

דף נוסחאות לקורס מבוא לסטטיסטיקה - 094423 - מבחן

רווח סמך: רווח סמך קטן הוא יותר מדויק (יותר נתונים מקטינים את ר"ס)
משמעות: לר"ס ברמת סמך 1 - alpha עבור פרמטר כלשהו, יש הסתברות של 1 - alpha שהרווח
שיתקבל ע"ס מדגם מקרי יכלול בתוכו את הערך האמיתי של הפרמטר.

התפלגות נורמאלית N(mu, sigma^2)
ר"ס לתוחלת, שונות ידועה: X-bar +/- Z_1-alpha/2 * sigma/sqrt(n)
סימטרי סביב ממוצע המדגם, מיקומו אקראי
רווח-סמך הקצר ביותר
אורכו קבוע מראש ושווה ל 2Z_1-alpha/2 * sigma/sqrt(n)

אם נחפש שאורך רווח הסמך לא יעלה על 2d: 2Z_1-alpha/2 * sigma/sqrt(n) <= 2d -> n >= (Z_1-alpha/2 * sigma/d)^2
בניית רווח-סמך כללי:
1. אנ"מ לפרמטר theta
2. בחירת פונקציה פיבולית q(theta, theta-hat)
3. מציאת השיבוטים X_p1, X_p2 כך ש p2 - p1 = 1 - alpha
4. P(X_p1 <= Q <= X_p2) = 1 - alpha
את theta נבחר כך שיהיה פונקציה של ס"מ.
אם נרצה רווח-סמך חד-צדדי נבחר p2 = 1 או p1 = 0
ר"ס לתוחלת, שונות איננה ידועה: X-bar +/- S/sqrt(n) * t_{(n-1), 1-alpha/2}

רווח סמך הקצר ביותר
מרכז רווח הסמך הוא ממוצע המדגם
אורך רווח הסמך הוא משתנה מקרי: 2t_{(n-1), 1-alpha/2} * S/sqrt(n)
לכן אין אפשרות לחשב מראש את גודל המדגם שיבטיח רווח באורך רצוי מאחר ואורך הרווח
תלוי באומדן לסטיית התקן שיתקבל במדגם.

רווח סמך לפרופורציה X ~ Bin(n, p)
רווח סמך 1 - alpha בקירוב: p +/- Z_1-alpha/2 * sqrt(p(1-p)/n)
מבוסס על קירובים: משפט הגבול המרכזי ואמידת השונות.
אורך רווח הסמך: 2Z_1-alpha/2 * sqrt(p(1-p)/n)
d מסוים ברמת הביטחון הנתונה היא: n >= (Z_1-alpha/2 * sqrt(p(1-p)/d))^2
גודל המדגם הגרוע ביותר: p(1-p) <= 0.25
ר"ס לשונות, תוחלת ידועה: [sum_{i=1}^n (X_i - mu)^2 / (n-1), sum_{i=1}^n (X_i - mu)^2 / n]
ר"ס לשונות, תוחלת איננה ידועה: [(n-1)S^2 / chi^2_{(n-1), 1-alpha/2}, (n-1)S^2 / chi^2_{(n-1), alpha/2}]
איננו האופטימאלי מבחינת אורכו
עבור סטיית התקן - פשוט שורש

רווח סמך להפרש תוחלות: (X-bar - Y-bar) +/- Z_1-alpha/2 * sqrt(sigma^2 * (1/n + 1/m))
אם השונות אינה ידועה (שונות שוות): S_p^2 = ((n-1)S_x^2 + (m-1)S_y^2) / (n+m-2)
שונות לא שוות: Q = ((X-bar - Y-bar) - (mu_x - mu_y)) / sqrt(S_x^2/n + S_y^2/m) ~ t(df)
ר"ס ליחס שונות: lambda = sigma_x^2 / sigma_y^2
היות ש F_{(n-1, m-1)} ~ S_x^2 / lambda ~ F_{(n-1, m-1)}
התפלגות ברנולי Ber(p_x), Ber(p_y)
ר"ס ל p_x - p_y: (p-hat_x - p-hat_y) +/- Z_1-alpha/2 * sqrt(p-hat_x(1-p-hat_x)/n + p-hat_y(1-p-hat_y)/m)
התפלגויות ממוגות (X_1, Y_1), (X_2, Y_2), ..., (X_n, Y_n)
ר"ס להפרש התוחלת: mu_x - mu_y
גדיר D_i = X_i - Y_i אזי D_1, ..., D_n ושווי התפלגות וכן: E(D_i) = E(X_i) - E(Y_i) = mu_x - mu_y = mu_d
נניח D_i ~ N(mu_d, sigma_D^2) ולכן: D-bar +/- t_{(n-1), 1-alpha/2} * S_D/sqrt(n)
לא ניתן להתעלם מהתלות
רווח הסמך השגוי (אם נתעלם מהתלות) יהיה רחב יותר מרווח הסמך הממוזג כי 2Cov(X, Y) > 0 | Var(X-bar - Y-bar) = Var(X-bar) + Var(Y-bar) - 2Cov(X-bar, Y-bar)

סימולציה ליצירת X ~ exp(lambda)
U = F_X(x) = 1 - e^{-lambda*x} -> F^{-1}(u) = -ln(1-u)/lambda -> X = -ln(1-U)/lambda
התפלגות אחידה X ~ Uni(a, b)
P(c <= x <= d) = (d-c)/(b-a), f_X(x) = 1/(b-a) a <= x <= b, else
E(X) = (a+b)/2, Var(X) = (b-a)^2/12
X-a/b-a ~ Uni(0,1)
פונקציה יוצרת מומנטים: M_X(t) = E(e^{tX})
iid: M_{X+Y} = M_X * M_Y
M_X(0) = 1
M_{aX+b}(t) = e^{bt} M_X(at)

שגיאת דגימה שגיאה כתוצאה מכך שנאספו נתונים על מדגם שהנו חלק מכלל האוכלוסייה.
גודל השגיאה תלוי בשיטת הדגימה ובגודל המדגם.
שגיאה שאינה שגיאת דגימה נובעת משימוש בשיטת דגימה לא נכונה.
דגימה הסתברותית דגימה שבה ידועה ההסתברות של כל פרט האוכלוסייה להיבחר.
דגימה מקרית פשוטה דגימה הסתברותית כך שכל שתי קבוצות שוות גודל הן בעלות ההסתברות שווה להיכלל במדגם. הנחה דגימה עם החזרה (לשם נוחות).

