

Matematisk statistik V
 Föreläsningsanteckningar
 Tommy Norberg
 22 april 2004

F8-F9: Ch 7 Skattningsteori

Nyckelord: punktskattning, väntevärdesriktighet, konsistens, konfidensintervall, t -fördelningen, prediktionsintervall

Citat ur Devore & Farnum, s 271.

The general objective of statistical inference is to use sample information as a basis for drawing various types of conclusions. In an estimation problem, we want to make an educated guess about the value of some population characteristic or parameter, such as the population mean battery lifetime μ , the proportion π of all components of a certain type that need service while under warranty, or the difference $\mu_1 - \mu_2$ between the population mean lifetimes for two different types of batteries. The simplest type of estimate is a point estimate, a single number that represents our best guess for the value of the parameter. Thus we might report a point estimate of 758 hours for the population mean lifetime of all brand X 100-watt lightbulbs; we are not saying that $\mu = 758$, only that sample data suggests 758 as a very plausible value for μ .

7.1 Punktskattning (s 272-275)

Vi ska tänka oss att vi har en parameter, θ , som vi är intresserade utav, och att vi har observationer x_1, \dots, x_n av en variabel x , som vi vill använda till att skatta eller göra en begåvad gissning av θ :s värde. Vi låter X beteckna motsvarande stokastiska variabel, och antar att den har väntevärdet μ och standardavvikelsen σ .

I det som följer kommer vi alltsomoftast tyst utgå ifrån att mätmetodiken är väntevärdesriktig, d.v.s. att $\theta = \mu$. Då låter vi istället μ beteckna paramatern vi observerar.

Vi ska även studera observationsmetodikens precision. I sådana fall fokuserar vi intresset på variansen σ^2 och/eller standardavvikelsen σ .

Exempel 1 Vi mäter en teknisk storhet (hållfasthet, sträckgräns, e dyl) som vi just nu kan kalla μ , och antar att variabeln vi observerar, x , är normalfördelad med parametrar μ och okänd standardavvikelse σ .

Exempel 2 Vi vill bestämma en viss mätutrustnings noggranhets, så vi mäter på en ”normal” som vi vet är μ enheter stor. Vi skulle kunna anta att variabeln vi observerar, x , är normalfördelad med parametrar μ och σ .

Exempel 3 Vi vill bestämma olycksfrekvensen för en viss typ av vägavsnitt i Sverige. Den här typen av vägar är relativt vanlig, så det kostar för mycket att analysera alla. Därför väljer man ut n st och bestämmer för var och en av dessa frekvensen olyckor under mars månad ett visst år. Här kan det vara rimligt att tänka sig att variabeln vi mäter, x , är Poissonfördelad med parameter λ (enhet: olyckor/månad). Uppgiften är att bestämma λ med så stor noggranhets som möjligt.

Exempel 4 Vi vill bestämma värdet av sannolikheten $p = P(A)$ att en viss händelse A inträffar då vi gör ett slumpmässigt försök. Vi utför därför n oberoende upprepningar av försöket och låter $x_i = 1$ om A inträffar i det i :te försöket och 0 annars.

Definition 1 (s 272) En punktskattning av en parameter θ är ett reellt tal $\hat{\theta}$, beräknat ur data x_1, \dots, x_n , som vi kan använda som en (mer eller mindre begåvad) gissning av θ :s värde.

En punktskattning $\hat{\theta}$ är alltså en funktion av observationerna. Detta betyder att den är en stickprovsvariabel. Vi använder samma beteckning för mptsvarande stokastiska variabel.

