

תיקון רציפות

- ✓ כאשר משתמשים במשפט הגבול המרכזי עבור אוכלוסיה המקיימת התפלגות **בדידה** (עוברים מ- m ל- m בדיד למ- m רציף) עבור $n \leq 100$ ✓

$$P(X \leq k) = \varphi\left(\frac{k+0.5-\mu}{\sigma}\right)$$

$$P(X < k) = P(X \leq k-1) = \varphi\left(\frac{k-0.5-\mu}{\sigma}\right)$$

$$P(X \geq k) = 1 - P(X < k) = 1 - \varphi\left(\frac{k-0.5-\mu}{\sigma}\right)$$

$$P(X > k) = 1 - P(X \leq k) = 1 - \varphi\left(\frac{k+0.5-\mu}{\sigma}\right)$$

$$P(X = k) = \varphi\left(\frac{k+0.5-\mu}{\sigma}\right) - \varphi\left(\frac{k-0.5-\mu}{\sigma}\right)$$

$$P(a \leq X \leq b) = \varphi\left(\frac{b+0.5-\mu}{\sigma}\right) - \varphi\left(\frac{a-0.5-\mu}{\sigma}\right)$$

- לממוצע לא עושים תיקון רציפות, לסכום כן!

קירוב נורמלי להתפלגות הבינומית

יהי X_i מ"מ ברנולי (מ"מ המייצג תוצאה של ניסוי בודד). הערך 1 מייצג "הצלחה" והערך 0 "כישלון".

X_i	1	0
$P(X_i)$	p	$1-p \equiv q$

$$E(X_i) = p$$

$$V(X_i) = pq$$

יהי Y מ"מ בינומי, סכום של n משתני ברנולי

$$Y \sim B(n, p) = \sum_{i=1}^n X_i$$

לפי מג"מ:

$$Y \sim B(n, p) = \sum_{i=1}^n X_i \sim N(np, npq)$$

$$\bar{X} = \frac{\sum_{i=1}^n x_i}{n} = \frac{X}{n} = \hat{p} \sim N\left(p, \frac{pq}{n}\right)$$

- \bar{X} הוא גם פרופורציית (אחוז) ההצלחות שהתקבלו ב- n הניסויים, ותוחלתו p ($\sum_{i=1}^n x_i$ זהו מספר ההצלחות ב- n ניסויי ברנולי)
- ע"מ להשתמש בקירוב זה, יש לדרוש את התנאי: $np \geq 10$ $nq \geq 10$
- ההתפלגות הבינומית היא בדידה. לכל n קטן מ-100 נעשה תיקון רציפות.

פונקציית ההתפלגות של המדגם

יהי X_1, \dots, X_n מדגם מקרי בת"ש מתוך אוכלוסיה בעלת התפלגות $f(x)$. **התפלגות המדגם** (=פונקציית ההתפלגות המשותפת) היא:

$$\prod_{i=1}^n f(x_i) / \prod_{i=1}^n P(X = x_i)$$

ההסתברות של הריאליזציה, ההסתברות לקבל צירוף מסוים של x_i כתוצאה מהמדגם.

התפלגות המדגם עבור מדגם מקרי בגודל n שנלקח מהאוכל' הבאות:

א. ברנולי

$$p(x) = p^x (1-p)^{1-x}, \mu = p, \sigma^2 = p * q$$

$$\prod p^x (1-p)^{1-x} = p^{\sum x_i} (1-p)^{n-\sum x_i}$$

ב. אחידה

$$\prod \frac{1}{b-a} = \left(\frac{1}{b-a}\right)^n$$

$$f(x) = \frac{1}{b-a}$$

ג. פואסוני

$$\prod \frac{e^{-\lambda} \lambda^{x_i}}{x_i!} = \frac{e^{-n\lambda} \lambda^{\sum x_i}}{\prod x_i!}$$

$$p(y) = \frac{e^{-\lambda} \lambda^y}{y!}$$

ד. נורמלית

$$f(x) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

$$\prod \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x_i-\mu)^2}{2\sigma^2}} = (2\pi\sigma^2)^{-\frac{n}{2}} * e^{-\frac{\sum (x_i-\mu)^2}{2\sigma^2}}$$

ה. גיאומטרית

$$p(x) = p(1-p)^{x-1}$$

$$\prod p(1-p)^{x_i-1} = p^n (1-p)^{\sum x_i - n}$$

ו. אקספוננציאלית

$$\prod \lambda e^{-\lambda x_i} = \lambda^n e^{-\lambda \sum x_i}$$

$$f(x) = \lambda e^{-\lambda x}$$

אמידה נקודתית ($\hat{\theta}$ - האומד הנקודתי ל- θ)

אמידה בשיטת המומנטים:

$$E(X^k) = \mu_k$$

$$\frac{\sum_{i=1}^n x_i^k}{n} = m_k$$

השיטה מתבססת על השוואת μ_k ל- m_k .

המשוואות המתקבלות עבור המומנט הראשון והשני:

$$E(X) = \bar{X}$$

$$V(X) = \frac{\sum_{i=1}^n x_i^2}{n} - (\bar{X})^2, E(X^2) = \frac{\sum_{i=1}^n x_i^2}{n}$$

$$\left[\frac{\sum_{i=1}^n (x_i - \bar{X})^2}{n} = \frac{\sum_{i=1}^n x_i^2}{n} - \bar{X}^2, E(X^2) = \sigma^2 + \mu^2 \right]$$

התפלגות אחידה בשיטת המומנטים $x \sim u(a, b)$

$$\hat{b} = \bar{X} + \sqrt{3\left(\frac{\sum_{i=1}^n x_i^2}{n} - (\bar{X})^2\right)} \quad \hat{a} = \bar{X} - \sqrt{3\left(\frac{\sum_{i=1}^n x_i^2}{n} - (\bar{X})^2\right)}$$

התפלגות המשתנה המקרי

$$F(x) \equiv \int_{-\infty}^x f(X=x) dx \equiv P(X < x)$$

תוחלת $E(X)$

$$\mu = \frac{\sum x_i}{N}$$

$$E(X) = \sum x_i \cdot P(x_i) = \mu$$

$$E(X) = \int_{-\infty}^{\infty} x f(x) dx$$

$$\sigma = \sqrt{\frac{\sum (x_i - \mu)^2}{N}}$$

תכונות:

$$E(a) = a$$

$$E(a + bX) = a + bE(X)$$

$$E(X^2 - X) = E(X^2) - E(X), E(\sum_{i=1}^n X_i) = \sum_{i=1}^n E(X_i)$$

$$E(X - \mu) = E(X) - \mu = 0$$

$$E(X^2) = V(X) + [E(X)]^2$$

$$E(X_1 \cdot X_2) = E(X_1) \cdot E(X_2)$$

$$E(aX + b) = aE(X) + b$$

$$E(G(X)) = \sum G(x_i) \cdot P(x_i)$$

$$E(G(X)) = \int_{-\infty}^{\infty} G(x) \cdot f(x) dx$$

שונות $V(X)$

$$\sigma^2 = V(X) = E((X - \mu)^2) = E(X^2) - \mu^2 = E(X^2) - [E(X)]^2$$

סטיות תקן

$$\sigma = \sqrt{V(X)} = \sqrt{\sigma^2}$$

תכונות של שונות וסטיות תקן:

$$V(X) = E(X - \mu)^2 = E(X^2) - \mu^2$$

$$\sigma_X \geq 0, V(X) \geq 0$$

$$\sigma_a = 0, V(a) = 0$$

$$V(bX) = b^2 V(X)$$

$$V(a + bX) = b^2 V(X)$$

$$\sigma_{a+bX} = |b| \sigma_X$$

$$V(\sum_{i=1}^n X_i) = \sum_{i=1}^n V(X_i)$$

התפלגות המקסימום $Y = \max X_i$

$$F_Y(Y) = P(Y \leq y) = P(\max X_i \leq y) = P(x_1 \leq y \cap \dots \cap x_n \leq y)$$

$$\prod_{i=1}^n P(x_i \leq y) = \prod_{i=1}^n F_{x_i}(y) = [F_X(y)]^n$$

$$f_Y(y) = \frac{F_Y(y)}{dy} = n \cdot [F_X(y)]^{n-1} \cdot f_X(y)$$

התפלגות המינימום $Y = \min X_i$

$$F_Y(Y) = P(Y \leq y) = P(\min X_i \leq y) = 1 - P(\min X_i > y) =$$

$$= 1 - P(x_1 > y, \dots, x_n > y) = 1 - \prod_{i=1}^n P(x_i > y) =$$

$$= 1 - \prod_{i=1}^n (1 - F_{x_i}(y)) = 1 - [1 - F_X(y)]^n$$

$$f_Y(y) = n \cdot [1 - F_X(y)]^{n-1} \cdot f_X(y)$$

סטטיסטיקה תיאורית - מושגים חשובים

פרמטר - מדד מתוך התפלגות האוכלוסיה. גודל קבוע, לא ידוע!

