גבול בין כל עמודה, ע"פ קוזלובסקי הרווח הוא לא נכון.

סיגמה  $\Sigma$

סיגמה היא סכום של אברים בתחום מסוים. למשל:

$$\Sigma x_i = X_1 + X_2 + X_3 \dots + X_n$$

כשהטווח הוא  $i=1, i=n$

n מסמלת את המשתנה האחרון.

כל אחד מהמשתנים הם בלתי תלויים. מיקומם ברשימה הוא לא קבוע או דירוגי. כלומר אם

רוצים לעשות  $\Sigma$  של  $1X$  ושל  $3X$  משנים את מיקומם ברשימה כך שיהיו עוקבים, 1 ו 2.

**חוקים**

חוק 1 - ישנם משתנים קבועים. C מסמן קבוע ו N הוא מספר החזרות של הקבוע

$$i=N$$

$$\Sigma c = Nc$$

$$i=1$$

אם c שווה 4 ו N שווה 5 אז התוצאה היא 20.

חוק 2 - נתן להוציא את הקבוע אל מחוץ ל

$$\Sigma c X_i = c \Sigma X_i$$

— ממוצע X

השימוש הוא במילה mean

$$\text{Mean} = \Sigma X_i / N$$

m הוא ממוצע של אוכלוסייה.

## הרצאה 4 - 7/11/00

כל התפלגות נורמלית בנויה בצורה דומה, ללא שום קשר לדבר הנמדד.

נוסחתו של גאוס -  $\Sigma e^{-x^2/2\sigma^2}$

Variance 1 - שונות (של מדגם):  $S^2 = \Sigma (x - \bar{x})^2 / N - 1$

Standard deviation - (סטית התקן): שורש של  $S^2$

Variance 2 - שונות (של אוכלוסייה) -  $\sigma^2 = \Sigma (x - \mu)^2 / N$

היחס בין השונות הוא ש expected של  $s^2$  הוא  $\sigma^2$

כשאר אנו עושים דגימה מתוך אוכלוסיה עלינו להשתמש רק במונה של  $N-1$ , רק כך ממוצע השוניות יהיה שווה לשוניות האוכלוסיה.  
אם משהו הוא קבוע מבחינת הממוצע (הממוצע הוא ידוע ומחייב) הוא "כופה" את נתוני ההתפלגות. חלק מהנתונים הוא ידוע.  
בהתפלגות נורמלית, סטית תקן אחת מעל/ מתחת לממוצע כוללת 34% מהאוכלוסיה בתחום שהיא תוחמת. סטית התקן השניה מכילה כ 13%. כך בהתפלגות נורמלית של כל דבר.

## הרצאה 6 - 11/12/00

### הסתברות

Mutually exclusive, מאורעות חריגים.  
 $P(A \cup B) = P(A) + P(B) - P(A \cap B)$  הסתברות האיחוד של A ו B היא ההסתברות של A ועוד ההסתברות של B פחות החיתוך של ההסתברות של A ו B.  
union - U, איחוד.  
intersection, overlap,  $\cap$  - חיתוך.  
למשל סטודנט הלומד בבר אילן וסטודנט הלומד בטכניון. הסיכוי שהוא יהיה גם פה וגם שם הם אפסיים (להלן מאורע חריג), אין חיתוך בין המיקרים ולכן  $P(A \cap B)$  שווה לאפס.  
Independent events, מאורעות בלתי תלויים.  
2 דגמאות הסתברותיות שאפשריות או לא אפשריות ביחד.  
 $P(A \cap B) = P(A) * P(B)$   
למשל :  
 $P(A) = .12$ , תלמידים בבר אילן.  
 $P(B) = .05$ , אכילת פלאפל.  
 $P(A \cap B) = .006$   
 $P(A \cup B) = P(A) + P(B) - .006 = .12 + .11 - .006 = .164$ , ההסתברות להיות בבר אילן או לאכול פלאפל או לעשות את שניהם.  
Dependent event, מאורע תלוי- יש השפעה של נתון אחד על השני.  
במקרה זה החיתוך אינו ידוע ויש לערוך תצפיות/ מחקר כדי לגלות מהו.  
למשל :  
 $P(A) = .20$  ארועי גשם בשנה.  
 $P(B) = .15$  שימוש במטריה בשנה.  
אם החיתוך היה שווה 11. אז התוצאה היא 24..  
◀ מקסימום החפיפה/ חיתוך האפשרית היא 15. מאחר ולא יכול להיות חיתוך גדול יותר מהערך הקטן מבין ההסתברויות.  
◀ מקסימום האיחוד הוא קרוב ל 35. מאחר ואם היה בדיוק 35 אז הסוג היה אירוע חריג.

**Permutation (סידורים?)**

הפרמוטציה של  $n$  איברים הוא  $n$  פקטוריאלי. (כמו  $n$  עצרת)

$$n! = (n-1)*(n-2)*(n-3)*\dots*1$$

**הרצאה 7 - 18/12/00**

הנוסחה:  $n! / (n-r)!$

למשל סדרה של מספרים: A, B, C. אם ה  $n$  שווה ל 3 וה  $r$  שווה ל 2 אנו מדברים רק על קומבינציה של 2 מספרים. הסדר כאן חשוב! האפשרויות הן: AC, BA, BC, CA, CB, AB. באותה סידרה, ה  $n$  שווה עדיין ל 3 (כמות האברים) וה  $r$  הפעם שווה ל 3. הסדרות הן: ABC, BAC, BCA, CAB, CBA, ACB.

