

Punktskattning

En funktion av en samling stok. var., $\eta = f(\xi_1, \dots, \xi_n)$, är en ny stok. var. som kallas för en statistika. Ex: $f(\xi_1, \dots, \xi_n) = \min\{\xi_1, \dots, \xi_n\}$ och $f(\xi_1, \dots, \xi_n) = \frac{1}{n} \sum_{i=1}^n \xi_i$.

- En statistika $\hat{\theta} = \hat{\theta}(\xi_1, \dots, \xi_n)$ som används för att skatta en parameter θ kallas för en skattare/estimator. När erhållna realiseringar/effall $x_1 = \xi_1(w), \dots, x_n = \xi_n(w)$ (data) av ξ_1, \dots, ξ_n sätts in i skattaren får vi en skattning/ett estimat $\hat{\theta}(x_1, \dots, x_n)$ av θ .

- En bra skattare $\hat{\theta}$ för en parameter θ uppfyller att medeltvadratfellet/mean squared error $MSE(\hat{\theta}) = E[(\hat{\theta} - \theta)^2]$ är litet. Vi vill alltså använda $\hat{\theta}$ s.a. $MSE(\hat{\theta})$ är litet. Två komponenter påverkar $MSE(\hat{\theta})$:

$$MSE(\hat{\theta}) = E[(\hat{\theta} - \theta)^2] = \underbrace{\text{Var}(\hat{\theta} - \theta)}_{E[\xi^2] = \text{Var}(\xi) + E[\xi]^2, \xi = \hat{\theta} - \theta} + E[\hat{\theta} - \theta]^2 =$$

$$= \text{Var}(\hat{\theta}) + (E[\hat{\theta}] - \theta)^2 = SE(\hat{\theta})^2 + \text{bias}(\hat{\theta})^2 \text{ där}$$

standardfel är $SE(\hat{\theta}) = \sqrt{\text{Var}(\hat{\theta})}$ (stand. avgv. för $\hat{\theta}$) och bias/systematiska fel är

$$\text{bias}(\hat{\theta}) = E[\hat{\theta}] - \theta.$$

Ofta börjar man med att försöka
sög om att $\hat{\theta}$ är väntevärdesriktig, dvs

$$\text{bias}(\hat{\theta}) = 0 \Leftrightarrow E[\hat{\theta}] = \theta \quad (\text{så att } \text{MSE}(\hat{\theta}) = \text{Var}(\hat{\theta}))$$

och har man två konkurrerande väntevärdes-
riktiga skattare $\hat{\theta}_1$ och $\hat{\theta}_2$ så visar man
den effektivaste, dvs den med lägst
standardfel/varians:

$$\hat{\theta}_1 \text{ mer effektiv än } \hat{\theta}_2 \Leftrightarrow \text{Var}(\hat{\theta}_1) \leq \text{Var}(\hat{\theta}_2)$$

$$(\Leftrightarrow)$$

dvs relativa effektiviteten uppfyller

$$\frac{\text{Var}(\hat{\theta}_2)}{\text{Var}(\hat{\theta}_1)} \geq 0.$$

En väntevärdesriktig skattare genererar
skattningar som i snedeltal, givet upp-
rapade forsök, ger den samma parametern.

Dock kan det svåra ofta att ange från
forsök till forsök, såvida inte $\text{SE}(\hat{\theta})$ är känd
(varvid $\text{MSE}(\hat{\theta})$ blir känd).

