



## Föreläsning 05

Statistisk inläring och dataanalys (Kungliga Tekniska Högskolan)

# Föreläsning 5

## 5.1 Maximum-likelihoodmetoden

Medan momentmetoden är snabb och enkel ger den inte alltid de bästa skattningarna. En annan metod är baserad på *likelihoodfunktionen*:

$$\begin{aligned} L(\boldsymbol{\theta}|\mathbf{x}) &= L(\boldsymbol{\theta}|x_1, \dots, x_n), \\ &= f_{\mathbf{X}}(x_1, \dots, x_n|\boldsymbol{\theta}), \\ &= \prod_{i=1}^n f_{X_i}(x_i|\boldsymbol{\theta}), \end{aligned}$$

för stickprovet  $x_1, \dots, x_n$  från en population  $f_X(x|\boldsymbol{\theta})$ . Likelihoodfunktionen är bara den simultana täthetsfunktionen (eller sannolikhetsfunktionen) men där värdena  $x_1, \dots, x_n$  är fixa och  $\boldsymbol{\theta}$  varierar. I diskreta fallet ger likelihoodfunktionen sannolikheten av stickprovet för olika värden av parametrarna. I detta sätt svarar likelihoodfunktionen frågan om hur troligt är datan för varje givet val av  $\boldsymbol{\theta}$ .

Låt  $\hat{\boldsymbol{\theta}}$  vara valet för  $\boldsymbol{\theta}$  som maximeras sannolikheten att se datan  $\mathbf{x}$ ; det vill säga

$$\hat{\boldsymbol{\theta}} = \operatorname{argmax}_{\boldsymbol{\theta} \in \Omega} L(\boldsymbol{\theta}|\mathbf{x}).$$

Funktionen  $\hat{\boldsymbol{\theta}}$  är en funktion  $\hat{\boldsymbol{\theta}} : \mathbb{R}^n \rightarrow \mathbb{R}^d$  som skickar möjliga observationer till parametervärden. När vi tillämpar funktionen till stickprovet  $\mathbf{x}$  får vi den *maximum-likelihoodskattning* (ML-skattning)  $\hat{\boldsymbol{\theta}}$  av  $\boldsymbol{\theta}$ .

Maximum-likelihoodskattningen är ofta unik; det vill säga, det finns bara ett val för  $\boldsymbol{\theta}$  som maximeras  $L(\boldsymbol{\theta}|\mathbf{x})$ . Detta behöver dock inte fallet. (ML-skattningen behöver egentligen inte existera!) Det vore mer korrekt att kalla  $\hat{\boldsymbol{\theta}}$  en ML-skattning men det är brukligt att kalla den ML-skattningen. ML-skattningar är ganska lätta att definiera men de kan vara svåra att beräkna i praktiken. Dock eftersom ML-skattningen maximerar en funktion har vi ett generellt recept för att beräkna dem:

1. Beräkna de partiella derivaten  $\frac{\partial}{\partial \theta_j} L(\boldsymbol{\theta}|\mathbf{x})$  där  $\boldsymbol{\theta} = (\theta_1, \dots, \theta_d)$ . Det är vanligt att beräkna de partiella derivaten av *log-likelihoodfunktionen*  $\log L(\boldsymbol{\theta}|\mathbf{x})$  istället. Båda metoderna ger det samma svar eftersom log är en strikt växande funktion.
2. Beräkna kritiska punkter av  $L(\boldsymbol{\theta}|\mathbf{x})$  (eller  $\log L(\boldsymbol{\theta}|\mathbf{x})$ ); dvs, lösningarna till systemet av ekvationerna ges igenom sätter den partiella derivaten lika med 0.
3. Bestäm vilka kritiska punkter är (lokala) maxima igenom att beräkna matrisen

$$\left[ \frac{\partial^2}{\partial \theta_i \partial \theta_j} L(\boldsymbol{\theta}|\mathbf{x}) \right]_{i,j=1}^d$$

och bestämma för vilka kritiska punkter matrisen är negativt definit (alla egenvärden negativa). (Detta går också att göra med  $\log L(\boldsymbol{\theta}|\mathbf{x})$ ). Ibland kan vi använda algebraisk metoder för att verifiera att en kritisk punkt är en lokal maximum istället.

4. Jämför värdet av  $L(\boldsymbol{\theta}|\mathbf{x})$  för alla lokal maxima och gränspunkter (inklusive oändlighet) och ta punkten med det största värdet.