**Combination צירופים**

כאן הסדר לא חשוב.

הנוסחה:  $n! / r! (n - r)!$

פפפ עעע  
פפע עעפ  
פעפ עפע  
פעפ עפפ

אם ה  $r$  וה  $n$  שווים יש רק אפשרות אחת.

בשלוש זריקות של מטבע האפשרויות הן:

כלומר 8 אפשרויות

מה ההסתברות לקבל רק פעמיים פ? -  $3/8$

מה ההסתברות לקבל לפחות פעמיים פ? -  $4/8$ , או  $1/2$

מה ההסתברות לא לקבל בכלל פ? -  $1/8$ .

בינום binomial (של ברנולי).

ההסתברות לקבל  $p$ .

הנוסחה:  $(n r) p^r (1 - P)^{n - r}$

להוציא מלך אחד מחבילה כשכול פעם אני מחזיר את הקלף לחבילה:

$n = 2$  מספר המטלות/ מספר הניסויים

$1 = r$  מספר ההצלחות

$P = 4/54$

$$(2 1) (4/54) (48/52) = 24/169$$

משולש פסקל כאשר  $p = q = 1/2$  עם מטבע.

				1		זריקות
				1		1
		1		2	1	2
	1		3	3	1	4
1	4		6	4	1	5

## הרצאה 11 - 8/1/01

### תיאורית בייס Baye's theory

הסתברות מותנית -  $P(A1 / B1) = P(A1 \cap B1) / P(B1)$

		A3	A2	A1	
		ברגל	ציבורי	פרטי	
B1	זכר	0.2	0.1	0.3	0.6
B2	נקבה	0.1	0.5	0.25	0.4

$$P(\text{פרטי/זכר}) = P(\text{פרטי} \cap \text{זכר}) / P(\text{זכר}) \quad 0.5 = 0.6 / 0.3$$

השלמת נתוני טבלה תלויה בידיעה שהקטגוריות הקיימות הן כול הקטגוריות שיש. כלומר הסה"כ הוא ידוע (או 1) ואז אפשר להשלים את הטבלה.  
בנוסף אפשר לדעת ע"פ הנתונים בטבלה מה הסיכויים של חיתוך נתונים הטבלה להתקיים, ע"פ אינפורמציה חלקית בטבלה.

בינום : מה הממוצע של פלי מ 100 זריקות של מטבע?  $50 =$

$$M = n * p \quad 100 * 0.5 = 50$$

מה ה SD של פלי ב 100 זריקות?

$$S = n * p * q \quad \text{שורש של } 100 * 0.5 * 0.5 = 5$$

שונות בבינום היא אותה נוסחה כמו סטית התקן רק ללא השורש.

לדוגמא : לקבל לפחות 7 פעמים פלי בתוך 10 הטלות מטבע.

$$\text{ממוצע} = n * p = 0.5 * 10 = 5$$

$$\text{סטית התקן} = \text{שורש של } n * p * q = \text{שורש של } 10 * \frac{1}{2} * \frac{1}{2} = 1.58$$

$$\text{בינום} = C(10, 7) * 0.57 * 0.53 + C(10, 8) * 0.58 * 0.52 + C(10, 9) * 0.59 * 0.51 + C(10, 10) * 0.510 * 0.50 = 0.172$$

ככול שה  $n$  שואף לאין סוף, ההתפלגות שואפת להיות נורמלית וכך ניתן להשתמש בהתפלגות זו כדי לקבל את הסיכויים לתוצאה.

הבעיה היא שהנתון 7 הוא בדיד והגרף הוא רציף. לכן ניקח את הנתון המקורב - 6.5

$$z = \frac{x - \bar{x}}{SD} = \frac{6.5 - 5}{1.58} = 0.949$$

$$z = \frac{5.5 - 5}{1.58}$$

כניסה עם ה 2 לגרף נותנת את התוצאה 0.171 שהיא מקורבת למדי בכדי להסתמך עליה כעל

סיכויים. הקירוב יתאפס כשה  $n$  יהיה שווה לאין סוף.

תוצאה מדויקת יכולה להתקבל רק בחישוב בינום.

## 15/1/1 - 12 הרצאה

### מתאם קורלציה

טווח ה  $r$  נע בין 1 ל -1

$$r = \text{CoV}(x, y) / \sqrt{S_x^2 * S_y^2}$$

נוסחת הקורלציה

דוגמא 1:

נתונים סדרת נתונים:

ממוצעים:  $\bar{x} = 9$

$$\bar{y} = 800$$

חישובים ע"פ הנוסחה, המונה:

$$8 - 9 \quad 600 - (-1) * (-$$

$$700 \quad 100)$$

$$9 - 9 \quad 700 - 0$$

$$700$$

$$10 - 9 \quad 800 - (1) * (100$$

$$700 \quad )$$

$$200$$

$$\text{המכנה: } 200 / \sqrt{2 * 20,000}$$

התוצאה היא  $1 = 200 / 200$  כלומר מתאם חיובי מלא.

דוגמא נוספת בה המתאם  $r = .40$ , ממנה נוציא את  $r^2 = .16$  או 16%. נתון זה אומר ש 16%

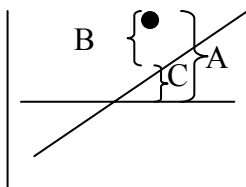
מהשונות מוסבר ע"י המתאם.  $r^2$  הוא השונות המוסברת.

לדוגמא: 4 נבדקים מסודרים על גרף. ביניהם נמתח קו אשר נחשב לטוב ביותר, לכול הנבדקים יש

את המרחק המינימלי מהקו.

