

# דף נוסחאות / סיכום

## מבוא לסטטיסטיקה קורס 094480

### יש להדפיס במדפסת לייזר ברזולוציה של (לפחות) DPI 600

דף עזר מספר 1:

#### תורת הסתברות

##### הגדרות:

- ניסוי מקרי - הוא ניסוי שאין לדעת בוודאות מראש את תוצאתו אלא קיים מרחב של תוצאות אפשריות.
- מרחב מדידים - קבוצה חסומה - אוסף כל התוצאות האפשריות של הניסוי המקרי.
- מאורע פשוט,  $\omega$  - תוצאה בודדת של מרחב המדידים (איבר ב- $\Omega$ ).
- מאורע  $A, B, C, \dots$  - אוסף של מאורעות פשוטים (קבוצה חלקית ל- $\Omega$ ).
- קבוצת חלקית (מובלת) -  $A \subset B$  - תת-קבוצה של  $B$  (מוכלת ב- $B$ ) אם כל איבר של  $A$  הוא איבר של  $B$ .
- אומרים שמאורע  $A$  מוביל את  $A$  (הוא תוצאת הניסוי שהיא  $\omega$  מקודמת המרחב) שייכת ל- $A$ .
- מאורע משלים למאורע  $A$  -  $A^c$  - כל האיברים ב- $\Omega$  שאינם שייכים ל- $A$ .
- מאורע ריק  $\emptyset$  (קבוצה ריקה) - אינו כולל אף איבר.

##### מעולות בין מאורעות:

- איחוד -  $A \cup B$  - כל האיברים שהם או  $A$  או  $B$  או בשתי הקבוצות.
- חיתוך -  $A \cap B$  - כל האיברים שהם גם  $A$  וגם  $B$ .

##### אם החיתוך בין הקבוצות הוא ריק, בלומר אין להן איברים משותפים, אזי הקבוצות נקראות זרות:

- מאורע  $A$  והמאורע המשלים לו  $A^c$  הם מאורעות זרים:  $A \cap A^c = \emptyset$
- האיחוד בין מאורעות  $A$  ו- $A^c$  הוא מרחב המדידים:  $A \cup A^c = \Omega$

- הפרש -  $A \setminus B$  - כל האיברים שהם  $A$  אך לא  $B$ :  $A \setminus B = A \cap B^c$

##### כללי מעולות בין קבוצות:

חוק החילוף:  $A \cup B = B \cup A$   
 $A \cap B = B \cap A$

חוק הפילוג:  $A \cap (B \cup C) = (A \cap B) \cup (A \cap C)$   
 $A \cup (B \cap C) = (A \cup B) \cap (A \cup C)$

##### חוקי דה-מורג:

$$(A \cap B)^c = A^c \cup B^c$$
$$(A \cup B)^c = A^c \cap B^c$$

#### מבוא להסתברות וסטטיסטיקה

##### דף עזר 3

##### מונקצית התפלגות ברנולי עם פרמטר $p$

מבצעים ניסוי ברנולי עם הסתברות  $p$  להצלחה.  $X$  מקבל 1 אם הניסוי הצליח ו-0 אחרת. אזי:

$$X \sim \text{Ber}(p) \quad p_X(x) = \begin{cases} p & x=1 \\ q & x=0 \\ 0 & \text{otherwise} \end{cases}$$

##### מונקצית התפלגות גאומטרית עם פרמטר $p$

$X$  מייצג מספר הניסויים עד להצלחה הראשונה. סדרת ניסויי ברנולי בלתי תלויים עם הסתברות  $p$  להצלחה, אזי:

$$X \sim \text{Geo}(p) \quad p_X(x) = \begin{cases} q^{x-1}p & x=1,2,3,\dots \\ 0 & \text{otherwise} \end{cases}$$

אם  $X \sim \text{Geo}(p)$  (גאומטרי) אז התוחלת שלו היא:  $E(X) = \frac{1}{p}$

השונות שלו היא:  $\text{Var}(X) = \frac{q}{p^2}$

##### תכונות חוסר יזכרון, התפלגות גאומטרית

אם ידוע כי לא הייתה הצלחה עד לניסוי ה- $n$  אזי: ההסתברות שהצלחה הראשונה תתרחש בניסוי ה- $n+k$  היא:  $q^{k-1}p$

$$\text{Var}(X) = E(X^2) - E^2(X)$$

ניסוח העבודה:

##### תכונות של שונות:

$$\text{Var}(X) \geq 0$$

$$\text{Var}(a) = 0$$

$$\text{Var}(a \cdot X) = a^2 \cdot \text{Var}(X)$$

$$\text{Var}(a + X) = \text{Var}(X)$$

$$\text{Var}(X + Y) = \text{Var}(X) + \text{Var}(Y)$$

$$\text{Var}(X - Y) = \text{Var}(X) + \text{Var}(Y)$$

$$\text{Var}(X + Y) = \text{Var}(X) + \text{Var}(Y) + 2\text{Cov}(X, Y)$$

$$\text{Var}(X - Y) = \text{Var}(X) + \text{Var}(Y) - 2\text{Cov}(X, Y)$$

##### סיימת תכונ

$$\sigma(X) = \sqrt{\text{Var}(X)}$$

#### מבוא להסתברות וסטטיסטיקה

##### דף עזר 4

נניח כי יש לנו מדידים  $X_1, \dots, X_n$ .