**Exempel 5.1.** Låt  $X_1, \dots, X_n$  vara oberoende och  $\text{Ber}(p)$ -fördelade. Så har vi likelihood-funktionen

$$L(p|\mathbf{x}) = L(p|x_1, \dots, x_n) = p^{\sum_{i=1}^n x_i} (1-p)^{n-\sum_{i=1}^n x_i}.$$

För att förenkla beräkningen ska vi använda log-likelihoodfunktionen

$$\log L(p|\mathbf{x}) = \log(p) \left( \sum_{i=1}^n x_i \right) + \log(1-p) \left( n - \sum_{i=1}^n x_i \right).$$

Vi har att

$$\frac{d}{dp} \log L(p|\mathbf{x}) = \frac{1}{p} \left( \sum_{i=1}^n x_i \right) - \frac{1}{1-p} \left( n - \sum_{i=1}^n x_i \right).$$

För att bestämma de kritiska punkterna sätter vi detta lika med 0 och löser för  $p$ :

$$\begin{aligned} 0 &= \frac{1}{p} \left( \sum_{i=1}^n x_i \right) - \frac{1}{1-p} \left( n - \sum_{i=1}^n x_i \right), \\ &= \frac{n}{p} \bar{x} - \frac{n}{1-p} (1 - \bar{x}), \\ pn(1 - \bar{x}) &= n\bar{x}(1 - p), \\ p &= \bar{x}. \end{aligned}$$

Det finns bara en kritisk punkt,  $p = \bar{x}$ . Denna punkt är bara möjlig när  $0 < n\bar{x} < n$ . Annat skulle en term i ovanstående ekvationen försvinna och lämna oss med en ekvation utan en lösning. Så om vi anta att  $0 < n\bar{x} < n$  ser vi att  $p = \bar{x}$  är en lokal maximum eftersom

$$\frac{d^2}{dp^2} \log L(p|\mathbf{x}) = -\frac{1}{p^2} n\bar{x} - \frac{1}{(1-p)^2} (n - n\bar{x}) < 0$$

där vi använder att  $n\bar{x} > 0$  och  $n - n\bar{x} > 0$ . För  $n\bar{x} = 0$  respektiv  $n\bar{x} = n$  är log-likelihoodfunktionerna strikt avtagande och strikt växande. Därför får vi ML-skattningarna 0 respektiv 1 vilka är lika med  $\bar{x}$  i båda fall. Det följer att ML-skattningen för  $p$  är  $\hat{p} = \bar{X}$ .

Vi kan också använda det samma metoden med modellen som beror på mer än en parameter.

**Exempel 5.2.** Låt  $X_1, \dots, X_n$  vara oberoende och  $N(\mu, \sigma^2)$ -fördelade med båda  $\mu$  och  $\sigma^2$  okända. Likelihoodfunktionen är

$$L(\mu, \sigma^2|\mathbf{x}) = \frac{1}{(2\pi\sigma^2)^{\frac{n}{2}}} \exp \left( -\frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2 \right),$$

och så har vi log-likelihoodfunktionen

$$\log L(\mu, \sigma^2|\mathbf{x}) = -\frac{n}{2} \log(2\pi) - \frac{n}{2} \log(\sigma^2) - \frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2.$$

Det följer att de partiella derivaten är

$$\begin{aligned} \frac{\partial}{\partial \mu} \log L(\mu, \sigma^2|\mathbf{x}) &= \frac{1}{\sigma^2} \sum_{i=1}^n (x_i - \mu), \\ \frac{\partial}{\partial \sigma^2} \log L(\mu, \sigma^2|\mathbf{x}) &= -\frac{n}{2\sigma^2} + \frac{1}{2\sigma^4} \sum_{i=1}^n (x_i - \mu)^2. \end{aligned}$$

(Notisera att den andra derivatan är med avseende på  $\sigma^2$  (inte  $\sigma$ ).) Det följer att

$$\begin{aligned} \hat{\mu} &= \frac{1}{n} \sum_{i=1}^n x_i = \bar{x}, \\ \hat{\sigma}^2 &= \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2 = \frac{n-1}{n} s^2 \end{aligned}$$

är den kritiska punkten. Vi behöver att verifiera att denna punkt är en lokal maximum. Istället för att beräkna andra derivat notiserar vi att för alla  $\mu$  vi har

$$\sum_{i=1}^n (x_i - \mu)^2 > \sum_{i=1}^n (x_i - \bar{x})^2.$$

Så har vi att

$$L(\mu, \sigma^2 | \mathbf{x}) = \frac{1}{(2\pi\sigma^2)^{\frac{n}{2}}} \exp\left(-\frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2\right) \leq \frac{1}{(2\pi\sigma^2)^{\frac{n}{2}}} \exp\left(-\frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \bar{x})^2\right) = L(\bar{x}, \sigma^2 | \mathbf{x})$$

för alla  $\sigma^2 > 0$ . Så är det tillräckligt att visa likelihoodfunktionen maximerar på  $\sigma^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2$  när  $\mu = \bar{x}$ . Den andra derivatan med avseende på  $\sigma^2$  är

$$\begin{aligned} \frac{d^2}{d(\sigma^2)^2} \log L(\bar{x}, \sigma^2 | \mathbf{x}) &= \frac{n}{2\sigma^4} - \frac{1}{\sigma^6} \sum_{i=1}^n (x_i - \bar{x})^2, \\ &= -\frac{2}{\sigma^6} \left( \sum_{i=1}^n (x_i - \bar{x})^2 - \frac{n\sigma^2}{2} \right). \end{aligned}$$

Eftersom faktoren

$$\left( \sum_{i=1}^n (x_i - \bar{x})^2 - \frac{n\sigma^2}{2} \right)$$

är positiv när  $\sigma^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2$  vet vi att den andra derivatan är negativ. Därför har vi en lokal maximum. Vi har också att  $L(\bar{x}, \sigma^2 | \mathbf{x}) \rightarrow 0$  när  $\sigma^2 \rightarrow \infty$ . Så har vi en global maximum. Det har redan konstaterats att  $\hat{\mu}$  ger en global maximum för  $\mu$  för vilken  $\sigma^2$  som helst. Därför är ML-skattningen av  $(\mu, \sigma^2)$  lika med  $(\hat{\mu}, \hat{\sigma}^2) = (\bar{x}, \frac{n-1}{n} s^2)$ .