הנקודה היא  $y'$  היא הקשר הקיים בין  $X$  ל  $Y$ . היא מתארת הכי טוב את

הנקודות. הקו האופקי הוא ממוצע  $Y$ .





$r^2$  השונות המוסברת של  $X$ ,  $Y$  עם הקשר ביניהם. שונות של  $Y$  המוסברת ע"י  $X$ . לרוב  $Y$  הוא המשתנה התלוי (הביטוי העתיד לבוא).

$B$  הוא המרחק של הנקודה מהקו,  $C$  היא מראה את ניבוי הניסוי. ככול שהנקודה קרובה יותר לקו יש דיוק גדול יותר. אם הכול נמצא על הקו, השונות מוסברת באופן מושלם. כך יתאפשר ניבוי טוב יותר. כשמרחק זה יהיה שווה לאפס, הנקודה היא על הקו.  $A$  מיצג את השונות, התפזרות תוצאות המחקר.

### spearman / Rank correlation

שימוש בנוסחה זו יעשה בסולמות מסוג סדר. למשל תחרות יופי:

מועמדת	שופט 1	שופט 2	d הפרש	d <sup>2</sup>
A	3	2	1	1
B	1	1	0	0
C	2	3	-1	1
D	4	4	0	0
$\Sigma d^2 =$				2

כדי לדעת מה המתאם בין הדירוגים הנ"ל יש להשתמש במתאם ספירמן.

$$N=4$$

$$1 - [6 * (2) / 60] = 1 - 0.2 = 0.8$$

קיים קשר גבוה בין דירוגי השופטים.

### הרצאה 13 - 26/2/01

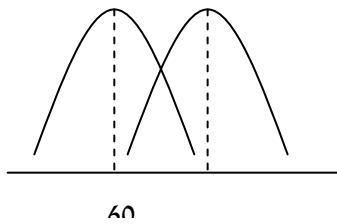
### Standard error of estimate

$$S_{y.x} = S_y \sqrt{1 - r^2}$$

זוהי סטית התקן של המנובאים. אם  $r = 1$  אזי הניבוי הוא מדויק. אם  $r = 0$  אז אין שום ניבוי. זה כאילו  $X$  לא קיים וסטית התקן של הניבוי היא סטית התקן של  $y$ .

### Hypotheses השערות

קיימת תרופה מסויימת שיעילותה מוכחת. אם תצא תרופה חדשה, יעילותה הטובה יותר היא בחזקת השערה. זוהי אמונה שלא נבדקה עדיין בשטח.



דוגמא נוספת היא ציוני פסיכומטרי בת"א מול אלו של ב"א. ההשערה היא שב"א גבוה יותר מזה של ת"א. יותר מ 600 אך לא ידוע בכמה.

בכדי לקבוע שאכן הציון גבוה יותר אנו נוקטים בשיטה של decision theory. מה הקשר בין הדיווח לאמת. דוגמת האדם שיושב מול מכ"ם ורואה ציפור. האפשרויות שעומדות לפניו הן:

	ציפור	מטוס
ציפור	1	3
מטוס	4	2

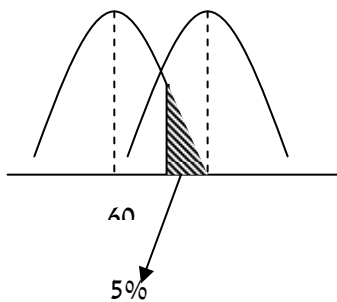
אם דיווח ציפור ובאמת זה ציפור - יצא טוב

אם דיווח מטוס וזה באמת מטוס - יצא טוב

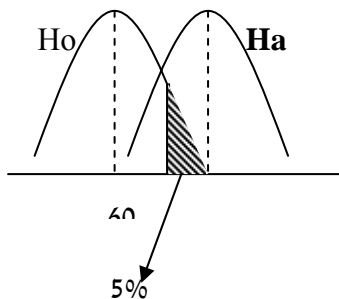
אם דיווח ציפור וזה מטוס - יצא לא טוב

אם דיווח מטוס וזה באמת ציפור - יצא לא טוב.

החוקר את ציוני הפסיכומטרי רוצה להיות בטוח שמה שיצא לו הוא נכון. לכן יש גבול שאם הציון יוצא מעליו הוא בטוח כמעט בודאות. ערך שהוא ערך קריטי. הוא נקבע להיות מעל לשטח 5% של הציון הגבוה של אותם ערכים שצריך להיות מעליהם.



### הרצאה 14 - 5/3/1



גרף  $H_0$  הוא הגרף הראשון שממנו מתחילים את השוואה לגרף  $H_a$ . גרף  $H_0$  הוא גרף המיצג את האוכלוסיה, הממוצע שלו הוא ממוצע של אוכלוסיה,  $\mu = 550$ . גרף  $H_a$  הוא גרף המיצג את התפלגות האוכלוסיה המשוער ע"פ מדגם שנערך ע"י החוקר.  $\hat{X} = 560$ . כדי לקבוע שה"מדגם" גבוה מהאוכלוסיה על ציון הממוצע שלו להיות מעל לערך הקריטי שנקבע.

יתכן שהקביעה ש  $\hat{X}$  היא שגויה ולמעשה הוא נמצא בכול זאת מעבר לערך הקריטי בתוך תחום

ההתפלגות של  $H_0$ , טעות זו נקראת Type one/ $\alpha$  error.