Ex ξ_1, ξ_2, ξ_3 observerade och tillsädd, med $E[\xi_i] = \mu, \text{Var}[\xi_i] = \sigma^2$
Två skattare för μ : $\hat{\mu}_1 = \frac{\xi_1 + 2\xi_2 + 3\xi_3}{6}, \hat{\mu}_2 = \frac{\xi_1 + 4\xi_2 + \xi_3}{6}$

$$\Rightarrow \text{bias}(\hat{\mu}_1) = E\left[\frac{\xi_1 + 2\xi_2 + 3\xi_3}{6}\right] - \mu = \frac{E[\xi_1] + 2E[\xi_2] + 3E[\xi_3]}{6} - \mu = \frac{6\mu}{6} - \mu = 0$$

$$\text{bias}(\hat{\mu}_2) = E\left[\frac{\xi_1 + 4\xi_2 + \xi_3}{6}\right] - \mu = \frac{E[\xi_1] + 4E[\xi_2] + E[\xi_3]}{6} - \mu = \frac{6\mu}{6} - \mu = 0$$

Relativ effektskattet: ber. med samma varians

$$\frac{SE(\hat{\mu}_2)^2}{SE(\hat{\mu}_1)^2} = \frac{\text{Var}(\hat{\mu}_2)}{\text{Var}(\hat{\mu}_1)} = \frac{\text{Var}\left(\frac{\xi_1 + 4\xi_2 + \xi_3}{6}\right)}{\text{Var}\left(\frac{\xi_1 + 2\xi_2 + 3\xi_3}{6}\right)} = \frac{\frac{\sigma^2 + 4\sigma^2 + \sigma^2}{6^2}}{\frac{\sigma^2 + 2^2\sigma^2 + 3\sigma^2}{6^2}} = \frac{18\sigma^2}{14\sigma^2} = \frac{9}{7} \approx 1.29 \approx 2$$

bias = 0

$\Rightarrow \hat{\mu}_1$ är en bättre skattare. ◻

$$S^2 = \frac{1}{n-1} \sum_{i=1}^n (\xi_i - \bar{\xi})^2, \quad \bar{\xi} = \frac{1}{n} \sum_{i=1}^n \xi_i, \quad \text{är en vartervärdesvitrig}$$

Skattare för $\sigma^2 = \text{Var}(\xi_i)$

$$(n-1) S^2 = \sum_{i=1}^n ((\xi_i - \mu) + (\mu - \bar{\xi}))^2 = \sum_{i=1}^n (\xi_i - \mu)^2 - 2(\bar{\xi} - \mu) \sum_{i=1}^n (\xi_i - \mu) + n(\mu - \bar{\xi})^2$$

beror ej på μ

$$= \sum_{i=1}^n (\xi_i - \mu)^2 - n(\bar{\xi} - \mu)^2$$

$$= n(\bar{\xi} - \mu)$$

$$\Rightarrow (n-1) \mathbb{E}[S^2] = \sum_{i=1}^n \underbrace{\mathbb{E}[(\xi_i - \mu)^2]}_{= \text{Var}(\xi_i) = \sigma^2} - n \underbrace{\mathbb{E}[(\bar{\xi} - \mu)^2]}_{= \text{Var}(\bar{\xi}) = \frac{1}{n^2} \text{Var}(\sum_{i=1}^n \xi_i) = \frac{n\sigma^2}{n^2} = \frac{\sigma^2}{n}}$$

$$+ \text{ty } \mathbb{E}[\bar{\xi}] = \mu$$

$$= n\sigma^2 - n\frac{\sigma^2}{n} = (n-1)\sigma^2$$

$$\Rightarrow \mathbb{E}[S^2] = \sigma^2$$
◻

I ett binomialexperiment har vi $\sum_{i=1}^n \xi_i$, där $\xi_i \sim \text{Bin}(1, p)$, med medelvärdet $\hat{P} = \frac{\xi}{n}$ och $\mathbb{E}[\hat{P}] = \frac{\mathbb{E}[\xi]}{n} = \frac{np}{n} = p$, $\text{Var}(\hat{P}) = \frac{\text{Var}(\xi)}{n^2} = \frac{np(1-p)}{n^2} = \frac{p(1-p)}{n}$. N.a.s. är \hat{P} vartervärdsitlig för p där $\xi \sim \text{Bin}(n, p)$ med p okänd.