סטטיסטי - מדד המתקבל מתוך המדגם (פונ' של התצפיות)

$\hat{\theta}$ - אומד לפרמטר, הינו סטטיסטי. הערך המספרי של $\hat{\theta}$ נקרא **אומדן**.

- כשמבקשים למצוא **אומדן** לפרמטר, מוצאים את האומד ומציבים את תוצאות המדגם!

סטטיסטי ממוצע המדגם $\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$

$$\sigma_{\bar{x}} = \frac{\sigma}{\sqrt{n}}, V(\bar{x}) = \frac{\sigma^2}{n}, E(\bar{x}) = \mu$$

(1) **סכום הסטיות** של התצפיות מהממוצע הינו אפס

$$\sum_{i=1}^n (x_i - \bar{x}) = 0$$

(2) **סכום ריבועי הסטיות** של התצפיות מהממוצע $\sum_{i=1}^n (x_i - \bar{x})^2$ אינו אפס, אך הוא הקטן ביותר מכל מדי המיקום. (כאשר מציבים את הממוצע זה הערך הכי נמוך שניתן לקבל מסכום ריבועי הסטיות)

התפלגות המדגם המקרי

$$f_{X_1, X_2, \dots, X_n}(x_1, x_2, \dots, x_n) = \prod_{i=1}^n f_X(x_i)$$

משפט הגבול המרכזי

$$E(X_i) = \mu$$

$$V(X_i) = \sigma^2$$

אזי, עבור $n > 30$:

$$\sum_{i=1}^n X_i \sim N(n\mu, n\sigma^2)$$

$$\bar{X} \sim N\left(\mu, \frac{\sigma^2}{n}\right)$$

הערה: אם X_1, \dots, X_n מתפלגים נורמלית אז המשפט מתקיים לכל n .

אמידה בשיטת הנראות המקסימלית:

פונ' הנראות מתארת את ההסתברות לקבלת ריאליזציה מסוימת במדגם כתלות בפרמטר הלא ידוע θ . נשים לב שזו בדיוק **התפלגות המדגם**.

אנ"מ $\hat{\theta}$ הוא הממקסם את פונ' הנראות: $\hat{\theta} = \operatorname{argmax}_{\theta} [L(\theta)]$
שיטת העבודה:

- נרשום את פונ' הנראות כפונקציה של θ
במקרה הבדיד: $L(\theta) = \prod_{i=1}^n P_{\theta}(X_i = x_i)$
במקרה הרציף: $L(\theta) = \prod_{i=1}^n f_{\theta}(x_i)$
- בד"כ נוציא \ln מהפונקציה.
- נגזור לפי θ ונמצא מקסימום (נציין כי היינו גוזרים פעם שנייה **ובדקים שזהו המקסימום**). אם יש יותר מאומדן אחד, נגזור את פונ' הנראות לפי כל אומדן.

אנ"מים מוכרים

- אחידה רציפה:** כאשר $U(0, \theta)$ אז $\hat{\theta} = \max(X_i)$
כאשר $U(a, b)$ אז $\hat{a} = \min(X_i)$, $\hat{b} = \max(X_i)$
- פואסונית:** $\hat{\lambda} = \bar{X}$
- בינומית:** אם $X \sim B(n, p)$ האנ"מ להסתברות p הוא $\hat{p} = \frac{X}{n}$
- מעריכית:** $\hat{\lambda} = \frac{1}{\bar{X}}$
- נורמלית:** כש- μ לא ידוע: $\hat{\mu} = \frac{\sum X_i}{n} = \bar{X}$, $\hat{\sigma}^2 = \frac{\sum (X_i - \bar{X})^2}{n}$
כש- μ ידוע: $\hat{\sigma}^2 = \frac{\sum (X_i - \mu)^2}{n}$

חוקי ln

$$\ln\left(\frac{x}{y}\right) = \ln x - \ln y, \quad \ln(xy) = \ln x + \ln y, \quad \ln(x^k) = k \ln x$$

$$\frac{d}{dx} \ln(x) = \frac{1}{x}, \quad \ln\left(\prod_{i=1}^n x_i\right) = \sum_{i=1}^n \ln x_i$$

אומד חסר הטייה לשונות האוכלוסייה על סמך מדגם בגודל n

$$s^2 = \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n-1} = \frac{\sum_{i=1}^n x_i^2 - n \cdot \bar{x}^2}{n-1}$$

- על מנת לאמוד את סטיית התקן באוכלוסייה (σ) משתמשים באומד s (מוציאים שורש).
- s אינו חסר-הטייה לסטיית התקן.

משפט: פונקציה של אומד נראות מקסימלית

נתונה פונ' $y = g(\theta)$ ומבקשים למצוא אנ"מ ל- y .

אם $\hat{\theta}$ הוא אנ"מ ל- θ , אז האנ"מ ל- y הינו: $\hat{y} = g(\hat{\theta})$.

לכן, אם $X \sim B(n, p)$ **אנ"מ להסתברות p הוא $\hat{p} = \frac{X}{n}$**

ממוצע ריבועי השגיאות (MSE)

$$MSE = E\left[\left(\hat{\theta} - \theta\right)^2\right] = V(\hat{\theta}) + \left[E(\hat{\theta}) - \theta\right]^2$$

- $V(\hat{\theta})$ השונות של האומד $\hat{\theta}$
- האיבר $\left[E(\hat{\theta}) - \theta\right]$ מכונה ההטיה ($Bias$) של האומד $\hat{\theta}$
- בד"כ נושאו בין אומדים ע"פ ערך ה-MSE שלהם, ברוב המקרים חוסר הטיה משפר MSE.**

תכונות אומדים

א. **חוסר הטיה:** $\hat{\theta}$ יקרא אומד חסר הטיה כאשר מתקיים: $E(\hat{\theta}) = \theta$

$$MSE \equiv V(\hat{\theta}) \leftarrow E(\hat{\theta}) - \theta = 0 \quad \leftarrow \text{ההטיה שווה ל-0}$$

* על מנת שיתקיים חוסר הטיה נדרש להוכיח $E(\hat{\theta}) = \theta$

* **ראינו בכיתה שלכל n ולכל התפלגות מתקיים: $E(\bar{X}) = E(X_i)$**

* ולכן, ממוצע המדגם \bar{X} הוא א.ח.ה לתוחלת ההתפלגות $E(\bar{X}) = \mu$

* s^2 א.ח.ה לשונות ההתפלגות $E(s^2) = \sigma^2$

* **טענה:** עבור פונ' ליניארית, אם $\hat{\theta}$ א.ח.ה ל- θ ו-

$$G(\theta) = a + b \cdot \theta, \quad \text{אזי } a + b \cdot \hat{\theta} \text{ א.ח.ה ל-} G(\theta)$$

* כשמבקשים אח"ה "על סמך n תצפיות" זה רומז ללכת לממוצע.

א. **עקיבות:** $\hat{\theta}$ יקרא אומד עקיב אם הוא אומד חסר הטיה וכן

$$\lim_{n \rightarrow \infty} V(\hat{\theta}) = 0 \quad (\text{תנאי מספיק ולא הכרחי})$$

* כדי לבדוק עקיבות: 1. נבדוק חוסר הטיה. 2. נבדוק את שונות האומד חוסר הטיה לעומת עקיבות

אין קשר בין 2 התכונות!

עקיבות מתייחסת לדיוק המדגם כתלות בגודל המדגם, כלומר מה קורה לאומדן כשדוגמים פרטים רבים באותו מדגם. **הטייה** מתייחסת לתוחלת האומד, כלומר מה קורה לאומדן כשמבצעים מדגמים רבים.