1. חצי - הערך שמשתנה המקריים קטנים ממנו או שווים לו:

$$\text{עבור } n=1: \text{MED} = X_{(\frac{n+1}{2})}$$

$$\text{עבור } n=2k: \text{MED} = \frac{X_{(k)} + X_{(k+1)}}{2}$$

2. רבעון ראשון-תחתון - הערך שרבע מהמקריים קטנים ממנו או שווים לו:

$$Q_1 = X_{(\lfloor \frac{n}{4} \rfloor)}$$

3. רבעון שני - חצינו:

$$Q_2 = X_{(\lfloor \frac{n}{2} \rfloor)}$$

4. רבעון שלישי-עליון - הערך שלשלוש רבעים מהמקריים קטנים ממנו או שווים לו:

$$Q_3 = X_{(\lfloor \frac{3n}{4} \rfloor)}$$

5. תחום בין-רבעוני - בו מרוכזות 50% מהתצפיות המרכזיות:  $IQR = Q_3 - Q_1$ .

10. השברון ה- $p$  - הערך שכרומרייה  $p$  מהתצפיות קטנות ממנו או שווים לו:

$$X_p = X_{(\lfloor np \rfloor)}$$

##### משתנה מקרי רציף:

##### מונקצית צפיפות

משתנה מקרי רציף מאופיין ע"י מונקצית צפיפות  $f_X(x)$  המקיימת:

$$f_X(x) \geq 0$$

$$\int_{-\infty}^{\infty} f_X(t) dt = 1$$

##### עבור משתנה מקרי רציף

$$P(X = x) = 0$$

$$P(a \leq X \leq b) = \int_a^b f_X(x) dx$$

##### מונקצית התפלגות מצטברת

יהי  $X$  מ"מ. פונקצ. התפלגות מצטברת של  $X$  היא מונקציה  $F_X: R \rightarrow [0,1]$  מוגדרת ע"י

$$F_X(x) = P(X \leq x)$$

תכונות חוסר יזכרון של משתנה מקרי גאומטרי:

#### מבוא לסטטיסטיקה והסתברות

##### דף עזר 2

##### משתנים מקריים

**משתנה מקרי** - נתון מרחב ההסתברות. הפונקציה  $X = x(w)$  המתאימה לכל נקודה  $w$  במרחב המדידים  $\Omega$  נקראת מ"מ (משתנה מקרי)

**משתנה מקרי בדיד** - מ"מ  $X$  יקרא מ"מ בדיד אם הוא מקבל סדרה סופית או בת מניה של ערכים.

##### מונקציות הסתברות עבור מ"מ בדיד:

המונקציה  $P_X(x) = P(X = x)$  נקראת מונקצית הסתברות עבור מ"מ  $X$ .

$$\sum_x P_X(x) = 1$$

##### מונקציות התפלגות מצטברת עבור מ"מ בדיד:

$$F_X(x) = \sum_{x_i \leq x} P_X(x_i) = P(X \leq x)$$

(הסתברות) שמ"מ  $X$  קטן או שווה למספר  $x$ , נקראת מונקצית התפלגות מצטברת של מ"מ  $X$  בנק  $x$ .

##### תכונות מונקציות התפלגות:

$$0 \leq F_X(x) \leq 1 \quad \forall x \in \mathcal{R}$$

2. הפונקציה  $F_X(x)$  מונוטונית לא יורדת ב- $x$ , כלומר:  $x_1 < x_2$  אז

$$F_X(x_1) \leq F_X(x_2)$$

$$\lim_{x_0 \rightarrow -\infty} F_X(x_0) = 0, \quad \lim_{x_0 \rightarrow \infty} F_X(x_0) = 1$$

**הערה:** אם נתונה  $F_X(x)$  ניתן לחשב ממנה את  $F_X(x)$ .

אם נתונה  $F_X(x)$ :  $F_X(x) = F_X(x) - F_X(x-1)$ , כלומר, מונקצית ההסתברות שווה למונקצית התפלגות פחות מונקציות התפלגות בנקודה שלפני.

##### מונקציות התפלגות היפרגאומטרית עם פרמטרים $N, R, n$

נתונה אוכלוסיה בת  $N$  איברים שמתוכם  $R$  מיוחדים. בוחרים מהאוכלוסיה  $n$  איברים ללא החזרה.  $X$  - מ"מ המונה את מספר האיברים המיוחדים שנבחרו.

$$P_X(x) = \begin{cases} \frac{\binom{R}{x} \binom{N-R}{n-x}}{\binom{N}{n}} & 0 \leq x \leq n \\ 0 & \text{otherwise} \end{cases}$$

והמסגרים:  $X \sim HG(N, R, n)$

##### מונקציות התפלגות פואסונית עם פרמטר $\lambda$

$X$  - מ"מ פואסוני מונה אירועים לאורך זמן:

$$P_X(x) = \begin{cases} \frac{e^{-\lambda} \lambda^x}{x!} & x = 0, 1, 2, \dots \\ 0 & \text{otherwise} \end{cases}$$

אם  $X \sim \text{Pois}(\lambda)$  (פואסוני) אז התוחלת שלו היא:  $E(X) = \lambda$

השונות שלו היא:  $\text{Var}(X) = \lambda$

##### תכונות יחידה $\text{Unif}(a, b)$ עבור מ"מ בדיד

ישנם  $N$  ערכים  $b-1, b-2, \dots, a+1, a$  שזוהו מקבלים מ"מ בדיד. ההסתברות לקבל ערך כלשהו מתוך ה- $n$  ערכים האלה היא

$$P_X(x) = \begin{cases} \frac{1}{N} & x = a, a+1, a+2, \dots, b-1, b \\ 0 & \text{otherwise} \end{cases}$$