Konfidensintervall / Intervallskattning

Hittills har vi förtat på punktskattningar av en parameter av intresset, θ . Men beröende på det observerade stickprovet kommer skattningen $\hat{\theta} = \hat{\theta}(x_{1:n})$ att ligga olika långt från det sanna värdet θ (Om $MSE(\hat{\theta})$ är litet vet vi att detta avstånd tenderar att vara litet). Men kan man på något sätt ge en skattning som inkorporerar osäkerheten som följer av att vi använder oss av ett enda observerat stickprov x_1, \dots, x_n från ξ_1, \dots, ξ_n ?

Ett konfidensintervall $[L, U] = [L(\xi_1, \dots, \xi_n), U(\xi_1, \dots, \xi_n)]$ för θ är ett slokastiskt interval som uppfyller

$$P(\theta \in [L, U]) = 1 - \alpha, \quad \alpha \in (0, 1)$$

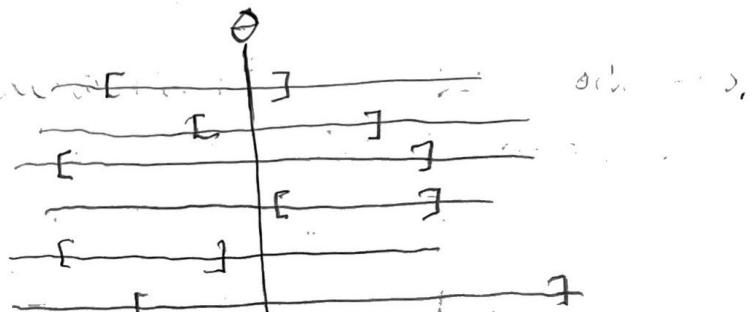
där $(1 - \alpha)100\%$ kallas konfidensgraden.

Om vi gör om ett experiment 1000 ggr får vi 1000 st observerade stickprov $x_j = \{x_{ij}\}_{i=1}^{n_j}, j=1, \dots, 1000$, från $\xi_1, \dots, \xi_n \Rightarrow$ intervallet $I_j = [L(x_{j1}, \dots, x_{jn_j}), U(x_{j1}, \dots, x_{jn_j})], j=1, \dots, 1000$, där ungefärligt $(1 - \alpha) \cdot 1000$ av dem kommer att innehålla den sanna men okända parametern θ .

Vi vill så klart ha α så litet som möjligt men $[L, U] \rightarrow (-\infty, \infty)$ när $\alpha \rightarrow 0$;

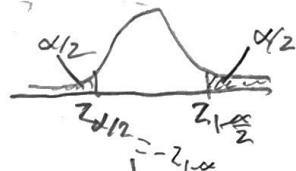
$P(\theta \in (-\infty, \infty)) = 1$ är inte så användbart.

Vad som är " tillräcklig säkerhet" varierar man oftast beroende på $1 - \alpha = 0.95$



Men hur kan man skapa sättiga konfidensintervall?

Om $\xi_1, \dots, \xi_n \sim N(\mu, \sigma^2)$ vet vi att



$$1-\alpha = P\left(-z_{\alpha/2} \leq \frac{\bar{\xi}-\mu}{\sigma/\sqrt{n}} \leq z_{\alpha/2}\right) = P\left(-2z_{\alpha/2} \frac{\sigma}{\sqrt{n}} \leq \bar{\xi}-\mu \leq 2z_{\alpha/2} \frac{\sigma}{\sqrt{n}}\right)$$
$$= P\left(\bar{\xi} - 2z_{\alpha/2} \frac{\sigma}{\sqrt{n}} \leq \mu \leq \bar{\xi} + 2z_{\alpha/2} \frac{\sigma}{\sqrt{n}}\right)$$

$$\Rightarrow L_2(\xi_1, \dots, \xi_n) = \bar{\xi} - z_{\alpha/2} \frac{\sigma}{\sqrt{n}}, U = U(\xi_1, \dots, \xi_n) = \bar{\xi} + 2z_{\alpha/2} \frac{\sigma}{\sqrt{n}}$$

funktion av ξ_1, \dots, ξ_n

Vilket man gett data x_1, \dots, x_n , med $\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$, uttrycker som:

