

SANNOLIKHETSLÄRA OCH STATISTIK

FÖRELÄSNING 8

Mattias Villani

**Avdelningen för Statistik och Maskininlärning
Institutionen för datavetenskap
Linköpings universitet**



ÖVERSIKT

- ▶ Punktskattning
- ▶ Samplingfördelning
- ▶ Konfidensintervall
- ▶ Konfidensintervall för populationsväntevärden
- ▶ Konfidensintervall för proportioner

PUNKTSKATTNING

- ▶ Grundproblem: sannolikhetsmodeller har **okända parametrar**, θ .
 - ▶ Ex medelinkomsten i Sverige: populationens väntevärde μ
 - ▶ Ex andelen defekta komponenter i produktionen av en produkt p .
 - ▶ Ex spamfilter: β_0 och β_1 är parametrar

$$\Pr(\text{Spam}|\text{antal\$}) = \frac{\exp(\beta_0 + \beta_1 \cdot \text{antal\$})}{1 + \exp(\beta_0 + \beta_1 \cdot \text{antal\$})}$$

- ▶ Vi vill använda (tränings)data för att bestämma värden för dessa parametrar.
- ▶ **Punktskattning**: vår **bästa gissning** utifrån data.

MOMENTMETODEN

- ▶ Ex $X_1, \dots, X_n | \mu \stackrel{iid}{\sim} \text{Poisson}(\mu)$. $E(X) = \mu$.
- ▶ Rimlig punktskattning av populationväntevärdet μ :

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

- ▶ **Moment** av ordningen $k = 1, 2, \dots$

$$\mu_k = \mathbb{E}(X^k)$$

- ▶ **Samplemoment** av ordningen k

$$m_k = \frac{1}{n} \sum_{i=1}^n X_i^k$$

- ▶ Ex: $k = 1$: $\mu_1 = \mu = \mathbb{E}X$ och $m_1 = \frac{1}{n} \sum_{i=1}^n X_i = \bar{X}$.

MOMENTMETODEN

- ▶ Momentmetoden för att skatta k modellparametrar $\theta_1, \dots, \theta_k$: Lös följande ekvationssystem m a p $\theta_1, \dots, \theta_k$:

$$\mu_1 = m_1$$

$$\mu_2 = m_2$$

$$\vdots$$

$$\mu_k = m_k$$

- ▶ Notera att μ_1, \dots, μ_k är funktioner av $\theta_1, \dots, \theta_k$. Mer korrekt:

$$\mu_1(\theta_1, \dots, \theta_k) = m_1$$

$$\mu_2(\theta_1, \dots, \theta_k) = m_2$$

$$\vdots$$

$$\mu_k(\theta_1, \dots, \theta_k) = m_k$$

MOMENTMETODEN

- ▶ Ibland mer praktiskt att jobba med centralmoment.
- ▶ **Centralmoment** av ordningen $k = 2, 3, \dots$

$$\mu'_k = \mathbb{E} (X - \mu_1)^k$$

- ▶ **Samplecentralmoment** av ordningen k

$$m'_k = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^k$$

- ▶ Notera att $\mu'_2 = \text{Var}(X)$ och

$$m'_2 = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2 \neq \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2$$

MOMENTMETODEN - BETA EXEMPEL

- Ex $X_1, \dots, X_n | \alpha, \beta \stackrel{iid}{\sim} \text{Beta}(\alpha, \beta)$:

$$\mu_1 = \mathbb{E}X = \frac{\alpha}{\alpha + \beta}$$

$$\mu_2 = \mathbb{E}(X^2) = \frac{\alpha(\alpha + 1)}{(\alpha + \beta)(\alpha + \beta + 1)}$$

- Momentskattningar: lös för α och β

$$\frac{\alpha}{\alpha + \beta} = m_1$$

$$\frac{\alpha(\alpha + 1)}{(\alpha + \beta)(\alpha + \beta + 1)} = m_2$$

$$\text{ger } \hat{\alpha} = \frac{m_1(m_2 - m_1)}{m_1^2 - m_2} \text{ och } \hat{\beta} = \frac{(m_1 - m_2)(m_1 - 1)}{m_1^2 - m_2}.$$

MAXIMUM LIKELIHOOD-METODEN

- ▶ **Maximum likelihood (ML) estimatorn:** Välj det θ som maximerar sannolikheten för data:

$$\hat{\theta} = \operatorname{argmax}_{\theta \in \Theta} P(x_1, \dots, x_n | \theta)$$

- ▶ Kontinuerliga fallet:

$$\hat{\theta} = \operatorname{argmax}_{\theta \in \Theta} f(x_1, \dots, x_n | \theta)$$