אם  $X \sim \text{Unif}(a, b)$  (אחד) אז התוחלת ושונות שלו היא:

$$E(X) = \frac{(a+b)}{2}$$
$$\text{Var}(X) = \frac{(b-a)(b-a+1)}{12}$$

במילים אחרות זה מרחב המדידים שווה להסתברות

מקרה פרטי כאשר  $X$  מקבל ערכים  $\{1, 2, \dots, N\}$ .

$$P_X(x) = \begin{cases} \frac{1}{N} & x = 1, 2, \dots, N \\ 0 & \text{otherwise} \end{cases}$$

$$E(X) = \frac{(N+1)}{2}$$
$$\text{Var}(X) = \frac{(N^2-1)}{12}$$

$$P(X > t + s | X > t) = P(X > s) = e^{-\lambda s}$$

כלומר, אם נתון כי ידע רגע  $t$  (כולל לא אחיז מופע, אזי ייתחילים מחדשה), ולכן זמן עד המופע הבא הוא  $\exp(\lambda)$ .

##### התפלגות נורמלית:

$$X \sim N(\mu, \sigma^2)$$

$$f_X(x) = \frac{1}{\sqrt{2\pi} \cdot \sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}} \quad -\infty < x < \infty$$

$$E(X) = \mu$$

$$\text{Var}(X) = \sigma^2$$

מקרה פרטי של התפלגות נורמלית - מ"מ נורמלי סטנדרטי:

$$Z \sim N(0,1) \quad \text{מ"מ מתפלג נורמלית עם תוחלת } \mu = 0 \text{ ושונות } \sigma^2 = 1$$

##### חישוב הסתברות בהתפלגות נורמלית (תיקונון משתנה נורמלי - מעבר להתפלגות נורמלית סטנדרטית):

$$Z = \frac{X - \mu}{\sigma} \sim N(0,1) \quad \text{אזי } X \sim N(\mu, \sigma^2)$$

##### חישוב הסתברות בהתפלגות נורמלית סטנדרטית:

$$P(Z < z) = \Phi(z)$$

מאחר שההסתברות סימטרית סביב ה-0:

$$P(Z > z) = 1 - P(Z < z) = 1 - \Phi(z) = \Phi(-z) = P(Z < -z)$$

**סימון:**

$z_\alpha$  - ערך של  $Z$  שמסתאוולו שטח (הסתברות) של  $\alpha$

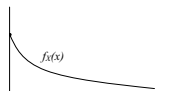
**הערה:**

מתקיים בגלל הסימטריה של התפלגות נורמלית סטנדרטית  $z_\alpha = -z_{1-\alpha}$

**חישוב ערך  $\alpha$  כשער  $\alpha$  ידוע:**  $X \sim N(\mu, \sigma^2)$

$$Z = \frac{X - \mu}{\sigma} \sim N(0,1)$$

$$X_\alpha = \mu + z_\alpha \cdot \sigma$$



סימון:  $X \sim \text{exp}(\lambda)$  (לוקר ייקבובי של התפלגות)

$$f_X(x) = \begin{cases} \lambda e^{-\lambda x} & x \geq 0 \\ 0 & \text{otherwise} \end{cases}$$

$$E(X) = \frac{1}{\lambda} \quad \text{Var}(X) = \frac{1}{\lambda^2}$$

$$P(X > a) = \int_a^\infty f_X(x) dx = \int_a^\infty \lambda e^{-\lambda x} dx = e^{-\lambda a}$$

תכונות חוסר יזכרון של משתנה מקרי גאומטרי:

דף עור 7

אמידה נקודתית

אובלסטיה - איסק פריטיס עליהם מעוניינים במידע מסוים.
דוגם - יוסף חקיק של פריטיס מהאוכלוסייה.
פרמטר - מספר קבוע המאפיין את האוכלוסייה (theta).
סטטיסטיק - ערך הנקבע לחישוב מתוכם נתונים המדגם.
אמד - סטטיסטי המשמש לאמידת פרמטר (theta-hat).
אומדן - הערך המספר של האמד.

אמד theta-hat ייקרא אומד חסר הטעה למפרט theta, אם קיים:
E(theta-hat) = theta

אחרת, הוא אומד הטעה והטיות ייתנה עיי:
Bias(theta-hat) = E(theta-hat) - theta

עדיף אומד חסר הטעה על אומד מוטע.

שגיאה ריבועית ממוצעת (MSE):

MSE(theta-hat) = E((theta-hat - theta)^2) = Var(theta-hat) + [Bias(theta-hat)]^2

עדיף הוא האמד בעל MSE מינימלית.