שונות S^2 = 1/(n-1) * sum_{i=1}^n (X_i - X-bar)^2 = 1/(n-1) * sum_{i=1}^n (X_i^2 - 2X_i * X-bar + X-bar^2)
חציין הערך שמחצית גדולים או שווים לו וחציית קטנים או שווים לו.
אינו רגיש לערכים חריגים (בניגוד לממוצע).
S^2 = 1/(n-1) * [sum_{i=1}^n X_i^2 - 2 * (sum_{i=1}^n X_i) * X-bar + n * X-bar^2]
E(S^2) = sigma^2
טווח Range = X_{(n)} - X_{(1)}

רבעון עליון הערך 1/4 גדולים או שווים לו ו 3/4 קטנים או שווים לו Q_3
רבעון תחתון הערך 3/4 גדולים או שווים לו ו 1/4 קטנים או שווים לו Q_1
תחום בין רבעוני התחום בו מצויות 50% מהתצפיות המרכזיות: IR = [Q_1, Q_3]
מדד פיזור שאינו רגיש לערכים חריגים.
שיברון p באוכלוסיה ערך epsilon כך שמתקיים P(X <= epsilon) = p
שיברון p במדגם ערך x שפרופורציה p מתוך n ערכי המדגם קטנים-שווים לו.
ככל שהמדגם גדל, יתקרב השיברון המדגמי לתיאורטי.
step = 3/2 * (Q_3 - Q_1), LF = Q_1 - step, UF = Q_3 + step
ה whisker הוא התצפית הראשונה מכל כיוון ללא חריגה מ LF/UF

לכל מספר ממשי פונקציית ההתפלגות האמפירית מוגדרת: F_n(x) = 1/n * sum_{i=1}^n I(X_i <= x)
סטטיסטי הוא פונקציה של התצפיות במדגם - של המ"מ, ואינו תלוי בפרמטר לא ידוע.
סטטיסטי המשמש לאמידת פרמטר נקרא אומדן עבור הפרמטר. הערך הספציפי שהתקבל מהמדגם נקרא אומדן.

יהי X מ"מ. המומנט ה-r של X מוגדר mu_r = E(X^r) אם התוחלת מוגדרת.
יהי X_1, ..., X_n מדגם מקרי מהתפלגות F_X(.). המומנט המדגמי r- מוגדר:
mu_r = 1/n * sum_{i=1}^n X_i^r
mu_1 = E(X)
mu_2 = Var(X) + mu_1^2
E(M_r) = 1/n * sum_{i=1}^n E(X_i^r) = mu_r
לפי חוק המספרים הגדולים mu_r -> M_r
חיסרון: לא לוקחים בחשבון אילוצים על תחומים. יתרון: פשוט
אומדן בשיטת המומנטים עבור theta ~ Uni[0, theta] הוא 2X-bar (חסר הטיה, עקיב).
sum_{i=1}^n (X_i - X-bar)^2 = sum_{i=1}^n (X_i - mu)^2 - n * (X-bar - mu)^2

פונקציית הנראות של n מ"מ X_1, ..., X_n מוגדרת להיות הצפיפות (או ההסתברות) המשותפת: L(theta) = L(theta; x_1, ..., x_n) = f_{X_1, ..., X_n}(x_1, ..., x_n; theta)
אומדן נראות רביית (אנ"מ) הוא הערך theta המביא למקסימום את L(theta).
אינו בהכרח יחיד
קל לעבוד עם ln(theta)
אם theta אנ"מ ל theta אי לכל פונקציה r(theta), אנ"מ לה הוא r(theta).
אנ"מ ל p כאשר Ber(p), theta = p, הוא X-bar.
אנ"מ עבור (mu, sigma^2), N(mu, sigma^2) נקבל: theta = (mu, sigma^2)
אנ"מ עבור theta ~ Uni[0, theta] הוא X_{(n)}, theta = X_{(n)} (מוטה, עקיב).
משפור: theta* = (n+1)/n * X_{(n)} (עקיב).

יהי T_n = T(X_1, ..., X_n) סטטיסטי המבוסס על מדגם בגודל n והנו אומדן ל theta(theta). ההטיה b_{T_n}(theta) = E(T_n) - theta(theta) מוגדרת על-ידי
b_{T_n}(theta) = 0 אם T_n הוא אומדן חסר הטיה (אנ"מ)
אם נוציא הרבה מאוד מדגמים, ועל-סמך כל מדגם נחשב את הערך של האומדן, אזי הממוצע של כל האומדנים יהיה שווה בקירוב לערך הפרמטר אותו רוצים לאמוד.

השגיאה הריבועית הממוצעת (MSE) של אומדן T_n ביחס ל theta(theta) מוגדרת כ MSE_{T_n}(theta) = E[(T_n - theta(theta))^2]
אם MSE = 0 או בהכרח T_n = tau(theta)
אם MSE_{T_n}(theta) <= MSE_{S_n}(theta) לכל theta in theta אז נעדיף את T_n משיקולי MSE.
MSE(T) = E[(T_n - tau(theta))^2] = Var(T) + (E(T) - tau(theta))^2
ריבוע הטיה
אם T_1, ..., T_n סדרת אומדים ל theta(theta) כך ש T_n = T_n(X_1, ..., X_n) נאמר כי סדרת האומדים עקיבה ל theta(theta) אם לכל epsilon > 0 מתקיים:
lim_{n->inf} P(|T_n - tau(theta)| < epsilon) = 1
משפט אם lim_{n->inf} MSE(T) -> 0 אזי הסדרה עקיבה ביחס ל theta(theta).
אם {T_n} סדרת אומדים עקיבה ל theta(theta) רציפה אזי {tau(theta)} סדרת אומדים עקיבה ל theta(theta).
אם T_n אח"ה ל theta(theta) לניאירית אז tau(T_n) אח"ה ל theta(theta).
מומנטים הם עקיבים וחי"