Ett $(1-\alpha) \cdot 100\%$ konfidensintervall för väntevärdelet $\mu = E[\xi]$ ges av $\bar{x} \pm z_{\alpha/2} \frac{\sigma}{\sqrt{n}}$

Notera att ett observerat interval antingen innehåller den sanna ständla parametern θ eller inte. Tolkningen är att vi med sannolikhet α har fått ett så "ovanligt" observerat stickprov att det erhållna intervallet ej innehåller θ .

Ex Observerar man temperaturer i 49 komponenter antas normalfördelad, med $\sigma = 14^\circ\text{C}$. Medeltemp. uppmäts till $\bar{x} = 68^\circ\text{C}$ \Rightarrow 99%igt konfidensintervall för μ ges av

$$\bar{x} \pm z_{\alpha/2} \frac{\sigma}{\sqrt{n}} = 68 \pm 2.58 \frac{14}{\sqrt{49}} \Leftrightarrow [62.84, 73.16]$$

Eftersom $\alpha = 0.01 \Rightarrow z_{\alpha/2} = z_{0.995} = 2.58$ ($0.995 = 2.58$ från tabellen)

Allmänt gäller för en skattare $\hat{\theta}$ som har en $N(\theta, \text{Var}(\hat{\theta}))$ -fördelning att ett $(1-\alpha) \cdot 100\%$ igt konfidensintervall ges av $\hat{\theta} \pm z_{\alpha/2} \sqrt{\text{Var}(\hat{\theta})}$
 $= SE(\hat{\theta})$

I praktiken är det dock (nästan) aldrig fallet att σ är känd och ofta är antagandet att ξ_1, \dots, ξ_n kommer från en normalfördelning inte sann.

I det senare fallet kan vi använda CGS-approximationen och utgå från att $\frac{\bar{\xi} - E(\bar{\xi})}{\sqrt{Var(\bar{\xi})}} \stackrel{approx}{\sim} N(0,1)$ för att konstruera $[L, U]$. T.ex. om $\xi_i \sim \text{Bin}(n, p)$ där vi vill intervallskatta $P \in (0,1)$, förutsatt att σ är känd:

$$1 - \alpha \approx P\left(\bar{\xi}_{\alpha/2} \leq \frac{\hat{P} - P}{\sqrt{\frac{P(1-P)}{n}}} \leq \bar{\xi}_{1-\alpha/2}\right) = \dots = P\left(\hat{P} - Z_{1-\alpha/2} \sqrt{\frac{P(1-P)}{n}} \leq P \leq \hat{P} + Z_{1-\alpha/2} \sqrt{\frac{P(1-P)}{n}}\right)$$

$$\hat{P} = \frac{\bar{\xi}}{n} = \frac{1}{n} \sum_{i=1}^n \xi_i \stackrel{approx}{\sim} N(P, \frac{P(1-P)}{n}) \quad \text{Var.}$$

$$P\left(\hat{P} - Z_{1-\alpha/2} \sqrt{\frac{\hat{P}(1-\hat{P})}{n}} \leq P \leq \hat{P} + Z_{1-\alpha/2} \sqrt{\frac{\hat{P}(1-\hat{P})}{n}}\right) = L = U$$

\Rightarrow intervallet ges av $\hat{P} \pm Z_{1-\alpha/2} \sqrt{\frac{\hat{P}(1-\hat{P})}{n}}$