Resultatet i Exempel 5.2 generaliserar även till multivariatnormalfördelningar. Matrisen  $\hat{\Sigma}$  i följande satsen kallas *stickprovscovariansmatrisen* och den är garanterat att vara positivt definit.

**Sats 5.1.** Låt  $\mathbf{X}_1, \dots, \mathbf{X}_n$  vara oberoende och  $N(\boldsymbol{\mu}, \Sigma)$ -fördelade. ML-skattningen för  $(\boldsymbol{\mu}, \Sigma)$  är

$$\hat{\boldsymbol{\mu}} = \frac{1}{n} \sum_{i=1}^n \mathbf{X}_i = \bar{\mathbf{X}}, \quad \hat{\Sigma} = \frac{1}{n} \sum_{i=1}^n (\mathbf{X}_i - \bar{\mathbf{X}})(\mathbf{X}_i - \bar{\mathbf{X}})^T.$$

Nu kanske börjar vi märka att det verkar finnas en koppling mellan ML-skattningar och tillräckliga statistikor. Exempelvis, såg vi att  $\bar{X}$  respektive  $S^2$  är tillräckliga statistikor för  $\mu$  och  $\sigma^2$  när  $X_1, \dots, X_n$  är  $N(\mu, \sigma^2)$ -fördelade och vi såg även att ML-skattningarna för  $\mu$  respektive  $\sigma^2$  är  $\bar{X}$  och  $\frac{n-1}{n} s^2$ . Denna koppling är inte en slump. Med hjälp av Sats 4.4 vet vi att

$$f_{\mathbf{X}}(\mathbf{x} | \boldsymbol{\theta}) = h(\mathbf{x})g(T(\mathbf{x}) | \boldsymbol{\theta})$$

för några funktioner  $h$  och  $g$  när  $T(\mathbf{x})$  är en tillräcklig statistika för  $\boldsymbol{\theta}$  där  $\mathbf{x}$  är ett stickprov från en  $f_{\mathbf{X}}(\mathbf{x} | \boldsymbol{\theta})$ -population. Det följer att  $L(\boldsymbol{\theta} | \mathbf{x})$  maximera på  $\hat{\boldsymbol{\theta}}$  om och endast om  $g(T(\mathbf{x}) | \boldsymbol{\theta})$  maximera på  $\hat{\boldsymbol{\theta}}$ . Det vill säga att ML-skattningen kan alltid vara skriven som en funktion av tillräckliga statistikor.

### 5.1.1 Invariants egenskap hos ML-skattningar

En användbar egenskap av ML-skattningar är att de är invariant under transformationer. In synnerhet är ML-skattningen av  $\tau(\boldsymbol{\theta})$  lika med  $\tau(\hat{\boldsymbol{\theta}})$  när  $\tau: \mathbb{R}^d \rightarrow \mathbb{R}^d$  är injektiv och  $\hat{\boldsymbol{\theta}}$  är ML-skattningen för  $\boldsymbol{\theta}$ . Om  $\tau$  är inte injektiv (och därmed inte inverterbar) är likelihoodfunktionen  $L(\tau(\boldsymbol{\theta}) | \mathbf{x})$  inte väldefinierad. Istället definierar vi *inducerad likelihoodfunktionen*

$$L^*(\boldsymbol{\eta} | \mathbf{x}) = \sup_{\boldsymbol{\theta}: \tau(\boldsymbol{\theta}) = \boldsymbol{\eta}} L(\boldsymbol{\theta} | \mathbf{x}).$$

På ett sätt säger detta att från alla  $\boldsymbol{\theta}$  som kunde ge oss  $\boldsymbol{\eta}$  tar vi den som ger oss den största likelihooden. Valet  $\hat{\boldsymbol{\eta}}$  som maximerar inducerad likelihoodfunktionen definieras som ML-skattningen för  $\boldsymbol{\eta}$ . Medan det kan låta svårt att bestämma sådana ML-skattningar går det lättare med följande:

**Sats 5.2.** Om  $\hat{\theta}$  är ML-skattningen av  $\theta$  så har vi att  $\tau(\hat{\theta})$  är ML-skattningen för  $\tau(\theta)$  för vilken funktion  $\tau : \mathbb{R}^d \rightarrow \mathbb{R}^\ell$  som helst.

*Bevis.* Låt  $\hat{\eta}$  maximera  $L^*(\eta|\mathbf{x})$ . Vi skulle vilja visa att

$$L^*(\hat{\eta}|\mathbf{x}) = L^*(\tau(\hat{\theta})|\mathbf{x}).$$

Då gäller det att  $\tau(\hat{\theta})$  maximera också  $L^*(\eta|\mathbf{x})$  och därför är en ML-skattning.