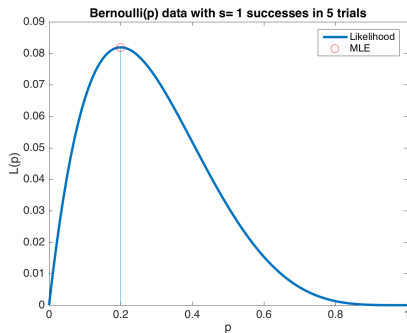
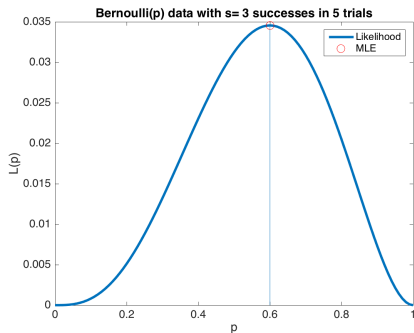
- ▶ **Likelihoodfunktionen** är sannolikheten för stickprovet sett som en funktion av parametern

$$L(\theta) = P(x_1, \dots, x_n | \theta)$$

- ▶ ML-estimatorn maximerar $L(\theta)$.
- ▶ Ex data från Bernoulli med sannolikhet p :
 $X_1 = 0, X_2 = 1, X_3 = 1, X_4 = 0, X_5 = 1$.

$$L(p) = (1 - p)pp(1 - p)p = p^3(1 - p)^2$$

MAXIMUM LIKELIHOOD - BERNOULLIEXEMPEL



MAXIMUM LIKELIHOOD-METODEN

- ▶ Vi kan hitta **ML-skattningen** analytiskt: Lös m a p θ

$$\frac{\partial L(\theta)}{\partial \theta} = 0.$$

- ▶ Oftast enklare att **maximera log-likelihoodfunktionen**

$$\frac{\partial \ln L(\theta)}{\partial \theta} = 0$$

- ▶ Ex Bernoulli:

$$\frac{\partial \ln L(p)}{\partial p} = \frac{\partial}{\partial p} (s \ln p + f \ln(1 - p)) = \frac{s}{p} + f \frac{-1}{1 - p} = \frac{s}{p} - \frac{f}{1 - p} = 0$$

vilket ger lösningen $\hat{p} = \frac{s}{s+f} = \frac{s}{n}$.

- ▶ Kontrollera \hat{p} är ett maximum - andraderivatan är negativ i $p = \hat{p}$

$$\frac{\partial^2 \ln L(\theta)}{\partial \theta^2} = -\frac{s}{p^2} - \frac{f}{(1 - p)^2} < 0$$

för alla $p \in [0, 1]$, inklusive $p = \hat{p}$.

ML-METODEN

- Notera att oberoende data är praktiskt

$$L(\theta) = \prod_{i=1}^n f(x_i|\theta)$$

så log-likelihooden blir en summa som är lättare att derivera

$$\ln L(\theta) = \sum_{i=1}^n \ln f(x_i|\theta).$$

- Ex: $X_1, \dots, X_n | \lambda \stackrel{iid}{\sim} \text{Exp}(\lambda)$ ger

$$L(\lambda) = \prod_{i=1}^n \lambda e^{-\lambda x_i} = \lambda^n e^{-\lambda \sum_{i=1}^n x_i}$$

och

$$\frac{\partial \ln L(\lambda)}{\partial \lambda} = \frac{n}{\lambda} - \sum_{i=1}^n x_i = \frac{n}{\lambda} - n\bar{x},$$

och därmed $\hat{\lambda} = 1/\bar{x}$.

SAMLINGFÖRDELNINGEN

- ▶ Hur bra är en **estimator** $\hat{\theta}$?
- ▶ **Väntevärdesriktig**? $\mathbb{E}\hat{\theta} = \theta$.
- ▶ $\text{Bias}(\hat{\theta}) = \mathbb{E}(\hat{\theta}) - \theta$.
- ▶ **Samplingfördelningen** beskriver variationen i $\hat{\theta}$ över alla stickprov av en viss storlek n .
- ▶ **Standardfelet** för $\hat{\theta}$ är $\sqrt{\text{Var}(\hat{\theta})}$, dvs standardavvikelsen för $\hat{\theta}$ över alla stickprov av storleken n .
- ▶ **Mean Squared Error** (MSE):

$$MSE(\hat{\theta}) = \mathbb{E} (\hat{\theta} - \theta)^2 = \text{Var}(\hat{\theta}) + [\text{Bias}(\hat{\theta})]^2$$