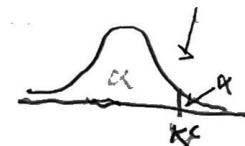
Till fallet där σ är okänd och $\xi_1, \dots, \xi_n \sim N(\mu, \sigma^2)$ (och även om ξ_1, \dots, ξ_n inte är normalford. men σ är känd):

$$\begin{aligned} \bar{\xi} - \mu &\stackrel{approx}{\sim} t(v) \quad t(v)-fördel. med v=n-1 "frihetsgrader" \\ S/\sqrt{n} &, S^2 = \frac{1}{n-1} \sum_{i=1}^n (\xi_i - \bar{\xi})^2 \quad \text{"standard "df"} \end{aligned}$$

$t(v)$ -fördelningen är symmetrisk och ser ut ungefärlig som en $N(0,1)$ -fördelning men har mer massa för stora och små värden ("tjocka svansar"); vi kan $\lim_{v \rightarrow \infty} t(v) = N(0,1)$ (eftersom S^2 går mot σ^2 när $n \rightarrow \infty$).

Tätheten $f_{t(v)}(x)$ (och därför fördeln. funkt.) har en ganska komplex form (se sid. 268) men värden för fördelningstakten finns i Appendix I i boken.

$$\xi \sim t(v) \Rightarrow E[\xi] = 0, \quad Var(\xi) = \frac{v}{v-2} \text{ för } v \geq 3.$$



Vrafat här
 $1-\alpha = P\left(-x_{\alpha/2}^c \leq \frac{\bar{z}-\mu}{S/\sqrt{n}} \leq x_{\alpha/2}^c\right) = \dots = P\left(\bar{z} - \frac{x_{\alpha/2}^c S}{\sqrt{n}} \leq \mu \leq \bar{z} + \frac{x_{\alpha/2}^c S}{\sqrt{n}}\right)$
 (skrivs ofta $t_{1-\frac{\alpha}{2}}(v) = t_{1-\frac{\alpha}{2}}(n-1)$)

⇒ Ett $(1-\alpha) \cdot 100\%$ igt konfidensintervall för $\mu = \mathbb{E}[\xi_i]$
 när $\sigma^2 = \text{Var}(\xi_i)$ är okänd och skattas med s^2 ges
 av $\bar{x} \pm t_{1-\frac{\alpha}{2}}(n-1) \frac{s}{\sqrt{n}}$

- Ex. Längderna på 15 stugor till bostar antas normalfördelade och vi vet att

$$\sum_{i=1}^{15} x_i = 199, \quad \sum_{i=1}^n x_i^2 = 2787$$

$$\Rightarrow \bar{x} = \frac{199}{15} \approx 13.267 \text{ och } s^2 = \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n-1} = \frac{\sum_{i=1}^n x_i^2 - 2\bar{x}\sum_{i=1}^n x_i - n\bar{x}^2}{n-1} =$$

$$= \frac{\sum_{i=1}^n x_i^2 - \frac{(\sum_{i=1}^n x_i)^2}{n}}{n-1} = \frac{2787 - \frac{199^2}{15}}{14} \Rightarrow s \approx 3.24$$

Antal frihetsgrader $v = \text{df}$ är $n-1=14$ och
 vi vill ha ett 95% igt konfidensintervall

$$\Rightarrow t_{1-\frac{\alpha}{2}}(v) = t_{1-\frac{0.05}{2}}(14) = 2.15 \text{ enl. tabell}$$

⇒ konfidensintervalllet ges av

$$\bar{x} \pm t_{1-\frac{\alpha}{2}}(n-1) \frac{s}{\sqrt{n}} = 13.267 \pm 2.15 \frac{3.24}{\sqrt{15}} = 13.267 \pm 1.799$$

~~Förslag om konfidensintervall för varianser~~

Eftersom kvarteren av $Z \sim N(0,1)$ har att $\mathbb{E}[Z^2] = \text{Var}(Z) + \mathbb{E}[Z]^2$

och $0 < \sigma^2$ så är $\sum_{i=1}^n (\xi_i - \mu)^2 = \sum_{i=1}^n Z_i^2$ vettig för att hitta

eftersom konfidensintervall för σ^2 ges av $\sigma^2 \sim \chi^2_{n-1}(\sigma^2)$