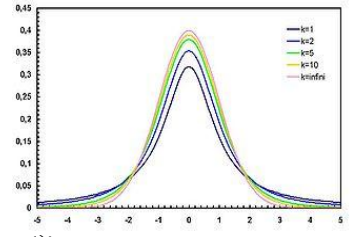
סטטיסטי מספיק היא פונקציה של המדגם הנוצרת בתוכה את כל האינפורמציה שיש במדגם ביחס לפרמטר.
סטטיסטי S = S(X_1, ..., X_n) יקרא סטטיסטי מספיק אם ההתפלגות המותנה של X_1, ..., X_n בהינתן S = s ב"ת theta לכל ערך s.
ס"מ אינו יחיד וכל פונקציה חח"ע של ס"מ היא גם ס"מ.
T*(X_1, ..., X_n) = r(T(X_1, ..., X_n))
עבור X_1, ..., X_n ~ Ber(p) ב"ת S = sum_{i=1}^n X_i ס"מ ל p וגם X-bar.
עבור X_1, ..., X_n ~ N(mu, sigma^2) ב"ת S = sum_{i=1}^n X_i, T = sum_{i=1}^n X_i^2 ס"מ ל (mu, sigma^2) וגם (X-bar, sum_{i=1}^n X_i^2) ס"מ
אם השונות אינה ידועה והתוחלת ידועה אזי ס"מ לשונות הוא sum_{i=1}^n (X_i - mu)^2
אם התוחלת אינה ידועה והשונות ידועה, אזי ס"מ לתוחלת הוא X-bar
משפט אנ"מ תלוי במדגם מקרי דרך ס"מ

משפט ראובל-בלקוול אומדן שהנו פונקציה של ס"מ הוא עם MSE קטן יותר בהשוואה ל MSE של אומדן שאינו פונקציה של ס"מ.
X_1, ..., X_n מדגם מקרי וי S ס"מ. T אומדן ל theta(theta) נגדיר T* = E(T|S)
1. T* הינו סטטיסטי פונקציה של S
2. E(T*) = E(T)
3. לכל theta in theta: MSE_{T*}(theta) <= MSE_T(theta)
* לא ניתן לשפר שוב

התפלגות t נראה מאוד דומה לנורמלי. יהי $Z \sim N(0,1)$ ו $W_k \sim \chi^2(k)$ ב"ת אז

$$T = \frac{Z}{\sqrt{\frac{W_k}{k}}} \sim t(k), \quad f_T(t) = \frac{\Gamma\left(\frac{k+1}{2}\right)}{\Gamma\left(\frac{k}{2}\right)\sqrt{k\pi}} \left(1 + \frac{t^2}{k}\right)^{-\frac{k+1}{2}}, \quad k > 0, -\infty < t < \infty$$

$$E(T) = 0, k > 1, \text{Var}(T) = \frac{k}{(k-2)}, k > 2$$



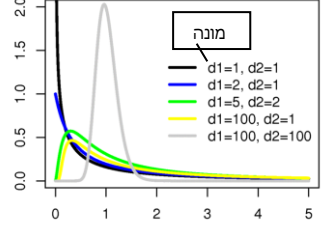
$\Gamma(n) = (n-1)!$
 $\Gamma(0.5) = \sqrt{\pi}$
 סימטרי סביב 0
 פיי"מ לא קיימת
 זנבות עבים יותר משל נורמלי

יהיו $X_1, \dots, X_n \sim N(\mu, \sigma^2)$ אז מתקיים $T = \frac{\bar{X} - \mu}{\frac{s}{\sqrt{n}}} \sim t_{(n-1)}$

התפלגות F יהיו $U_1 \sim \chi^2(k_1)$ ו $U_2 \sim \chi^2(k_2)$ ב"ת

$$F = \frac{U_1/k_1}{U_2/k_2} \sim F(k_1, k_2), \quad f_F(y) = \frac{\Gamma\left(\frac{k_1+k_2}{2}\right)}{\Gamma\left(\frac{k_1}{2}\right)\Gamma\left(\frac{k_2}{2}\right)} \left(\frac{k_1}{k_2}\right)^{\frac{k_1}{2}} \frac{y^{\frac{k_1}{2}-1}}{\left(1 + \frac{k_1 y}{k_2}\right)^{\frac{k_1+k_2}{2}}}, \quad y > 0$$

$$E(F) = \frac{m}{m-2}, m > 2, \text{Var}(F) = \frac{2m^2(k+m-2)}{k(m-2)^2(m-4)}, m > 4$$



$\lim_{m \rightarrow \infty} kF \sim \chi^2(k)$
 $\frac{1}{F} \sim F(m, k)$
 אם $X \sim t_{(d)}$ אז $X^2 \sim F_{(1,d)}$
 אם $Y = \frac{1}{F} P(X \leq q_X(p)) = p$
 אם $q_X(p) = \frac{1}{q_Y(1-p)}$ ו $P(Y \leq q_Y(1-p)) = 1-p$

משפט הפירוק סטטיסטי $S(x_1, \dots, x_n)$ יהיו מספיק ביחס לפרמטר θ אם לכל x_1, \dots, x_n מתקיים (צפיפות או הסתברות):

$$f_{X_1, \dots, X_n}(x_1, \dots, x_n; \theta) = h(x_1, \dots, x_n) \cdot g(S(x_1, \dots, x_n); \theta)$$

כאשר h פונקציה אי-שלילית ואינה פונקציה של θ . g אי-שלילית ותלויה במשתנים שלה רק דרך S .