Eftersom  $\hat{\eta}$  maximerar  $L^*(\eta|\mathbf{x})$  har vi

$$L^*(\hat{\eta}|\mathbf{x}) = \sup_{\eta} L^*(\eta|\mathbf{x}) = \sup_{\eta} \sup_{\theta: \tau(\theta)=\eta} L(\theta|\mathbf{x})$$

som följer direkt från definitionen av  $L^*(\eta|\mathbf{x})$ . Att maximera  $L(\theta|\mathbf{x})$  först över en delmängd  $\{\theta : \tau(\theta) = \eta\}$  och då över alla  $\eta$  är ekvivalent till att maximera över alla  $\theta$ . Så har vi

$$L^*(\hat{\eta}|\mathbf{x}) = \sup_{\theta} L(\theta|\mathbf{x}) = L(\hat{\theta}|\mathbf{x})$$

där  $\hat{\theta}$  är ML-skattningen för  $\theta$ . Vi har även att

$$L(\hat{\theta}|\mathbf{x}) = \sup_{\theta: \tau(\theta)=\tau(\hat{\theta})} L(\theta|\mathbf{x}) = L^*(\tau(\hat{\theta})|\mathbf{x}).$$

□

**Exempel 5.3.** Vi såg i Exempel 5.2 att

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

är ML-skattningen för  $\sigma^2$  när  $X_1, \dots, X_n$  oberoende och  $N(\mu, \sigma^2)$ -fördelade. Eftersom  $\sigma = \tau(\sigma^2)$  där  $\tau(\theta) = \sqrt{\theta}$  har vi att  $\hat{\sigma} = \sqrt{\hat{\sigma}^2}$  är ML-skattningen för standardavvikelsen  $\sigma$ .

**Exempel 5.4.** Låt  $\tau(\theta) = \theta(1 - \theta)$  och låt  $X_1, \dots, X_n$  vara oberoende och  $\text{Ber}(p)$ -fördelade. Variansen av  $X \sim \text{Ber}(p)$  är  $p(1 - p)$  och vi såg i Exempel 5.1 att  $\hat{p} = \bar{X}$  är ML-skattningen för  $p$ . Så följer det från Sats 5.2 att  $\tau(\bar{X}) = \bar{X}(1 - \bar{X})$  är ML-skattningen för variansen av modellen.

Vi avslutar diskussionen om ML-skattningar med en observation om deras relation till momentmetodens punktskattningar. För ett  $N(\mu, \sigma^2)$ -fördelat stickprov såg vi att

$$(\hat{\mu}, \hat{\sigma}^2) = \left( \bar{x}, \frac{n-1}{n} s^2 \right) = (\tilde{\mu}, \tilde{\sigma}^2).$$

Det vill säga att båda punktskattningar är det samma. Det är också snabbt att verifiera att ML-skattningen för  $p$  är det samma som momentmetodens punktskattningen när  $x_1, \dots, x_n$  är ett stickprov från en  $\text{Ber}(p)$ -population. Medan det kan ofta hända dessa är dessa två punktskattningar inte alltid det samma.

**Exempel 5.5.** Låt  $X_1, \dots, X_n$  vara obereonde och  $U(0, \theta)$ -fördelade. Momentmetodens punktskattningen för  $\theta$  beräknas som

$$\frac{\theta}{2} = E[X|p] = \frac{1}{n} \sum_{i=1}^n X_i = \bar{X},$$

och därför är  $\tilde{\theta} = 2\bar{X}$ . Å andra sidan har vi att likelihoodfunktionen för  $\theta$  är

$$L(\theta|\mathbf{X}) = \frac{1}{\theta^n}$$

vilken är en strikt avtagande funktion av  $\theta$ . Det följer att ML-skattningen för  $\theta$  är  $\hat{\theta} = \max_{i \in [n]} (X_i)$ . Generellt sett är  $\tilde{\theta} \neq \hat{\theta}$ . Exemplevis, om vi har ett stickprov

$$x_1 = 0.28, x_2 = 0.43, x_3 = 0.88, x_4 = 0.48, x_5 = 0.2$$

sa får vi att  $\tilde{\theta} = 0.908$  och  $\hat{\theta} = 0.88$ .

Nu har vi olika metoder för att bestämma en punktskattning för parametrarna som kan ge oss olika skattningar. Frågan blir "vilken punktskattning borde vi använda?"

## 5.2 Utvärdering av skattningar

För att bestämma oss vilken skattning är bättre behöver vi metoder för att utvärdera skattningar. Utvärdering av skattningar använder vanligtvis funktioner som kallas förlustfunktioner. Studiet av utförande och optimalitet av skattningar som utvärderas igenom förlustfunktioner kallas *beslutsteori*. Given data  $x_1, \dots, x_n$  från en  $f_X(x|\theta)$ -population (där  $\theta$  är okänt) kan vi fatta ett val (eller *beslut*) om  $\theta$ . Exempelvis, vi fattar beslutet att  $\theta = \hat{\theta}$ . *Handlingsrummet*  $\mathcal{A}$  är mängden av alla tillåtna beslut. Handlingsrummet är vanligtvis lika med parameterrummet; dvs,  $\mathcal{A} = \Omega$  men det behöver inte vara så.