הערה: אם theta-hat הוא אומד חסר הטעה למפרט theta, קיים:
MSE(theta-hat) = Var(theta-hat) = Bias(theta-hat) = 0

סטיית התקן של הממוצע כאשר השונות של X ידועה:

SE\_X-bar = sigma\_Y / sqrt(n)

סטיית התקן של הממוצע כאשר השונות של X לא ידועה:

SE\_X-bar = S / sqrt(n)

S^2 = sum\_{i=1}^n (x\_i - X-bar)^2 / (n-1) = sum\_{i=1}^n x\_i^2 - nX-bar^2 / (n-1)

רוח סמן להפחש תוחלת כאשר המדגמים מוונים:

f\_{alpha/2}^{(n-1)} <= d-delta <= f\_{1-alpha/2}^{(n-1)}

כאשר:

s\_d = sqrt(1/(n-1) \* sum\_{i=1}^n (d\_i - d-bar)^2) SE(d) = s\_d / sqrt(n)

מכאן:

delta in [d-bar +/- t\_{1-alpha/2}^{(n-1)} SE(d-bar)]

רוח סמן להפחש פרמטריזציה:

באותו אופן כמו בהפחש תוחלות:

p\_1 - p\_2 in [p-hat\_1 - p-hat\_2 +/- z\_{1-alpha/2} \* sqrt(p-hat\_1(1-p-hat\_1)/n\_1 + p-hat\_2(1-p-hat\_2)/n\_2)]

מתרון תרגיל כיתה 8

רוח סמן להפחש:

S^2(n-1) = sum\_{i=1}^n (X\_i - X-bar)^2 / sigma^2 ~ chi^2(n-1)

S^2(n-1) / sigma^2 <= S^2(n-1) / sigma^2 <= S^2(n-1) / sigma^2

רוח סמן להפחש תוחלת במדגמים בלתי-תלויים:

X-bar ~ N(mu, sigma^2/n), Y-bar ~ N(mu, sigma^2/n)

1. כאשר שונות של שתי האוכלוסיות ידועה:

X-bar - Y-bar ~ N(mu\_1 - mu\_2, sigma^2/n\_1 + sigma^2/n\_2) ~ N(0,1)

ולכן:

mu\_1 - mu\_2 in [X-bar - Y-bar +/- z\_{1-alpha/2} \* sqrt(sigma^2/n\_1 + sigma^2/n\_2)]

2. כאשר שונות לא ידועה:

נבנה לשתי האוכלוסיות אותה שונות. את השונות נאמד על-ידי:

S\_p^2 = ((n\_1 - 1)S\_1^2 + (n\_2 - 1)S\_2^2) / (n\_1 + n\_2 - 2)

כאשר:

S^2 = sum\_{i=1}^n (x\_i - X-bar)^2 / (n-1), S\_p^2 = sum\_{i=1}^n (y\_i - Y-bar)^2 / (n\_p - 1)

מתקיים:

(X-bar - Y-bar) - (mu\_1 - mu\_2) ~ t\_{(n\_1+n\_2-2)}

ולכן:

mu\_1 - mu\_2 in [X-bar - Y-bar +/- t\_{1-alpha/2}^{(n\_1+n\_2-2)} \* S\_p^2 \* sqrt(1/n\_1 + 1/n\_2)]

- לכל סוג של השערה אלטרנטיבית נקבע אזור דחייה שאם הממוצע המדגמי נופל בו, המסקנה תהיה לדחות את H\_0.
האזור המשלים נקרא אזור הקבלה.
אם הממוצע המדגמי נופל בו המסקנה תהיה לא לדחות את H\_0.

- עבור השערה חד-צדדית ימנית נא דחה עבור ערכים גדולים של הממוצע.
עבור השערה חד-צדדית שמאלית נא דחה עבור ערכים קטנים של הממוצע.
עבור השערה דו-צדדית נא דחה עבור ערכים קטנים או גדולים של הממוצע.

ערכים קריטיים לאלטרנטיבות שונות:

מפתח השערות:

אזור דחייה:

Z = (X-bar - mu\_0) / (sigma\_0 / sqrt(n)) > z\_{1-alpha} H\_0: mu = mu\_0 H\_1: mu > mu\_0

Z = (X-bar - mu\_0) / (sigma\_0 / sqrt(n)) < z\_alpha H\_0: mu = mu\_0 H\_1: mu < mu\_0

|Z| = |(X-bar - mu\_0) / (sigma\_0 / sqrt(n))| > z\_{1-alpha/2} H\_0: mu = mu\_0 H\_1: mu != mu\_0

מבחן השערות על תוחלת mu כאשר השונות איננה ידועה:

T = (X-bar - mu\_0) / (S / sqrt(n)) > t\_{1-alpha}^{n-1} H\_0: mu = mu\_0 H\_1: mu > mu\_0

T = (X-bar - mu\_0) / (S / sqrt(n)) < t\_alpha^{n-1} H\_0: mu = mu\_0 H\_1: mu < mu\_0

|T| = |(X-bar - mu\_0) / (S / sqrt(n))| > t\_{1-alpha/2}^{n-1} H\_0: mu = mu\_0 H\_1: mu != mu\_0

דף עור 7

אמידה מרווחית

רמת הסמן (רמת הביטחון) 1-alpha, כאשר alpha - רמת המבחנות

רוח בר-סמן לתוחלת mu של אוכלוסייה כאשר שונות sigma^2 ידועה:

X-bar ~ N(mu, sigma^2/n) => Z = (X-bar - mu) / (sigma / sqrt(n)) ~ N(0,1)

ולכן,

mu in [X-bar +/- z\_{1-alpha/2} \* sigma / sqrt(n)]

זהו רווח בר-סמן לתוחלת mu ברמת ביטחון 1-alpha.

