

EX | 7.2.1

Antag att vi har X_1, \dots, X_n från $\text{Bin}(20, p)$. Vi vill skatta p

$$E(X) = 20p = M_1 = \frac{1}{n} \sum_{i=1}^n X_i = \bar{X} \rightarrow 20p = \bar{X} \rightarrow p = \frac{\bar{X}}{20}$$

→ $\boxed{p = \frac{\bar{X}}{20}}$ för binomialfördelning

EX | $X \sim N(\mu, \sigma^2)$. Vi vill skatta μ och σ

$\left. \begin{array}{l} E(X) = M_1 \\ E(X^2) = M_2 \end{array} \right\}$ Ekvationssystemet vi ska lösa!

$\left. \begin{array}{l} \mu = \bar{X} \\ \sigma^2 + \mu^2 = \frac{1}{n} \sum_{i=1}^n X_i^2 \end{array} \right\} \rightarrow \left. \begin{array}{l} \mu = \bar{X} \\ \sigma^2 = \frac{1}{n} \sum_{i=1}^n X_i^2 - \mu^2 \end{array} \right\} \rightarrow$

$$\rightarrow \left. \begin{array}{l} \mu = \bar{X} \\ \sigma^2 = \frac{1}{n} \sum_{i=1}^n X_i^2 - \bar{X}^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2 \end{array} \right\}$$

MAXIMUM LIKELIHOOD METODEN (ML)

$$f_{X_1, \dots, X_n}(x_1, \dots, x_n) = \prod_{i=1}^n \underbrace{f_{X_i}(x_i)}_{f(x_i)} = \prod_{i=1}^n \underbrace{f(x_i)}_{\text{Likelihood}}$$

DEFINITION: $L(\theta) = \prod_{i=1}^n f(x_i)$ kallas likelihood

DEFINITION: $\hat{\theta}_{ML} = \arg\max L(\theta)$ är ML skattaren.

Det är ofta lättare att beräkna $\ln L(\theta) = \sum_{i=1}^n \ln f(x_i)$

$\boxed{\frac{\partial}{\partial \theta} \ln L(\theta) = 0}$ → maximum likelihood metoden

EX1 Antag X_1, \dots, X_n är ett stickprov från $X \sim \text{Po}(\mu)$

$$\begin{aligned} \ln L(\theta) &= \sum_{i=1}^n \ln \left(\underbrace{\frac{\mu^{x_i}}{x_i!} \cdot e^{-\mu}}_{P_X(x_i)} \right) = \sum_{i=1}^n (x_i \ln \mu - (-\mu) - \ln(x_i!)) = \\ &= \ln \mu \sum_{i=1}^n x_i - n\mu - \sum_{i=1}^n \ln(x_i!) \end{aligned}$$

$$\frac{\partial}{\partial \mu} \ln L(\mu) = \frac{1}{\mu} \sum_{i=1}^n x_i - n = 0 \rightarrow \mu = \frac{1}{n} \sum_{i=1}^n x_i = \bar{X}$$

stickprovsmedelvärdet brukar vara en bra skattare.

EX1 Antag X_1, \dots, X_n stickprov från $X \sim N(\mu, \sigma^2)$

$$\ln L(\mu, \sigma) = \sum_{i=1}^n \ln \left(\frac{1}{\sigma \sqrt{2\pi}} \cdot e^{-\frac{1}{2\sigma^2} (x_i - \mu)^2} \right) = \dots = -n\sqrt{2\pi} - n \ln \sigma - \frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2$$

$$\textcircled{1} \frac{\partial}{\partial \mu} \ln L(\mu) = -\frac{1}{\sigma^2} \sum_{i=1}^n (x_i - \mu) = 0 \quad \textcircled{2} \frac{\partial}{\partial \sigma} \ln L(\sigma) = -\frac{n}{\sigma} + \frac{1}{\sigma^3} \sum_{i=1}^n (x_i - \mu)^2 = 0$$

Lös ut μ och σ^2

$$\textcircled{1} \frac{\partial}{\partial \mu} \ln L(\mu) = 0 \rightarrow \mu = \frac{1}{n} \sum_{i=1}^n x_i = \bar{X}$$

$$\text{Stoppa in } \textcircled{1} \text{ i } \textcircled{2} \rightarrow \frac{\partial}{\partial \sigma} \ln L(\sigma) = 0 = -\frac{n}{\sigma} + \frac{1}{\sigma^3} \sum_{i=1}^n (x_i - \bar{X})^2 \rightarrow$$

$$\rightarrow \sigma^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{X})^2$$

TVÅSIDIGT KONFIDENSINTERVALL

DEFINITION: Låt X_1, \dots, X_n vara slumpvariabler med fördelningen som har en parameter θ med ett sant okänt värde θ_0 .

Ett konfidensintervall för θ med konfidensgrad $1-\alpha$ är:

$[a(X_1, \dots, X_n), b(X_1, \dots, X_n)]$ så att $P(a < \theta_0 < b) = 1-\alpha$

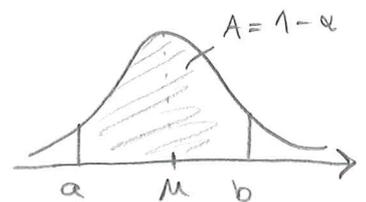
KONFIDENSINTERVALL FÖR μ I $N(\mu, \sigma^2)$ DÅ σ^2 ÄR KÄND

$X_1, \dots, X_n \sim N(\mu, \frac{\sigma^2}{n}) \xrightarrow{*} \bar{X} \sim N(\mu, \frac{\sigma^2}{n})$ om σ^2 är känd.

kom ihåg: σ^2 -regeln

Om $Y \sim N(\mu, \sigma^2) \rightarrow P(\mu - z_{\alpha/2} \sigma < Y < \mu + z_{\alpha/2} \sigma) = 1-\alpha$

Summan av normalfördelade variabler är också normalfördelad.*



SATS

Om $X_i \sim N(\mu, \sigma^2)$ och σ^2 är känd så är:

$I_\mu = (\bar{X} - z_{\alpha/2} \frac{\sigma}{\sqrt{n}}, \bar{X} + z_{\alpha/2} \frac{\sigma}{\sqrt{n}})$ ett konf-intervall för μ med konf-grad $1-\alpha$

BÄVIS

$P(\mu - z_{\alpha/2} \frac{\sigma}{\sqrt{n}} < \bar{X} < \mu + z_{\alpha/2} \frac{\sigma}{\sqrt{n}}) = 1-\alpha$ vilket är samma som \rightarrow

$\rightarrow P(-\bar{X} - z_{\alpha/2} \frac{\sigma}{\sqrt{n}} < -\mu < -\bar{X} + z_{\alpha/2} \frac{\sigma}{\sqrt{n}}) = 1-\alpha \quad \parallel \quad \rightarrow$

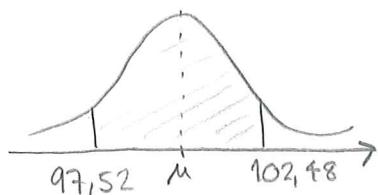
$\rightarrow P(\underbrace{\bar{X} - z_{\alpha/2} \frac{\sigma}{\sqrt{n}}}_{\text{punkt a}} < \mu < \underbrace{\bar{X} + z_{\alpha/2} \frac{\sigma}{\sqrt{n}}}_{\text{punkt b}}) = 1-\alpha$

punkt a

punkt b.

EX Antag $X \sim N(\mu, 16)$ $n=10$ $\bar{X}=100$ $\alpha=0,05 \rightarrow Z_{\alpha/2}=1,96$

$$I_{\mu} = (100 - 1,96 \times \frac{4}{\sqrt{100}}, 100 + 1,96 \times \frac{4}{\sqrt{100}}) = (97,52, 102,48)$$



- * ju mindre α desto större slh att intervallet täcker det sanna värdet
- * ju mindre α desto bredare intervall

KONFIDENSINTERVALL FÖR μ I $N(\mu, \sigma^2)$ DÅ σ^2 OKÄND

kom ihåg: Om $X \sim N(\mu, \sigma^2) \rightarrow \frac{X - \mu}{\sigma} \sim N(0,1)$

SATS: Antag $X_1, \dots, X_n \sim N(\mu, \sigma^2)$ Då är $\frac{\bar{X} - \mu}{S/\sqrt{n}} \sim t(n-1)$

Då $n \rightarrow \infty$ i T-fördelningen går $t(n-1) \rightarrow N(\mu, \sigma)$

↑
T-fördelad med $n-1$ frihetsgrader

Täthetsfunktionen för T-fördelning

$$f(x) = \frac{\Gamma\left(\frac{n+1}{2}\right)}{\Gamma\left(\frac{n}{2}\right)} \left(1 + \frac{x^2}{L}\right)^{-\frac{n+1}{2}} \quad \Gamma = \text{gammafunktioner}$$

Vi kommer inte behöva kunna dessa. Men bra att se

SATS Om $X_1, \dots, X_n \sim N(\mu, \sigma^2)$ där σ^2 är okänd så är:

$$I_{\mu} = \left(\bar{X} - t_{\alpha/2}(n-1) \frac{S}{\sqrt{n}}, \bar{X} + t_{\alpha/2}(n-1) \frac{S}{\sqrt{n}} \right) \text{ ett konfidensintervall m. konfgrad } 1-\alpha$$

$\alpha/2$ kvantiler för $t(n-1)$ fördelningen
Anns i tabeller i boken.