Definition 2 (s 274) En punktskattning $\hat{\theta}$ av θ är väntevärdesriktig (”unbiased”) om

$$E[\hat{\theta}] = \theta$$

Definition 3 (s 275) En punktskattning $\hat{\theta}$ av θ är konsistent, om

$$P(|\hat{\theta} - \theta| > \epsilon) \rightarrow 0 \text{ då } n \rightarrow \infty$$

Man kan visa att punktskattningen $\hat{\theta}$ av θ är konsistent, om dess varians går mot noll då antalet observationer går mot oändligheten, d v s om

$$\text{Var}[\hat{\theta}] \rightarrow 0 \text{ då } n \rightarrow \infty$$

Exempel 5 Medelvärdet \bar{x} är en väntevärdesriktig och konsistent skattning av väntevärdet μ , ty

$$E[\bar{x}] = \mu$$

och

$$\text{Var}[\bar{x}] = \frac{\sigma^2}{n}$$

Att $\text{Var}[\bar{x}] = \sigma^2/n \rightarrow 0$ innebär alltså att

$$P(|\bar{X} - \mu| > \epsilon) \rightarrow 0$$

då $n \rightarrow \infty$ för varje $\epsilon > 0$. Men ϵ/δ -jargon från matematiken gäller alltså för alla $\epsilon > 0$ och för varje $\delta > 0$ att det finns ett heltal n_0 , sådant att

$$n \geq n_0 \Rightarrow P(|\bar{X} - \mu| > \epsilon) \leq \delta$$

Alltså, vill vi ha ett maximalt absolut fel $\epsilon > 0$, så kan vi för hur litet $\delta > 0$ som helst försäkra oss om att

$$P(|\bar{X} - \mu| \leq \epsilon) \geq 1 - \delta$$

om bara antalet observationer är tillräckligt många. En svårighet är att det är sällan som man på förhand vet hur många observationer man behöver göra.

Exempel 6 Stickprovsvariancen s^2 är en skattning av variansen σ^2 och s är en skattning av σ . Senare i kurserna ska vi se om någon av dessa är väntevärdesriktig.

7.2-7.3 Konfidensintervall (stora stickprov) (s 276-282, 285-291)

Vi börjar med fallet då vi har n oberoende mätningar x_1, \dots, x_n av en variabel x och vi är intresserade utav x :s väntevärde μ och vi låter σ beteckna x :s standardavvikelse.

Vi punktskattar μ väntevärdesriktigt och konsistent med medelvärdet \bar{x} . Alltså,

$$\hat{\mu} = \bar{x}$$

Vi ska nu se hur vi kan m.h.a. centrala gränsvärdessatsen dra slutsaser om hur bra vår punktskattning av μ är.

Vi noterar först att enl cgs är medelvärdet \bar{x} approximativt normalfördelat med parametrar μ och σ/\sqrt{n} . Approximationen blir bättre ju fler observationer som görs. Man bör nog ha minst ca 30 observationer.

Vi ska gå igenom ett antal faktarutor i läroboken.

Den första faktarutan ger formeln för ett konfidensintervall för μ med den approximativa konfidensgraden 95%. Här är $z_c = 1.96 \approx 2$.

(Sida 278) **Konfidensintervall för μ med approximativ konfidensgrad 95%**

$$\mu = \bar{x} \pm 1.96 s / \sqrt{n}$$

Regeln förutsätter att $n \geq 30$.

Vi generalisering nu detta resonemang till godtyckliga konfidensgrader.

(Sida 281) **Konfidensintervall för μ med approximativ konfidensgrad 90%, 95% eller 99%**

$$\mu = \bar{x} \pm z_c s / \sqrt{n}$$

där $z_c = 1.645, 1.96$ eller 2.575 beroende på om konfidensgraden är 90%, 95% eller 99%. Regeln förutsätter att $n \geq 30$.

Proportioner är ju också medelvärden, så resonemanget fungerar även för dem. Byt ut medelvärdet \bar{x} mot relativ frekvensen f/n . Observera att

$$\mu_{f/n} = p$$

$$\sigma_{f/n}^2 = \frac{p(1-p)}{n}$$

och att

$$s^2/n \approx \frac{\frac{f}{n}(1 - \frac{f}{n})}{n}$$

Nu bör vi kunna förstå faktaruta på s 286 längst ned.