רוח בר-סמן לתוחלת mu של אוכלוסייה כאשר שונות sigma^2 איננה ידועה:

נאמד את השונות בעזרת סטטיסטי S^2 המוחשב מהמדגם:

S^2 = sum\_{i=1}^n (x\_i - X-bar)^2 / (n-1) = sum\_{i=1}^n x\_i^2 - nX-bar^2 / (n-1)

מתקיים:

T = (X-bar - mu) / (S / sqrt(n)) ~ t\_{(n-1)}

כאשר (n-1) - מספר דרגות חופש - פרמטר בהתפלגות t\_{(n-1)} (סטודנט) i. הוא נגזר המדגם שעליו מבוסס הממוצע.

ולכן:

mu in [X-bar +/- t\_{1-alpha/2} \* S / sqrt(n)]

זהו רווח בר-סמן לתוחלת mu ברמת ביטחון 1-alpha.

רוח סמן ל-p - פרובינציה באוכלוסייה:

p-hat = x/n

הפרופורציה במדגם x - הנו מספר בעלי התכונה במדגם נגזר n.

עבור n-1 מספיק גדולים מתקיים בקירוב:

p-hat ~ N(p, p(1-p)/n) => (p-hat - p) / sqrt(p(1-p)/n) ~ N(0,1)

ולכן:

p in [p-hat +/- z\_{1-alpha/2} \* sqrt(p-hat(1-p-hat)/n)]

דף עור 9

בדיקת השערות:

המטרה: רוצים להחליט על נכונותה או אי-נכונותה של השערה מסוימת לנבי פרמטרים של התפלגות.

ההחלטה לא לירות או לא לדחות את השערה נעשית על סמן המדגם, ובעזרת כלל הכרעה קובעו כלל תוצאת צריכות מדגם אפשרית לירות או לא לירות את השערה.

השערות:

- השערת האפס H\_0 - השערה אותה מעוניינים לבדוק.
השערה האלטרנטיבית H\_1 - השערה המנוגדת.
האפשרויות צריכות להיות ממצות ומוצאות.

סוגי השערות:

- השערה משוטה - מגדירה באופן חד ערכי את ההתפלגות.
השערה מורכבת - כוללת הרבה אפשרויות.
H\_0: mu = 800
H\_1: mu < 800
השערה מורכבת-חד-צדדית (דו-צדדית):
H\_0: mu = 800
H\_1: mu != 800
השערה מורכבת-דו-צדדית (דו-צדדית):
H\_0: mu = 800
H\_1: mu != 800

החלטות האפשריות:

- דחיית H\_0 - זוהי החלטה האומרת שהשערת האפס איננה נכונה.
אי דחיית H\_0 (קבלה, H\_1) - החלטה האומרת שהשערת האפס איננה ניתנת וזוהי נכונה.

שגיאות:

alpha - הסתברות לטעות מסוג ראשון או רמת מובהקות
alpha = P(reject the null hypothesis | the null hypothesis is correct)
beta - הסתברות לטעות מסוג שני
beta = P(don't reject the null hypothesis | the alternative hypothesis is correct)
pi = 1 - beta - עוצמת המבחן

pi = P\_{mu\_1} (Reject H\_0) = P\_{mu\_1} (accept H\_0) alpha = P\_{mu\_0} (Reject H\_0)

דף עור 11

מבחני טיב התאמה:

המטרה לבדוק האם המדגם שייך להתפלגות מסוימת.
H\_0: X ~ F
H\_1: otherwise

סטטיסטי המבחן עבור מבחני טיב התאמה הינו:

chi^2\_p = sum\_{i=1}^k (n\_i - np\_0)^2 / (np\_0) = sum\_{i=1}^k (O\_i - E\_i)^2 / E\_i

- n\_i - מספר התצפיות בקטגוריה i
k - מספר הקטגוריות

אזור דחייה ברמת מובהקות alpha

דחה את H\_0 אם chi^2\_p > chi^2\_{1-alpha}(k-1)

- 1 - מספר הפרמטרים שאומדו

בדיקת אי תלות:

לקח מדגם מסוים מתוך אוכלוסייה המסווגת לפי שני משתנים איכותיים (קטגוריאליים). המטרה היא לבדוק האם שני המשתנים חיים בצלם תלויים.

A - משתנה קטגוריאלי ראשון (שורות)
B - משתנה קטגוריאלי שני (עמודות)

H\_0: P(A\_i, B\_j) = P(A\_i) \* P(B\_j)

H\_1: otherwise

סטטיסטי המבחן עבור מבחני אי תלות:
chi^2\_r = sum\_{i=1}^k sum\_{j=1}^l (n\_{ij} - np\_{ij})^2 / (np\_{ij}) = sum\_{i=1}^k sum\_{j=1}^l (O\_{ij} - E\_{ij})^2 / E\_{ij}

p\_i = n\_i / n

p\_j = n\_j / n

- n\_i - מספר התצפיות בעלות ערך i במשתנה A
n\_j - מספר התצפיות בעלות ערך j במשתנה B

אזור דחייה ברמת מובהקות alpha

דחה את H\_0 אם chi^2\_r > chi^2\_{1-alpha}(r-1)(c-1)

- r - מספר הקטגוריות במשתנה A
c - מספר הקטגוריות במשתנה B

דף עור 6

אמידה נקודתית

אובלסטיה - איסק פריטיס עליהם מעוניינים במידע מסוים.
דוגם - יוסף חקיק של פריטיס מהאוכלוסייה.
פרמטר - מספר קבוע המאפיין את האוכלוסייה (theta).
סטטיסטיק - ערך הנקבע לחישוב מתוכם נתונים המדגם.
אמד - סטטיסטי המשמש לאמידת פרמטר (theta-hat).
אומדן - הערך המספר של האמד.