SAMLINGFÖRDELNINGEN

- ▶ Ex. Poisson. ML-estimator för μ : \bar{X} .
- ▶ Väntevärdesriktig: $\mathbb{E}(\hat{\mu}) = \mu$ och $\text{Var}(\hat{\mu}) = \frac{\sigma^2}{n} = \frac{\mu}{n}$.
- ▶ Notera att $\text{Var}(\hat{\mu}) = \frac{\mu}{n}$ beror på den okända parametern μ . Lösning: sätt $\mu = \hat{\mu} = \bar{x}$.
- ▶ Tekniker för att härleda samplingfördelningen för en estimator $\hat{\theta}$:
 - ▶ Om X_1, \dots, X_n är iid från $N(\mu, \sigma^2)$ så är $\hat{\theta} = \bar{X} \sim N(\mu, \sigma^2/n)$, exakt.
 - ▶ **Centrala gränsvärdessatsen:** $\hat{\theta} \stackrel{\text{approx}}{\sim} N[\theta, \text{Var}(\hat{\theta})]$.
 - ▶ **Bootstrapsimulering.**
- ▶ **Bootstrap:**
 - ▶ Skapa N **bootstrapstickprov** $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(N)}$ av samma storlek som det ursprungliga stickprovet genom dragning **med återläggning**.
 - ▶ Beräkna estimatet $\hat{\theta}(\mathbf{x}^{(1)}), \dots, \hat{\theta}(\mathbf{x}^{(N)})$ för var och ett av dessa N stickprov.
 - ▶ Den empiriska fördelningen för $\hat{\theta}(\mathbf{x}^{(1)}), \dots, \hat{\theta}(\mathbf{x}^{(N)})$ (tänk histogram) är en approximation av samplingfördelningen för $\hat{\theta}$.

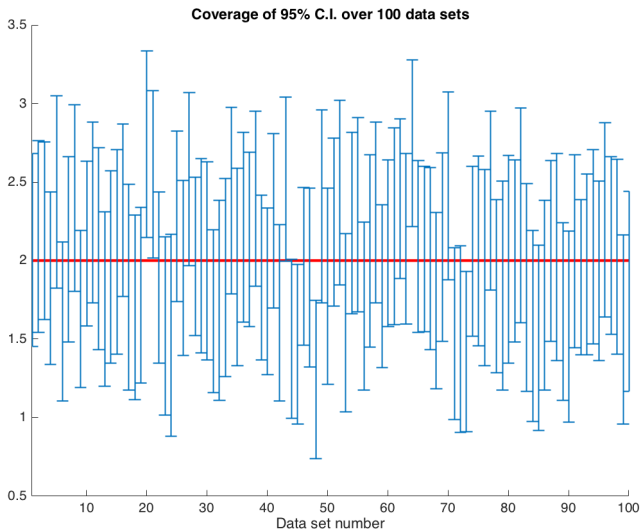
KONFIDENSINTERVALL (K.I.)

- ▶ Punktskattning ger bara en bästa gissning för θ . Konfidensintervall är ett försök att beskriva osäkerheten om θ .
- ▶ **95%-igt konfidensintervall** för θ är ett intervall $[a, b]$ sådant att

$$P\{a \leq \theta \leq b\} = 0.95.$$

- ▶ Viktigt: det är **intervallet** som är **slumpmässigt**. Parametern θ är en fix konstant.
- ▶ **Tolkning**: ett 95%-igt konfidensintervall $[a, b]$ kommer att **täcka** ($\theta \in [a, b]$) parametervärdet θ i 95% av alla möjliga stickprov.
- ▶ Man kan naturligtvis ha andra **konfidensnivåer** än 95%. 90%, 95% och 99% är vanligast. Se den lite klumpiga allmänna notationen $(1 - \alpha) \cdot 100\%$ -igt konfidensintervall i Baron.

KONFIDENSINTERVALL



KONFIDENSINTERVALL - STANDARDPROCEDUR

- ▶ Antag normalfördelad väntevärdesriktig estimator $\hat{\theta}$, t ex \bar{X} vid normalfördelade data (eller CLT argument). Då gäller

$$Z = \frac{\hat{\theta} - \theta}{\sigma(\hat{\theta})} \sim N(0, 1)$$

- ▶ Låt z_α vara $(1 - \alpha)\%$ percentilen i $N(0, 1)$ -fördelningen. Det värde som klipper av ytan α till **höger**. Tabell A4 i Baron ger att $z_{0.025} = 1.96$.
- ▶ Då gäller att

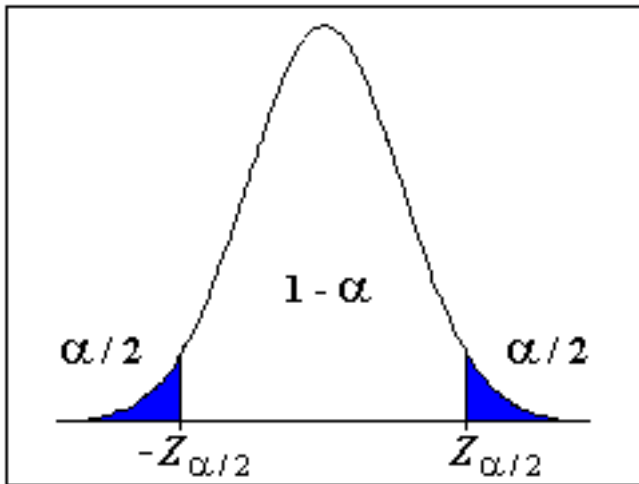
$$\mathbf{P} \left\{ -z_{0.025} \leq \frac{\hat{\theta} - \theta}{\sigma(\hat{\theta})} \leq z_{0.025} \right\} = 0.95$$