(Sida 286) **Konfidensintervall för p med approximativ konfidensgrad 90%, 95% eller 99%**

$$p = \hat{p} \pm z_c \sqrt{\hat{p}(1 - \hat{p})/n}$$

där $\hat{p} = f/n$ är relativ frekvensen och $z_c = 1.645, 1.96$ eller 2.575 beroende på om konfidensgraden är 90%, 95% eller 99%. Regeln ska endast användas då man är hyfsat övertygad om att både $np \geq 5$ och att $n(1 - p) \geq 5$ gäller.

Vi generalisering nu resonemanget till en situation där vi har två stickprov och vi är intresserade utav differensen av väntevärdena.

Stickprov nr 1 består av n_1 mätningar av en variabel med väntevärde μ_1 och standardavvikelse σ_1 . Stickprovets medelvärde betecknas \bar{x}_1 och dess standardavvikelse s_1 . Analoga beteckningar har vi för stickprov nr 2.

(Sida 289) **Fakta om fördelningen för en differens av medelvärden**

1. $E[\bar{X}_1 - \bar{X}_2] = \mu_1 - \mu_2$
2. $\text{Var}[\bar{X}_1 - \bar{X}_2] = \frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}$
3. $\bar{X}_1 - \bar{X}_2$ är normalfördelad om de två stickproven består av normalfördelade observationer
4. $\bar{X}_1 - \bar{X}_2$ är approximativt normalfördelad om stickprovsstorlekarna n_1, n_2 ej är för små

(Sida 290) **Konfidensintervall för en differens av väntevärden**

$$\mu_1 - \mu_2 = \bar{x}_1 - \bar{x}_2 \pm z_c \sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}$$

där $z_c = 1.645, 1.96$ eller 2.575 beroende på om konfidensgraden är approximativt 90%, 95% eller 99%.

7.4-7.5 Konfidensintervall i normalfördelningsmodellen (s 294-299, 302-306)

Sats 1 (s 294) Om vårt stickprov består av normalfördelade observationer, så är det standardiserade medelvärdet,

$$Z = \frac{\bar{X} - \mu}{\sigma/\sqrt{n}}$$

exakt normalfördelat med parametrar 0 och 1. Om man i det byter den teoretiska standardavvikelsen σ mot den uppmätta s , och får

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

har vi istället en t -fördelad variabel med $n - 1$ frihetsgrader.

Obs. att för fixt $x > 0$ gäller att

$$P(-x \leq Z \leq x) < P(-x \leq T \leq x) \searrow P(-x \leq Z \leq x)$$

då antalet observationer $n \rightarrow \infty$. Härur fälder att om z_α och $t_\alpha(n-1)$ uppfyller $P(Z > z_\alpha) = P(T > t_\alpha(n-1)) = \alpha$, så gäller att

$$z_\alpha < t_\alpha(n-1) \searrow z_\alpha$$

(sida 296) **Konfidensintervall för μ med exakt konfidensgrad $1 - \alpha$**

$$\mu = \bar{x} \pm t_{\alpha/2} s / \sqrt{n}$$

där $t_{\alpha/2}$ uppfyller $P(T > t_{\alpha/2}) = \alpha/2$ för $T \sim t(n-1)$. Här förutsätts att stickprovet består av $N(\mu, \sigma^2)$ -fördelade observationer.

Sats 2 (s 298) Antag att \bar{x}, s vara medelvärde och standardavvikelse av n oberoende observationer x_1, \dots, x_n av $X \sim N(\mu, \sigma^2)$. Låt X_{n+1} vara nästa ännu ej sedda observation. Då

$$\frac{X_{n+1} - \bar{x}}{s \sqrt{1 + 1/n}} \Big| x_1, \dots, x_n \sim t(n-1)$$

(Sida 298) **Prediktionsintervall** Givet att de n första observationerna är x_1, \dots, x_n , gäller

$$P\left(\bar{x} - t_{\alpha/2} s \sqrt{1 + 1/n} \leq X_{n+1} \leq \bar{x} + t_{\alpha/2} s \sqrt{1 + 1/n}\right) = 1 - \alpha$$

Läs ej avsnittet om toleransintervall.