יתכן וקבענו ש  $\hat{X}$  נמצא מצידו השני של הערך הקריטי לכיוון ממוצע  $\mu$  וקביעה זו היא שגיאה והערך האמיתי של  $H_a$  נמצא מהצד האחד של הערך הקריטי, לכיוון ממוצע  $\hat{X}$ , טעות זאת

נקראת Type two/ $\beta$  error.

היחס בין הטעויות הוא הפוך, ככול שגדל הסיכוי לטעות מסוג  $\alpha$  קטן הסיכוי לטעות מסוג  $\beta$ .

הרצאה 15 - 19/3/1רווח בר סמך Confidence interval

טווח של 95% יתואר בצורה הבאה :

$$\hat{X} - 1.96 S_{\hat{x}} \leq \mu \leq \hat{X} + 1.96 S_{\hat{x}}$$

למשל נתונים של מחקר מסויים (השפעת אקמול) הם  $\hat{X} = 40$   $S = 10$   $n = 25$ .

לגבי הקבוצה, לא נשתמש ב  $S$  של האוכלוסיה אלא בשגיאת התקן של הממוצע  $\hat{x} = 10/5 = 2$  כך שלגבי נתוני האוכלוסיה :

$$40 - 1.96 * (2) \leq \mu \leq 40 + 1.96 * (2)$$

התוצאה היא  $36.08 \leq \mu \leq 43.92$  כלומר, 95% מהאוכלוסיה יהיו בטווח התוצאות הללו ע"פ תוצאות המדגם.

ככול ש  $n$  גדול יותר, המדגם מיצג יותר והטווח של  $\mu$  יהיה קטן יותר, אנו מתקרבים יותר מ 2 הקצוות אל הממוצע.

**Power** - הסיכויים/ הסתברות לדחות את  $H_0$  למעשה זה  $1 - \beta$ . אם  $\alpha$  עולה גם ה  $\alpha$  power עולה ותחום  $H_1$  גדל. תוצאה שיצאה בתוך  $\alpha$  נקראת significant result, אם לא היא נקראת not significant result.

T test

$$t = (\hat{X} - \mu) / (S / \sqrt{n})$$

סטית התקן כאן מגיעה ממדגם ולא מאוכלוסיה. כול הנתונים מגיעים מהמדגם. הנתון היחיד שהוא לא ממדגם הוא ה  $\mu$ , זהו נתון קיים.

למשל: מדגם שנתונים הם  $\hat{x} = 520$   $S = 90$   $n = 36$   $\mu = 500$   
 $520 - 500 / (90 / \sqrt{36}) = 20 / 15 = 1.667$

את הנתון הזה משווים לנתון שבטבלה של t test. בעמודת ה df (דרגות חופש) הנתון הוא  $n - 1$ . כך שבמדגם של 36 נבדקים הנתון הוא 35. מאחר ובעמודה שלנו אין 35 נקח את הנתון הנמוך יותר שמופיע והוא 30. בעמודה של 5% הנתון הוא 2.042. היות והנתון שיצא לנו הוא 1.667, הוא נמוך יותר מהערך הקריטי ובמקרה זה נדחה את  $H_0$

אם ה  $n$  שואף ל אין סוף, ה  $t$  יצא 1.96. כמו בהשערה סטטיסטית של אוכלוסיה מאחר ובמקרה זה המדגם כבר שואף לאוכלוסיה.

הרצאה 16 - 26/3/1

במבחן  $t$  סטית התקן היא של המדגם. הגבול כאן הוא ע"פ הטבלה:  $t = (\hat{X} - \mu) / (S / \sqrt{n})$

במבחן  $z$  סטית התקן היא של אוכלוסיה והגבולות כאן הם  $\pm 1.96$  :

$$Z = \frac{\bar{X} - \mu}{\sigma / \sqrt{n}}$$

### מבחן t של שני מדגמים בלתי תלויים

יש לנו 2 מדגמים ואין נתוני בסיס של אוכלוסיה (למשל 2 תרופות חדשות שיצאו לשוק ולא נבדקו באוכלוסיה עדיין).

$$t = \frac{\bar{X}_1 - \bar{X}_2}{S_{\bar{X}_1 - \bar{X}_2}}$$

השערת  $H_0$  במקרה הזה היא  $\mu_1 = \mu_2$  או  $\mu_1 - \mu_2 = 0$ .

השערת  $H_1$  היא  $\mu_1 \neq \mu_2$

אנו משתמשים ב  $\mu$  ולא ב  $\bar{X}$  כי הם מיצגים אוכלוסיה.

דרך חישוב המכנה היא :

$$S^2 = \frac{S_1^2 (N_1 - 1) + S_2^2 (n_2 - 1)}{n_1 + n_2 - 2}$$

נוסחה זו מבטאת את הממוצע של כול המדגמים.

השלב הבא בחישוב המכנה הוא לקחת את ה  $S^2$  שיצא וכדי לקבל את מה שמופיע במכנה לבצע

את החישוב הבא :

$$\sqrt{S^2 / N_1 + S^2 / N_2}$$

שלבם אלו מוצאים  $S^2$  משותף אשר כולל את שני ה S שבמחקר.

את ה t שיצא כאן אני בודק בטבלת מבחן ה t ע"פ דרגות חופש של  $(N_1 - 1) + (N_2 - 1)$ . את הערך

הקריטי שהתקבל אני משווה ל t וממשיך כרגיל עם קבלה או דחייה של  $H_0$ .