אמד theta-hat ייקרא אומד חסר הטעה למפרט theta, אם קיים:
E(theta-hat) = theta

אחרת, הוא אומד הטעה והטיות ייתנה עיי:
Bias(theta-hat) = E(theta-hat) - theta

עדיף אומד חסר הטעה על אומד מוטע.

שגיאה ריבועית ממוצעת (MSE):

MSE(theta-hat) = E((theta-hat - theta)^2) = Var(theta-hat) + [Bias(theta-hat)]^2

עדיף הוא האמד בעל MSE מינימלית.

הערה: אם theta-hat הוא אומד חסר הטעה למפרט theta, קיים:
MSE(theta-hat) = Var(theta-hat) = Bias(theta-hat) = 0

סטיית התקן של הממוצע כאשר השונות של X ידועה:

SE\_X-bar = sigma\_Y / sqrt(n)

סטיית התקן של הממוצע כאשר השונות של X לא ידועה:

SE\_X-bar = S / sqrt(n)

S^2 = sum\_{i=1}^n (x\_i - X-bar)^2 / (n-1) = sum\_{i=1}^n x\_i^2 - nX-bar^2 / (n-1)

רוח סמן להפחש תוחלת כאשר המדגמים מוונים:

f\_{alpha/2}^{(n-1)} <= d-delta <= f\_{1-alpha/2}^{(n-1)}

כאשר:

s\_d = sqrt(1/(n-1) \* sum\_{i=1}^n (d\_i - d-bar)^2) SE(d) = s\_d / sqrt(n)

מכאן:

delta in [d-bar +/- t\_{1-alpha/2}^{(n-1)} SE(d-bar)]

רוח סמן להפחש פרמטריזציה:

באותו אופן כמו בהפחש תוחלות:

p\_1 - p\_2 in [p-hat\_1 - p-hat\_2 +/- z\_{1-alpha/2} \* sqrt(p-hat\_1(1-p-hat\_1)/n\_1 + p-hat\_2(1-p-hat\_2)/n\_2)]

- $-1 \leq r \leq 1$
  - $r = 0$ : אין קשר ליניארי בין  $x$  ל- $y$ .
  - $0 < r \leq 1$ : קשר ליניארי חיובי בין  $x$  ל- $y$ . (כש  $x$  עולה,  $y$  עולה, ושלח).
  - $-1 \leq r < 0$ : קשר ליניארי שלילי בין  $x$  ל- $y$ . (כש  $x$  עולה,  $y$  יורד).
- $F^2$  - פרופורצית הישגות המוסברות על-ידי  $X$  (מקדם טרמינטיבי)
- $$0 \leq r^2 \leq 1$$

$Y$  - משתנה תלוי / משתנה מוסבר.

$X$  - משתנה בלתי-תלוי / משתנה מסביר.

לעמיתים יש מספר משתנים מסבירים (מן משתנים מסבירים) ואם הם סטוכסטיים:  $X_1, X_2, \dots, X_p$

**המטרה** רוצים להסביר את  $Y$  בעזרת המשתנה/ים המוסבר/ים ולתת תחזית לציון  $Y$  על סמך ציון  $X$ .

**שליבים:**

- ציון דיאגרמת פיזור: מדיאגרמת הפיזור נלמד האם קיים קשר ליניארי בין  $X$  ל- $Y$ .
- מציאת קו הרגרסיה:  $Y = aX + b$  כאשר  $a$  - שיפוע הקו ו- $b$  - תוחך הקו
- חיזוי  $Y$  על-פי  $X$  ו-על-פי קו הרגרסיה.

**סימונים:**

$$S_{XY} = \sum (X_i - \bar{X})(Y_i - \bar{Y}) = \sum XY - n\bar{X}\bar{Y}$$

$$S_{XX} = S_X^2(n-1) = \sum (X_i - \bar{X})^2 = \sum (X_i^2) - n(\bar{X})^2$$

$$S_{YY} = S_Y^2(n-1) = \sum (Y_i - \bar{Y})^2 = \sum (Y_i^2) - n(\bar{Y})^2$$

**שיפוע הקו:**

$$a = \frac{S_{XY}}{S_{XX}} = r \sqrt{\frac{S_{YY}}{S_{XX}}}$$

**התחנה:**

$$b = \bar{Y} - a\bar{X}$$

**שונות הרגרסיה**

$$\hat{\sigma}^2 = \frac{\sum (Y_i - (aX_i + b))^2}{n-2} = \frac{S_{YY} - a^2 S_{XX}}{n-2}$$

**מקדם המתאם:**

$$r = \frac{S_{XY}}{\sqrt{S_{XX}S_{YY}}} = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\left(\sum_{i=1}^n (x_i - \bar{x})^2\right)\left(\sum_{i=1}^n (y_i - \bar{y})^2\right)}} = \frac{\sum_{i=1}^n (x_i \cdot y_i) - n \cdot \bar{x} \cdot \bar{y}}{\sqrt{\left(\sum_{i=1}^n x_i^2 - n\bar{x}^2\right)\left(\sum_{i=1}^n y_i^2 - n\bar{y}^2\right)}}$$