טרנספורמציה חד-חד-ערכית של מ"מ רציף:

X מ"מ רציף (דיפרנציאבילי) המקבל ערכים בקטע $[a, b]$ מ $\mathbb{R} \rightarrow \mathbb{R}$ פונקציה גזירה ועולה (יורדת) ממש, $Y = g(X)$

$$f_Y(y) = \begin{cases} f_X(g^{-1}(y)) \cdot \left| \frac{\partial}{\partial y} g^{-1}(y) \right| & , g(a) \leq y \leq g(b) \\ 0 & , \text{else} \end{cases}$$

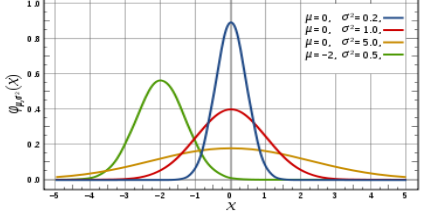
$$Y = a + bX \rightarrow f_Y(y) = \frac{1}{|b|} f_X\left(\frac{y-a}{b}\right)$$

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

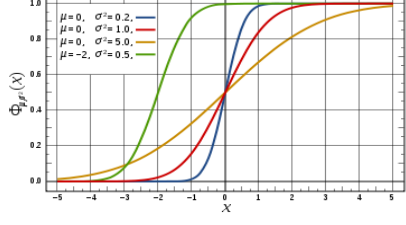
$$E(X) = \mu, \text{Var}(X) = \sigma^2$$

$$f(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}}, \quad \bar{X} \sim N\left(\mu, \frac{\sigma^2}{n}\right)$$

כל קומבינציה ליניארית של משתנים נורמליים הוא נורמלי.
 פונקציות הצפיפות:



פונקציית הצטברות:



אקספוננציאלי מוח

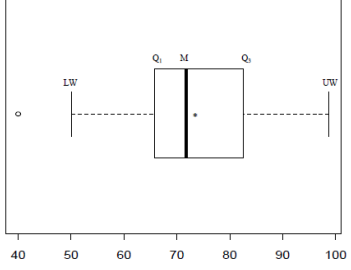
$$X = Y + \theta \sim \exp(\lambda)$$

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

סיכום קטן לאמדים

ס"מ	אנ"מ	מומנטים אמד	N
$(\bar{X}, \sum X_i^2)$ ס"מ ל $\theta = (\mu, \sigma^2)$	$\hat{\mu} = \bar{X}$ $\hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2$	$\hat{\mu} = \bar{X}$ $\hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2$	N
$\sum X_i$	$\hat{\lambda} = \bar{X}$	$\hat{\lambda} = \bar{X}$	Pois
$\sum X_i$	$\hat{\theta} = \frac{1}{\bar{X}}$	$\hat{\theta} = \frac{1}{\bar{X}}$	exp
$\sum X_i$	$\hat{p} = \bar{X}$	$\hat{p} = \bar{X}$	Bin
$\sum X_i$	$\hat{p} = \bar{X}$	$\hat{p} = \bar{X}$	Ber

דיאגרמת קופסא



בדיאגרמת קופסא ניתן לראות את מרבית הסטטיסטים: חציון, רבעון עליון ותחתון, Whiskers, טווח הנתונים וטווח בין-רבעוני.
 יש תצפית אחת הקטנה מהגדר התחתונה ולכן זאת תצפית הרגילה.
 ההתפלגות אינה סימטרית. חוסר הסימטריה מתבטא בעיקר במחצית המרכזית של הנתונים, כלומר בקופסא. הממוצע גבוה מהחציון, ומרחקו של הרבעון העליון מהחציון גדול יותר מאשר מרחקו של הרבעון התחתון מהחציון. להתפלגות נטייה מימנה.
 ניתן ללמוד על פזור הנתונים סביב הממוצע והחציון.

Pvalue הנה ההסתברות לקבל תוצאה קיצונית לפחות כמו תוצאה שקיבלנו בניסוי, בהנחה שהשערת האפס נכונה.
 אם $\alpha \geq Pvalue$ נדחה את השערת האפס בר"מ α
 אחרת לא נדחה את השערת האפס בר"מ α
הגדרה 1 ההסתברות לקבל תוצאה מדגמית כפי שקיבלנו או קיצונית ממנה בכיוון H_1 .
הגדרה 2 מקסימאלית עבודה עדיין מקבלים את השערת H_0 .
או מינימאלית שממנה מתחילים לדחות את H_0 .

מקדם המתאם הליניארי בין שני משתנים מקריים מוגדר $\rho = \frac{Cov(X,Y)}{\sigma_X \sigma_Y}$ ($-1 \leq \rho \leq 1$)
 $Cov(X,Y) = E\{(X - E(X))(Y - E(Y))\} = E(XY) - E(X)E(Y)$

האומד למקדם המתאם הליניארי (מדגמי) מוגדר להיות $R = \frac{S_{XY}}{\sqrt{S_{XX}S_{YY}}} = \hat{\beta}_1 \sqrt{\frac{S_{XX}}{S_{YY}}}$

הערך של R אינו תלוי באיזה משתנה נקרא Y ואיזה X.
 הערך של R אינו תלוי ביחידות המדידה של המשתנים.
 תמיד בין מינוס 1 ל 1. אם הוא קרוב ל 0 אז הוא מעיד על חוסר קשר ליניארי בין המשתנים.