Konfidensintervall för varianser

Givet $\xi_1, \dots, \xi_n \sim N(\mu, \sigma^2)$ så gäller att

$$\sum_{i=1}^n \frac{(\xi_i - \mu)^2}{\sigma^2} = \sum_{i=1}^n z_i^2 \sim \chi^2(n) \quad \text{"chi-två"-fördelad med } n \text{ frihetsgrader}$$

$\chi^2(v)$ -fördelningen är skew med en ganska komplex fäthetsfunktion (sid. 303).

Vi har att $\xi \sim \chi^2(v) \Rightarrow E[\xi] = v, \text{Var}(\xi) = 2v$.

Om $v=2$ så gäller $\chi^2(2) = \chi^2(2) = \text{Exp}(z)$.

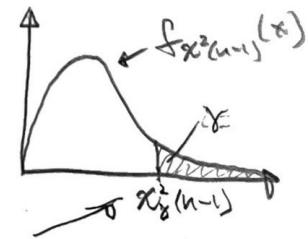
Eftersom $\xi_i - \mu \sim N(0, \sigma^2)$ så skulle vi kunna använda $\frac{1}{\sigma^2} \sum_{i=1}^n (\xi_i - \mu)^2 \sim \chi^2(n)$ för att härleda ett konfidensintervall för σ^2 . Men, vi känner ej till μ i praktiken. Lösning: Använd $\frac{1}{\sigma^2} \sum_{i=1}^n (\xi_i - \bar{\xi})^2$, men vad har detta statistika för fördelning?

Man kan visa att $\chi^2(n-1) \sim \frac{1}{\sigma^2} \sum_{i=1}^n (\xi_i - \bar{\xi})^2 = \frac{n-1}{\sigma^2} \sum_{i=1}^{n-1} (\xi_i - \bar{\xi})^2 = S^2(n-1)$

För $\xi \sim \chi^2(n-1)$, låt $\gamma = P(\xi > \chi^2_{\gamma}(n-1))$, $\gamma \in (0, 1)$.

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

$$= P\left(\underbrace{\frac{(n-1)S^2}{\chi^2_{\alpha/2}(n-1)}}_L \leq \underbrace{\frac{\sigma^2}{S^2}}_U \leq \underbrace{\frac{(n-1)S^2}{\chi^2_{1-\alpha/2}(n-1)}}_R\right), \text{ dvs}$$



Finn nu i tabell

ett $(1-\alpha) \cdot 100\%$ rgt konfidensintervall för $\sigma^2 = \text{Var}(\xi_i)$, där $\xi_1, \dots, \xi_n \sim N(\mu, \sigma^2)$ med både μ och σ^2 okända.

$$\left[\frac{(n-1)S^2}{\chi^2_{\alpha/2}(n-1)}, \frac{(n-1)S^2}{\chi^2_{1-\alpha/2}(n-1)} \right], \text{ För intervallet för } \sigma^2, \text{ ta roten av intervallets gränserna}$$

Ex $n=15$ huvudarkstabsletter där
koncentrationen av alla ingredienser
antas normalford. Givet punktskattningen
0,8 för stand. avv., finn ett 95%igt konf. int. f.v.

hosni Tabell: $\chi^2_{\alpha/2}(n-1) = \chi^2_{0.025}(14) \stackrel{n-1}{=} 26.119$

$$\chi^2_{1-\alpha/2}(n-1) = \chi^2_{0.975}(14) = 5.629$$

$\kappa=15$

$s=0,8$

$$\Rightarrow \left[\sqrt{\frac{(n-1)s^2}{\chi^2_{\alpha/2}(n-1)}}, \sqrt{\frac{(n-1)s^2}{\chi^2_{1-\alpha/2}(n-1)}} \right] = [0,5857, 1,2617]$$