א. **יעילות:** עבור שני אומדים **חסרי הטיה** $\hat{\theta}_1$ ו- $\hat{\theta}_2$: אם מתקיים

$$V(\hat{\theta}_1) < V(\hat{\theta}_2) \quad \text{נאמר ש-} \hat{\theta}_1 \text{ יעיל יותר מ-} \hat{\theta}_2.$$

* אח"ה יעיל יותר \leftrightarrow MSE קטן יותר

התפלגויות שימושיות**התפלגות נורמלית**

אם $X \sim N(\mu, \sigma^2)$, אז $Y = X + \text{const}$ מתפלג $Y \sim N(\mu + \text{const}, \sigma^2)$
אם $X \sim N(\mu, \sigma^2)$, אז המ"מ $Y = cX$ (קבוע) מתפלג $Y \sim N(c\mu, c^2\sigma^2)$
אם $X \sim N(\mu, \sigma^2)$, אז המ"מ $Y = \frac{X - \mu}{\sigma}$ מתפלג $Y \sim N(0, 1)$
סכום מ"מ נורמליים הוא מ"מ נורמלי עם סכום התוחלות וסכום השונות.
הפרש מ"מ נורמליים הוא מ"מ נורמלי עם הפרש התוחלות וסכום השונות.

התפלגות חי-בריבוע χ^2 (לא סימטרית)

התפלגות סכום הריבועים של מ"מ נורמליים סטנדרטיים.

אם $X_i \sim N(\mu, \sigma_i^2)$ בלתי תלויים, אז $\sum_{i=1}^n \left(\frac{X_i - \mu_i}{\sigma_i}\right)^2 \sim \chi^2(n)$
אם $X_i \sim \chi^2(k_i)$ בלתי תלויים, אז $\sum_{i=1}^n X_i \sim \chi^2(\sum_{i=1}^n k_i)$

התפלגות t (סימטרית)

אם $Z \sim N(0, 1)$ ו- $U \sim \chi^2(k)$ בלתי תלויים אז $\frac{Z}{\sqrt{U/k}} \sim t(k)$

התפלגות של מספר סטטיסטיים חשובים

נתונים מ"מ ב"ת ושוי התפלגות $x_i \sim N(\mu, \sigma^2)$ $i = 1, \dots, n$ אז:

$$\sum_{i=1}^n \left(\frac{X_i - \mu}{\sigma}\right)^2 \sim \chi^2(n)$$

אם התוחלת לא ידועה:

$$\sum_{i=1}^n \left(\frac{X_i - \bar{X}}{\sigma}\right)^2 \sim \chi^2_{(n-1)}$$

אם σ לא ידועה: $\frac{\bar{X} - \mu}{\sigma/\sqrt{n}} \sim t_{(n-1)}$ לעומת $\frac{\bar{X} - \mu}{\sigma/\sqrt{n}} \sim N(0, 1)$

תפלגות F

אם $U \sim \chi^2_m$ ו- $V \sim \chi^2_n$ ב"ת, אזי $\frac{U/m}{V/n} \sim F(m, n)$

אם $X \sim F(m, n)$, אזי $\frac{1}{X} \sim F(n, m)$

רב"ס

$P(T_1 < \theta < T_2) = 1 - \alpha$	(T_1, T_2)	רב"ס דו"צ
$P(T_3 < \theta) = 1 - \alpha$	(T_3, ∞)	רב"ס ח"צ תחתון
$P(T_4 > \theta) = 1 - \alpha$	$(-\infty, T_4)$	רב"ס ח"צ עליון

רב"ס עבור תוחלת μ על סמך מדגם של n תצפיות

הנחות ומידע	רב"ס דו צדדי	רב"ס חד צדדי תחתון	רב"ס חד צדדי עליון
σ^2 ידוע. התפלגות נורמלית של האוכלוסייה או $n > 30$	$\mu \in \left(\bar{X} - Z_{1-\frac{\alpha}{2}} \frac{\sigma}{\sqrt{n}}, \bar{X} + Z_{1-\frac{\alpha}{2}} \frac{\sigma}{\sqrt{n}}\right)$	$\mu \in \left(\bar{X} - Z_{1-\alpha} \frac{\sigma}{\sqrt{n}}, \infty\right)$	$\mu \in \left(-\infty, \bar{X} + Z_{1-\alpha} \frac{\sigma}{\sqrt{n}}\right)$
σ^2 אינו ידוע, $n > 30$.	$\mu \in \left(\bar{X} - Z_{1-\frac{\alpha}{2}} \frac{s}{\sqrt{n}}, \bar{X} + Z_{1-\frac{\alpha}{2}} \frac{s}{\sqrt{n}}\right)$	$\mu \in \left(\bar{X} - Z_{1-\alpha} \frac{s}{\sqrt{n}}, \infty\right)$	$\mu \in \left(-\infty, \bar{X} + Z_{1-\alpha} \frac{s}{\sqrt{n}}\right)$
σ^2 אינו ידוע. $n \leq 30$ וכן האוכלוסייה מתפלגת נורמלית.	$\mu \in \left(\bar{X} - t_{1-\frac{\alpha}{2}}^{n-1} \frac{s}{\sqrt{n}}, \bar{X} + t_{1-\frac{\alpha}{2}}^{n-1} \frac{s}{\sqrt{n}}\right)$	$\mu \in \left(\bar{X} - t_{1-\alpha}^{n-1} \frac{s}{\sqrt{n}}, \infty\right)$	$\mu \in \left(-\infty, \bar{X} + t_{1-\alpha}^{n-1} \frac{s}{\sqrt{n}}\right)$

אורך רווח הסמך (רלוונטי רק עבור רב"ס דו-צדדי!)

- הרב"ס הינו סימטרי, כאשר במרכזו נמצא ממוצע המדגם.
- המרחק מהממוצע לכל אחד מהגבולות זהה (כי הרב"ס סימטרי).
- מרחק זה מכונה **טעות הדגימה** וגודלו $\varepsilon = Z_{1-\frac{\alpha}{2}} \frac{\sigma}{\sqrt{n}}$
- אורך הרב"ס (מרחק מגבול עליון לתחתון): $L = 2 \cdot \varepsilon = 2 \cdot Z_{1-\frac{\alpha}{2}} \frac{\sigma}{\sqrt{n}}$
- גדלים אלו **אינם תלויים בתצפיות** שהתקבלו בפועל במדגם!
- אורך הרב"ס גדל כאשר:
- 1. רמת הסמך $(1 - \alpha)$ גדלה. 2. השונות גדלה. 3. גודל המדגם קטן
- גודל המדגם המינימלי n המבטיח טעות דגימה שאינה עולה על d ברמת סמך $1 - \alpha$: $Z_{1-\frac{\alpha}{2}} \frac{\sigma}{\sqrt{n}} \leq d \Rightarrow n \geq \left(\frac{Z_{1-\frac{\alpha}{2}} \sigma}{d}\right)^2$

התמודדות עם טענות "מישהו טוען ש..." באמצעות רב"ס

- מישהו טוען שהפרמטר שווה לערך מסוים C - מקבלים את הטענה אם רב"ס דו"צ כולל את הערך.
- מישהו טוען שהפרמטר נמוך מערך מסוים C - מקבלים את הטענה אם רב"ס ח"צ עליון אינו כולל את הערך (= כל הרב"ס נמוך מ- C).
- מישהו טוען שהפרמטר גבוה מערך מסוים C - מקבלים את הטענה אם רב"ס ח"צ תחתון אינו כולל את הערך (= כל הרב"ס גבוה מ- C).

בניית רב"ס לפרופורציה p

נתונה אוכלוסייה עם פרופורציה p. במדגם של n תצפיות מהאוכלוסייה, כל תצפית מתפלגת ברנולי עם סיכוי p להצלחה:

$$X_i \sim B(1, p) \Rightarrow \sum_{i=1}^n X_i \equiv X \sim B(n, p)$$

$$\hat{p} = \frac{\sum_{i=1}^n X_i}{n} = \frac{X}{n}, \quad \hat{q} = 1 - \hat{p}$$

האומדים הנקודתיים ל-p ול-q: עבור $np \geq 10, nq \geq 10$ מתקיים משפט הגבול המרכזי ואז: $\hat{p} \sim N\left(p, \frac{pq}{n}\right)$ אז הרב"ס ל-p באמינות $1 - \alpha$:

$$\begin{aligned} \text{רווח סמך דו צדדי} & p \in \left(\hat{p} - Z_{1-\frac{\alpha}{2}} \sqrt{\frac{\hat{p}\hat{q}}{n}}, \hat{p} + Z_{1-\frac{\alpha}{2}} \sqrt{\frac{\hat{p}\hat{q}}{n}}\right) \\ \text{רווח סמך חד צדדי תחתון} & p \in \left(\hat{p} - Z_{1-\alpha} \sqrt{\frac{\hat{p}\hat{q}}{n}}, 1\right) \\ \text{רווח סמך חד צדדי עליון} & p \in \left(0, \hat{p} + Z_{1-\alpha} \sqrt{\frac{\hat{p}\hat{q}}{n}}\right) \end{aligned}$$

גודל המדגם המינימלי n המבטיח טעות דגימה שאינה עולה על d ברמת סמך $1 - \alpha$: $Z_{1-\frac{\alpha}{2}} \sqrt{\frac{\hat{p}\hat{q}}{n}} \leq d \Rightarrow n \geq \left(\frac{Z_{1-\frac{\alpha}{2}}}{d}\right)^2 \cdot \hat{p}\hat{q}$

לצורך מציאת הגודל n

- אם קיימת הערכה מוקדמת לפרופורציה האמיתית לפיה $p \leq C$, נבחר את הערכים $\hat{p} = C, \hat{q} = 1 - C$
- אם אין הערכה מוקדמת לפרופורציה האמיתית ניקח את הערכים $p = q = 0.5$ (בערכים אלו הפונ' pq מקבלת מקסימום בתחום (0,1)).
- הערה: האומד p נמצא במרכז הרב"ס.