הסתברות חשובות

$$Cov(X, Y + Z) = Cov(X, Y) + Cov(X, Z)$$

$$Var(Y) = E(Var(Y|X)) + Var(E(Y|X))$$

אם $X_i \sim Pois(\lambda)$ אז $P(\sum X_i = s) = \frac{(n\lambda)^s e^{-n\lambda}}{s!}$ וגם $(X_i | \sum X_i = s) \sim Bin\left(s, \frac{1}{n}\right)$
חיבור פואסון $X + Y \sim Pois(\mu + \lambda)$

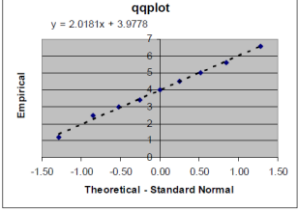
מג"מ X_1, \dots, X_n מ"מ ב"ת ש"ה $E(X_i) = \mu$ ו $Var(X_i) = \sigma^2$ אז $N(0,1) \rightarrow \frac{\bar{X}_n - \mu}{\frac{\sigma}{\sqrt{n}}}$
WALD אם X_i ב"ת ש"ה ו N בדיד אז $E(X_i) = \mu$ ו $Var(X_i) = \sigma^2$
אינדיקטורים: $E(I) = p$ וגם $E(I^2) = p^2$ וגם $Cov(X, Y) = 0$ וגם $Var(X) = pq$
מרקוב: $P(X \geq a) \leq \frac{E(X)}{a}$

ציבשב: $P(|X - \mu| \geq a) < \frac{Var(X)}{a^2}$ ו $E(X) = \mu$ ו $Var(X) = \sigma^2$
מינימום של אקספוננט הוא אקספוננט $X \sim \exp(\lambda), Y \sim \exp(\mu)$ אז $\min(X, Y) \sim \exp(\lambda + \mu)$

וגם **תחרות בין אקספוננטים** $P(X < Y) = \frac{\lambda}{\lambda + \mu}$

סכום גאמות אם $X \sim Gamma(r, \lambda)$ ו $Y \sim Gamma(s, \lambda)$ אז $X + Y \sim Gamma(r + s, \lambda)$
אקספוננציאלי $F_X(a) = P(X \leq a) = 1 - e^{-\lambda a}$

QQPLOT



i	1	2	3	4
$\frac{i}{n+1}$	0.1	0.2	0.3	0.4
שברנים תיאוריים ע"פ התפלגות נורמלית סטנדרטית	-1.282	-0.842	-0.524	-0.253
שברנים אמפיריים (התצפיות מסודרות מוקטנת גודלה)	1.2	2.5	3	3.4

זנבנות
 מדגם מקרי X_1, \dots, X_n מהתפלגות $Ber(p)$ נתעניין בבדיקת השערות $H_0: p = p_0$ ו $H_1: p = p_1$

מבחן זנב: נדחה את השערת האפס אם $\lambda(x) = \frac{\binom{n}{x} p_0^x (1-p_0)^{n-x}}{\binom{n}{x} p_1^x (1-p_1)^{n-x}} = \frac{p_1^x (1-p_0)^{n-x}}{p_0^x (1-p_1)^{n-x}} = \left(\frac{p_1(1-p_0)}{p_0(1-p_1)}\right)^x \left(\frac{1-p_1}{1-p_0}\right)^{n-x} > k_\alpha$

אם $p_1 > p_0$ אז $\ln\left(\frac{p_1(1-p_0)}{p_0(1-p_1)}\right) > \ln k_\alpha^* = k_\alpha^*$ ו $\ln\left(\frac{p_1(1-p_0)}{p_0(1-p_1)}\right) < \ln k_\alpha^* = k_\alpha^*$ אז $p_1 < p_0$

לכן, נדחה את השערת האפס אם $X = \sum_{i=1}^n X_i > C_\alpha$

לכן, נדחה את השערת האפס אם $X = \sum_{i=1}^n X_i < C_\alpha$

נמצא את הערך של C_α :

נניח $p_1 > p_0$ אזי המבחן יהיו $\sum_{i=1}^n X_i > C_\alpha$ דחה את השערת האפס אם $\sum_{i=1}^n X_i = C_\alpha$ ודחה את השערת האפס בהסתברות α אם $p_1 = 3/4, p_0 = 1/2, n = 10, \alpha = 0.05$

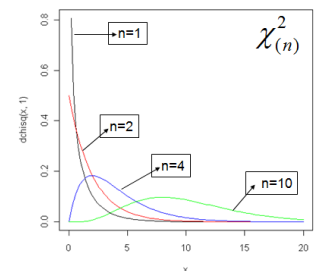
$$P_{H_0}(X = 10) + P_{H_0}(X = 9) = 0.5^{10} + 10 \cdot 0.5^{10} = 0.01074 < 0.05$$

תוצאות בדיקת השערות לדוגמא

מדגם מקרי $X_1, \dots, X_n \sim N(\mu, \sigma^2)$ והשונות ידועה.
 $H_0: \mu = \mu_0$
 $H_1: \mu = \mu_1$
 אם $\mu_0 < \mu_1$ דחה אם $\bar{X} > C_\alpha$
 $C_\alpha = \frac{\sigma}{\sqrt{n}} z_{1-\alpha} + \mu_0$
 $\pi = \Phi\left(\frac{\mu_1 - \mu_0}{\frac{\sigma}{\sqrt{n}}} - z_{1-\alpha}\right)$
 - פונקציית העוצמה עולה ב $\mu_1 - \mu_0$
 - פונקציית העוצמה עולה ב n
 - עבור $\mu_1 = \mu_0$ נקבל $\pi = \alpha$
 - גודל המדגם המינימאלי הדרוש לעוצמה ור"מ נתונות: $n = \frac{(z_{1-\alpha} + z_{1-\beta})^2 \sigma^2}{(\mu_1 - \mu_0)^2}$
 לסיכום נדחה אם $\bar{X} < \mu_0 - z_{1-\alpha} \frac{\sigma}{\sqrt{n}}$
 אם $\mu_0 > \mu_1$ דחה אם $\bar{X} < C_\alpha$, זאת אומרת $\frac{\bar{X} - \mu_0}{\frac{\sigma}{\sqrt{n}}} < -z_{1-\alpha} = z_\alpha$