T-fördelningsskvantilen

$$t_{0,025}(\infty) = Z_{0,025} = 1,96 \quad \text{alltså } t_{0,025} \rightarrow Z_{0,025} \text{ då } n \rightarrow \infty$$

$$t_{0,025}(1) = 12,71$$

↑
extremt dåligt fall eftersom $n-1=1$ ger $n=2!$

CENTRALA GRÄNSVÄRDES SATSEN (CGS)

$\bar{X} \in N(\mu, \frac{\sigma^2}{n})$ detta ger då ett approximativt konf, intervall för μ
då:

För σ känt

$$I_\mu = (\bar{X} - z_{\alpha/2} \frac{\sigma}{\sqrt{n}}, \bar{X} + z_{\alpha/2} \frac{\sigma}{\sqrt{n}})$$

För σ okänd

$$I_\mu = (\bar{X} - z_{\alpha/2} \frac{s}{\sqrt{n}}, \bar{X} + z_{\alpha/2} \frac{s}{\sqrt{n}})$$

KONFIDENSINTERVALL FÖR σ^2

Kom ihåg: $\frac{\bar{X} - \mu}{s/\sqrt{n}} \rightarrow$ Vi behöver alltså hitta ett motsvarande uttryck för att få σ^2

DEFINITION χ^2 -fördelning (chi-två fördelning)

Om $z_i \sim N(0,1)$ så gäller $\sum_{i=1}^n z_i^2 \sim \chi^2(n)$

χ^2 -fördelad med n frihetsgrader.

Täthetsfunktionen för χ^2 -fördelning

$$f(x) = \begin{cases} \frac{x^{\frac{n}{2}-1} \cdot e^{-x/2}}{2^{n/2} \cdot \Gamma(\frac{n}{2})} & x \geq 0 \\ 0 & x < 0 \end{cases}$$

Behöver inte kunna denna men bra att se.

$$z_i = \frac{X_i - \mu}{\sigma} \sim N(0,1) \text{ om } X_i \sim N(\mu, \sigma^2)$$

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

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

vi tappar en frihetsgrad eftersom vi skattar μ med \bar{X}

Låt $\chi^2_\alpha(n)$ vara α -kvantilen i $\chi^2(n)$ -fördelningen.

$$P\left(\chi^2_{1-\alpha/2}(n-1) \leq \frac{(n-1)s^2}{\sigma^2} \leq \chi^2_{\alpha/2}(n-1)\right) = 1-\alpha$$

Vi vill lösa ut σ^2 för att få ut konf.-intervallet för σ^2

$$P\left(\underbrace{\frac{(n-1)s^2}{\chi^2_{\alpha/2}(n-1)}} \leq \sigma^2 \leq \underbrace{\frac{(n-1)s^2}{\chi^2_{1-\alpha/2}(n-1)}}\right) = 1-\alpha$$

våra gränsvärden i konf. intervallet

SATS Om $X_i \sim N(\mu, \sigma^2)$ så är

$$I_{\sigma^2} = \left(\underbrace{\frac{(n-1)s^2}{\chi^2_{\alpha/2}(n-1)}}_{av}, \underbrace{\frac{(n-1)s^2}{\chi^2_{1-\alpha/2}(n-1)}}_{av} \right) \text{ ett konf. intervall för } \sigma^2$$

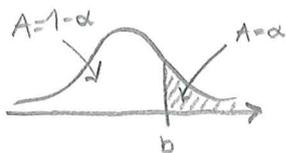
Man tar bara roten ur gränserna för att få σ

$$I_\sigma = \left(\sqrt{\frac{(n-1)s^2}{\chi^2_{\alpha/2}(n-1)}}, \sqrt{\frac{(n-1)s^2}{\chi^2_{1-\alpha/2}(n-1)}} \right) \text{ konf. intervall för } \sigma$$

Värden för χ^2 -kvantiler finns i boken.

NOTERA! Dessa konfidensintervall för σ^2 måste vara normalfördelade.

ENSIDIGT KONFIDENSINTERVALL

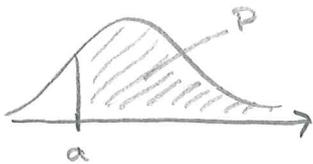


Tex. hur stor/liten kan fördelningen vara
Alltså finna övre/undre gräns

EX]

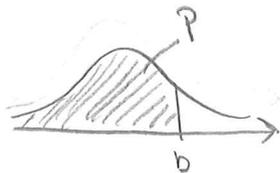
Ett undre begränsat konfidensintervall för θ ges av:

$$(a(\bar{X}_1, \dots, \bar{X}_n), \infty) \text{ Där } P(a < \theta_0) = 1 - \alpha$$



Ett övre begränsat konfidensintervall för θ ges av:

$$(-\infty, b(\bar{X}_1, \dots, \bar{X}_n)) \text{ Där } P(\theta_0 < b) = 1 - \alpha$$



EXI Ett ensidigt kont, intervall för μ i $N(\mu, \sigma^2)$

Med känt σ^2 och övre begränsat

$$I_{\mu} = (-\infty, \bar{X} + t_{\alpha}(n-1) \frac{s}{\sqrt{n}})$$

Med okänt σ^2 och undre begränsat

$$I_{\mu} = (\bar{X} - z_{\alpha} \frac{\sigma}{\sqrt{n}}, \infty)$$

Skillnaden mellan en- och tvåsidigt intervall är att man kvantiler från $\alpha/2$ till α för ensidigt och man sätter undre gränsen till $-\infty$ om man söker ett övre begränsat ensidigt intervall, och tvärtom för undre begränsat.

HYPOTESPRÖVNING

EXEMPEL: Hypotesprövningsträgeställningar:

- * Gör ett nytt läkemedel någon effekt?
- * Dör rökare än icerökare?
- * Har mätinstrument ett systematiskt fel?

Vi vill testa nollhypotesen $H_0: \theta = \theta_0$ mot mothypotesen $H_1: \theta \neq \theta_0$.

- Om testet tex bekräftar ett systemfel så förkastar vi H_0 till förmån för $H_1 \rightarrow \theta$ är signifikant skild från θ_0 .

○ EX 1.] En nyanställd mäter klorhalten 5 ggr

56,7 61,3 58,5 62,1 59,5

Sann klorhalt är 60 mg/l

Kan vi säga att den nyanställda systematiskt mäter fel?

$$\bar{x} = 59,62 \quad s^2 = 4,692$$

Antag $x_i \sim N(\mu, \sigma^2)$

Vi vill testa: $H_0: \mu = 60$ $H_1: \mu \neq 60$

KONFIDENSINTERVALLMETODEN

Om vi bildar ett 95% konfidensintervall I_θ för θ kan vi direkt använda detta för att göra hypotestestet.

- Eftersom intervallet är beräknat så att det ska täcka det sanna värdet på θ : 95% av fallen så kan vi förkasta H_0 om intervallet inte innehåller θ_0 .
- Om intervallet däremot innehåller θ_0 kan vi inte förkasta H_0 dvs påstå $\theta \neq \theta_0$.

Konfidensgraden vi använder när vi beräknar konfidensintervall är också konfidensgraden för hypotestestet vi gör.

- Det har blivit standard inom många områden att använda 95% konfidensgrad, men vi kan också använda högre konfidensgrad (tex 99%) eller en lägre konfidensgrad (tex 90%)

FORTS Ex 1

EH 95% konfidensintervall $\rightarrow \alpha = 0,05$

$$I_{\mu} = (\bar{x} \pm t_{\alpha/2} (n-1) \frac{s}{\sqrt{n}}) = (56,93 \quad 62,81)$$

$60 \in I_{\mu} \rightarrow$ vi kan ej förkasta H_0

DEFINITION Felrisk/signifikansnivå Typ 1. fel

Felrisken eller signifikansnivån definieras som:

$$\alpha = P(H_0 \text{ förkastas} \mid H_0 \text{ sann})$$

Översatt till konfidensintervall är signifikansnivån:

$$\alpha = P(\theta \notin I_{\theta} \text{ om } \theta = \theta_0) = P(\theta_0 \notin I_{\theta})$$

Alltså fås signifikansnivån som 1-konfidensgraden.