En *förlustfunktion* är en funktion som kvantifiera vårt beslut  $\mathbf{a} \in \mathcal{A}$  relativt till  $\theta$ . Om  $\mathbf{a}$  är långt borta från  $\theta$  så ska förlusten (dvs, värdet av förlustfunktionen) vara stor. Samtidigt om  $\mathbf{a}$  är nära till  $\theta$  så ska förlusten vara liten. I detta sätt ger förlustfunktioner oss ett mått på hur mycket vi förlorar när vi skatta  $\theta$  med  $\mathbf{a}$ .

Några vanliga förlustfunktioner inkluderar

- *absolutfel*:  $L(\theta, a) = |a - \theta|$ ,
- *kvadratfel*:  $L(\theta, a) = (a - \theta)^2$ .

Det går att generalisera dessa till flera parameter  $\theta = (\theta_1, \dots, \theta_d)$  och beslut  $\mathbf{a} = (a_1, \dots, a_d)$  som

$$L(\theta, \mathbf{a}) = \sum_{i=1}^n |a_i - \theta_i| \quad \text{respektiv} \quad L(\theta, \mathbf{a}) = \sum_{i=1}^n (a_i - \theta_i)^2.$$

Eftersom  $L(\theta, \mathbf{a}) = 0$  om och endast om  $\mathbf{a} = \theta$  är vår förlust 0 om och endast om vårt beslut är korrekt. Olika val av förlustfunktioner kan ge oss naturligtvis olika förlust för det samma beslut. Exempelvis, absolutfel ger en större förlust för liten fel relativt till kvadratfel. Å andra sidan ger kvadratfel oss en större förlust för större fel relativt till absolutfel. Vi kan även definiera förlustfunktioner som ger oss mindre fel när vi överskattar  $\theta$  än när vi underskattar  $\theta$ :

$$L(\theta, a) = \begin{cases} (a - \theta)^2 & a < \theta, \\ 10(a - \theta)^2 & a \geq \theta. \end{cases}$$

Anta att vi beräknar en skattning  $W(\mathbf{X})$  för parametrarna  $\theta$  från ett stickprov  $\mathbf{X}_1, \dots, \mathbf{X}_n$  dras från en  $f_X(\mathbf{x}|\theta)$ -population. Given en förlustfunktion kvantifieras kvaliteten av en skattning  $W(\mathbf{X})$  med hjälp av en *riskfunktion*:

$$R(\theta, W(\mathbf{X})) = E[L(\theta, \mathbf{a})].$$

Noterar att väntevärdet tas mot avseende på  $\mathbf{X}$ ; dvs,

$$R(\theta, W(\mathbf{X})) = \int_{\mathbf{X}} L(\theta, \mathbf{a}) f_{\mathbf{X}}(\mathbf{x}|\theta) d\mathbf{x}.$$

Värdet  $R(\theta, W(\mathbf{X}))$  för en given  $\theta$  är den genomsnittliga förlusten som vi ska ha när vi använder  $W(\mathbf{X})$  som en punktskattning för  $\theta$ . Eftersom  $\theta$  är okänt stämmer det att använda en punktskattning  $W(\mathbf{X})$  för vilken  $R(\theta, W(\mathbf{X}))$  är liten för alla  $\theta$ . Ett sätt för att jämföra olika punktskattningar är att jämföra deras risk. Exempelvis, om  $W_1(\mathbf{X})$  och  $W_2(\mathbf{X})$  är olika punktskattningar för  $\theta$  där

$$R(\theta, W_1(\mathbf{X})) < R(\theta, W_2(\mathbf{X})), \quad \text{för alla } \theta \in \Omega$$

så är  $W_1(\mathbf{X})$  den föredragna punktskattningen.

Generellt sett kan det vara svårt att beräkna risken för en given punktskattning och förlustfunktion men för kvadratfel har vi en formel som gör det lättare:

$$E[(W(\mathbf{X}) - \theta)^2] = \text{Var}[W(\mathbf{X})] + (E[W(\mathbf{X})] - \theta)^2.$$

Formellen stämmer eftersom

$$\begin{aligned} E[(W(\mathbf{X}) - \theta)^2] &= E[(W(\mathbf{X}) - E[W(\mathbf{X})] + E[W(\mathbf{X})] - \theta)^2], \\ &= E[(W(\mathbf{X}) - E[W(\mathbf{X})])^2] + E[(W(\mathbf{X}) - E[W(\mathbf{X})])(E[W(\mathbf{X})] - \theta)] + E[(E[W(\mathbf{X})] - \theta)^2], \\ &= \text{Var}[W(\mathbf{X})] + (E[W(\mathbf{X})] - E[W(\mathbf{X})])(E[W(\mathbf{X})] - \theta) + (E[W(\mathbf{X})] - \theta)^2, \\ &= \text{Var}[W(\mathbf{X})] + (0)(E[W(\mathbf{X})] - \theta) + (E[W(\mathbf{X})] - \theta)^2, \\ &= E[(W(\mathbf{X}) - \theta)^2] = \text{Var}[W(\mathbf{X})] + (E[W(\mathbf{X})] - \theta)^2. \end{aligned}$$