התפלגות חי-בריבוע

יהיו $Z_1, \dots, Z_n \sim N(0,1)$ מקבל ערכים חיוביים בלבד - צפיפות לא סימטרית עם זנב ימני - עבור $n > 50$ ההתפלגות קרובה מאוד לנורמלית - סכום של מ"מ ב"ת המתפלגים חי-בריבוע, מתפלג חי-בריבוע
 אם $X_1, \dots, X_n \sim N(\mu, \sigma^2)$ אז $W_n = \frac{1}{\sigma^2} \sum_{i=1}^n (X_i - \mu)^2 \sim \chi^2_{(n)}$
 אם התוחלת איננה ידועה אז $W_n = \frac{1}{\sigma^2} \sum_{i=1}^n (X_i - \bar{X}_n)^2 \sim \chi^2_{(n-1)}$
 היותו \bar{X}_n ו $X_1 - \bar{X}_n$ **בלתי מתואמים** ומתפלגים נורמאליים, הם גם ב"ת ולפיכך **בלתי תלויים**, $\bar{X}_n, \sum_{i=1}^n (X_i - \bar{X}_n)^2$



$f_X(x) = \frac{1}{\Gamma(\frac{n}{2})} \left(\frac{x}{2}\right)^{\frac{n}{2}-1} e^{-\frac{x}{2}} I_{(0,\infty)}$
 $E(X) = n$, $Var(X) = 2n$, $E(e^{tX}) = \left(\frac{1}{1-2t}\right)^{\frac{n}{2}}$, $t < \frac{1}{2}$
 אם ידוע כי $Z_i \sim \exp(\lambda)$ מ"מ ב"ת אז $\sum_{i=1}^n Z_i \sim \text{Gamma}(n, \lambda)$ וגם:
 $2\lambda \sum_{i=1}^n Z_i \sim \text{Gamma}(n, \lambda) \sim \text{Gamma}\left(\frac{n}{2}, \frac{\lambda}{2}\right) \sim \chi^2_{(2n)}$, $2 \sum_{i=1}^n Z_i \sim \text{Gamma}\left(\frac{n}{2}, \frac{\lambda}{2}\right)$

בדיקת ההשערות $H_0: \theta \in W, H_1: \theta \in W^c$ כך ש $W \cup W^c = \Theta$
יחס נראות מוכלל מוגדר להיות:

אזי המבחן יהיה: דחה את השערת האפס אם $k_\alpha < \lambda$ כאשר k_α נקבע כך שלמבחן תהיה ר"מ α .
 $\Lambda(x) = \frac{\sup_{\theta \in W} f_\theta(x)}{\sup_{\theta \in \Theta} f_\theta(x)}$
 $0 \leq \Lambda(x) \leq 1$
 - בד"כ מבחן טוב, אך יש מקרים בהם יש טובים ממנו.
אם X_1, \dots, X_n מדגם מקרי, אזי $\Lambda^* = -2 \log \Lambda \sim \chi^2_{(k-1)}$
מוכלל זהו מבחן יחס נראות כאשר יש צורך לאמוד פרמטרים.
 במבחן יחס נראות מוכלל המונה והמכנה מתחלפים, וגם הסימון.

טיב התאמה למשתנה K קטגוריות (ערכים) A_1, \dots, A_K
 בהסתברויות p_1, \dots, p_K בהתאמה $(\sum_{i=1}^K p_i = 1)$.
 יש n תצפיות, כל תצפית מסווגת לאחת הקטגוריות.
 X_j מספר התצפיות בקטגוריה j אזי X_j מולטינומי.
לבדיקת ההשערות:
 כאשר $\vec{p}_0 = (p_{10}, \dots, p_{K0})$ ידועות
 וכן $\sum_{j=1}^K p_{j0} = 1$ תחת השערת האפס:
 $E(X_j) = np_{j0}$

קטגוריה	A_1	A_2	\dots	A_K	סה"כ
תצפיות	X_1	X_2	\dots	X_K	n
שכיחות צפויה	np_{10}	np_{20}	\dots	np_{K0}	n

$\sum_{j=1}^K E(X_j) = \sum_{j=1}^K np_{j0} = n \sum_{j=1}^K p_{j0} = n$

סטטיסטי המבחן **לבדיקת טיב התאמה** של המדגם למודל ההסתברותי מוגדר להיות:
 $\chi^2 = \sum_{j=1}^K \frac{(X_j - np_{j0})^2}{np_{j0}}$
טענה עבור מדגם "מספיק גדול" אם השערת האפס נכונה, אזי:
 $\chi^2 \sim \chi^2_{(k-1)}$, $[X_j | H_0] \sim \text{Pois}(np_{j0})$
 כאשר $k-1$ בגלל שהשכיחות בקטגוריה האחרונה נקבעת ע"ס כל הקודמות לה.
 מקובל להניח שהקירוב לעיל מספיק טוב אם $p_{j0} \geq 5$ לכל j .
 אם כלל האצבע לא מתקיים, נאחד קטגוריות סמוכות.
 ניתן להשתמש ב $\Lambda^* = -2 \log \Lambda \sim \chi^2_{(k-1)}$
הערה אם לצורך חישוב ההסתברויות תחת השערת האפס יש לאמוד r פרמטרים, אזי נאמד כל פרמטר על ידי אג"מ וד"ח יהיו $r-1$.
 כל שימוש באג"מ מוריד דרגת חופש.

ההחלטה	בדיקת השערות	
	דחיית H_0	אי דחיית H_0
טעות מסוג I	V	טעות מסוג II
טעות מסוג II	V	טעות מסוג I

טעות מסוג I: דחיית השערת האפס כאשר השערת האפס נכונה.
טעות מסוג II: אי-דחיית השערת האפס כאשר השערת האפס אינה נכונה.
 $\alpha = P(\text{type I error}) = P_{H_0}(\text{reject } H_0)$
 $\beta = P(\text{type II error}) = P_{H_1}(\text{accept } H_0)$
 בד"כ הקטנת ההסתברות לטעות מסוג ראשון גוררת הגדלת ההסתברות לטעות מסוג שני וההפוך.
השערה פשוטה היא השערה הקובעת באופן יחיד את ההתפלגות של האוכלוסייה שממנה נלקח המדגם. בכל מקרה אחר, ההשערה נקראת **השערה מורכבת** (מורכבת = לא שווה, גדול מ, קטן מ).
פונקציית המבחן לכל כלל הכרעה D מותאמת פונקציה $\phi_D(x) = \phi_D(x_1, \dots, x_n)$ (נקראת המבחן) $x \in X$ (אוסף התוצאות האפשריות) כך ש $0 \leq \phi_D(x) \leq 1$ קובעת את ההסתברות בה נדחה את השערות האפס עבור התוצאה x .
 $\alpha = E_{H_0}(\phi_D(X))$, $\beta = E_{H_1}(1 - \phi_D(X))$
רמת המובהקות (ר"מ) של המבחן היא ההסתברות המירבית המותרת לטעות מסוג ראשון. נקבעת לפי מידת הנוק - נוק רב גורר ר"מ נמוכה.
העוצמה של המבחן מוגדרת להיות ההסתברות לדחיית השערת האפס כאשר האלטרנטיבה נכונה. $\pi = 1 - \beta$.