DEFINITION Typ 2 fel

Under hypotesprövning säger vi att vi gör ett typ 2 fel om vi inte förkastar H_0 trots att H_1 är sann. Sannolikheten för att göra ett typ 2 fel brukar betecknas med β

$$\beta = P(H_0 \text{ förkastas inte} \mid H_1 \text{ sann})$$

OLIKA UTFALL I HYPOTESTEST:

- ① H_0 sann och hypotestestet förkastar inte H_0
- ② H_0 sann och ——— " ——— förkastar H_0
Detta är ett typ 1 fel och har enligt konstruktionen SLH α att inträffa
- ③ H_1 är sann och hypotestestet förkastar H_0
- ④ H_1 är sann och hypotestestet förkastar inte H_0 .
Detta är ett typ 2 fel och SLH för att detta inträffar brukar betecknas med β

I boken kallar dom hypotestest med konfidensintervall för hypotesprövning och det med teststorhet för signifikansprövning

DEFINITION Teststorheter

En teststorhet $T = T(X_1, \dots, X_n)$ är en funktion av observationerna, och alltså en slumpvariabel. $T_{obs} = T(x_1, \dots, x_n)$ är ett observerat värde av teststorheten för givna observationer.

I allmänhet innehåller teststorheten en skattning $\hat{\theta}^*$ av parametern θ . Tex för test av väntevärdet i en normalfördelning med:

Okänd varians:

$$T = \frac{\bar{X} - \mu_0}{s/\sqrt{n}}$$

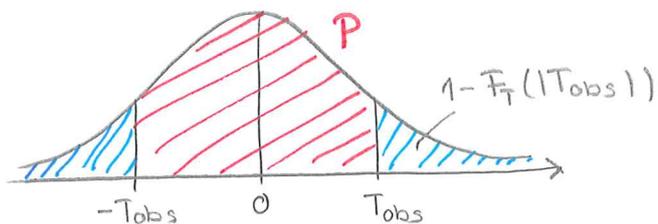
känd varians:

$$T = \frac{\bar{X} - \mu_0}{\sigma/\sqrt{n}}$$

DEFINITION P-värde

P-värdet eller signifikanssannolikheten definieras som sannolikheten under nullhypotesen att vi får ett värde $|T|$ som är lika stort eller större än det observerade värdet $|T_{obs}|$

Symmetrisk fördelning



$$P = 2(1 - F_T(|T_{obs}|))$$

EX]

- Om $0,01 < P < 0,05$ kallas det enstjärnigsignifikant (*)
- Om $0,001 < P < 0,01$ ——— " ——— två ——— " ——— (**)
- Om $P < 0,001$ ——— " ——— tre ——— " ——— (***)

Ju högre "stjärnfaktor" desto större SLH att vi förkastar rätt

FORT EX 1]

$$T = \frac{\bar{X} - b_0}{s/\sqrt{n}} \sim t(n-1)$$

$$T_{obs} = \frac{59,62 - 60}{\sqrt{4,692}/\sqrt{5}} = -0,3933$$

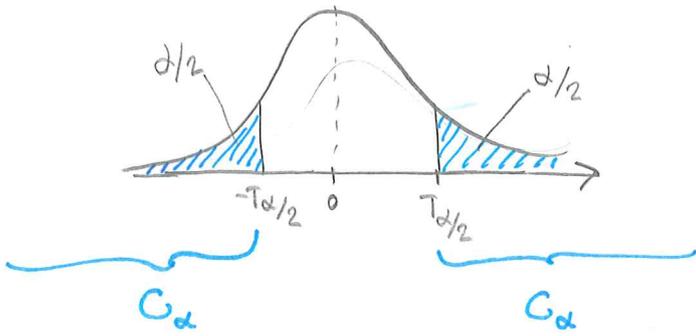
$$P = 2(1 - F_{t(4)}(0,3933)) = 0,71 \rightarrow \text{minsta felrisken vi kan förkasta } H_0 \text{ på är 71\%}$$

$$P > 0,05 \rightarrow \text{förkastar ej } H_0$$

DEFINITION Kritiskt område

Givet en signifikansnivå α definierar vi det kritiska området C_α som de värden på teststorheten T som leder till att man förkastar H_0 på nivå α .

Förkastas på nivå α $p < \alpha \rightarrow 1 - F_T(|T_{\text{obs}}|) < \frac{\alpha}{2}$



C_{α} = kritiska området. Befinner sig T_{obs} inom detta området så förkastas H_0 vilket är ekvivalent med att säga att om vi hamnar utanför konfidensintervallet så förkastas H_0

KRITISKT OMRÅDE FÖR NORMALFÖRDELNING

Med test för väntevärdet för normalfördelade data får vi

- Om variansen är känd: T under H_0 är $N(0,1)$ -fördelad förkastas H_0 på nivån α om $|T| > z_{\alpha/2}$
- Om variansen är okänd: T under H_0 är $t(t)$ -fördelad, förkastas H_0 på nivå α om $|T| > t_{\alpha/2}(t)$

$$C_{\alpha} = \{T: |T| > z_{\alpha/2}\} \rightarrow \text{Förkastas } H_0 \text{ om } |T| > z_{\alpha/2}$$

$$|T| = \left| \frac{\bar{x} - \mu_0}{\sigma/\sqrt{n}} \right| \rightarrow \left| \frac{\bar{x} - \mu_0}{\sigma/\sqrt{n}} \right| > z_{\alpha/2} \iff (\bar{x} - \mu_0) > z_{\alpha/2} \cdot \frac{\sigma}{\sqrt{n}} \iff$$

$$\iff \left\{ \begin{array}{l} \mu_0 < \bar{x} - z_{\alpha/2} \frac{\sigma}{\sqrt{n}} \\ \mu_0 > \bar{x} + z_{\alpha/2} \frac{\sigma}{\sqrt{n}} \end{array} \right\} \iff \mu \notin I_{\mu} = \left(\bar{x} \pm z_{\alpha/2} \frac{\sigma}{\sqrt{n}} \right)$$

Här ifrån kommer kopplingen mellan hypotesprövningen och konfidensintervall

FORTS EX 1

$$C_{\alpha} = \{T: |T| > t_{\alpha/2}(4)\} \quad \text{så} \quad \alpha = 0,05 \rightarrow C_{\alpha} = \{T: |T| > 2,77\}$$

ENSIDIGTTEST

Nollhypotesen $H_0: \theta = \theta_0$ testas mot antingen

Mothypotesen $H_1: \theta > \theta_0$ eller $H_1: \theta < \theta_0$.

Använder vi konfidensintervall så:

• Om $H_1: \theta > \theta_0$ $I_\theta = \{a, \infty\}$

• Om $H_1: \theta < \theta_0$ $I_\theta = \{-\infty, b\}$

Förkastar H_0 om intervallet inte täcker θ_0

Använder vi teststorhet så:

• Om $H_1: \theta > \theta_0$ $P = P_{H_0}(T \geq T_{obs})$

• Om $H_1: \theta < \theta_0$ $P = P_{H_0}(T \leq T_{obs})$

Förkastar sedan H_0 om $p \leq \alpha$

KRITISKT OMRÅDE - ENSIDIGT

Om $H_1: \theta > \theta_0$

• Om σ är känd \rightarrow förkasta H_0 på nivå α om $T > z_\alpha$

• Om σ är okänd \rightarrow förkasta H_0 på nivå α om $T > t_\alpha(f)$

Om $H_1: \theta < \theta_0$

• Om σ är känd \rightarrow _____ " _____ $T < -z_\alpha$

• Om σ är okänd \rightarrow _____ " _____ $T < -t_\alpha(f)$

DEFINITION styrkefunktion

Vi vill testa $H_0: \theta = \theta_0$ och har bestämt en teststorhet T och en signifikansnivå α . Testets styrkefunktion $\pi(\theta)$ definieras som:

$$\pi(\theta) = P(T \in C_\alpha | \theta) = P(H_0 \text{ förkastas} | \text{parametervärdet är } \theta)$$

SATS

Antag att vi har n observationer av en normalfördelning $N(\mu, \sigma^2)$ med känd varians och att vi vill testa $H_0: \mu = \mu_0$

Testets styrka ges då av

Om $H_1: \mu > \mu_0$

$$\pi(\mu) = 1 - \Phi\left(z_\alpha - \frac{(\mu - \mu_0)\sqrt{n}}{\sigma}\right)$$

Om $H_1: \mu < \mu_0$

$$\pi(\mu) = \Phi\left(z_\alpha - \frac{(\mu - \mu_0)\sqrt{n}}{\sigma}\right)$$

Om $H_1: \mu \neq \mu_0$

$$\pi(\mu) = 1 + \Phi\left(-z_{\alpha/2} - \frac{(\mu - \mu_0)\sqrt{n}}{\sigma}\right) - \Phi\left(z_{\alpha/2} - \frac{(\mu - \mu_0)\sqrt{n}}{\sigma}\right)$$

SATS

Antag att vi har normalfördelade observationer $N(\mu, \sigma^2)$ och önskar testa $H_0: \mu = \mu_0$. Om vi vill att testet ska ge utslag för avvikelser $\mu - \mu_0$ med sln $1 - \beta$ ska antal observationer väljas som:

* känd varians:

Ensidigt $n = \frac{\sigma^2}{(\mu - \mu_0)^2} (z_\alpha + z_\beta)^2$

Tvåsidigt $n = \frac{\sigma^2}{(\mu - \mu_0)^2} (z_{\alpha/2} + z_\beta)^2$

* Om variansen skattas med s^2

Ensidigt $n \approx \frac{\sigma^2}{(\mu - \mu_0)^2} (z_\alpha + z_\beta)^2 + \frac{z_\alpha^2}{2}$

Tvåsidigt $n \approx \frac{\sigma^2}{(\mu - \mu_0)^2} (z_{\alpha/2} + z_\beta)^2 + \frac{z_{\alpha/2}^2}{2}$

"JÄMFÖRELSE"

NOLLHYPOTES: $H_0: \mu = \theta$ eller $H_0: \mu \geq \theta$ osv.