Risken  $E[(W(\mathbf{X}) - \theta)^2]$  kallas vanligtvis *medelkvadratfel*. Med hjälp av ovanstående formel kan vi jämföra våra två punktskattningar för parametern  $\theta$  i Exempel 5.5:

**Exempel 5.6.** Låt  $X_1, \dots, X_n$  vara oberoende och  $U(0, \theta)$ -fördelade. Vi bestämde att  $\tilde{\theta} = 2\bar{X}$  respektiv  $\hat{\theta} = \max_{i \in [n]}(X_i)$  var den momentmetodens punktskattningen och ML-skattningen i Exempel 5.5. Medelkvadratfelet för  $\tilde{\theta}$  beräknas som

$$\begin{aligned} E[(2\bar{X} - \theta)^2] &= \text{Var}[2\bar{X}] + (E[2\bar{X}] - \theta)^2, \\ &= \frac{4}{n^2} \sum_{i=1}^n \text{Var}[X_i] + \left( \frac{2}{n} \sum_{i=1}^n E[X_i] - \theta \right)^2, \quad (X_1, \dots, X_n \text{ oberoende}) \\ &= \frac{4}{n^2} \frac{n\theta^2}{12} + \left( \frac{2}{n} \frac{n\theta}{2} - \theta \right)^2, \\ &= \frac{\theta^2}{3n}. \end{aligned}$$

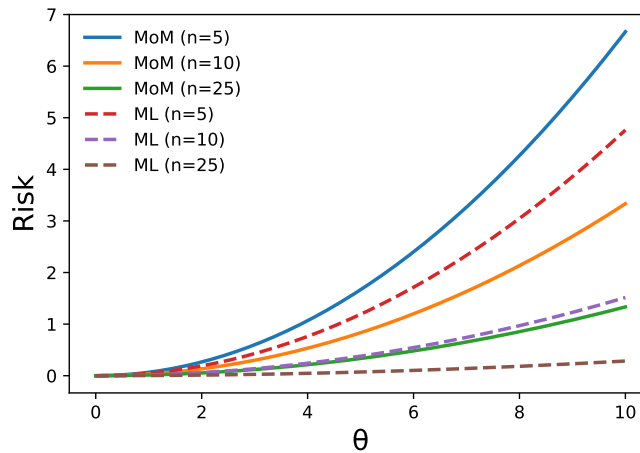
För att beräkna medelkvadratfelet för  $\hat{\theta}$  behöver vi  $\text{Var}[\hat{\theta}]$  och  $E[\hat{\theta}]$ . Det kan visas att samplingfördelningen för  $\hat{\theta}$  är

$$f_{\hat{\theta}}(w) = \frac{1}{\theta^n} n w^{n-1}.$$

Med hjälp av denna formel kan vi beräkna  $E[\hat{\theta}] = \frac{n}{n+1}\theta$  och  $E[\hat{\theta}] = \frac{n}{n+2}\theta^2$ . Det följer att

$$\begin{aligned} E[(\hat{\theta} - \theta)^2] &= \text{Var}[\hat{\theta}] + (E[\hat{\theta}] - \theta)^2, \\ &= \frac{n\theta^2}{n+2} - \left( \frac{n\theta}{n+1} \right)^2 + \left( \frac{\theta}{n+1} \right)^2, \\ &= \frac{2}{(n+2)(n+1)} \theta^2. \end{aligned}$$

Eftersom båda  $\frac{1}{3n}$  och  $\frac{2}{(n+2)(n+1)}$  går till 0 som  $n \rightarrow \infty$  blir risken för båda punktskattningar mindre för större stickprover. Risken för ML-skattningen blir dock mindre snabbare än momentmetodens punktskattningen. Det visas i följande



Så ser ML-skattningen ut som en bättre punktskattning för  $\theta$ .

Å andra sidan såg vi i Exempel 5.6 att  $E[\tilde{\theta}] = \theta$  medan  $E[\hat{\theta}] < \theta$  för alla  $n$ . Så i genomsnitt är momentmetodens punktskattningen mer exakt än ML-skattningen. För att formalisera denna idé definiera vi *systematiskt felet* av en punktskattning  $W(\mathbf{X})$  för  $\theta$  som

$$\text{Bias}[W(\mathbf{X})] = E[W(\mathbf{X})] - \theta.$$

Så har vi att

$$E[(W(\mathbf{X}) - \theta)^2] = \text{Var}[W(\mathbf{X})] + \text{Bias}[W(\mathbf{X})]^2.$$

Generellt sätt mäter systematiskt fel *riktighet* (hur långt borta är punktskattningen från målet) medan variansen mäter *precision* (punktskattningens konsekvens). Vi vill vänligtvis ha båda små riktighet och precision eftersom det

menar att risken är låg. Ibland accepterar vi en liten systematiskt fel om det låter oss reducera variansen (eller vice versa). Det gjorde vi, exempelvis, i Exempel 5.6.

Vi har också sett i Exempel 5.6 att det kan ta mycket arbete för att beräkna båda variansen och systematiskt felet även för enkla modeller. Ett sätt för att bestämma bra punktskattningar i allmänhet är att kräver att systematiskt felet är 0; dvs,  $E[W(\mathbf{X})] = \theta$ . I sådant fall kallas  $W(\mathbf{X})$  en *väntevärdesriktig skattning*.