(Sida 303) **Tvåstickprovsfallet**

Sats 3 Låt \bar{x}_1, s_1 och \bar{x}_2, s_2 vara medelvärde och standardavvikelse av två oberoende stickprov av storlekarna n_1 resp. n_2 av $X_1 \sim N(\mu_1, \sigma_1)$ och $X_2 \sim N(\mu_2, \sigma_2)$. Då gäller att

$$t = \frac{\bar{x}_1 - \bar{x}_2 - (\mu_1 - \mu_2)}{\sqrt{s_1^2/n_1 + s_2^2/n_2}} \stackrel{\text{ap}}{\sim} t(\nu)$$

där

$$\nu \approx \frac{s_1^2/n_1 + s_2^2/n_2}{\frac{(s_1^2/n_1)^2}{n_1 - 1} + \frac{(s_2^2/n_2)^2}{n_2 - 2}}$$

(avrunda nedåt). Härur följer att

$$\mu_1 - \mu_2 = \bar{x}_1 - \bar{x}_2 \pm t_{\alpha/2} \sqrt{s_1^2/n_1 + s_2^2/n_2}$$

är ett konfidensintervall för differensen $\mu_1 - \mu_2$ med den approximativa konfidensgraden $1 - \alpha$.

(Sida 305) **Parade (bivariata) data**

Låt \bar{d}, s_d beteckna medelvärde och standardavvikelse för ett stickprov bestående av n N-fördelade differenser $D = X - Y$. Då är

$$\mu_d = \bar{d} \pm t_{\alpha/2} s_d / \sqrt{n}$$

ett konfidensintervall för $\mu_d = \mu_X - \mu_Y$ med konfidensgraden $1 - \alpha$. Antalet frihetsgrader i $t_{\alpha/2}$ är $n - 1$.

7.6 Mer skattningsteori (s 309-313)

I detta avsnittet läser vi bara om

Trolighetsmetoden ("Maximum likelihood estimation") (s 309-313)

Exempel 7.16 (s 311) 12 st livstider

$$\begin{array}{ccccccc} 10\,502 & 9\,560 & 11\,671 & 12\,825 & 8\,987 & 7\,924 \\ 9\,508 & 8\,875 & 14\,439 & 11\,320 & 6\,549 & 10\,654 \end{array}$$

Modell: exponentialfördelning med parameter λ .

Problem: Skatta λ .

Teoretisk lösning: Antag att du har n oberoende observationer t_1, \dots, t_n av $T \sim \text{Exp}(\lambda)$. Exponentialfördelningens täthetsfunktion är

$$f(t) = \lambda e^{-\lambda t} \quad t > 0$$

Det följer att stickprovets täthet är

$$f(t_1, \dots, t_n) = \prod_i f(t_i) = \prod_i \lambda e^{-\lambda t_i} = \lambda^n e^{-\sum_i t_i}$$

Observera att nu har vi observerat t_1, \dots, t_n , så ovanstående täthet beror bara av λ , som är okänd. Stickprovets *trolighetsfunktion* eller bara *trolighet* är

$$L(\lambda) = \lambda^n e^{-\sum_i t_i}$$

Observera att detta är en funktion av λ .

ML-skattningen (*trolighetsskattningen*, "maximum likelihood-skattningen") av λ är det värde $\hat{\lambda}$ som maximerar troligheten $L(\lambda)$. Man brukar skriva