MOTHYPOTES: $H_1: \mu \neq \theta$ eller $H_1: \mu < \theta$

Sannolikheten att inte göra ett Typ 2 fel kallas testets styrka och ges av $1 - \beta$ där $\beta = P(H_0 \text{ förkastas ej} | H_1 \text{ sann})$

P-värdet ger minsta signifikansnivå som vi kan förkasta H_0 på.

OBEROENDE STICKPROV

Antag att vi har:

$$X_{11}, \dots, X_{1n_1} \sim N(\mu_1, \sigma_1^2)$$

Vi vill testa om $\mu_1 = \mu_2!$

$$X_{21}, \dots, X_{2n_2} \sim N(\mu_2, \sigma_2^2)$$

$$\bar{X}_1 = \frac{1}{n_1} \sum_{i=1}^{n_1} X_{1i} \quad S_1^2 = \frac{1}{n_1 - 1} \sum_{i=1}^{n_1} (X_{1i} - \bar{X}_1)^2$$

$$\bar{X}_2 = \frac{1}{n_2} \sum_{i=1}^{n_2} X_{2i} \quad S_2^2 = \frac{1}{n_2 - 1} \sum_{i=1}^{n_2} (X_{2i} - \bar{X}_2)^2$$

Definiera $\theta = \mu_1 - \mu_2$

Skattning av $\theta \rightarrow \hat{\theta} = \bar{X}_1 - \bar{X}_2$

Vi vill testa: $H_0: \theta = 0$ ($\mu_1 = \mu_2$) $H_1: \theta \neq 0$ ($\mu_1 \neq \mu_2$)

I regel har vi tre olika fall:

- 1) σ_1 och σ_2 kända
- 2) $\sigma_1 = \sigma_2 = \sigma$ okända
- 3) $\sigma_1 \neq \sigma_2$ samt okända

HUR GÖRS TESTET I DE OLIKA FALLEN?

1) σ_1 och σ_2 okända

$$V(\theta^*) = \frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2} \quad (\text{ty oberoende})$$

$$I_\theta = (\bar{X}_1 - \bar{X}_2 \pm z_{\alpha/2} \sqrt{V(\theta^*)})$$

↑ Eftersom de var normalfördelade

För hypotestest

$$T = \frac{\bar{X}_1 - \bar{X}_2}{\sqrt{V(\theta^*)}} \sim N(0,1) \text{ om } H_0 \text{ är sann}$$

Ex 1

Antag att vi har 2 stickprov

$$\bar{X}_1 = 2,45 \quad n_1 = 10$$

$$\bar{X}_2 = 3,23 \quad n_2 = 20$$

Vi känner till att $\sigma_1 = \sigma_2 = 1/2$

$$V(\theta^*) = \frac{(1/2)^2}{10} + \frac{(1/2)^2}{20} = 0,1936^2$$

$$I_\theta = (2,45 - 3,23 \pm 1,96 \times 0,1936) = (-1,9, -0,4)$$

$$T_{\text{obs}} = \frac{2,45 - 3,23}{0,1936} \approx -4$$

$$P = 2(1 - \Phi(4)) = 5 \times 10^{-5}$$

H_0 förkastas på väldigt hög signifikansnivå

2) $\sigma_1 = \sigma_2 = \sigma$ okänd

Om σ_1 och σ_2 är okända, skatta $V(\theta^*)$ med:

$$S^2 = \frac{S_1^2}{n_1} + \frac{S_2^2}{n_2}$$

Eftersom $\sigma_1 = \sigma_2$ är S_1 och S_2 2 olika skattningar av samma varians.

SATS

En väntevärdesriktig skattning av σ^2 från k normalfördelningar $N(\mu, \sigma^2)$ fås som

$$S_p^2 = \frac{(n_1-1) \cdot S_1^2 + (n_2-1) \cdot S_2^2 + \dots + (n_k-1) \cdot S_k^2}{(n_1-1) + (n_2-1) + \dots + (n_k-1)}$$

↑
står för polad

Dessutom gäller att $\frac{(N-k)S_p^2}{\sigma^2} \sim \chi^2(N-k)$ totala antalet observationer $N = \sum_{i=1}^k n_i$

BEVIS

Visa att $E(S_p^2) = \sigma^2$ följer av att $E(S_i^2) = \sigma^2$

Vi har att $\frac{(n_i-1)S_i^2}{\sigma^2} \sim \chi^2(n_i-1)$

$\chi^2(N-k)$ följer av additionssatsen för χ^2 -fördelningen

Satsen säger att $S_p^2 = \frac{(n_1-1)S_1^2 + (n_2-1)S_2^2}{(n_1-1) + (n_2-1)}$

Vi skattar $V(\theta^*)$ med: $S_p^2 \cdot \left(\frac{1}{n_1} + \frac{1}{n_2}\right)$

$T = \frac{\bar{X}_1 - \bar{X}_2}{S_p \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}}$ är t -fördelad ty skattad $t(n_1+n_2-2)$ under H_0

EX 2. | samma stickprov som i Ex 1 men vi vet även:

$$S_1^2 = 0,57^2 \quad S_2^2 = 0,48^2 \quad \text{vi vet att } \sigma_1 = \sigma_2 = \sigma$$

Beräkna: $S_p^2 = \frac{9 \times 0,57^2 + 19 \times 0,48^2}{9 + 19} = 0,5107^2$

$$V(\theta^*) = 0,5107^2 \left(\frac{1}{10} + \frac{1}{20} \right) = 0,1978^2$$

$$I_\theta = \bar{X}_1 - \bar{X}_2 \pm t_{\alpha/2(28)} \sqrt{V(\theta^*)} = (-1,18, -0,38)$$

$$T_{\text{obs}} = \frac{2,45 - 3,23}{0,1978} = -3,94$$

$$P = 2(1 - F_{t(28)}(|T_{\text{obs}}|)) = 0,0004 \quad \text{väntevärdena är skilda}$$

Vi kan förkasta H_0 på en hög signifikansnivå.

3.) $\sigma_1 \neq \sigma_2$ och okända

$$V(\theta^*) = \frac{S_1^2}{n_1} + \frac{S_2^2}{n_2}$$

Vi kan ej pola stickprovsvariansen
ty $\sigma_1 \neq \sigma_2$

SATS

Om σ_1 och σ_2 okända så gäller att:

$$T = \frac{\bar{X}_1 - \bar{X}_2}{\sqrt{V(\theta^*)}} \quad \text{är approximativt } t\text{-fördelad} \\ \text{med } f \text{ frihetsgrader } t_{(f)}$$

$$f = \frac{\left(\frac{S_1^2}{n_1} + \frac{S_2^2}{n_2} \right)^2}{\frac{\left(\frac{S_1^2}{n_1} \right)^2}{n_1 - 1} + \frac{\left(\frac{S_2^2}{n_2} \right)^2}{n_2 - 1}} = \text{skattning av antalet} \\ \text{frihetsgrader}$$

* Om $\frac{S_1^2}{n_1} \gg \frac{S_2^2}{n_2}$ blir antalet frihetsgrader $n_1 - 1$

* Om $S_1 \approx S_2$ och $n_1 \approx n_2 \approx n$ blir $f \approx 2(n - 1)$

EX 3. | Samma värden som i EX 2 fast $\sigma_1 \neq \sigma_2$

Skatta: $V(\theta^*) \approx \frac{0,57^2}{10} + \frac{0,48^2}{20} = 0,2098^2$

$$f = \frac{\left(\frac{0,57^2}{10} + \frac{0,48^2}{20}\right)^2}{\frac{\left(\frac{0,57^2}{10}\right)^2}{9} + \frac{\left(\frac{0,48^2}{20}\right)^2}{19}} = 15,5$$

Alltså hälften så många frihetsgrader som i EX 2.