קיימת דרך נוספת לחשב את המכנה והיא לעשות ישר  $\sqrt{S^2 / N_1 + S^2 / N_2}$

אבל אנו לא עושים זאת מאחר ואין כאן חלוקה בשונות משותפת. ההנחה היא

שהמדגמים דומים והאוכלוסיות הן שוות (כך גם היא הנחת  $H_0$ ) ולכן אנו נעשה שונות משותפת.

בדוגמה זו השונויות הן נפרדות, שונות אחת מהשניה.

### הרצאה 17 - 2/4/01

דוגמא ל t test של 2 מדגמים לא תלויים :

$$n_2 = 36 \quad n_1 = 2 \quad S_2^2 = 3 \quad S_1^2 = 2 \quad \bar{X}_2 = 7 \quad \bar{X}_1 = 4$$

הנוסחה היא :  $t = \frac{\bar{X}_1 - \bar{X}_2}{S_{\bar{X}_1 - \bar{X}_2}}$  המכנה שווה ל-  $\sqrt{S^2 / N_1 + S^2 / N_2}$

חישוב המונה :  $4 - 7 = -3$

$$S^2 = \frac{(n_1 - 1) * S_1^2 + (n_2 - 1) * S_2^2}{N_1 + N_2 - 2} = \frac{48 + 105}{59} = 2.6$$

$$S_2 = \frac{4 - 7}{\sqrt{2.6/25 + 2.6/36}} = \frac{-3}{0.4} = -7.5$$

יש לבדוק את הערך הקריטי בטבלת t עם דרגת חופשיות מתאימה שהיא  $n_1 + n_2 - 2$  כלומר 59.

הנתון המתאים ביותר בטבלה הוא 40. הערך הקריטי הוא 2.021, נתון שהוא ערך מוחלט. ה t

שיצא הוא תוצאה מובהקת.

סוגי מבחני השערות סטטיסטיות :

Z test

t test של מדגם אחד.

t test של 2 מדגמים בלתי תלויים

t test מסוג before after / match

t test מסוג "r"

מחקר לפני אחרי before after או match

מדגם זה הוא לקבוצה שנדגמה בזמן מסויים בעבר ואח"כ נדגמה שוב כדי לבדוק קשר לשינוי מסויים שהתרחש במהלך הזמן שעבר. אפשרות נוספת היא דגימה של 2 קבוצות זהות מכול הבחינות פרט להבדל אחד, מתוך מטרה לבדוק את השפעת הגורם. למשל עישון.

$$t = D^{\wedge} / S_{D^{\wedge}} \quad \text{הנוסחה היא :}$$

למשל שיפור מבחני בגרות של נבדקים

$D - D^{\wedge}$	D (הפרש B-A)	B	A	
$(1 - (-1))^2 = 4$	+1	8	7	
$(-1 - (-1))^2 = 0$	-1	9	10	4 נבדקים שניגשו פעמיים למבחן בגרות
$(0 - (-1))^2 = 1$	0	8	8	מסויים. B הוא ציון ה "לפני" ו A הוא
$(-4 - (-1))^2 = 9$	-4	6	10	ציון ה"אחרי"
14	הממוצע הוא -1			הסה"כ הוא 14

חישוב המונה  $D^{\wedge}$  : עמודה D כלומר ממוצע ההפרשים שווה -1.

$$S_{D^{\wedge}} = S_D / \sqrt{n} \quad \text{חישוב המכנה :}$$

$$4 = \sqrt{n} \quad \text{היות ויש לנו 4 נבדקים בטבלה.}$$

$$S_D = \sqrt{\sum (D_i - D^{\wedge})^2 / (N - 1)} = \sqrt{(14 / 3)} = 2.16$$

$$S_{D^{\wedge}} = S_D / \sqrt{n} = 2.6 / 2 = 1.08$$

$$t = D^{\wedge} / S_{D^{\wedge}} = -1 / 1.08 = -0.97 \quad \text{חישוב כול הנוסחה :}$$

הערך הקריטי הוא לפי אותה טבלת דרגות חופשיות דו זנבית. חישוב דרגות החופשיות הוא מספר הנבדקים פחות 1, כלומר 3. הערך הקריטי בטבלה הוא 3.182. היות והתוצאה היא -0.97 התוצאה אינה מובהקת ואנו לא דוחים את  $H_0$

הרצאה 18 - 16/4/01

פתרון תרגיל סמסטר ב'.

א. בדיקת ההשערה תעשה ע"י השוואת הממוצעים של 2 הקבוצות

ב. מדובר בהשוואה בין 2 מדגמים לא תלויים אחד בשני. בכדי לבדוק האם ההבדל הוא

מובהק יש לבצע השערה סטטיסטית :

a.  $H_0$  - ממוצע מספר הפעמים שנשים אומרות תודה שווה לממוצע מספר הפעמים שגברים אומרים תודה.

$H_1$  - ממוצע מספר הפעמים שנשים אומרות תודה שונה ממוצע מספר הפעמים שגברים אומרים תודה.