בניית רב"ס לשונות / סטיית התקן

התוחלת ידועה. התפלגות נורמלית של האוכלוסייה:

רב"ס דו צדדי	רב"ס חד צדדי תחתון	רב"ס חד צדדי עליון
$\sigma^2 \in \left(\frac{\sum_{i=1}^n (X_i - \mu)^2}{\chi^2_{\frac{\alpha}{2}, n}}, \frac{\sum_{i=1}^n (X_i - \mu)^2}{\chi^2_{1-\frac{\alpha}{2}, n}}\right)$	$\sigma^2 \in \left(\frac{\sum_{i=1}^n (X_i - \mu)^2}{\chi^2_{1-\alpha, n}}, \infty\right)$	$\sigma^2 \in \left(-\infty, \frac{\sum_{i=1}^n (X_i - \mu)^2}{\chi^2_{\alpha, n}}\right)$

התוחלת איננה ידועה. התפלגות נורמלית של האוכלוסייה:

רב"ס דו צדדי	רב"ס חד צדדי תחתון	רב"ס חד צדדי עליון
$\sigma^2 \in \left(\frac{(n-1)s^2}{\chi^2_{\frac{\alpha}{2}, n-1}}, \frac{(n-1)s^2}{\chi^2_{1-\frac{\alpha}{2}, n-1}}\right)$	$\sigma^2 \in \left(\frac{(n-1)s^2}{\chi^2_{1-\alpha, n-1}}, \infty\right)$	$\sigma^2 \in \left(-\infty, \frac{(n-1)s^2}{\chi^2_{\alpha, n-1}}\right)$

- לעיתים שימושי: $(n-1)s^2 = \sum_{i=1}^n (X_i - \bar{X})^2$
- מכיוון שהתפלגות χ^2 אינה סימטרית, האומד הנקודתי לא ימצא במרכז הרווח. כמו כן, הנוסחאות לגבי אורך הרב"ס אינן רלוונטיות.
- אם מבקשים רב"ס ל-s, עושים רב"ס ל-s² ומוציאים שורש.

בדיקת השערות

- השערת האפס - H_0 (ברירת המחדל, המצב הקיים)
- השערת האלטרנטיבה - H_1 (המידע החדש המועלה כנגד H_0)
- אזור דחייה (C) - תחום תוצאות המדגם שעבורו מחליטים לדחות את H_0 (ולקבל את H_1).
- אזור הקבלה (\bar{C}) - תחום תוצאות המדגם שעבורו מחליטים לא לדחות את H_0 .

טעות מסוג ראשון - דחיית H_0 כאשר היא נכונה: $P_{H_0}(C/H_0) \equiv \alpha$
טעות מסוג שני - אי דחיית H_0 למרות ש- H_1 נכונה: $P_{H_1}(\bar{C}/H_1) \equiv \beta$
עוצמת המבחן - ההסתברות של דחייה מוצדקת של H_0 , $P_{H_0}(C/H_1) \equiv 1 - \beta$
 כלומר קבלת H_1 כאשר H_1 אכן נכונה:

H_1 נכונה	H_0 נכונה
אנחנו צודקים!	ביצענו שגיאה מסוג 1
ביצענו שגיאה מסוג 2	אנחנו צודקים!

- נרצה למזער את α ו- β (כש- α קטנה β בהכרח גדל)
- ככל שאזור הדחייה C קטן יותר \leftarrow עוצמת המבחן קטנה יותר
- טעות מסוג ראשון חמורה יותר מאשר טעות מסוג שני.
- לכן, נקבע רמת מובהקות (חסם עליון ל- α) הנדרשת מהמבחן (נהוג 5%), ותחת חסם זה נחפש מבחן בעל עוצמה מקסימאלית (β מיני).
- אם דוחים את H_0 ברמת מובהקות α , אז נדחה את H_0 גם עבור רמת מובהקות $\tilde{\alpha} < \alpha$.
- אם לא דוחים את H_0 ברמת מובהקות α , אז לא נדחה את H_0 גם עבור כל רמת מובהקות $\tilde{\alpha} > \alpha$.

מבחן חד-צדדי

- מבחן חד-צדדי ימני: $H_0: \mu \leq \mu_0$ או $H_1: \mu > \mu_0$
- מבחן חד-צדדי שמאלי: אי השוויון של H_1 הפוך.

בדיקת השערות במודל נורמלי

בדיקת השערות על תוחלת

מודל נורמלי - מתקיים כאשר סטטיסטי המבחן (ממוצע המדגם) מתפלג נורמלית, כלומר כאשר מתקיים לפחות אחד מהשניים:

- האוכלוסייה מתפלגת נורמלית
- מתקיימים התנאים לשימוש במג"מ ($n > 30$ או $nq, np > 10$)
- כשנתון משקל בד"כ נניח לגבי התפלגות נורמלית

מבחן דו צדדי	מבחן חד צדדי	מערכת ההשערות
$H_0: \mu = \mu_0$ $H_1: \mu \neq \mu_0$	$H_0: \mu \geq \mu_0$ $H_1: \mu < \mu_0$	$H_0: \mu \leq \mu_0$ $H_1: \mu > \mu_0$
$\bar{x} > \mu_0 + z_{1-\frac{\alpha}{2}} \cdot \frac{\sigma}{\sqrt{n}}$ או: $\bar{x} < \mu_0 - z_{1-\frac{\alpha}{2}} \cdot \frac{\sigma}{\sqrt{n}}$	$\bar{x} < \mu_0 - Z_{1-\alpha} \cdot \frac{\sigma}{\sqrt{n}}$	$\bar{x} > \mu_0 + Z_{1-\alpha} \cdot \frac{\sigma}{\sqrt{n}}$ הערך הקריטי k
$n \geq \left[\frac{(z_{1-\frac{\alpha}{2}} + z_{1-\beta}) \cdot \sigma}{\mu_0 - \mu_1}\right]^2$	$n \geq \left[\frac{(z_{1-\alpha} + z_{1-\beta}) \cdot \sigma}{\mu_0 - \mu_1}\right]^2$	גודל מדגם מיני המבטיח $1 - \beta$ רצויים
לאחר שבוצע מדגם והתקבלה תוצאה מדגמית $\bar{X} = C$ אזי עבור תוצאה זו:		
$2 \cdot \left(1 - \phi\left(\frac{ C - \mu_0 }{\sigma/\sqrt{n}}\right)\right)$	$\phi\left(\frac{C - \mu_0}{\sigma/\sqrt{n}}\right)$	$1 - \phi\left(\frac{C - \mu_0}{\sigma/\sqrt{n}}\right)$
P-value (מובהקות התוצאה)		

הטבלה היא עבור שונות ידועה. אם השונות לא ידועה אומדים אותה באמצעות s² ומחליפים את ההתפלגות הנורמלית בהתפלגות t⁽ⁿ⁻¹⁾.

מובהקות התוצאה P-value

ההסתברות לקבל תוצאה כמו שהתקבלה במדגם בפועל או תוצאה קיצונית ממנה (בכיוון האלטרנטיבה), בהנחה שהשערת האפס נכונה.

* אם התוצאה שהתקבלה במדגם בפועל (C) היא הערך הקריטי של המבחן (K) אז P-value היא בדיוק ההסתברות לטעות מסוג ראשון!

* כלל ההחלטה: אם ערך PV קטן מ- α בו אנו מעוניינים (נאמר שתוצאת הניסוי היא "מובהקת") נדחה את H_0 , אחרת לא נדחה את H_0 .

P-value היא רמת המובהקות המינימאלית שעבורה יוחלט לדחות את H_0 .
 * אם על סמך מדגם ניתן לדחות את H_0 עבור רמת מובהקות α , תידחה גם עבור כל רמת מובהקות $\tilde{\alpha} < \alpha$ (כי $PV \leq \alpha < \tilde{\alpha}$).