**אמידת קו רגרסיה:**

$$S_{XY} = \sum (X_i - \bar{X})(Y_i - \bar{Y}) = \sum XY - n\bar{X}\bar{Y}$$

$$S_{XX} = S_X^2(n-1) = \sum (X_i - \bar{X})^2 = \sum (X_i^2) - n(\bar{X})^2$$

$$S_{YY} = S_Y^2(n-1) = \sum (Y_i - \bar{Y})^2 = \sum (Y_i^2) - n(\bar{Y})^2$$

$$\hat{\alpha} \sim N\left(\alpha, \frac{\sigma^2 \sum X_i^2}{nS_{XX}}\right)$$

$$\hat{\beta} \sim N\left(\beta, \frac{\sigma^2}{S_{XX}}\right)$$

$$\hat{\beta} = \frac{S_{XY}}{S_{XX}} = r \sqrt{\frac{S_{YY}}{S_{XX}}} = r \sqrt{\frac{(n-1)S_Y^2}{(n-1)S_X^2}} = r \frac{S_Y}{S_X}$$

$$\hat{\alpha} = \bar{Y} - \hat{\beta}\bar{X}$$

$$\hat{\sigma}^2 = \frac{\sum e_i^2}{n-2} = \frac{\sum (Y_i - (\hat{\alpha} + \hat{\beta}X_i))^2}{n-2} = \frac{S_{YY} - \hat{\beta}^2 S_{XX}}{n-2} = \frac{SS_E}{n-2} = MS_E$$

**ANOVA**

$$SS_T = SS_B + SS_E = S_{YY} = \hat{\beta}^2 S_{XX} + \sum e_i^2 = S_Y^2(n-1) = \sum (Y_i - \bar{Y})^2$$

$$SS_B = SS_{reg} = \hat{\beta}^2 S_{XX} = S_{YY} r^2$$

$$SS_E = \sum e_i^2 = S_{YY} - \hat{\beta}^2 S_{XX} = S_{YY}(1 - r^2)$$

$$MS_E = \frac{SS_E}{n-2} = \frac{\sum e_i^2}{n-2} = \hat{\sigma}^2$$

$$MS_E = \frac{SS_B}{df_B}$$

$$\frac{SS_B}{SS_T} + \frac{SS_E}{SS_T} = 1 = r^2 + (1 - r^2) = \frac{\hat{\beta}^2 S_{XX}}{S_{YY}} + \frac{\sum e_i^2}{S_{YY}}$$

$$F = \frac{MS_B}{MS_E} = \frac{\hat{\beta}^2 S_{XX}}{\hat{\sigma}^2} = T^2 = \left[ \frac{\hat{\beta}}{SE(\hat{\beta})} \right]^2 = \frac{r^2(n-2)}{1-r^2}$$

**סכום של משתנים מקריים**

השונות של סכומם של  $m$  משתנים  $Y$  ו- $X$  נתונה ע"י:

$$Var(X+Y) = Var(X) + Var(Y) + 2Cov(X, Y)$$

כאשר  $Cov(X, Y)$  היא השונות המשותפת של שני המשתנים:

$$Cov(X, Y) = E[(X - E(X))(Y - E(Y))] = E[(X - Y) - E(X) - E(Y)]$$

$$E[(X - Y) - E(X) - E(Y)] = \sum_{i,j} \sum_{k,l} (ij - kl) P_{ij,kl}(x, y)$$

תכונות השונות המשותפת:

- $-\infty \leq Cov(X, Y) \leq +\infty$
- $Cov(X, Y) > 0 \Leftrightarrow$  קיים מתאם חיובי בין שני המשתנים
- $Cov(X, Y) = 0 \Leftrightarrow$  המשתנים בלתי מתואמים
- $Cov(aX + b, cY + d) = a \cdot c \cdot Cov(X, Y) \Leftrightarrow$
- $Cov(X, X) = Var(X)$

-  $Var(X+Y) = Var(X) + Var(Y)$  אם  $X$  ו- $Y$  בלתי מתואמים אוי

-  $X$  ו- $Y$  בלתי-תלויים  $\Leftrightarrow Cov(X, Y) = 0$  (אם  $X$  ו- $Y$  בלתי מתואמים, הכיוון ההפוך לא תמיד נכון)

**מקדם המתאם:**

$$\rho(X, Y) = \frac{Cov(X, Y)}{\sqrt{Var(X)Var(Y)}}$$

תכונות מקדם המתאם:

- $-1 \leq \rho(X, Y) \leq 1$
- $\rho(X, Y) = 0 \Leftrightarrow$  המשתנים בלתי מתואמים
- $\rho(X, Y) = 1 \Leftrightarrow Y = aX + b, a > 0$
- $\rho(X, Y) = -1 \Leftrightarrow Y = aX + b, a < 0$

	SS	df	MS	F
רגרסיה	$SS_B$	1	$MS_B = SS_B/df_B$	$MS_B/MS_E = F_{(1, n-2)}$
שארית	$SS_E$	n-2	$MS_E = SS_E/df_E = \hat{\sigma}^2$	
סה"כ	$SS_T = SS_B + SS_E = S_{YY}$	n-1		

$$SS_B + SS_E = SS_T = S_{YY} = S_Y^2(n-1) = \sum (Y_i - \bar{Y})^2 = \sum (Y_i - \hat{Y}_i + \hat{Y}_i - \bar{Y})^2$$

$$SS_B = \sum e_i^2 = \sum (Y_i - \hat{Y}_i)^2 = S_{YY}(1 - R^2)$$

$$SS_E = SS_T - SS_B = S_{YY} - R^2 S_{YY} = a^2 S_{XX}$$

$$R^2 = \frac{SS_B}{SS_T} = \frac{SS_B}{SS_T}$$

**קשר בין מקדם הרגרסיה לשיפוע הקו:**

$$a = r \sqrt{\frac{S_{YY}}{S_{XX}}} = r \cdot \frac{s_y}{s_x} = r \cdot \frac{\sqrt{\sum_{i=1}^n (y_i - \bar{y})^2}}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2}}$$