b. הגדרת ה  $\alpha$  ואיזורי הדחיה וקבלה:

$$\alpha = 0.05 \text{ למבחן דו זנבי}$$

$$df = 8 \text{ הגדרת דרגות החופש}$$

$$t_c = 2.306 \text{ הגדרת } t \text{ קריטי}$$

$$-2.306 \leq t_c \leq 2.306 \text{ איזור קבלה}$$

$$t_c > 2.306 \text{ וגם } t_c < -2.306 \text{ איזור דחיה}$$

c. חישוב הסטטיטי

$$t = \frac{\bar{X}_1 - \bar{X}_2}{\sqrt{S^2/N + S^2/N}} \text{ ע"פ הנוסחה}$$

תחילה יש לחשב את השונות המשותפת המופיעה בנוסחה כ  $S^2$ . חישובה הוא

$$S^2 = S_1^2 (N-1) + S_2^2 (N-1) / N_1 + N_2 - 2$$

$$\text{כלומר } 8 = 34 + 30 / 8$$

$$t = 5 - 3 / \sqrt{8/5 + 8/5} = 1.12 \text{ חישוב ערך ה } t$$

d. המסקנה הסטטיסטית היא שהתוצאה אינה מובהקת.  $t < 2.306$

ג. עכשיו מדובר בשני מדגמים שכן תלויים אחד בשני. הנוסחה המתאימה היא

$$t = \frac{\bar{D} - \bar{D}_0}{SD / \sqrt{n}} \text{ חישוב ה } SD \text{ בנוסחה זו הוא ע"פ חישוב סטית תקן של הפרשי}$$

$$SD = \sqrt{\sum (D_i - \bar{D})^2 / N - 1} \text{ המדדים. הנוסחה היא}$$

במקרה זה ההתייחסות היא לא אל כול מדד בפני עצמו אלא אל זוגות של מדדים. לכן

$$df = 4 \text{ כלומר } t_c = 2.776$$

$$SD = 44/4 = \sqrt{11} = 3.32$$

$$t = 2 / 3.32 / \sqrt{5} = 2 / 1.48 = 1.35$$

$$t < 2.776 \text{ התוצאה אינה מובהקת.}$$

ד. במקרה זה ההשוואה היא בין מדגם הגברים האמריקאים לבין אוכלוסית הגברים

הישראליים. כלומר השוואה בין מדגם לאוכלוסיה. לכן גם נתונים סטית התקן של

$$\sigma = 1 \text{ וממוצע האוכלוסיה } \mu = 2 \text{ יש להשתמש במבחן } Z$$

$$Z = \frac{\bar{X} - \mu}{\sigma / \sqrt{n}} \text{ הנוסחה היא}$$

$$3 - 2 / 1 / 2.24 = 2.24 \text{ התוצאה היא}$$

ה  $Z$  הקריטי הוא 1.96 כך שתוצאה זו היא מובהקת.

מבחן r

לבדיקה האם מתאם שהתקבל במדגם הוא מובהק או אם לאו.

$$t = r \sqrt{N - 2} / \sqrt{1 - r^2}$$

ההשערות במבחן זה הן

$$H_0 : \rho = 0$$

$$H_1 : \rho \neq 0$$

למשל במחקר יצאו התוצאות הבאות :  $r = 0.6$  ו  $n = 60$ .

לאחר הצבה פשוטה בנוסחה התוצאה היא 6. את הנתון בודקים בטבלת t לפי דרגות חופשיות של

$n - 2$ . מאחר ו  $t_c$  שווה ל 2 התוצאה במקרה זה היא מובהקת.

הרצאה 19 - 23/4/01הנחיות לביצוע מבחן t

1. על ההתפלגות להיות נורמלית
  2. מספר המבדקים צריך להיות שווה ב 2 הקבוצות  $n_1 = n_2$
  3. Homoskedascity שוויון בין השונותיות  $s_1 = s_2$ . כדי לבדוק זאת ישנו מבחן פשוט מאוד המחלק את השונות הגבוהה יותר בנמוכה יותר. נקרא מבחן F.
- מבחן זה הוא גם סוג של השערה סטטיסטית. השערת ה  $H_0$  כאן היא ש F הוא לא מובהק כלומר נמוך מערך שיופיע בטבלה. אם אכן כך ולא דחינו את  $H_0$  אז  $s_1 = s_2$  וניתן לערוך את המבחן.
- את התוצאה בודקים בתוך טבלה של F כשלכול n נכנסים עם דרגות החופש שלו. כלומר אם ה n של המכנה שווה ל 9, נכנס לטבלה עם המספר 8. התוצאה בטבלה היא של 2 מספרים. המספר העליון יותר מתיחס ל 0.05 והנמוך יותר מתיחס ל 0.01.

מבחן  $\chi^2$  - chi square

$$\chi^2 = \sum [(O_i - E_i)^2 / E_i]$$

לדוגמא : הטלת מטבע 100 פעמים. ה O מציין את המצוי observed וה

E מציין את צפוי/ רצוי expected :

חישוב הנוסחה :

	Oi	Ei
פלי	55	50
עץ	45	50
	100	100

$$\chi^2 = \sum [(O_i - E_i)^2 / E_i] = (55 - 50)^2 / 50 + (45 - 50)^2 / 50 = 1$$

ההשערות במקרה זה הן :

$$H_0 : O_i = E_i$$

$$H_1 : O_i \neq E_i$$

הסכומים  
תמיד יהיו  
שווים

עם הערך 1 נכנסים לטבלה של ערכים קריטיים של  $\chi^2$  כשאת חישוב ה df עושים בדרך שונה מבדרך כלל. כאן מחשבים את מספר הקטגוריות שיש מינוס אחד. כלומר במקרה של מטבע יש 2 קטגוריות ולכן ה df יהיה שווה ל 1.

התוצאה בטבלה היא 3.84 כלומר לא מובהקת, לא דוחים את  $H_0$ . המטבע הוא תקין.