יחס נראות $\lambda(x) = \frac{f_1(x)}{f_0(x)}$ (צפיפות או הסתברות).
הלמה של Neyman - Pearson (מבחן יחס-נראות) לבדיקת ההשערות הפשוטות
 $H_0: X \sim f_0$
 $H_1: X \sim f_1$
 ברמת מובהקות הקטנה או שווה ל α , מבחן בעל עוצמה מירבית הינו פונקציה יורדת של x נותן מבחן $x < C$

עבור k_α אי-שלילי, $0 \leq \tau_\alpha \leq 1$, וכן k_α ו τ_α נקבעים כך שיתקיים:
 $\alpha = P_{H_0}(\text{reject } H_0) = E_{H_0}(\phi(X))$
 מבין כל המבחנים עם ר"מ α (או קטנים יותר), מבחן י"נ הוא מבחן עם עוצמה מירבית.
טענה אם $f_0(\cdot) = f(\cdot, \theta_0)$, $f_1(\cdot) = f(\cdot, \theta_1)$ אזי מבחן י"נ תלוי ב X רק דרך ס"מ.

תהי A משפחה של התפלגויות של X ועלינו לבדוק את ההשערות
 $H_0: F = F_0$
 $H_1: F \in A$
 אם לכל $F_1 \in A$ המבחן בעל עוצמה מירבית בר"מ α של F_0 מול F_1 הוא אותו מבחן, אזי נאמר כי מבחן זה הוא בעל **עוצמה מירבית במידה שווה** ביחס ל A , בר"מ α .
 $X_1, \dots, X_n \sim N(\mu, \sigma^2)$ - $H_0: \mu = \mu_0, H_1: \mu > \mu_0$ השונות ידועה.
 מבחן בר"מ α בעל עוצמה מירבית הינו:
 $\bar{X} > C_\alpha = z_{1-\alpha} \frac{\sigma}{\sqrt{n}} + \mu_0$
 (אינו תלוי בערך של $\mu_0 > \mu_0$ ולכן זהו מבחן בעל עוצמה מירבית במידה שווה בר"מ α).
 - אם נשתמש במג"מ המבחן לא יהיה בעל עוצמה מירבית.

הקשר בין שני הסטטיסטים. בעזרת $(-1)^{n+1} \frac{\epsilon^n}{n!}$
 $\log(1 + \epsilon) = \epsilon - \frac{\epsilon^2}{2!} + \frac{\epsilon^3}{3!} - \dots$
 $\Lambda^*(x) = -2 \sum_{j=1}^k X_j \log \left(\frac{np_{0j} - 1 + X_j}{np_{0j}} \right) = -2 \sum_{j=1}^k X_j \left[\frac{np_{0j} - 1}{X_j} - \frac{1}{2} \left(\frac{np_{0j} - 1}{X_j} \right)^2 + \dots \right]$
 $-2 \sum_{j=1}^k (np_{0j} - X_j) + \sum_{j=1}^k \frac{1}{X_j} (np_{0j} - X_j)^2$
 $\sum_{j=1}^k \frac{1}{np_{0j}} (np_{0j} - X_j)^2 = \chi^2$
 $E_{H_0}(X_j) = np_{0j}$

אמידה
 $E(Y|X = x_0) = \beta_0 + \beta_1 x_0 \equiv \eta(x_0)$
 $\hat{\eta}(x_0) = \hat{\beta}_0 + \hat{\beta}_1 x_0 \sim N\left(\beta_0 + \beta_1 x_0, \sigma^2 \left[\frac{1}{n} + \frac{(x_0 - \bar{X})^2}{S_{XX}} \right]\right)$
 מאחר ו $(\hat{\beta}_0, \hat{\beta}_1)$ ב"ת $\hat{\sigma}^2$ אזי
 $\frac{\hat{\beta}_0 + \hat{\beta}_1 x_0 - (\beta_0 + \beta_1 x_0)}{\hat{\sigma} \sqrt{\frac{1}{n} + \frac{(x_0 - \bar{X})^2}{S_{XX}}}} \sim t_{(n-2)}$
 ולכן ר"ס עבור $\eta(x_0)$ (הקצר ביותר יתקבל בנקודה \bar{X}) הינו:
 $\hat{\beta}_0 + \hat{\beta}_1 x_0 \pm \hat{\sigma} \sqrt{\frac{1}{n} + \frac{(x_0 - \bar{X})^2}{S_{XX}}} \cdot t_{(n-2, 1-\frac{\alpha}{2})}$
ניבוי/תחזית נקבל תוספות בשונות עבור התצפית החדשה ולכן ר"ס:
 $\hat{\beta}_0 + \hat{\beta}_1 x_0 \pm \hat{\sigma} \sqrt{1 + \frac{1}{n} + \frac{(x_0 - \bar{X})^2}{S_{XX}}} \cdot t_{(n-2, 1-\frac{\alpha}{2})}$
 רווח סמך רחב יותר, והאורך **אינו** שואף לאפס כאשר גודל המדגם שואף לאינסוף.