Jämför med den polade stickprovsvariansskattningen $S_p^2 = 0,1978^2$, vilket är nära

$$I_\theta = (\bar{X}_1 - \bar{X}_2 \pm t_{\alpha/2}(15) \cdot 0,2098) = (-1,23, -0,133)$$

$$T_{\text{obs}} = -3,718$$

$$P = 2(1 - F_{t(15)}(3,718)) = 0,002 \quad \text{kan förtasta } H_0 \text{ på hög signifikansnivå}$$

JÄMFÖRELSEAVVIK

JÄMFÖRELSE AV VARIANSEN

Vi vill testa $H_0: \sigma_1^2 = \sigma_2^2$ $H_1: \sigma_1^2 \neq \sigma_2^2$

Inför teststorheten: $T = \frac{S_1^2/\sigma_1^2}{S_2^2/\sigma_2^2} \quad (*)$

Vi vet att $(n-1)s^2/\sigma^2$ är $\chi^2(n-1)$ -fördelad
Alltså måste T-kvoten (*) vara en kvot mellan 2 χ^2 -fördelningar.

DEFINITION

Om vi har $V_1 \sim \chi^2(f_1)$ så är $\frac{V_1/f_1}{V_2/f_2} \sim F(f_1, f_2)$ F-fördelad

tätthetsfunktion finns på slidsen.

F-fördelningstabell finns i boken ∇

EX 4

Samma data som i EX3.

$$I_{\sigma_1^2/\sigma_2^2} = \left[\frac{0,57^2/0,48^2}{F_{\alpha/2}(9,19)}, \frac{0,57^2/0,48^2}{F_{1-\alpha/2}(9,19)} \right] = [0,48; 5,19]$$

Alltså kvoten kan vara "lite varsom helst"

$I_{\sigma_1^2/\sigma_2^2}$ Är mycket känsligt! Måste vara normalfördelat

$1 \in I_{\sigma_1^2/\sigma_2^2} \rightarrow$ kan ej förkasta $H_0 \rightarrow$

\rightarrow Rimligt att använda polad stickprovsvarians

Om ej Normalfördelat \rightarrow vi måste veta i vilket fall (1,2,3) vi befinner oss i.

STICKPROV I PAR

Antag:

Hör ihop $\overline{X}_{11}, \dots, \overline{X}_{1n} \sim N(,)$
Hör ihop $\overline{X}_{21}, \dots, \overline{X}_{2n} \sim N(,)$

Samman antal

Alltså, vi har förvandlat 2 stickprov till 1
Bilda $D_i = \overline{X}_{1i} - \overline{X}_{2i} \sim N(\Delta, \sigma^2)$ ty samma variation

Vi vill testa $H_0: \Delta = 0$ $H_1: \Delta \neq 0$

tex. mäter vikt före/etter banthingskur hos n antal personer.

"JÄMFÖRELSE OCH INFERENS"

PUNKTSKATTNING

$$\theta^* \pm \text{"Någon kvantil"} \cdot \sqrt{V(\theta^*)}$$

POOLAD SKATTNING AV GEMENSAM VARIANS:

Om vi antar att variansen är lika kan vi vinna en del precision på att använda data från båda stickproven för att skatta den gemensamma variansen. En skattning när vi använder data från flera stickprov för att skatta variansen kallas "pooled estimate"

SATS

En väntevärdesriktig skattning av ett gemensamt σ^2 från k olika normalfördelningar $N(\mu_j, \sigma^2)$ fås genom den sammanvägda skattningen.

$$s_p^2 = \frac{(n_1-1)s_1^2 + (n_2-1)s_2^2 + \dots + (n_k-1)s_k^2}{(n_1-1) + (n_2-1) + \dots + (n_k-1)}$$

Dessutom gäller att $(N-k)s_p^2/\sigma^2$ är $\chi^2(N-k)$ -fördelad där $N = \sum_{i=1}^k n_i$

STICKPROV I PAR

Vanliga situationer då mätningar uppkommer i par:

- * Studera hur mkt personer ökar i vikt då de slutar röka
- * Systematiska skillnader mellan två mätmetoder och använder varje metod på var och ett av ett antal prover och jämför de två metoderna för varje prov

För varje prov kan vi bilda differensen D_i som antas vara $N(\mu, \sigma^2)$

$$D_i = x_i - y_i \sim N(\Delta, \sigma^2)$$

Vi vill testa $\Delta=0$, vilket är vanligt hos normalfördelning

FÖR- OCH NACKDELAR M. STICKPROV I PAR

Ofta är det mer effektivt att använda stickprov i par än oberoende stickprov, särskilt om variationen mellan mätningarna är stor.

Vi kan dela upp variationen i D_i som $\sigma^2 = \sigma_0^2 + \sigma_\Delta^2$ där σ_0^2 beskriver variationen mellan objekten.

Vid stickprov i par har vi:

$$V(\bar{D}) = 2\sigma_\Delta^2/n$$

Vid oberoende stickprov har vi:

$$V(\bar{X} - \bar{Y}) = 2(\sigma_0^2 + \sigma_\Delta^2)/n$$

Fördelar

Vi vinner på metoden om $\sigma_0^2 > 0$

Nackdelar

Vi förlorar frihetsgrader
→ större osäkerhet.

INFERENS FÖR DISKRETA DATA Kap 9

BINOMIALFÖRDELNINGEN

Antag att vi observerar $X \sim \text{Bin}(n, p)$ ^{antal försök}
← Parametern vi söker

$p^* = \frac{x}{n}$ är en naturlig skattning av p .

$E(X) = n \cdot p \Rightarrow E(p^*) = \frac{np}{n} = p$ alltså en väntevärdes-
riktig skattning

$V(X) = n \cdot p(1-p) \Rightarrow V(p^*) = \frac{p(1-p)}{n}$

Väldigt svårt att ta fram ett exakt konfidensintervall för Bin. Vanligt att normalapproximera, om n är stort

Tumregel: Om $np(1-p) > 10$ så är normalapproximation ok!

P^* är approximativt $N\left(p, \frac{p(1-p)}{n}\right)$ -fördelad

$$T = \frac{P^* - p}{\sqrt{\frac{p(1-p)}{n}}} \in N(0,1) \quad (\text{approximativt standard normal-fördelad})$$

$I_p = \left[P^* \pm z_{\alpha/2} \sqrt{\frac{p(1-p)}{n}} \right]$ är ett konfidensintervall med approximativ konfidensgrad $1-\alpha$

ENSIDIGA INTERVALL:

$$I_p = \left[0, P^* + z_{\alpha} \sqrt{\frac{p(1-p)}{n}} \right]$$

$$I_p = \left[P^* - z_{\alpha} \sqrt{\frac{p(1-p)}{n}}, 1 \right] \quad (0 \leq p \leq 1)$$

Vi vill testa: $H_0: P = P_0$ mot $H_1: P \neq P_0$ eller
 $H_1: P < P_0$ eller $H_1: P > P_0$

Under H_0 gäller att $X \sim \text{Bin}(n, p_0)$

Så p-värdet beräknas som för:

$$H_1: P > P_0 \rightarrow p = P_{H_0}(X \leq x)$$

$$H_1: P < P_0 \rightarrow p = P_{H_0}(X \geq x)$$

$$H_1: P \neq P_0 \rightarrow \text{om } x \geq np_0 \rightarrow p = 2 \cdot P_{H_0}(X \geq x)$$

$$x \leq np_0 \rightarrow p = 2 \cdot P_{H_0}(X \leq x)$$

EX] $n=1000$ $x=2$

Antag defekta enheter max 0,1%

Vi vill testa: $H_0: p=0,001$ $H_1: p > 0,001$

$$\begin{aligned} p &= P(X \geq 2 | X \sim \text{Bin}(1000; 0,001)) = 1 - P(X \leq 1 | X \sim \text{Bin}(1000; 0,001)) = \\ &= 1 - P(X=1) - P(X=2) = 1 - \binom{1000}{1} 0,001^1 \cdot (1-p)^{999} - \binom{1000}{2} 0,001^2 \cdot (1-p)^{998} = 0,264 \end{aligned}$$

P-värdet är för stort. H_0 kan ej förkastas.