**Exempel 5.7.** Låt  $X_1, \dots, X_n$  vara oberoende och  $N(\mu, \sigma^2)$ -fördelade. Vi såg i Sats 4.1 att  $\bar{X}$  respektiv  $S^2$  är väntevärdesriktiga skattningar för  $\mu$  och  $\sigma^2$  eftersom

$$E[\bar{X}] = \mu \quad \text{och} \quad E[S^2] = \sigma^2.$$

Vi kan också beräkna deras varianser som

$$\text{Var}[\bar{X}] = \frac{\sigma^2}{n} \quad \text{och} \quad \text{Var}[S^2] = \frac{2\sigma^4}{n-1}.$$

Eftersom dessa skattningar är väntevärdesriktiga följer det att deras varianser är lika med deras medelkvadratfel.

Om vi tar istället ML-skattningen  $\hat{\sigma}^2 = \frac{n-1}{n} S^2$  för  $\sigma^2$  har vi systematisktfelet

$$\text{Bias}[\hat{\sigma}^2] = \frac{n-1}{n} \sigma^2 - \sigma^2 = -\frac{\sigma^2}{n}$$

och varians

$$\text{Var}[\hat{\sigma}^2] = \left(\frac{n-1}{n}\right)^2 \text{Var}[S^2] = \left(\frac{n-1}{n}\right)^2 \frac{2\sigma^4}{n-1} = \frac{2(n-1)}{n^2} \sigma^4.$$

Så är medelkvadratfelet

$$E[(\hat{\sigma}^2 - \sigma^2)^2] = \frac{2(n-1)}{n^2} \sigma^4 + \frac{\sigma^4}{n^2} = \frac{2n-1}{n^2} \sigma^4 < \frac{2}{n-1} \sigma^4 = E[(S^2 - \sigma^2)^2].$$

Så ML-skattningen har en mindre medelkvadratfel än  $S^2$  trots att ML-skattningen är inte väntevärdesriktig och  $S^2$  är väntevärdesriktig.

I våra exemplar hittills har ML-skattningen varit en bättre skattning än alternativen i den meningen att den har haft mindre medelkvadratfel. Det betyder inte nödvändigtvis att man ska alltid använda ML-skattningen. Exempelvis, för stora stickprover finns det inte så mycket skillnad mellan skattningarna. Så kan det vara lättare att använda, exempelvis, momentmetodens skattningen om det är lättare att beräkna. Vi också vet att ML-skattningen (i dessa exemplar) är inte väntevärdesriktig och underskattar parametern. Detta kan vara ett problem i vissa tillämpningar. Samtidigt vi har bara jämfört skattningarna igenom medelkvadratfelet som straffar båda underskattning och överskattning det samma. Detta stämmer nödvändigtvis inte för en skalaparameter som  $\sigma^2$  som är alltid positivt och kan potentiellt vara mer känsligt till underskattning. Sist kan vi se att medelkvadratfelet är en funktion  $\theta$  som kan variera över  $\Omega$ . Det är möjligt att det finns en skattning som fungerar bättre för vissa värdena av  $\theta$  än ML-skattningen men inte för andra.

## 5.3 Cramér-Rao olikheten

Nu fokuserar vi på väntevärdesriktiga skattningar  $W(\mathbf{X})$ . För sådana skattningar minimerar vi medelkvadratfel precis när vi minimerar variansen  $\text{Var}[W(\mathbf{X})]$ . Det finns en kraftfull sats som hjälper oss att hitta väntevärdesriktiga skattningar med den minsta möjliga variansen. Satsen ger oss en nedre gräns om variansen för sådana skattningar. Om vi når denna gräns vet vi att vi har en skattning med det minsta möjliga medelkvadratfelet.

**Sats 5.3.** Låt  $X_1, \dots, X_n$  vara oberoende och likafördelade med sannolikhetsfunktion (eller täthetsfunktion)  $f_X(x|\theta)$ . Låt  $W(\mathbf{X}) = W(X_1, \dots, X_n)$  vara en väntevärdesriktig skattning av  $g(\theta)$  och anta att  $f_X(x|\theta)$  tillhör en exponentialfamilj. Så gäller det att

$$\text{Var}[W(\mathbf{X})] \geq \frac{g'(\theta)^2}{I(\theta)},$$

där

$$I(\theta) = n E \left[ \left( \frac{\partial}{\partial \theta} \log f_X(X|\theta) \right)^2 \right] = -n E \left[ \frac{\partial^2}{\partial \theta^2} \log f_X(X|\theta) \right].$$



Antalet  $I(\theta)$  kallas *informationstalet* eller *Fisher informationen* och värdet  $\frac{1}{I(\theta)}$  kallas den *Cramér-Rao under gränsen*.