רווח סמך ל σ^2 : $\sigma^2 \in \left[\frac{\hat{\sigma}^2 (n-2)}{\chi^2_{(n-2, 1-\frac{\alpha}{2})}}, \frac{\hat{\sigma}^2 (n-2)}{\chi^2_{(n-2, \frac{\alpha}{2})}} \right]$

	SS	DF	MS	F
Model	$SSR = \hat{\beta}_1^2 S_{XX} = \beta^T (X^T Y) - n(\bar{Y})^2$	r	$\frac{SSR}{r}$	$\frac{\hat{\beta}_1^2 S_{XX}}{\hat{\sigma}^2} = \frac{MSR}{MSE}$
Error	$SSE = \sum e^2 = S_{YY} - SSR$	$n - 1 - r$	$\hat{\sigma}^2 = \frac{SSE}{n - 1 - r}$	
Total	$SST = S_{YY} = \sum Y_i^2 - n(\bar{Y})^2$	$n - 1$		

$\hat{\eta}(x_0) = x_0^T \hat{\beta} \sim N(x_0^T \beta, \sigma^2 x_0^T (X^T X)^{-1} x_0) \rightarrow \hat{\eta}(x_0) \pm \hat{\sigma} \sqrt{x_0^T (X^T X)^{-1} x_0} t_{(n-p-1), 1-\frac{\alpha}{2}}$
 $Y_0 - x_0^T \hat{\beta} \sim N(0, \sigma^2 [1 + x_0^T (X^T X)^{-1} x_0]) \rightarrow x_0^T \hat{\beta} \pm \hat{\sigma} \sqrt{1 + x_0^T (X^T X)^{-1} x_0} t_{(n-p-1), 1-\frac{\alpha}{2}}$
 הראשון תוחלת והשני תחזית, כאשר $x_0^T \hat{\beta}$

עידו שמוי אביב 2010, מראה: פרופ' מלכה גורפיץ

אי-תלות המסקנה צריכה לגרום לשינוי הרגלים, אנחנו לעולם לא יודעים את הגורם.

תצפית מסווגת לפי שני משתנים: הקטגוריות של משתנה אחד: A1, ..., Aj, ... An הקטגוריות של משתנה שני: B1, ..., Bj, ... Bn

נסכם את הנתונים בטבלת שכיחויות דו-מימדית: טבלת שכיחויות בתא ij, שכיחות בשורה i, שכיחות בעמודה j

המטרה: לבדוק האם יש קשר בין שני המשתנים. ניסוח ההשערה: H0: P(Ai ∩ Bj) = P(Ai)P(Bj) vs H1: otherwise

לצורך בדיקת ההשערות לעיל עלינו למצוא את השכיחות הצפויה בכל אחת מהתאים שטבלה בהנחה שהשערת האפס נכונה.

גודל המספרים: טבלת תאונה קטלנית/לא קטלנית לפי גודלה/ביטוייה/קטנה

לכן, אם המשתנים ב"ת נגיעי: P(A1 ∩ B1) = 0.45 * 0.325 = 0.14625, P(A1 ∩ B2) = 0.45 * 0.375 = 0.16875, P(A1 ∩ B3) = 0.45 * 0.3 = 0.135

ולכן טבלת השכיחויות הצפויות הנה: טבלת שכיחויות צפויות

רגרסיה ליניארית מרובה עבור p משתנים מסבירים: Yi = beta_0 + beta_1 Xi1 + ... + beta_p Xi p + ei

תפולגות רב-נומאלית כאשר R^m מ-1 חיונית לחיטוי Vm x m: W ~ MVN(mu, V)

אזי נגזור ונקבל p + 1 משוואות נורמאליות X^T Y = X^T X b + e

מטריצת השונות המשותפת של beta היא sigma^2 (X^T X)^-1

רווח סמך לכל מקדם בנפרד. סטטיסטי המבחן (עבור 0 המשמעות היא תרומת המקדם למודל): beta_j - beta_j^0 / (sigma^2 (X^T X)^-1)^0.5 ~ t_{(n-p-1), alpha/2}

Stepwise Regression Methods: Forward - הוספת משתנה אחת בכל צעד. Backward - מתחילים מכל המודלים ומוצאים את מי שמעלה את ה SSE הכי מעט.

רגרסיה ליניארית פשוטה הנחות המודל: epsilon מתפלג נורמאלי עם sigma^2

Y = beta_0 + beta_1 x + epsilon, epsilon ~ N(0, sigma^2). E(Y|X=x) = beta_0 + beta_1 x, Var(Y|X=x) = sigma^2

ה x-ים הם נתונים תמיד מבחינתם. צריך לאמוד את x, ההתפלגות המותנה של Y בהיתן x היא מודל הפיזור סביב הקו.

אומד ריבועים פחותים הוא הישר המביא למינימום את סכום ריבועי המרחקים האנכיים בין כל הנקודות לישר.

היתן לאמוד את התוחלת המותנה של Y בהיתן x. ייתן לבצע עבור תצפית חדשה בטוח: Y-hat = beta_0 + beta_1 x

אמידת השונות (SSE) הן השונות הלא מוסברות, המטרה היא להקטין אותן, SST (קבוע).

טענה: תחת הנחת הנורמאליות, beta_0, beta_1, sigma^2 הם אומדי נראות מירבית.

בדיקת השערות על השיפוע: סטטיסטי המבחן: T = (beta_1 - beta_1^0) / (sigma_hat / sqrt(S_XX)) ~ t_{(n-2)}

כלל הכרעה עבור מבחן ברמת מובהקות alpha: A: P-value = P(T >= t) vs B: P-value = P(T <= t) vs C: P-value = P(|T| >= |t|)

ר"ס: beta_hat +/- t_{(n-1), (1-alpha/2)} * sigma_hat / sqrt(S_XX)

בדיקת מובהקות הרגרסיה עבור H0: beta_0 = 0: T^2 = (beta_hat_1 - beta_1^0)^2 / (sigma_hat^2 / S_XX) ~ F_{(1, n-2)}

האומד לריבוע מקדם המתאם הליניארי מוגדר על ידי R^2 = SSR / SST = 1 - SSE / SST