JÄMFÖRELSE AV BINOMIALFÖRDELNINGEN

Antag $X_1 \sim \text{Bin}(n_1, p_1)$ $X_2 \sim \text{Bin}(n_2, p_2)$

Kan vi påstå att $p_1 = p_2$? Vi skattar $p_1 \neq p_2$ med $p_1^* \neq p_2^*$

$p_i^* = \frac{X_i}{n_i}$ Normalapproximation ok om $n_i p_i (1-p_i) > 10$

Vi skattar $p_1 - p_2$ med $p_1^* - p_2^*$

$$E(p_1^* - p_2^*) = p_1 - p_2 \quad V(p_1^* - p_2^*) = \frac{p_1(1-p_1)}{n_1} + \frac{p_2(1-p_2)}{n_2}$$

$$T = \frac{p_1^* - p_2^* - (p_1 - p_2)}{\sqrt{\frac{p_1(1-p_1)}{n_1} + \frac{p_2(1-p_2)}{n_2}}} \stackrel{\epsilon}{\sim} N(0,1)$$

$$I_{p_1 - p_2} = \left[p_1^* - p_2^* \pm z_{\alpha/2} \sqrt{\frac{p_1(1-p_1)}{n_1} + \frac{p_2(1-p_2)}{n_2}} \right]$$

Om vi vill testa $H_0: p_1 = p_2$, så kan vi använda poolad variansskattning, eftersom $V_{H_0}(p_1^* - p_2^*) = p(1-p)\left(\frac{1}{n_1} + \frac{1}{n_2}\right)$

EH gemensamt p skattas som:

$$p^* = \frac{X_1 + X_2}{n_1 + n_2}$$

Då blir teststorheten:

$$T = \frac{p_1^* - p_2^*}{\sqrt{p^*(1-p^*)\left(\frac{1}{n_1} + \frac{1}{n_2}\right)}} \stackrel{\epsilon}{\sim} N(0,1) \text{ under } H_0$$

EX | Antag $n_1 = 100$ $X_1 = 15$
 $n_2 = 200$ $X_2 = 18$

$$\text{Vi har } p_1^* = \frac{15}{100} = 0,15 \quad p_2^* = \frac{18}{200} = 0,09$$

Vi vill testa: $H_0: p_1 = p_2$

$$p^* = \frac{15 + 18}{100 + 200} = 0,11 \quad \longrightarrow \quad T = \frac{p_1^* - p_2^*}{\sqrt{p^*(1-p^*)\left(\frac{1}{n_1} + \frac{1}{n_2}\right)}} = 1,566$$

Vi har $|T_{obs}| < z_{\alpha/2} = 1,96$

Vi kan ej förkasta H_0

↑ kritiska värdet för $N(0,1)$

ICKEPARAMETRISKA METODER

Kom ihåg att medianen är ett tal M så att

$$P(X < M) = P(X > M) = 1/2$$

Antag att vi har ett stickprov av storleken n .

Vi vill testa: $H_0: M = M_0$ mot $H_1: M \neq M_0$ eller $H_1: M > M_0$

Låt $Q_+ = \#\{X_i > M_0\}$ $Q_- = \#\{X_i < M_0\}$ $H_1: M < M_0$

TECKENTEST FÖR MEDIANEN

Om H_0 är sann $\rightarrow Q_+ \sim \text{Bin}(n, 1/2)$ $Q_- \sim \text{Bin}(n, 1/2)$

H_0 förkastas om:

För $H_1: M < M_0 \rightarrow Q_+$ är för liten

För $H_1: M > M_0 \rightarrow Q_-$ är för liten

För $H_1: M \neq M_0 \rightarrow \min(Q_+, Q_-)$ är för liten.

EX | Antag $n=15$

Vi vill testa: $H_0: M=55$ mot $H_1: M < 55$

Räkna antalet större än 55 \rightarrow säg att vi får 3.

$$p = P(Q_+ \leq 3 | Q_+ \sim \text{Bin}(15, 1/2)) = P(X=0) + P(X=1) + P(X=2) + P(X=3) = 0,0176$$

Väldigt litet p -värde $\rightarrow H_0$ förkastas.

Det enda vi antar är symmetrisk fördelning \rightarrow testet tappar i styrka.

WILCOXONS RANKTEST (Ett starkare test)

- 1) Beräkna alla $|x - M_0|$
- 2) Ordna dessa i ökande storleksordning
- 3) Beräkna R_i som ranken för $|x_i - M_0|$ gånger $\text{sign}(x_i - M_0)$

Definiera $W_+ = \sum_{i: R_i > 0} R_i$ $|W_-| = \sum_{i: R_i < 0} |R_i|$

Teststorheten: är $W = \min(W_+, |W_-|)$ finns i tabell.

EX 8.1.2

X_i :	115,1	117,8	116,5	121,0
$ X_i - M_0 $:	4,9	2,2	3,5	1
Rank:	4	2	3	1
R_i :	-4	-2	-3	1

$$W_+ = 1 \quad |W_-| = |-2-3-4| = 9$$

$W = \min(1, 9) = 1$ slå upp kritiska värdet i tabell.

Det finns ett ekvivalent test för jämförelser \rightarrow Wilcoxons Rank-sumtest

MATEMATISK STATISTIK F12 ^{15/6} 2014

EXEMPEL: STICKPROV I PAR

I fallet med parade observationer $(X_1, Y_1), \dots, (X_m, Y_m)$ bildar vi först differenserna $Z_i = X_i - Y_i$ och använder sedan ett ranktest på differanserna för att testa om differansernas median är noll.

Från EX. 10, b.2 i boken \rightarrow

Vi vill testa hur mycket minne 2 statistiska paket använder vid analys av ett dataset. Vi gör mätningar:

Vi vill testa: $H_0: \mu_x = \mu_y$ $H_1: \mu_x \neq \mu_y$

Vi ordnar differanserna och beräknar rankerna:

$ D_i $:	0	10	12	25	25	40	50	100
Rank:	1	2	3	4,5	4,5	6	7	8
R_i :	-1	2	3	-4,5	4,5	6	7	8

Vi har nu

$$W_+ = 2 + 3 + 4,5 + 6 + 7 + 8 = 30,5$$

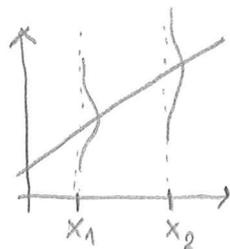
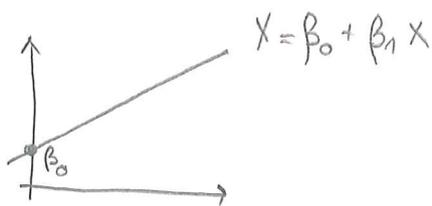
$$|W_-| = |-1 - 4,5| = 5,5$$

Eftersom vi har ett 2-sidigt test så använder vi:

$$W = \min(|W_-|, W_+) = 5,5$$

från tabell ser vi då att den kritiska punkten blir: 6 för ett 2-sidigt test med $\alpha = 0,1$.

$W = 5,5 < 6 \rightarrow H_0$ förkastas.



β_0 = interceptet

β_1 = lutningen

MINSTA KVADRAT METODEN (MK)

Antag den lite mer generella modellen att väntevärdet för y ges av en funktion f .

$$E(y_i) = f(\theta_1, \dots, \theta_k, x_i)$$

Parametrarna $\theta_1, \dots, \theta_k$ skattas enligt minsta kvadrat metoden genom att minimera kvadratfelet för den data vi har.

$$S(\theta_1, \dots, \theta_k) = \sum_{i=1}^k (y_i - f(\theta_1, \dots, \theta_k, x_i))^2 \quad \text{m.a.p } \theta_1, \dots, \theta_k$$

Lösningen brukar fås som lösningen till:

$$\frac{\partial S}{\partial \theta_i} = 0 \quad \text{för } i=1, \dots, k$$

Ekvationssystemet löses av $\theta_1^*, \dots, \theta_k^*$ som kallas M-k-skattningar av $\theta_1, \dots, \theta_k$

$$S(\beta_0, \beta_1) = \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i)^2$$

$$\frac{\partial S}{\partial \beta_0} = -2 \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i) = 0$$

$$\frac{\partial S}{\partial \beta_1} = -2 \sum_{i=1}^n x_i (y_i - \beta_0 - \beta_1 x_i) = 0$$

$$\left. \begin{array}{l} \sum y_i - n\beta_0 - \beta_1 \sum x_i = 0 \\ \sum x_i y_i - \beta_0 \sum x_i - \beta_1 \sum x_i^2 = 0 \end{array} \right\}$$

$$\left. \begin{array}{l} \sum y_i - n\beta_0 - \beta_1 \sum x_i = 0 \\ \sum x_i y_i - \beta_0 \sum x_i - \beta_1 \sum x_i^2 = 0 \end{array} \right\}$$

$$\left. \begin{array}{l} \beta_0 n + \beta_1 \sum x_i = \sum y_i \\ \beta_0 \sum x_i + \beta_1 \sum x_i^2 = \sum x_i y_i \end{array} \right\}$$

Multiplisera m.

$$-\frac{1}{n} \sum x_i$$



$$\beta_1 \left(\underbrace{\sum x_i^2 - \frac{(\sum x_i)^2}{n}}_{S_{xx}} \right) = \underbrace{\sum x_i y_i - \frac{\sum x_i \sum y_i}{n}}_{S_{xy}}$$

$$\beta_1 = \frac{S_{xy}}{S_{xx}}$$

$$S_{xx} = \sum x_i^2 - n\bar{x}^2 = \sum (x_i - \bar{x})^2 \quad \text{"variationen i x-led"}$$

$$S_{xy} = \sum x_i y_i - n\bar{x}\bar{y} = \sum (x_i - \bar{x})(y_i - \bar{y}) \quad \text{"samvariation"}$$

$$\beta_0 n + \beta_1 \sum x_i = \sum y_i \xrightarrow{\cdot \frac{1}{n}} \beta_0 + \beta_1 \bar{x} = \bar{y} \rightarrow \beta_0 = \bar{y} - \beta_1 \bar{x}$$