* אם על סמך מדגם ניתן לקבל את H_0 עבור רמת מובהקות α , תתקבל (לא תידחה) גם עבור כל רמת מובהקות $\tilde{\alpha} > \alpha$ (כי $PV \geq \alpha > \tilde{\alpha}$).

בדיקת השערות על פרופורציה

נדרש $nq_0 > 10, np_0 > 10$ ע"מ שיתקיים מג"מ, ואז מקבלים: $\hat{p} \sim N\left(p, \frac{pq}{n}\right)$

מבחן דו צדדי	מבחן חד צדדי	מערכת ההשערות
$H_0: p = p_0$ $H_1: p \neq p_0$	$H_0: p = p_0$ $H_1: p < p_0$	$H_0: p = p_0$ $H_1: p > p_0$
$\hat{p} > p_0 + z_{1-\frac{\alpha}{2}} \cdot \sqrt{\frac{p_0 q_0}{n}}$ או: $\hat{p} < p_0 - z_{1-\frac{\alpha}{2}} \cdot \sqrt{\frac{p_0 q_0}{n}}$	$\hat{p} < p_0 - z_{1-\alpha} \cdot \sqrt{\frac{p_0 q_0}{n}}$	$\hat{p} > p_0 + z_{1-\alpha} \cdot \sqrt{\frac{p_0 q_0}{n}}$
$n \geq \left[\frac{Z_{1-\frac{\alpha}{2}} \sqrt{p_0 q_0} + Z_{1-\beta} \sqrt{p_1 q_1}}{p_0 - p_1}\right]^2$	$n \geq \left[\frac{Z_{1-\alpha} \sqrt{p_0 q_0} + Z_{1-\beta} \sqrt{p_1 q_1}}{p_0 - p_1}\right]^2$	גודל מדגם מיני המבטיח $1 - \beta$ רצויים
לאחר שבוצע מדגם והתקבלה תוצאה מדגמית $\bar{X} = C$ אזי עבור תוצאה זו:		
$2 \cdot \left(1 - \phi\left(\frac{ C - p_0 }{\sqrt{p_0 q_0/n}}\right)\right)$	$\phi\left(\frac{C - p_0}{\sqrt{p_0 q_0/n}}\right)$	$1 - \phi\left(\frac{C - p_0}{\sqrt{p_0 q_0/n}}\right)$
P-value (מובהקות התוצאה)		

הנוסחאות בדף המבחן (II)

בדיקת השערות על שונות

* מתקיים רק כשהאוכלוסייה מתפלגת נורמלית
 * רמת המובהקות הגבולית: להתפלגות χ^2 אין נוסחה לחישוב P-value. נסתכל בטבלת ההתפלגות במספר דרגות החופש המתאימות ונראה עבור איזה ערך P (1- α) היינו דוחים את ההשערה. ר"ה הגבולית היא בין ערך ה- α הזה לזה שבא אחריו.

בדיקת השערות על שיוויון שונות באוכלוסיות נורמליות

- $S^2 = \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n-1} = \frac{\sum_{i=1}^n x_i^2 - n \cdot \bar{x}^2}{n-1}$
- בטבלת F נתונים לנו רק הערכים עבור $f_{0.95}$, כלומר עבור $\alpha = 10\%$. לכן, נשתמש תמיד במובהקות $\alpha = 10\%$.
- ע"מ למצוא את $f_{\frac{\alpha}{2}}$ ניעזר בזהות: $f_{\frac{\alpha}{2}}^{n_1-1, n_2-1} = \frac{1}{f_{1-\frac{\alpha}{2}}^{n_2-1, n_1-1}}$
- אם לא נדחה את H_0 בר"מ 10%, גם לא נדחה אותה בר"מ קטנה יותר.

בדיקת השערות על הפרש פרופורציות

המרה לח"צ ימני: $H_0: p_1 = p_2$ או $H_1: p_1 > p_2$
 דחה אם: $Z_{\hat{p}_1 - \hat{p}_2} > Z_{1-\alpha}$ (שמאלי: הופכים סימן)

מבחן חי בריבוע לאי תלות

בודק האם קיימת תלות בין שני משתנים (בין שתי אוכלוסיות, בין שתי תכונות באותה אוכלוסיה וכו').

נתונות n תצפיות בלתי תלויות, בכל תצפית נרשמו זוג ערכים (Y, X) . מסווגים את התצפיות בלוח שכיחויות דו ממדי בעל r שורות ו- c עמודות. O_{ij} - שכיחות התצפית המוגדרת ע"פ התא (i, j) .

בשולי הלוח - השכיחויות של הקטגוריות של כל משתנה בנפרד.

X, Y	Y ₁	Y ₂	...	Y _j	...	Y _c	סה"כ
X ₁	O ₁₁	O ₁₂	...	O _{1j}	...	O _{1c}	$f_{1.} = \sum_{j=1}^{j=c} O_{1j}$
X ₂	O ₂₁	O ₂₂	f _{2.}
...
X _i	O _{i1}	O _{ij}	...	O _{ic}	...
...
X _r	O _{r1}	O _{ij}	...	O _{re}	f _{r.}
סה"כ	$f_{.1} = \sum_{i=1}^{i=r} O_{i1}$	f _{2.}	f _{.c}	n

לא קיימת תלות בין משתנה השורות למשתנה העמודות $H_0: P(X_i \cap Y_j) = P(X_i) \cdot P(Y_j) \Leftrightarrow$

אחרת: H_1

$$\chi^2_{emp} = \sum_{i=1}^r \sum_{j=1}^c \frac{(O_{ij} - E_{ij})^2}{E_{ij}} \quad \text{סטטיסטי המבחן לאי תלות:}$$

$$E_{ij} = \frac{f_{i.} \cdot f_{.j}}{n} \quad \text{כאשר:}$$

ע"מ להשתמש במבחן נדרש שבכל תא (i, j) יתקיים $E_{ij} \geq 5$.

כלל ההחלטה ברמת מובהקות α : דחה אם $\chi^2_{emp} > \chi^2_{1-\alpha, (r-1)(c-1)}$

בגרסיה לינארית פשוטה

מאפשרת לחזות את ערכו של משתנה אחד ע"פ ערכו של משתנים אחרים.

Y - המשתנה התלוי / המנובא / המוסבר

X - המשתנה הבלתי תלוי / המנבא / המסביר

מודל הרגרסיה הלינארית הפשוטה

הערך הממוצע של Y (התוחלת): $E(Y/X = x_i) = \beta_0 + \beta_1 x_i$

עבור פרט מסוים באוכלוסיה מתקיים: $y_i = \beta_0 + \beta_1 x_i + \varepsilon_i$ כאשר ε_i

מיצג "רעש" אקראי קטן. מניחים שהרעשים ב"ת ונורמליים - $\varepsilon_i \sim N(0, \sigma^2)$

לכל i , ולכן מתקיים: $(y_i | x) \sim N(\beta_0 + \beta_1 x, \sigma^2)$

מודל הרגרסיה הלינארית במדגם

אומד ל- β_0 (החותך); אומד ל- β_1 (השיפוע). השיגאה ε_i לא ידועה.

קו הרגרסיה שיתקבל מהמדגם: $\hat{y}_i = b_0 + b_1 x_i$

מקדמי הרגרסיה - שיטת הריבועים הפחותים

הסטייה של תצפית במדגם ביחס לקו נתון \hat{y}_i : $e_i = y_i - \hat{y}_i$

$$b_1 = \frac{\sum_{i=1}^n (y_i - \bar{y})(x_i - \bar{x})}{\sum_{i=1}^n (x_i - \bar{x})^2} = \frac{SS_{xy}}{SS_x} = \frac{\sum_{i=1}^n x_i y_i - n \bar{x} \bar{y}}{\sum_{i=1}^n x_i^2 - n \bar{x}^2} \quad b_0 = \bar{y} - b_1 \bar{x}$$

תכונות קו הריבועים הפחותים

• הקו עובר דרך נקודת הממועים (\bar{x}, \bar{y})

• האומדים b_1, b_0 הם חסרי הטיה. כלומר: $E(b_1) = \beta_1$ $E(b_0) = \beta_0$

משפט פירוק השוניות - $SST = SSE + SSR$

$$SST \equiv SS_y = \sum_{i=1}^n (y_i - \bar{y})^2 = \sum_{i=1}^n y_i^2 - n \bar{y}^2$$

$$SSR \equiv \sum_{i=1}^n (\hat{y}_i - \bar{y})^2 = b_1^2 SS_x \quad (\text{ברל"פ})$$

$$SSE \equiv \sum_{i=1}^n e_i^2 = \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

• ככל שקו הרגרסיה מתאים יותר למודל, יותר גדול SSR ויותר קטן SSE .