$$\hat{\lambda} = \arg \max_{\lambda} L(\lambda)$$

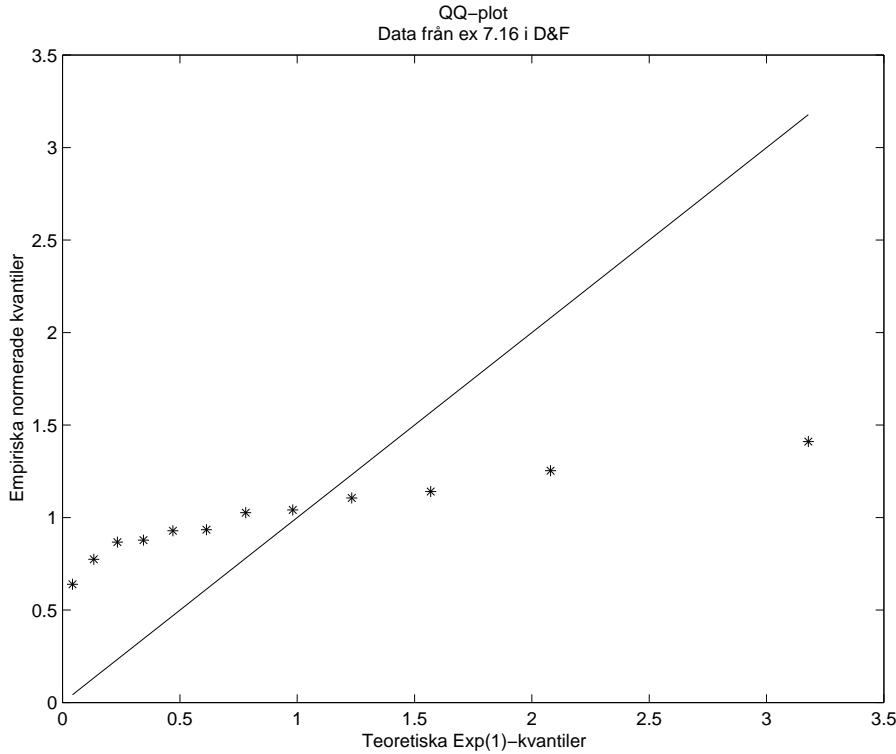
Löser man detta maximeringsproblem erhålls

$$\hat{\lambda} = \frac{1}{\bar{t}}$$

I Exponentialfördelningsmodellen gäller alltså att ML-skattningen av λ är ett genom medelvärdet.

Lösning av problemet: Vi har $n = 12$ observationer. Deras medelvärde är $\bar{t} = 10\,234.5$, så ML-skattningen av λ blir $\hat{\lambda} = 1/10\,234.5 \approx 9.771 \times 10^{-5}$.

Kontroll av modellen: Vi plottar det ordnade stickprovet normaliserat mot motsv. kvantiler i Exp(1)-fördelningen och erhåller



Exp-fördelningen verkar inte vara någon bra modell för data i detta exempel.

ML-metoden

Du har n oberoende observationer x_1, \dots, x_n av en stokastisk variabel X , vars täthet eller massfunktion är $f(x|\theta)$ (θ är en en- eller fler-dimensionell parameter och vår uppgift är att skatta $\vartheta = h(\theta)$).

1) Bilda trolighetsfunktionen

$$L(\theta) = \prod_i f(x_i|\theta)$$

2) Maximera denna, d.v.s. sök

$$\hat{\theta} = \arg \max_{\theta} L(\theta)$$

3) ML-skattningen av ϑ är

$$\hat{\vartheta} = h(\hat{\theta})$$

4) Tips: Låt $\mathcal{L}(\theta) = \ln L(\theta)$. Då gäller att

$$\hat{\theta} = \arg \max_{\theta} \mathcal{L}(\theta)$$

Det är i allmänhet betydligt enklare att maximera $\mathcal{L}(\theta)$.

5) Glöm inte att kontrollera modellen, t.ex. genom att plotta de empiriska kvantilerna mot de teoretiska.

ML-skattningens egenskaper (s 313)

- For large n , the sampling distribution of an MLE is approximately normal, and the estimator is nearly unbiased with a variance smaller than any other estimator.

Exempel 7 Om det vi mäter, x , är Bernoullifördelat med parameter p , så är ML-skattningen av p , lika med $\hat{p} = \sum_i x_i/n = \bar{x}$.