Mk-skattningar

$$\beta_1^* = \frac{S_{xy}}{S_{xx}} \quad \beta_0^* = \bar{y} - \beta_1^* \bar{x}$$

i boken skriver dom B_0 för β_0^* och B_1 för β_1^*

Mk-skattning av σ^2

Skattas av $s^2 = \frac{Q_0}{n-2}$ eftersom vi har n-2 frihetsgrader

$$Q_0 = \sum (y_i - \beta_0^* - \beta_1^* x_i)^2$$

Om vi räknar för hand kan vi:

$$Q_0 = S_{yy} - \beta_1^* S_{xy} = S_{yy} - \frac{S_{xy}^2}{S_{xx}}$$

SATS

$$\text{Vi har att: } E(\bar{Y}) = \beta_0 + \beta_1 \bar{X}$$

$$V(\bar{Y}) = \frac{\sigma^2}{n}$$

$$E(\beta_1^*) = \beta_1$$

$$V(\beta_1^*) = \frac{\sigma^2}{S_{XX}} = \frac{\sigma^2}{\sum (x_i - \bar{x})^2}$$

BEVIS

$$Y_i = \beta_0 + \beta_1 x_i + \varepsilon_i$$

$$\begin{aligned} E(\bar{Y}) &= E\left(\frac{1}{n} \sum Y_i\right) = E\left(\frac{1}{n} \sum (\beta_0 + \beta_1 x_i + \varepsilon_i)\right) \\ &= E\left(\beta_0 + \beta_1 \left(\frac{1}{n} \sum x_i\right) + \frac{1}{n} \sum \varepsilon_i\right) = \beta_0 + \beta_1 \bar{x} - \frac{1}{n} \sum E(\varepsilon_i) \end{aligned}$$

konstant, ε_i slumpmässig

\uparrow
 $N(0, \sigma^2)$

$$V(\bar{Y}) = V\left(\beta_0 + \beta_1 \bar{x} + \frac{1}{n} \sum \varepsilon_i\right) = V\left(\frac{1}{n} \sum \varepsilon_i\right) = \frac{1}{n^2} \sum \sigma^2 = \frac{\sigma^2}{n}$$

$$\begin{aligned} E(\beta_1^*) &= E\left(\frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2}\right) = \frac{\sum (x_i - \bar{x}) E(y_i - \bar{y})}{\sum (x_i - \bar{x})^2} \\ &= \frac{\sum (x_i - \bar{x}) \beta (x_i - \bar{x})}{\sum (x_i - \bar{x})^2} = \beta \end{aligned}$$

Eftersom det bara ε_i som är stokastisk

$$\text{Vi har att } \sum (x_i - \bar{x})(y_i - \bar{y}) = \sum (x_i - \bar{x}) y_i - \sum (x_i - \bar{x}) \bar{y} = 0$$

$$\rightarrow \sum (x_i - \bar{x}) y_i = \sum (x_i - \bar{x}) \bar{y}$$

Detta gör det enklare att beräkna $V(\beta_1^*)$

$$V(\beta_1^*) = V\left(\frac{\sum (x_i - \bar{x}) y_i}{\sum (x_i - \bar{x})^2}\right) = \frac{1}{\left(\sum (x_i - \bar{x})^2\right)^2} \cdot V\left(\sum (x_i - \bar{x}) y_i\right)$$

konstant

$\forall y_i$ oberoendety
 $\forall \varepsilon_i$ oberoende

$$= \frac{1}{\left(\sum (x_i - \bar{x})^2\right)^2} \sum (x_i - \bar{x})^2 \underbrace{V(x_i)}_{\sigma^2} =$$

$$= \frac{\sigma^2}{\sum (x_i - \bar{x})^2} = \frac{\sigma}{S_{XX}}$$

SATS

Om $\mu_{\bar{Y}}^*(x_0) = \beta_0^* + \beta_1^* x_0$ är den skattade linjens värde för ett fixt x_0 så är:

$$E(\mu_{\bar{Y}}^*(x_0)) = \beta_0 + \beta_1 x_0$$

Samt:

$$V(\mu_{\bar{Y}}^*(x_0)) = \sigma^2 \left(\frac{1}{n} + \frac{(x_0 - \bar{x})^2}{S_{xx}} \right)$$

BEVIS

$$\mu_{\bar{Y}}^*(x_0) = \beta_0^* + \beta_1^* x_0 = \bar{Y} + \beta_1^* (x_0 - \bar{x})$$

$$\begin{aligned} E(\mu_{\bar{Y}}^*(x_0)) &= E(\bar{Y}) + E(\beta_1^*) (x_0 - \bar{x}) = \beta_0 + \beta_1 \bar{x} + \beta_1 (x_0 - \bar{x}) = \\ &= \beta_0 + \beta_1 x_0 \end{aligned}$$

$$\begin{aligned} V(\mu_{\bar{Y}}^*(x_0)) &= \left\{ C(\bar{Y}, \beta_1^*) = 0 \right\} = V(\bar{Y}) + (x_0 - \bar{x})^2 \cdot V(\beta_1^*) = \\ &= \frac{\sigma^2}{n} + \frac{(x_0 - \bar{x})^2 \cdot \sigma^2}{S_{xx}} = \sigma^2 \left(\frac{1}{n} + \frac{(x_0 - \bar{x})^2}{S_{xx}} \right) \end{aligned}$$

SATS

Om ε_i är normalfördelad gäller att \bar{Y} , β_0^* , β_1^* och $\beta_0^* + \beta_1^* x_0$ också är normalfördelade

BEVIS

$$\bar{Y} = \frac{1}{n} \sum_1^1 y_i \leftarrow N(\beta_0 + \beta_1 x_i, \sigma^2)$$

$$\beta_1^* = \frac{\sum_1^1 (x_i - \bar{x})(y_i - \bar{y})}{S_{xx}} = \frac{\sum_1^1 (x_i - \bar{x}) \cdot \underbrace{(y_i - \bar{y})}_{\text{Normalfördelad}}}{S_{xx}} = \beta_1^* \in N(\beta_1, \sigma^2)$$

Normalfördelad ty konstant \cdot normalfördelad = normalfördelad

Samt \sum normalfördelad = normalfördelad

KONFIDENSINTERVALL OCH TEST

Låt θ vara någon av β_0, β_1 eller $\beta_0 + \beta_1 x_0$. Vi vet att dess skattningar är normalfördelade och vi har tagit fram variansen av skattningarna. Låt $d(\theta^*)$ vara standardavvikelsen för skattningen. Vi har då att:

$$T = \frac{\theta^* - \theta}{d(\theta^*)} \sim t(n-2)$$

Vilken på vanligt sätt används för att göra test och bilda konfidensintervall.

$$I_{\theta} = (\theta^* \pm t_{\alpha/2}(n-2) d(\theta^*))$$

För β_0

$$I_{\beta_0} = \left(\beta_0^* \pm t_{\alpha/2}(n-2) s \sqrt{\frac{1}{n} + \frac{\bar{x}^2}{s_{xx}}} \right)$$

För β_1

$$I_{\beta_1} = \left(\beta_1^* \pm t_{\alpha/2}(n-2) \frac{s}{\sqrt{s_{xx}}} \right)$$

För $\mu_Y(x_0) = \beta_0 + \beta_1 x_0$

$$I_{\mu_Y(x_0)} = \left(\beta_0^* + \beta_1^* x_0 \pm t_{\alpha/2}(n-2) s \sqrt{\frac{1}{n} + \frac{(x_0 - \bar{x})^2}{s_{xx}}} \right)$$

MATEMATISK STATISTIK F13

19/5 2014

• Han skrev inget på tavlan!

"TENTAGENOMGÅNG"

16 JANUARI 2014 (stefans.)

① $P(A)=0,2$ $P(B)=0,5$ $P((A \cup B)^c)=0,4$

a) om $P(A \cap B) = P(A)P(B) \rightarrow A$ och B oberoende

$$P(A \cup B) = P(A) + P(B) - P(A \cap B) = 1 - P((A \cup B)^c) = 0,6$$

$$\rightarrow P(A \cap B) = 0,5 + 0,2 - 0,6 = 0,1$$

$P(A) \times P(B) = 0,1 \rightarrow$ A och B är oberoende

b) $P(A|B)$?