**Exempel 5.8.** För  $X_1, \dots, X_n$  oberoende och  $\text{Po}(\lambda)$ -fördelade vet vi att

$$\mathbb{E}[\bar{X}] = \lambda$$

eftersom en  $\text{Po}(\lambda)$ -fördelad stokastisk variabel har väntevärdet lika med  $\lambda$ . (Här använder vi Sats 4.1.) Sannolikhetsfunktionen för en  $\text{Po}(\lambda)$ -fördelad stokastisk variabel är

$$f_X(x|\lambda) = \frac{\lambda^x e^{-\lambda}}{x!}.$$

Så har vi att Fisher informationen för  $\lambda$  är

$$\begin{aligned} I(\lambda) &= -n \mathbb{E} \left[ \frac{\partial}{\partial \lambda^2} \log f_X(X|\lambda) \right], \\ &= -n \mathbb{E} \left[ \frac{\partial}{\partial \lambda^2} (X \log(\lambda) - \lambda - \log(X!)) \right], \\ &= -n \mathbb{E} \left[ -\frac{X}{\lambda^2} \right], \\ &= \frac{n \mathbb{E}[X]}{\lambda^2}, \\ &= \frac{n}{\lambda}. \end{aligned}$$

Med hjälp av Sats 4.1 har vi att

$$\text{Var}[\bar{X}] = \frac{\text{Var}[X]}{n} = \frac{\lambda}{n} = \frac{1}{I(\lambda)}.$$

Det följer att  $\bar{X}$  har den minsta möjliga variansen (och därmed kvadratfelet) över alla väntevärdesriktiga skattningar för  $\lambda$ .

**Exempel 5.9.** Låt  $X_1, \dots, X_n$  vara oberoende och  $N(\mu, \sigma^2)$ -fördelade. Fisher informationen för  $\sigma^2$  är

$$\begin{aligned} I(\sigma^2) &= -n \mathbb{E} \left[ \frac{\partial^2}{\partial (\sigma^2)^2} \log f_X(X|\mu, \sigma^2) \right], \\ &= -n \mathbb{E} \left[ \frac{1}{2\sigma^4} - \frac{(X - \mu)^2}{\sigma^6} \right], \\ &= -n \left( \frac{1}{2\sigma^4} - \frac{\mathbb{E}[(X - \mu)^2]}{\sigma^6} \right), \\ &= -n \left( \frac{1}{2\sigma^4} - \frac{\text{Var}[X]}{\sigma^6} \right), \\ &= -n \left( \frac{1}{2\sigma^4} - \frac{1}{\sigma^4} \right), \\ &= \frac{n}{2\sigma^4}. \end{aligned}$$

Enligt Sats 5.3 har en väntevärdesriktig skattning  $W(\mathbf{X})$  för  $\sigma^2$

$$\text{Var}[W(\mathbf{X})] \geq \frac{2\sigma^4}{n}.$$

Sats 4.1 visar att  $S^2$  är en väntevärdesriktig skattning för  $\sigma^2$  men vi såg i Exempel 5.7 att

$$\text{Var}[S^2] = \frac{2\sigma^4}{n-1} > \frac{2\sigma^4}{n}.$$

Därför uppnås inte  $S^2$  den Cramér-Rao under gränsen.

I situationer när vi har inte hittas en skattning som uppnår den Cramér-Rao under gränsen är det oklart om det finns en skattning som gör det. Följande satsen ger oss ett sätt för att svara denna fråga.

**Sats 5.4.** Låt  $X_1, \dots, X_n$  vara oberoende och likafördelade med sannolikhetsfunktion (eller täthetsfunktion)  $f_X(x|\theta)$  som tillhör en exponentialfamilj. En väntevärdesriktig skattning  $W(\mathbf{X})$  för  $g(\theta)$  uppnår den Cramér-Rao under gränsen om och endast om

$$\frac{\partial}{\partial \theta} \log f_{\mathbf{X}}(\mathbf{x}|\theta) = a(\theta)(W(\mathbf{X}) - g(\theta))$$

för någon funktion  $a(\theta)$  där  $f_{\mathbf{X}}(\mathbf{x}|\theta)$  är simultan sannolikhetsfunktionen (eller täthetsfunktionen) för  $X_1, \dots, X_n$ .

**Exempel 5.10.** För normalfamiljen har vi simultan täthetsfunktionen

$$f_{\mathbf{X}}(\mathbf{x}|\mu, \sigma^2) = \frac{1}{(2\pi\sigma^2)^{\frac{n}{2}}} \exp\left(-\frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2\right).$$

Det följer att

$$\begin{aligned} \frac{\partial}{\partial \sigma^2} \log f_{\mathbf{X}}(\mathbf{x}|\mu, \sigma^2) &= -\frac{n}{2\sigma^2} + \frac{1}{2\sigma^4} \sum_{i=1}^n (x_i - \mu)^2, \\ &= \frac{n}{2\sigma^4} \left( \frac{1}{n} \sum_{i=1}^n (x_i - \mu)^2 - \sigma^2 \right). \end{aligned}$$

När vi tar  $a(\sigma^2) = \frac{n}{2\sigma^4}$  får vi

$$W(\mathbf{X}) = \frac{1}{n} \sum_{i=1}^n (X_i - \mu)^2.$$

Det följer att en skattning som uppnår den Cramér-Rao under gränsen kräver att vi vet  $\mu$ . Utan denna information kan vi inte uppnå gränsen.

I sådana fall som Exempel 5.10 säger vi att gränsen är *ouppnåeligt*.