דוגמא נוספת היא של קוביה:			
	$O_i$	$E_i$	$O_i - E_i$
1	10	10	0
2	20	10	$(20-10)^2/10 = 10$
3	5	10	$(5-10)^2/10 = 2.5$
4	5	10	$(5-10)^2/10 = 2.5$
5	10	10	0
6	10	10	0
	60	60	$\chi^2 = 15$
			df = 5

הערך המתקבל מהטבלה הוא 11.07. היות והתוצאה היא 15, הטבלה נמוכה יותר ולכן התוצאה היא מובהקת, דוחים את  $H_0$  והקוביה היא לא מאוזנת.

### מבחן מנדל Mendel

מנדל היה נזיר שערך ניסויים בוטנים באפונים. הוא גילה דרך לחשב בצורה דומה למבחן  $\chi^2$  ללא הנוסחה לעיל.

### הרצאה 20 - 30/4/01

המשך  $\chi^2$  -

נקרא גם goodness of fit, טיב ההתאמה בין מה שיצא לבין מה שצריך היה לצאת.

### מבחן לבדיקת תלות או אי תלות

למשל חלוקה של גברים ונשים במחלקות באוניברסיטה. המבחן יבדוק האם יש קשר בין חלוקת המינים, האם יש תלות ביניהם:

	כלכלה	פסיכולוגיה	מחשבים	סה"כ
גברים	50	40	50	140
נשים	50	80	30	160
	100	120	80	300

הנתונים הללו הם ה Observed. מהו ה expected?  
חישוב הצפוי הוא סה"כ הגבר כפול סה"כ כלכלה חלקי כול האוכלוסיה:

$47 \sim 140 * 100 / 300$ . נתון צפוי זה של גברים לכלכלה יכנס לטבלה בתא הרלוונטי בפינה הימנית עליונה. את הנתון של נשים בכלכלה ניתן לחשב כמו בטבלאות בייס ע"י חיסור של נתון הגברים מסה"כ נתון הכלכלה. כך שהטבלה תראה כך:

דרגות החופש כאן הן 2 - בכול

	כלכלה	פסיכולוגיה	מחשבים	סה"כ
גברים	47 50	56 40	37 50	140
נשים	53 50	64 80	43 30	160
	100	120	80	300

שורה יש 3 נתונים משתנים. מינוס אחד = 2.

בדוגמה זו  $\chi^2 = 18$ , ע"פ הטבלה הערך הקריטי הוא 5.99. ההשערות הסטטיסטיות כאן הן:

$$H_0 : O_i = E_i$$



$$H_1 : O_i \neq E_i$$

כך שהתוצאה היא מובהקת ויש תלות בין 2 המשתנים.

מאפיינים :

1. במבחן זה הנתונים הם מספריים ולא מדדים כמו בשאר המבחנים האחרים. לכן במבחן זה אין ממוצעים או סטיות תקן. מבחן זה לא מסתמך על התפלגות נורמלית. המספרים כאן הם על רצף, לא חלק מהתפלגות.
2. הסולמות כאן הם של הקטגוריות (מחשבים, פסיכולוגיה וכו') - הם שמייים.

### ניתוח שונות Analysis of variance

בדיקת סיבתיות בין נתונים מסויימים, למשל בין ציונים של תלמידים למורה שלהם :  
האם יש קשר בין ציונים התלמידים למוריהם?

ציונים					ממוצע
מורה א'	7	9	7	9	8
מורה ב'	10	9	9	10	9.5
מורה ג'	5	6	5	5	5.5

### הרצאה 21 - 7/5/01

דוגמא :

בדיקה של רמות כעס בקבוצת נבדקים, ב 4 רמות שונות של טמפרטורה - האם טמפרטורה משפיעה על רמת הכעס. התוצאות הן :

	A1	A2	A3	A4
	1	2	4	6
	1	1	3	6
	2	3	2	5
	4	6	3	3
X^=	2	3	3	5

סימון הטבלה יעשה באופן הבא :

	A1	A2	A3	A4
	X <sub>11</sub>			X <sub>41</sub>
	X <sub>21</sub>			X <sub>42</sub>
	X <sub>31</sub>			X <sub>43</sub>
	X <sub>41</sub>			X <sub>44</sub>
X <sup>^</sup> .=	X. <sub>1</sub>	X. <sub>2</sub>	X. <sub>3</sub>	X. <sub>4</sub>
X <sup>^</sup> =	3.25			

החישוב הוא שונות בין הקבוצות חלקי השונות בתוך הקבוצות. אם השונות "בין" גדולה מהשונות "בתוך" המשמעות היא שיש הבדלים ממשיים בין הקבוצות ולכן יש השפעה של הגורם (המשתנה הבלתי תלוי) על הקבוצה.

	Source	SS	df	MS	F
Between	B		$K - 1 = 3$		
Within	W		$K(n - 1)$		
Total	T		$(n - 1) = 15$		

הדבר הראשון שכדי לחשב במבנה זה הוא את דרגות החופש של ה Total. במקרה של ניסוי הטמפרטורה התוצאה היא  $15 - 1 = 16$ .  
נוסחת חישוב ה SS:

$$\sum \sum (X_{ij} - \bar{X}_{..})^2 = \sum \sum (X_{ij} - \bar{X}_{.j})^2 + \sum_{nj} (X_{.j} - \bar{X}_{..})^2$$

### הרצאה 22 - 14/5/01

המשך נושא קודם

לפי החישוב התוצאות הן:  $T = 47$ ,  $B = 19$

דרגות החופש של B הן 4 פחות 1 כלומר 3, df של W הם כמות כול הנתונים שיש פחות כמות הקבוצות, כלומר 12 וה df של T הם כמות כול הנתונים שיש פחות 1, כלומר 15.  
כך שטבלת הנתונים נראית כך:

	Sum square.		Mean square = שונות	
Source	SS	df	MS	F
B	19	3	$19/3=6.33$	$6.33/2.33=2.71$
W	28	12	$28/12=2.33$	
T	47	15	$47/15=3.13$	

חישוב הערך הקריטי הוא ע"פ נתוני ה df של W ושל B, כלומר ע"פ 3 במונה ו 12 במכנה. הערך הקריטי הוא 3.49. מאחר והתוצאה היא 2.71 התוצאה היא לא מובהקת ואין קשר. ניתן לנסח זאת או בניסוח  $P/n$  s (not significant) או לכתוב  $P > 0.05$ .