Två sätt: 1. Definition av betingad slh:

$$P(A|B) = \frac{P(A \cap B)}{P(B)} = \frac{0,1}{0,5} = 0,2$$

2. $P(A|B) = P(A)$ ty oberoende i a)

② $X \sim \text{Exp}(\lambda)$ $X_{\min} = \min(X_1, \dots, X_{10})$

a) $P(X > x) = 1 - P(X \leq x) = 1 - F(x) = 1 - (1 - e^{-\lambda x}) = e^{-\lambda x}$

b) $P(X_{\min} > x) = P(\min(X_1, \dots, X_{10}) > x) = P(X_1 > x) \cap P(X_2 > x) \cap \dots$
 $\dots \cap P(X_{10} > x) = P(X_1 > x) P(X_2 > x) \dots P(X_{10} > x) =$
 $= e^{-\lambda x} \cdot e^{-\lambda x} \cdot \dots \cdot e^{-\lambda x} = e^{-10\lambda x}$

c) $F_{X_{\min}}(x) = P(X_{\min} \leq x) = 1 - P(X_{\min} > x) = 1 - e^{-10\lambda x}$

d) $X_{\min} \sim \text{Exp}(10\lambda)$

③ Är kass enligt david!

④

	n	\bar{x}/\bar{y}	s
(x) Blyblandad	15	87	4
(y) Vanlig	10	90	5

a) Man får anta $X \sim N(\mu_1, \sigma_1^2)$ $Y \sim N(\mu_2, \sigma_2^2)$

Testa $H_0: \sigma_1 = \sigma_2$ $H_1: \sigma_1 \neq \sigma_2$

$$T_{\text{obs}} = \frac{s_x^2}{s_y^2} = \frac{25}{16}$$

Förkasta H_0 om $T_{\text{obs}} > F_{\alpha/2}(9, 14) = 2,54$

$T_{\text{obs}} < 2,54$ vi kan ej förkasta H_0

→ vi kan anta $\sigma_1 = \sigma_2$

I facit har dom antagit detta utan att testa

$$I_{\mu_1 - \mu_2} = \left(\bar{x} - \bar{y} \pm \underbrace{t_{\alpha/2}(23)}_{2,069} \times S_p \sqrt{\frac{1}{n_x} + \frac{1}{n_y}} \right)$$

$$S_p^2 = \frac{(n_x - 1)S_x^2 + (n_y - 1)S_y^2}{n_x + n_y - 2} = 19,52$$

Man kan använda poolad pga $\sigma_1 = \sigma_2$

$$\longrightarrow I_{\mu_1 - \mu_2} = (-3 \pm 2,069 \times 1,803) = [-6,73; 0,73]$$

b) statfan har inte specificerat $H_0 \rightarrow$ antar $H_0: \mu_1 = \mu_2$

$0 \in I_{\mu_1 - \mu_2} \rightarrow$ vi kan ej förkasta H_0 på $\alpha = 0,05$.

Vi kan inte heller förkasta på starkare nivå (lägre nivå) \rightarrow vi kan förkasta på $\alpha = 0,01$

5) Uppgiften visar en bra poäng.

	n	$\bar{\delta}$	S_{δ}
A	16	0,3	0,5
B	9	0,2	0,5

a) Vi vill testa $H_0: \delta = 0$ $H_1: \delta > 0$
 "T-test" om $N(\delta, \sigma^2)$

$$T = \frac{\bar{\delta}}{S_{\delta}/\sqrt{n}} \text{ under } H_0 \quad T \sim t(n-1)$$

För A

$$T_{\text{obs}} = \frac{0,3}{0,5/\sqrt{16}} = 2,4$$

$$t_{0,05} (16-1) = 1,75$$

$T_{\text{obs}} > 1,75 \rightarrow$ Förläsa H_0

För B

$$T_{\text{obs}} = \frac{0,2}{0,5/\sqrt{9}} = 1,2$$

$$t_{0,05} (9-1) = 1,86$$

$T_{\text{obs}} < 1,86 \rightarrow H_0$ kan ej förläsa

b) * pga olika n så har testen olika styrka

* Man borde göra ett test för skillnaden i väntevärde

$$\rightarrow H_0: \delta_A = \delta_B \quad H_1: \delta_A > \delta_B$$

Samma antaganden som i a). Använder poolad variansskattning

$$S_p^2 = 0,5^2$$

$$T = \frac{\bar{\delta}_A - \bar{\delta}_B}{S_p \sqrt{\frac{1}{16} + \frac{1}{9}}} = 0,24$$

$$t_{0,05} (23) = 1,71$$

$T < 1,71$ vi kan ej förläsa H_0

6. $H_0: \mu = 0$ $H_1: \mu > 0$

$$T = \frac{\bar{X}}{\sigma/\sqrt{n_1}} = z_{0,05} = 1,645 \rightarrow \frac{\bar{X}}{\sigma/\sqrt{n}} = 1,645 \rightarrow \frac{\bar{X}}{\sigma} = \frac{1,645}{\sqrt{20}}$$

Om vi skulle gjort ett 2-sidigt test:

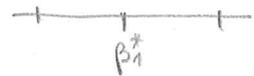
$$\frac{\bar{X}}{\sigma/\sqrt{n}} = z_{0,05/2} = 1,96 \rightarrow \sqrt{n} = \frac{1,96}{\bar{X}/\sigma} = \frac{1,96}{\frac{1,645}{\sqrt{20}}}$$

$$\rightarrow n = \left(\frac{1,96}{1,645/\sqrt{20}} \right)^2 = 28,4 \approx 29 \rightarrow n \geq 29$$

7. $\bar{X} = 2,325$ $\bar{Y} = 129,7$ $I_{\beta_1} = [6,29 ; 24,99]$
 $n = 16$

a) $M_X^*(x=0) = ?$

Hur ser I_{β_1} ut?



$$\beta_1^* = \frac{24,99 - 6,29}{2} = 15,64$$

$$\beta_0^* = \bar{Y} - \beta_1^* \cdot \bar{X} = 129,7 - 15,64 \times 2,325 = 93,337$$

$$M_Y^*(x=0) = \beta_0^* + \beta_1^* \cdot x = \beta_0^* = 93,337$$

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

$$T = \frac{\beta_1^* - \beta_1}{s/\sqrt{S_{xx}}} = \frac{15,64}{s/\sqrt{S_{xx}}}$$

$$I_{\beta_1} = \left(\beta_1^* \pm t_{\alpha/2}(n-2) \cdot \frac{s}{\sqrt{S_{xx}}} \right)$$

Längden på intervallet = $24,99 - 6,29 = 2 \cdot 2,145 \cdot \frac{s}{\sqrt{S_{xx}}} \rightarrow$

$$\frac{s}{\sqrt{S_{xx}}} = 4,36$$

$$\rightarrow T = \frac{15,64}{4,36} = 3,59$$

c), d) saknas info för att lösa detta

25 Maj 2012 (Aila)

Typ så här kommer dauids tentor se ut

$$\textcircled{2} \quad F(x) = \begin{cases} 0 & x < 0 \\ \frac{1}{4} \cdot x^3 & 0 \leq x < 1 \\ \frac{3}{4} \cdot x - \frac{1}{2} & 1 \leq x \leq 2 \\ 1 & x > 2 \end{cases}$$

a) $P(0,5 \leq X \leq 1,5) = ?$

Finns 2 sätt att lösa detta!

Krämligt sätt

$$= \int_{0,5}^{1,5} f(x) dx = \int_{0,5}^{1,5} \frac{\partial F(x)}{\partial x} dx$$

Lättare sätt

$$F(x) = P(X \leq x) \rightarrow$$

$$P(0,5 \leq X \leq 1,5) = P(X \leq 1,5) - P(X < 0,5) = F(1,5) - F(0,5) =$$

$$= \frac{3}{4} \times 1,5 - \frac{1}{2} - \frac{0,5^3}{4} = \frac{19}{32}$$

b) $P(X \geq 1,5) = 1 - F(1,5) = \frac{3}{8}$

$Y =$ "Antal $X_i > 1,5$ " $\sim \text{Bin}(10, \frac{3}{8})$

$$P(Y=2) = P_Y(2) = \binom{10}{2} \left(\frac{3}{8}\right)^2 \left(1 - \frac{3}{8}\right)^8 = 0,147$$

4. a) Minimera kvadratfel: $s = \sum_i (y_i - \beta_0 - \beta_1 x_i)^2$
m.a.p β_0 och β_1

b) Egentligen inga alls

c) $e_i = y_i - \beta_0^* - \beta_1^* x_i$: Modelvalidering

5. $\bar{X} = 54,11$ $s^2 = 207,19$

a) $X \sim N(\mu, \sigma^2)$ $H_0: \mu = 50$ $H_1: \mu < 50$

Z-test \rightarrow om σ är känd

t-test \rightarrow om σ är okänd och skattas av s

$$T = \frac{\bar{X} - \mu}{s/\sqrt{n}} = \frac{54,11 - 50}{\sqrt{207,19}/\sqrt{20}} = 1,277$$

$$t_\alpha(19) = 1,729$$

H_0 förkastas pga $T < t_\alpha$



c) Jmf tecken- och t-test. Teckentestet borde vara bättre eftersom det ej ser ut som en jdmn fördelning