**ביקת השערות לשיפוע הרגרסיה**

$$E(a) = \beta$$

$$Var(a) = \sigma_a^2 = \frac{\sigma^2}{S_{XX}}$$

$$\hat{Var}(a) = SE^2(a) = \hat{\sigma}_a^2 = \frac{\hat{\sigma}^2}{S_{XX}}$$

$$H_0: \beta = 0$$

$$H_1: \beta \neq 0$$

$$T = \frac{a - \hat{\beta}}{\hat{\sigma}_a} = \frac{a}{SE(a)} = \frac{a}{\hat{\sigma}_a} \sqrt{S_{XX}} \sim t_{n-2}$$

$$T^2 = \left[ \frac{a}{SE(a)} \right]^2 = \left[ \frac{a}{\hat{\sigma}_a} \right]^2 S_{XX} \sim F_{1, n-2}$$

$$a - t_{n-2}^{1-\alpha/2} \hat{\sigma}_a \leq \beta \leq a + t_{n-2}^{1-\alpha/2} \hat{\sigma}_a$$

רוח סמך עבור  $a$ :

	SS	df	MS	F
רגרסיה	$SS_B$	1	$MS_B$	$MS_B/MS_E$
שארית	$SS_E$	n-2	$MS_E$	
סה"כ	$SS_T$	n-1		

**מקדם המתאם - r**

$$r^2 = \frac{\hat{\beta}^2 S_{XX}}{S_{YY}} = \frac{SS_B}{SS_T} = \frac{S_{XY}^2}{S_{XX}S_{YY}}$$

$$r = \frac{S_{XY}}{\sqrt{S_{XX}S_{YY}}}$$

$$\rho = \frac{Cov(X, Y)}{\sqrt{Var(X)Var(Y)}}$$

**ביקת השערות**

$$H_0: \beta = 0 \quad H_1: \beta \neq 0$$

$$T = \frac{\hat{\beta} - \beta}{\hat{\sigma}_\beta} = \frac{\hat{\beta}}{SE(\hat{\beta})} = \frac{\hat{\beta}}{\hat{\sigma}_\beta} \sqrt{S_{XX}} \sim t_{n-2}^{1-\alpha/2}$$

$$T = \sqrt{\frac{r^2}{1-r^2}}(n-2)$$

$$T^2 = \left[ \frac{\hat{\beta}}{SE(\hat{\beta})} \right]^2 = \left[ \frac{\hat{\beta}}{\hat{\sigma}_\beta} \right]^2 S_{XX} = \frac{MS_B}{MS_E} \sim F_{1, n-2}^{1-\alpha}$$

**רוח סמך לשיפוע**

$$\hat{\beta} - t_{n-2}^{1-\alpha/2} \hat{\sigma}_\beta \leq \beta \leq \hat{\beta} + t_{n-2}^{1-\alpha/2} \hat{\sigma}_\beta$$

**התפלגויות משותפות**

יהיו  $X$  ו- $Y$  שני משתנים מקריים המוגדרים על אותו מרחב הסתברות  $\Omega: R^2 \rightarrow (X, Y)$

**פונקציית ההסתברות המשותפת** של הזוג  $(X, Y)$  מוגדרת ע"י:

$$P_{XY}(x, y) = P(X=x, Y=y)$$

אזי פונקציית ההסתברות השולית של מ"מ  $X$  ע"י:

$$P_X(x) = \sum_y P_{XY}(x, y)$$

בדומה, פונקציית ההסתברות השולית של מ"מ  $Y$  ע"י:

$$P_Y(y) = \sum_x P_{XY}(x, y)$$

**פונקציית ההסתברות המותנת** של  $X$  בהינתן  $Y=y$  מוגדרת ע"י:

$$P_{X|Y}(x|y) = \frac{P_{XY}(x, y)}{P_Y(y)} = \frac{P(X=x, Y=y)}{P(Y=y)} = P(X=x|Y=y)$$

**פונקציית ההסתברות המותנת** של  $Y$  בהינתן  $X=x$  מוגדרת ע"י:

$$P_{Y|X}(y|x) = \frac{P_{XY}(x, y)}{P_X(x)} = \frac{P(X=x, Y=y)}{P(X=x)} = P(Y=y|X=x)$$

**אי תלות**

שני משתנים מקריים נקראים **בלתי-תלויים** אם לכל זוג נקודות  $(x, y)$ :

$$P_{XY}(x, y) = P_X(x)P_Y(y)$$

אם  $X$  ו- $Y$  בלתי תלויים אזי:

$$P_{XY}(x|y) = \frac{P_{XY}(x, y)}{P_Y(y)} = \frac{P(X=x)P(Y=y)}{P(Y=y)} = P(X=x)$$

$$P_{Y|X}(y|x) = \frac{P_{XY}(x, y)}{P_X(x)} = \frac{P(X=x)P(Y=y)}{P(X=x)} = P(Y=y)$$

## נוצר ע"י אסף נתן מבחן אחרון לסיום התואר ! 8.3.2007

## בהצלחה לכולם, טוב לדעת שיש להנדסת חשמל סוף.