ההשערות ינוסחו כך:

$$\mu_1 = \mu_2 = \mu_3 = \mu_4 : H_0$$

$$\mu_1 \neq \mu_2 \neq \mu_3 \neq \mu_4 : H_1$$

### הנוסחה של Sheffe

$$(x_i - x_j)^2 / (S_w^2/n_1) + (S_w^2/n_2)$$

יעשה בה שימוש רק אם התוצאה היא מובהקת. בעזרת נוסחה זו ניתן לאתר בתוך הקבוצות איזה קבוצה היא זו המשפיעה על המובהקות והיטתה את התוצאה לכיוון זה. ניתן להשוות את הקבוצות, אחת לשניה כמה פעמים שרוצים, הדבר לא ישנה את אחוז המובהקות. אם למשל

היתה יוצאת תוצאה מובהקת בדוגמה הנ"ל, החישוב היה נראה כך : השוואה בין קבוצה A1 לבין קבוצה A4 -

$$5 - 2 / 2.33/4 + 2.33/4 = 9/1.2 = 7.72$$

את הערך הקריטי (3.49) יש להכפיל בדרגות החופש של B, כלומר ב 3. התוצאה היא 10.5, זהו הערך הקריטי המתקן שאליו משוים את התוצאה (7.72). במקרה זה היא לא מובהקת, כלומר יש סיכוי גדול יותר שהמובהקות נמצאת בין הקבוצות האחרות. כך שכדי להתחיל להשוות בקצוות.

## הרצאה 23 - 4/6/01

### 2 way anova ניתוח שונות דו כיווני

בטבלת הנתונים הבאה יש 2 משתנים בלתי תלויים ומשתנה אחד תלוי. המשתנים הבלתי תלויים הם הקטגוריות. קטגוריה אחת תהיה בטורים, היא תסומן ב A ותחולק (במקרה זה) ל 3 קבוצות (למשל תרופה עם 3 סוגי מינונים שונים). קטגוריה שניה תהיה בשורות, תסומן ב B ובמקרה זה תחולק ל 2 קבוצות (למשל עיתוי לקיחת התרופה). הנתונים שבתוך התאים מיצגים את הנבדקים באותה קבוצה, במקרה זה יש 2 נבדקים בכל תא. אלו הם המשתנים התלויים :

מינוני תרופות

		A1	A2	A3	
עיתוי	B1	2,4	3,5	1,3	$X_{.1.}=3$
לקיחת	B2	1,1	5,7	7,9	$X_{.2.}=5$
תרופה			$X_{..2}=5$		$X_{...}=4$

ע"פ הטבלה ניתן לבדוק :

הבדלים בין שורות (בין זמני התרופות)

הבדלים בין הטורים (הבדלים בין מינוני התרופות)

אינטראקציה - הבדלים בין התאים בטבלה (בין כול קבוצה וקבוצה).

הסימון :  $X_{irc}$

i - אינדיקסידואל, במקרה זה יכול להיות 1 או 2 מאחר ויש רק 2 בכול תא.

r - שורה, במקרה זה יכול להיות 1 או 2.

c - טור, יכול להיות 1, 2 או 3.

טבלת החישוב :

	source	SS	df	MS	F
between	row	12	1	12	12/1.7
	column	24	2	12	12/1.7
	interaction	32	2	16	16/1.7
	Within	10	6	1.7	

total	78	11
-------	----	----

נוסחת החישוב היא :

$$\sum \sum (X_{irc} - \hat{X}_{...})^2 = nC\sum (X_{.r} - \hat{X}_{...})^2 + nR\sum (X_{.c} - \hat{X}_{...})^2 + n\sum \sum (X_{rc} - \hat{X}_{.r} - \hat{X}_{.c} + \hat{X}_{...})^2 + \sum \sum \sum (X_{irc} - \hat{X}_{rc})^2$$

ה  $n$  הקטן מסמן את מספר הנבדקים בתא. הסימן  $X_{.r}$  מסמל ממוצע שורה, כך הסימן  $X_{.c}$  הוא ממוצע טור,  $X_{rc}$  הוא ממוצע תא ו  $X_{irc}$  הוא כול נתון בכול תא ותא.  $X_{...}$  הוא ממוצע הממוצעים.

ההשערות הן 3 קבוצות של השערות לכול שונות ושונות בטבלה :

שונות שורות :

$$\mu_{.1} = \mu_{.2} : H0$$

$$\mu_{.1} \neq \mu_{.2} : H1$$

שונות טורים :

$$\mu_{..1} = \mu_{..2} = \mu_{..3} : H0$$

$$\mu_{..1} \neq \mu_{..2} \neq \mu_{..3} : H1$$

שונות אינטראקציה :

$$\mu_{.11} = \mu_{.12} = \mu_{.13} : H0 \text{ וכו' } \dots$$

$$\mu_{.11} \neq \mu_{.12} \neq \mu_{.13} : H1 \text{ וכו' } \dots$$