• בהתאמה מושלמת, $SSR = SST, SSE = 0$

• "אחוז השונות המוסברת", "מקדם ההסבר": $R^2 \equiv \frac{SSR}{SST} = 1 - \frac{SSE}{SST}$

• המתאם בין שני משתנים מקריים במדגם: $\hat{\rho} \equiv r = \frac{SS_{xy}}{\sqrt{SS_x \cdot SS_y}}$

• בגרסיה לינארית פשוטה: $r^2 = R^2$

אמידת σ^2 (השונות באוכלוסיה) ושונויות האומדים b_1, b_0

$$s^2 = \hat{\sigma}^2 = MSE = \frac{SSE}{n-2} = \frac{SST - SSR}{n-2} = \frac{SS_y - b_1^2 SS_x}{n-2} \quad \text{א.ח.ה. ל- } \sigma^2$$

$$S_{b_1} = \frac{s}{\sqrt{SS_x}} \quad , \quad S_{b_0} = s \sqrt{\frac{1}{n} + \frac{\bar{x}^2}{SS_x}} \quad \text{א.ח.ה. לסטיית התקן של } b_1, b_0$$

בניית רווח סמך לפרמטרים שנאמדו

רב"ס דו-צדדי ל- β_1 ברמת סמך $1 - \alpha$: $\beta_1 \in \left(b_1 \pm \frac{s}{\sqrt{SS_x}} t_{1-\frac{\alpha}{2}}^{(n-2)} \right)$

- אם נמצא רב"ס הזה, לא ניתן להגיד ברמת סמך $1 - \alpha$ כי

הרגרסיה מובהקת = הקשר בין x ל- y הוא לינארי.

רב"ס דו-צדדי ל- β_0 ברמת סמך $1 - \alpha$: $\beta_0 \in \left(b_0 \pm S \sqrt{\frac{1}{n} + \frac{\bar{x}^2}{SS_x}} t_{1-\frac{\alpha}{2}}^{(n-2)} \right)$

רב"סים ח"צ: נחליף את $1 - \frac{\alpha}{2}$ ל- $1 - \alpha$ ונעדכן את הגבול המתאים ל- $\pm \infty$.

הנוסחאות בדף המבחן (IV)

בדיקת השערות על הפרש תוחלות

מערכת ההשערות: $H_0: \mu_1 - \mu_2 = d_0$

$H_1: \mu_1 - \mu_2 \neq d_0$

• נתונה הערכה לגבי גודל ההפרש d_0 , אחרת $d_0 = 0$

1. מערכת ההשערות - בדף הכלל הוא עבור דו"צ. עבור מבחן חד צדדי:

$H_1: \mu_1 - \mu_2 < 0$ - מבחן חד צ' שמאלי, דחה אם: $T_D < -t_{1-\alpha}^{(n_1+n_2-1)}$

$T_D > t_{1-\alpha}^{(n_1+n_2-1)}$ - מבחן חד צ' ימני, דחה אם:

2. תלות - האם המדגמים בלתי תלויים או תלויים (=מזווגים).

מדגמים הם תלויים כאשר זהות הפרטים שנבחרו במדגם אחד

מושפעת מזהות הפרטים שנבחרו במדגם השני. המקרים הנפוצים:

א. כאשר בשני המדגמים מופיעים אותם נבדקים (אנשים, מכונות...)

ב. כשיש התאמה מכוונת בין הנבדקים במדגמים (למשל: בגיל ומין).

3. שיוויון שונויות - אם השונויות בשתי האוכלוסיות אינן ידועות, האם ניתן

להניח שהן שוות? איך מחליטים: (1) זה נתון; (2) מבצעים בדיקת

השערות על שיוויון שונויות; (3) בודקים באמצעות רב"ס ליחס השונויות.

רב"ס להפרש תוחלות ברמת סמך $1 - \alpha$

הנחות מיוחדות	רב"ס דו צדדי	הנחות מיוחדות
מדגמים ב"ת, שונויות ידועות	$\mu_1 - \mu_2 \in \left(\bar{x} - \bar{y} - z_{1-\frac{\alpha}{2}} \cdot \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}} \right)$	$E(X_i) = \mu_1, V(X_i) = \sigma_1^2$ $E(Y_j) = \mu_2, V(Y_j) = \sigma_2^2$
$n_1, n_2 > 30$ או נורמליות	$\bar{x} - \bar{y} + z_{1-\frac{\alpha}{2}} \cdot \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}$	
מדגמים ב"ת, שונויות לא ידועות אך שוות בשתי האוכלוסיות.	$\mu_1 - \mu_2 \in \left(\bar{x} - \bar{y} - t_{1-\frac{\alpha}{2}}^{n_1+n_2-2} \cdot s_p \sqrt{\frac{1}{n_1} + \frac{1}{n_2}} \right)$	$s_p^2 = \frac{(n_1-1)s_1^2 + (n_2-1)s_2^2}{n_1+n_2-2}$
	$\bar{x} - \bar{y} + t_{1-\frac{\alpha}{2}}^{n_1+n_2-2} \cdot s_p \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}$	$X_1, \dots, X_{n_1} \sim N(\mu_1, \sigma^2)$ $Y_1, \dots, Y_{n_2} \sim N(\mu_2, \sigma^2)$
מדגמים ב"ת, שונויות לא ידועות ולא שוות.	$\mu_1 - \mu_2 \in \left(\bar{x} - \bar{y} - t'_{1-\frac{\alpha}{2}} \cdot \sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}} \right)$	$X_1, \dots, X_{n_1} \sim N(\mu_1, \sigma_1^2)$ $Y_1, \dots, Y_{n_2} \sim N(\mu_2, \sigma_2^2)$
	$\bar{x} - \bar{y} + t'_{1-\frac{\alpha}{2}} \cdot \sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}$	

רב"ס להפרש פרופורציות

נסתפק בתנאים $n_1 p_1 \geq 10, n_1 q_1 \geq 10, n_2 p_2 \geq 10, n_2 q_2 \geq 10$ וגם

$$p_1 - p_2 \in \left(\hat{p}_1 - \hat{p}_2 - Z_{1-\frac{\alpha}{2}} \sqrt{\frac{\hat{p}_1 \hat{q}_1}{n_1} + \frac{\hat{p}_2 \hat{q}_2}{n_2}} \right), \quad \hat{p}_1 - \hat{p}_2 + Z_{1-\frac{\alpha}{2}} \sqrt{\frac{\hat{p}_1 \hat{q}_1}{n_1} + \frac{\hat{p}_2 \hat{q}_2}{n_2}}$$

* הרב"סים סימטריים, וכל מה שלמדנו לגבי אורך רב"סים תקף גם כאן.

בניית רב"ס ליחס שונויות

דרך חלופית לבדיקת טענות שיוויון

אוכלוסיות: אם השונויות שוות היחס 1.

אם הערך 1 נמצא בתוך הרב"ס, נקבע

שהשונויות שוות ברמת מובהקות α .

מבחן חי בריבוע לטיב התאמה

האם לפי נתוני המדגם סביר שהאוכלוסיה מתפלגת בהתפלגות מסוימת.

השערת האפס היא כי האוכלוסיה מתפלגת בהתפלגות תיאורטית כלשהי.

השערת האלטרנטיבה היא כי התפלגות האוכלוסיה היא אחרת.

* במקרה זה האינפורמציה החדשה "תאושר" ע"י קבלת H_0 !

שלבי העבודה

1. מוודאים שהפרמטרים של ההתפלגות ידועים. אם לא, יש לאמוד אותם.

2. יש לנסח בבירור את ההשערות הנבדקות

3. מחלקים את נתוני המדגם לקבוצות (מספר הקבוצות מסומן ב- k):

• אם נתונה חלוקה לקבוצות, מוודאים שתוחלת מס' הפרטים הצפויים

בכל קבוצה i מקיים: $E_i \geq 5$. אם לא, מאחדים עם קבוצה סמוכה.

$$E_i = n \left[\sum_{x=a}^b P(X=x) \right], (R_{צ'י}), E_i = n[P(a \leq X \leq b)] \quad (B_{יד})$$

• אם לא נתונה חלוקה, ניצור אותה בעצמנו כך שהקבוצות יכסו את כל

התחום של ההתפלגות, ויתקיים כמו קודם: $E_i \geq 5$.