vilket kan skrivas om som

$$\mathbf{P} \{ \hat{\theta} - z_{0.025} \cdot \sigma(\hat{\theta}) \leq \theta \leq \hat{\theta} + z_{0.025} \cdot \sigma(\hat{\theta}) \} = 0.95$$

- ▶ Alltså: $[\hat{\theta} - z_{\alpha/2} \cdot \sigma(\hat{\theta}), \hat{\theta} + z_{\alpha/2} \cdot \sigma(\hat{\theta})]$ är ett $(1 - \alpha)\%$ -igt konfidensintervall för θ .

KOPPLINGEN MELLAN $(1 - \alpha)100\%$ -IGA INTERVALL OCH TABELLTALEN $z_{\alpha/2}$



KONFIDENSINTERVALL FÖR POPULATIONSVÄNTEVÄRDET

- ▶ $\theta = \mu$. $\hat{\theta} = \bar{X}$. $\sigma(\hat{\theta}) = \text{Std}(\bar{X}) = \sigma/\sqrt{n}$. σ antas känd.
- ▶ Centrala gränsvärdessatsen ger att $\hat{\theta} = \bar{X}$ är approximativt normalfördelad när n är stort ($n \geq 30$). Oavsett hur data är fördelade. Om data är normalfördelade är intervallet exakt.
- ▶ Alltså: $\bar{X} \pm z_{0.025} \cdot \frac{\sigma}{\sqrt{n}}$ är ett (approximativt) 95%-igt konfidensintervall för θ .
- ▶ **Bestämning av stickprovsstorlek n .** Vi kan bestämma n så att vi får ett konfidensintervall av given bredd.

KONFIDENSINTERVALL - OKÄNT STANDARDFEL

- ▶ I praktiken är $\sigma(\hat{\theta})$ inte känd utan måste skattas (estimeras) från data. Ex: $Std(\bar{X}) = \sigma/\sqrt{n}$ och σ är ofta okänd.
- ▶ Vid stort n får vi ett bra approximativt konfidsensintervall genom att ersätta $\sigma(\hat{\theta})$ med en skattning. T ex s/\sqrt{n} istället för σ/\sqrt{n} .
- ▶ Konfidsensintervall för populationsväntevärdet μ vid **små stickprov** en **normalfördelad population**:

$$t = \frac{\bar{X} - \mu}{s/\sqrt{n}} \sim t_{n-1}(0, 1)$$

- ▶ Så ett exakt 95%-igt konfidsensintervall för μ ges då av

$$\bar{X} \pm t_{0.025}(n-1) \frac{s}{\sqrt{n}}$$

där $t_{0.025}(n-1)$ är 97.5% percentilen i **t -fördelningen** med $\nu = n-1$ frihetsgrader. Läses av från Tabell A5 i Baron.

KONFIDENSINTERVALL FÖR EN ANDEL

- ▶ Ex. 196 av 2000 utfrågade svarar att de röstar på centerpartiet. Hur stor andel p röstar på centerpartiet i hela populationen?
- ▶ $\hat{p} = 0.098$ är ML-skattningen. Men hur säkra är vi?
- ▶ $\hat{p} = \frac{1}{n} \sum_{i=1}^n X_i$ där $X_i = 1$ om den i :te utfrågade person röstar på centerpartiet och $X_i = 0$ annars. Så \hat{p} är också ett medelvärde!
- ▶ Antag att $X_i \stackrel{iid}{\sim} \text{Bernoulli}(p)$. Då gäller $\mathbb{E}X_i = p$ och $\text{Var}(X_i) = p(1 - p)$. Alltså

$$\mathbb{E}\hat{p} = p \quad \text{Var}(\hat{p}) = \frac{1}{n^2} \sum_{i=1}^n \text{Var}(X_i) = \frac{1}{n^2} np(1 - p) = \frac{p(1 - p)}{n}.$$

- ▶ $\sigma(\hat{p})$ beror på p , som vi ersätter med en skattning:
 $s(\hat{p}) = \sqrt{\hat{p}(1 - \hat{p})/n}$.
- ▶ Centrala gränsvärdessatsen ger ett approximativt $(1 - \alpha)100\%$ -igt intervall

$$\hat{p} \pm z_{\alpha/2} \sqrt{\frac{\hat{p}(1 - \hat{p})}{n}} = (0.085, 0.111)$$