4. O_i - מס' הפרטים בכל קבוצה כפי שהתקבלו בפועל במדגם.

$$\chi^2_{emp} = \sum_{i=1}^k \frac{(O_i - E_i)^2}{E_i} \quad \text{חישוב סטטיסטי המבחן:}$$

עבור מדגם גדול, וכאשר מתקיים $E_i \geq 5$ הסטטיסטי מתפלג בקירוב χ^2

עם $(k-p-1)$ דרגות חופש, כש- p הוא מס' הפרמטרים שנאמדו. אם לא

נאמדו כל פרמטרים, $p=0$.

6. הגדרת המובהקות הרצויה α ומציאת הערך הטבלאי: $\chi^2_{1-\alpha, (k-p-1)}$

משמעות הערך הטבלאי: בהסתברות של $1 - \alpha$ נקבל את הערך הטבלאי

או ערך קטן ממנו, כאשר השערת האפס מתקיימת.

7. כלל החלטה: דחה את H_0 אם $\chi^2_{emp} > \chi^2_{1-\alpha, k-p-1}$

בדיקת השערות על ערכי β_0, β_1 ברמת מובהקות α

מבחן דו צדדי	מבחן חד צדדי	מבחן חד צדדי	מבחן חד צדדי
$H_0: \beta_0 = \mu_0$ $H_1: \beta_0 \neq \mu_0$	$H_0: \beta_0 = \mu_0$ $H_1: \beta_0 < \mu_0$	$H_0: \beta_0 = \mu_0$ $H_1: \beta_0 > \mu_0$	מערכת ההשערות
$ T_{b_0} = \left \frac{b_0 - \mu_0}{S_{b_0}} \right > t_{1-\frac{\alpha}{2}}^{n-2}$	$T_{b_0} = \frac{b_0 - \mu_0}{S_{b_0}} < -t_{1-\alpha}^{n-2}$	$T_{b_0} = \frac{b_0 - \mu_0}{S_{b_0}} > t_{1-\alpha}^{n-2}$	דחייה

עבור β_1 נציב בהתאם b_1, μ_1, S_{b_1} .

בדיקת השערות על מובהקות הרגרסיה

מערכת ההשערות: $H_0: \beta_1 = 0$ (אם נקבל את H_0 הרגרסיה לא מובהקת)
 $H_1: \beta_1 \neq 0$

דרך ראשונה: מבחן t על הפרמטר β_1

סטטיסטי המבחן: $T_{b_1} = \frac{b_1 - 0}{S_{b_1}} = \frac{b_1}{s/\sqrt{SS_x}} \sim t(n-2)$

כלל ההחלטה עבור רמת מובהקות α : דחה את H_0 אם $|T_{b_1}| > t_{1-\frac{\alpha}{2}}^{(n-2)}$

דרך שנייה: מבחן F

סטטיסטי המבחן: $F = \frac{MSR}{MSE} = \frac{SSR/1}{SSE/(n-2)}$

כלל ההחלטה עבור רמת מובהקות α : דחה את H_0 אם $F > f_{1-\alpha}^{1,(n-2)}$

עבור רגרסיה פשוטה בלבד מתקיים הקשר: $F = (T_{b_1})^2$

חיזוי ערכי המשתנה המוסבר

ר"ס לאומדן y_p בהינתן $X = x_p$ ברמת סמך $1 - \alpha$:

$$y_p \in \left((b_0 + b_1 \cdot x_p) \pm t_{1-\frac{\alpha}{2}}^{n-2} \cdot s \cdot \sqrt{1 + \frac{1}{n} + \frac{(x_p - \bar{x})^2}{SS_x}} \right)$$

ר"ס לתוחלת (המונתית) $E(y_p)$ בהינתן $X = x_p$ ברמת סמך $1 - \alpha$:
אותו דבר, רק בלי t בתוך השורש

רגרסיה ליניארית מרובה

- β_0 - החותך, β_i - השיפוע
- כמו ברל"פ, המקדמים β_i לא ידועים ומייצרים אומדי ריבועים פחותים b_i

בדיקת מובהקות מודל הרגרסיה השלם - מבחן F

מערכת ההשערות: $H_0: \beta_i = 0$
 $H_1: \text{else}$

סטטיסטי המבחן: $F = \frac{MSR}{MSE}$, דחה אם: $F_0 > F_{1-\alpha}^{(k, n-k-1)}$

הסקה על מקדמי הרגרסיה

1. מבחן t : זהו למבחן ברל"פ. מאפשר בחינה של כל פרמטר בנפרד.

אם המשתנה j -שווה ל-0, אין למשתנה j -השפעה לינארית על

המשתנה המוסבר.

$H_0: \beta_1 = 0$
 $H_1: \beta_1 \neq 0$

סטטיסטי המבחן: $t_0 = \frac{b_j}{s(b_j)}$, דחה אם $|t_0| > t_{1-\frac{\alpha}{2}}^{(n-k-1)}$

$s(b_j)$ - סטיית התקן של האומד למשתנה j -מופיע בפלט הרגרסיה.

רב"ס דו"צ לפרמטר β_j ברמת סמך $1 - \alpha$: $\beta_j \in [b_j \pm t_{1-\frac{\alpha}{2}}^{(n-k-1)} \cdot s(b_j)]$

2. מבחן F חלקי: בוחן את התרומה של קב' משתנים מסבירים להסבר

הרגרסיה (פרמטר יחיד או קבוצת פרמטרים). משווים את SSR בלי

קבוצת המשתנים לעומת הערך של רגרסיה שכוללת אותם.

סטטיסטי המבחן:

$$F_0 = \frac{(SSR_{full} - SSR_{partial}) / (df_{full} - df_{partial})}{SSE_{full} / (n - df_{full} - 1)} = \frac{(R_{full}^2 - R_{partial}^2) / (df_{full} - df_{partial})}{(1 - R_{full}^2) / (n - df_{full} - 1)}$$

דחה את H_0 בר"מ α אם: $F_0 > F_{1-\alpha}^{(df_{full} - df_{partial}, n - df_{full} - 1)}$

מקדם ההסבר ברגרסיה לינארית מרובה

מדד חלופי ל- R^2 : $R_{adj}^2 \equiv 1 - \frac{SSE / (n-k-1)}{SST / (n-1)}$

n - מס' תצפיות, k - מס' משתנים מסבירים.

תמיד מתקיים $R_{adj}^2 \leq R^2$. יש פחות תמריץ להוספת משתנים רבים למודל.

מוליטיקוליאריות

1. האומדים לשונות של מקדמי הרגרסיה (s_{b_i}) "מתנפחים": הרב"ס למקדמים האמיתיים (β_i) הרבה יותר רחבים מכפי שאמור להיות.

2. תוצאות מוזרות: סימני המקדמים הפוכים מהצפוי.

3. תוצאות מוזרות: מבחן F מעיד על מודל מובהק ומבחן t אינם מובהקים.

מבחנים למוליטיקוליאריות

1. בדיקת מתאם בין כל זוג משתנים מסבירים: $r_{x_j, x_k} = \frac{SS_{x_j x_k}}{\sqrt{SS_{x_j} \cdot SS_{x_k}}}$

אם המתאם גבוה (קרוב ל-1 בערך מוחלט) כדאי לוותר על אחד מהמשתנים.

2. מדד VIF: לכל משתנה מסביר x_j מחשבים את הערך $VIF_j = \frac{1}{1 - R_j^2}$

R_j^2 הוא מקדם ההסבר ברגרסיה בה המשתנה המוסבר הוא x_j והמשתנים המסבירים הם שאר המשתנים המסבירים. $VIF_j \geq 5$ - יש מוליטיקוליאריות

Regression Statistics						
Multiple R	$r = \frac{SS_{xy}}{\sqrt{SS_x SS_y}} = \frac{\sum x_i y_i - n \bar{x} \bar{y}}{\sqrt{SS_x SS_y}} = \sqrt{R^2}$		$\int u \cdot v' dx = u \cdot v - \int u' \cdot v dx$ - אינטגרציה בחלקים		מה שמודגש נכון רק לרגרסיה ליניארית פשוטה	
R Square	$R^2 = \frac{SSR}{SST} = 1 - \frac{SSE}{SST} = \frac{b_1^2 SS_x}{SS_y} = \frac{b_1^2 SS_x (\sum x_i^2 - n \bar{x}^2)}{\sum y_i^2 - n \bar{y}^2}$		<div><u>אריטמטיקה של נגזרות:</u> $(f \cdot g)'(x) = f'(x) \cdot g(x) + g'(x) \cdot f(x)$ $\left(\frac{f}{g}\right)'(x) = \frac{f'(x)g(x) - g'(x)f(x)}{g^2(x)}$</div>			
Adjusted R Square	לא נדרש ברגרסיה פשוטה.		<div>זהו P-value בשימוש במבחן F.</div>			
Standard Error	$S = \sqrt{S^2} = \sqrt{\hat{\sigma}^2} = \sqrt{\frac{SS_y - b_1^2 SS_x}{n-2}} = \sqrt{\frac{SSE}{n-2}}$					
Observations	n					
ANOVA						
	Df-degrees of freedom (דרגות חופש)	SS	MS	F	Significance F	
Regression	מס' המשתנים המסבירים k – (בפשוטה k=1)	SSR	$MSR = \frac{SSR}{k}$	$F_{stat} = \frac{MSR}{MSE}$	$P_v = P(F_{stat} < F_{\alpha})$ $F_{\alpha} = F_{1-\alpha}^{k, n-k-1}$	
Residual	n-k-1	SSE	$MSE = \frac{SSE}{n-k-1} = S^2$			
Total	n-1	SST				
	Coefficients	Standard Error	t Stat	P-value	Lower 95%	Upper 95%
Intercept	$b_0 = \bar{y} - b_1 \bar{x}$	$S_{b_0} = S \sqrt{\frac{1}{n} + \frac{\bar{x}^2}{SS_x}}$	$T_{b_0} = \frac{b_0}{S_{b_0}}$	$P_{value} = P(T_{b_0} < T_{\alpha})$ $T_{\alpha} = t_{1-\frac{\alpha}{2}}^{n-2}$	$b_0 - t_{1-\frac{\alpha}{2}}^{n-2} S_{b_0}$	$b_0 + t_{1-\frac{\alpha}{2}}^{n-2} S_{b_0}$
X	$b_1 = \frac{SS_{xy}}{SS_x}$	$S_{b_1} = \frac{S}{\sqrt{SS_x}}$	$T_{b_1} = \frac{b_1}{S_{b_1}}$	$P_{value} = P(T_{b_1} < T_{\alpha})$ $T_{\alpha} = t_{1-\frac{\alpha}{2}}^{n-2}$	$b_1 - t_{1-\frac{\alpha}{2}}^{n-2} S_{b_1}$	$b_1 + t_{1-\frac{\alpha}{2}}^{n-2} S_{b_1}$

זהו ערך סטטיסטי המבחן t של ההשערה $H_0: \beta_1 = 0$ כנגד $H_0: \beta_1 \neq 0$.

ברגרסיה פשוטה בלבד, מתקיים $T_{b_1} = \sqrt{F_{stat}}$

זהו P-value של ההשערה $H_0: \beta_1 = 0$ כנגד $H_0: \beta_1 \neq 0$.

עבור רגרסיה פשוטה בלבד, זה שווה ל-P-value בשימוש במבחן F.

התפלגויות בדידות מיוחדות (כאשר p החסתברות להצלחה בכל ניסיון)						
התפלגות	אחידה	ביונומית	היפר גיאומטרית	גאומטרית	ביונומית שלילית	פואסונית
תיאור	מ"מ בקפיצות של 1 עם הסתברות שווה לכל תוצאה	מספר הצלחות בסדרת n ניסויי ברנולי בִּית	מספר המיוחדים במדגם בגודל n (ללא החזרה) מאוכלוסייה בגודל N , שמתוכה D מיוחדים	מספר הניסיונות עד להצלחה הראשונה בסדרת ניסויי ברנולי בִּית	מספר הניסיונות עד להצלחה הראשונה בסדרת ניסויי ברנולי בִּית	מספר האירועים ביחידת זמן נתונה
סימון	$X \sim U(a, b)$	$X \sim B(n, p)$	$X \sim \text{HG}(N, D, n)$	$X \sim \text{G}(p)$	$X \sim \text{NB}(r, p)$	$X \sim P(\lambda)$
$P(X = k)$	$\begin{cases} \frac{1}{b-a+1}, & k = a, a+1, \dots, b \\ 0, & \text{אחרת} \end{cases}$	$\binom{n}{k} p^k q^{n-k}$ $k = 0, 1, \dots, n$	$\frac{\binom{D}{k} \binom{N-D}{n-k}}{\binom{N}{n}}$ $k = \max\{0, n - (N - D)\}, \dots, \min\{n, D\}$	$q^{k-1} p$ $k = 1, 2, \dots$	$\binom{k-1}{r-1} p^r q^{k-r}$ $k = r, r+1, \dots$	$\frac{e^{-\lambda} \lambda^k}{k!}$ $k = 0, 1, \dots$
$E(X)$	$\frac{b+a}{2}$	np	$\frac{nD}{N}$	$\frac{1}{p}$	$\frac{r}{p}$	λ
$V(X)$	$\frac{(b-a+1)^2 - 1}{12}$	npq	$\frac{nD}{N} \left(1 - \frac{D}{N}\right) \frac{N-n}{N-1}$	$\frac{q}{p^2}$	$\frac{r}{p^2}$	λ
תכונות	עבור המקרה הפרטי $X \sim U(a, b)$ $a = 1, b = N$ נקבל: $P(X) = \frac{1}{N}$	חיבור מ"מ בינומיים: $X \sim B(n, p), Y \sim B(m, p)$ אז P או $X + Y \sim B(n+m, p)$		$P(X > k) = q^k$; $P(X \geq k) = q^{k-1}$ $P(X > a + b \mid X > a) = P(X > b)$ $P(X = a + b \mid X > a) = P(X = b)$	תשוב: סכום של r מ"מ ב"ת שוויו התפלגות מתפלג בינומי שלילי עם פרמטרים r, p	מספר האירועים בקטעי זמן זרים בִּית $X \sim P(\lambda_1), Y \sim P(\lambda_2)$ $\Rightarrow X + Y \sim P(\lambda_1 + \lambda_2)$
התפלגויות רציפות מיוחדות						
התפלגות	אחידה	מעריכית	היפר גיאומטרית	גאומטרית	נורמלית	
תיאור	מיים שיכול לקבל כל ערך בהסתברות שווה		א. אם מספר האירועים ביחידת זמן מתפלג פואסונית עם קצב אירועים λ , אז הזמן הבין מופעי (בין אירועים עוקבים) מתפלג מעריכית עם פרמטר λ . ב. אורך חיים של תהליכים המקיימים חוסר זיכרון	מספר הניסיונות עד להצלחה הראשונה בסדרת ניסויי ברנולי בִּית	מספר הניסיונות עד להצלחה הראשונה בסדרת ניסויי ברנולי בִּית	מספר האירועים ביחידת זמן נתונה
סימון	$X \sim U(a, b)$	$X \sim \exp(\lambda)$				$X \sim N(\mu, \sigma^2)$
$f(x)$	$\begin{cases} \frac{1}{b-a}, & a \leq x \leq b \\ 0, & \text{אחרת} \end{cases}$	$\begin{cases} \lambda e^{-\lambda x}, & 0 < x \\ 0, & \text{אחרת} \end{cases}$				$\frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$
$F(x)$	$\begin{cases} 0, & x < a \\ \frac{x-a}{b-a}, & a \leq x < b \\ 1, & b \leq x \end{cases}$	$\begin{cases} 0, & x < a \\ \frac{x-a}{b-a}, & a \leq x < b \\ 1, & b \leq x \end{cases}$				$\Phi\left(\frac{x-\mu}{\sigma}\right)$
$E(X)$	$\frac{b+a}{2}$					μ
$V(X)$	$\frac{(b-a)^2}{12}$					σ^2
תכונות	במקרה הפרטי $X \sim U(0, 1)$ נקבל: $F(x) = \begin{cases} 0, & x < a \\ x, & a \leq x < b \\ 1, & b \leq x \end{cases}$			אם מיים מקיים את תכונת חוסר הזיכרון הוא מתפלג אקספוננציאלית ולחיפך. דוגמא: אם λ ממוצע תאונות דרכים בחדש אז $X \sim P(\lambda)$ מספר התאונות בחדש, ו $X \sim \exp(\lambda)$ זמן בין 2 תאונות. $P(X \leq t) = 1 - e^{-\lambda t}$, $P(X > t) = e^{-\lambda t}$	$f(\mu + x) = f(\mu - x)$, $Z_{(1-\alpha)} = -Z_\alpha$ $Z = \frac{X - \mu}{\sigma} \sim N(0, 1^2)$ $\Phi(-b) = 1 - \Phi(b)$, $p(x \geq a) = 1 - p(x \leq a)$ $P(a < Z < b) = \Phi(b) - \Phi(a)$ $P(-c < Z < c) = 2\Phi(c) - 1$	