

# Föreläsning 8, Matematisk statistik 7.5 hp för E

## Punktskattningar

Stas Volkov

# Matematisk statistik – slumpens matematik

**Sannolikhetsteori:** Hur beskriver man slumpen och slumpmässiga händelser?

- ▶ Slh. för **3** st **1**:or på **10** tärningslag?
- ▶ Givet fördelningen för vågor, hur höga/stora kan de **5 %** ”värsta” vågorna vara?
- ▶ Vi observerar ett radioaktivt material med känd halveringstid under **10** minuter; vilken fördelning kommer det observerade antalet sönderfall att följa?

**Statistikteori:** Vilka slutsatser kan man dra av ett datamaterial?

- ▶ Givet **3** st **1**:or på **10** tärningslag, är tärningen ”rättvis”?
- ▶ Givet **10** års mätningar av vågor, vad kan vi säga om fördelningen?
- ▶ Under **10** minuter observerar vi **5** sönderfall, vad är halveringstiden?

# Statistik

Från mätningar (insamlad data) dra slutsatser om verkligheten.

Vi behöver då en modell för våra mätningar!

Ofta innehåller vår modell okända parametrar samt ett antagande om fördelning för observationerna.

## Exempel: Kvalitetskontroll

Vi kontrollerar  $n$  st slumpmässigt utvalda komponenter från ett stort parti och ser om de fungerar.

Modell:  $X$  = antalet trasiga komponenter

$X \in \text{Bin}(n, p)$ , där  $p$  är andelen trasiga komponenter. Parametern  $p$  är okänd.

Möjliga frågeställningar:

1. Vad är en bra uppskattning av  $p$ ?
2. Hur stor är osäkerheten i uppskattningen?
3. Vilket intervall tror vi  $p$  ligger inom?
4. Hur stort måste  $n$  vara för att uppnå en ”tillräckligt liten” osäkerhet?

# Statistikteori — översikt

## Punktskattning

Hur gör man en bra gissning av en okänd storhet? Hur vet man att den är bra?

## Intervallskattning

Hitta istället ett intervall som täcker den okända storheten med en given (stor) sannolikhet.

## Hypotestest

Om gissningen blev **0.013**, kan rätt värde på den okända storheten ändå vara **0.01**?

## Regression

Sambandsanalys:

hur vet vi om två variabler påverkar varandra?

# Statistikteori, grundläggande begrepp

## Stickprov

Ett **stickprov**,  $x_1, x_2, \dots, x_n$ , är **observationer** av s.v.  $X_1, \dots, X_n$  från någon fördelning  $X_i \in F(\vartheta)$  där  $\vartheta$  är en okänd **parameter**.

## Skattning

En **skattning** av  $\vartheta$ ,  $\vartheta^*(x_1, \dots, x_n)$  är en observation av den s.v.  $\vartheta^*(X_1, \dots, X_n)$ . Båda betecknas oftast bara med  $\vartheta^*$ .

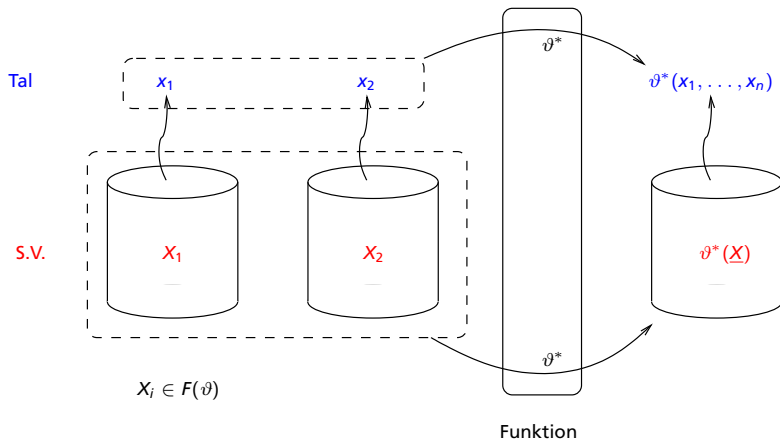
Bra egenskaper för en skattning är

Väntevärdesriktig :  $E(\vartheta^*) = \vartheta$ , inget systematiskt fel.

Effektiv: liten varians (osäkerhet)  $V(\vartheta^*)$ .

Konsistent:  $P(|\vartheta_n^* - \vartheta| > \varepsilon) \rightarrow 0$ ,  $n \rightarrow \infty$ , dvs "Blir bättre när vi får fler observationer",

# En skattning $\vartheta^*$ är både ett tal, en s.v. och en funktion



# Modell för mätning med slumpmässigt mätfel

Antag att vi vill mäta en storhet  $\mu$ . Om man gör  $n$  st mätvärden,  $x_1, \dots, x_n$  är dessa observationer av

$$X_j = \mu + \varepsilon_j = \text{"Rätt värde"} + \text{"Mätfel"}$$

där  $\varepsilon_j$  är ett slumpmässigt mätfel.  
Ofta antas att  $\varepsilon_j$  är oberoende och

$$\varepsilon_j \in N(\mathbf{0}, \sigma)$$

Detta ger att våra observationer blir

$$X_j \in N(\mu, \sigma)$$

Därför ser vi att väntevärdet är den storhet vi försöker mäta upp.



## Kom ihåg: väntevärde och varians

Väntevärdet anger **tyngdpunkten** för fördelningen

$$E(X) = \begin{cases} \int_{-\infty}^{\infty} x \cdot f_X(x) dx & \text{kontinuerlig} \\ \sum_k k \cdot p_X(k) & \text{diskret} \end{cases}$$

Variansen anger hur utspridd  $X$  är kring sitt väntevärde.

$$V(X) = E\left([X - E(X)]^2\right) = E(X^2) - E(X)^2 \geq 0.$$

Nuttiga formler:

$$E\left(\sum a_i X_i + b\right) = \sum a_i E(X_i) + b$$

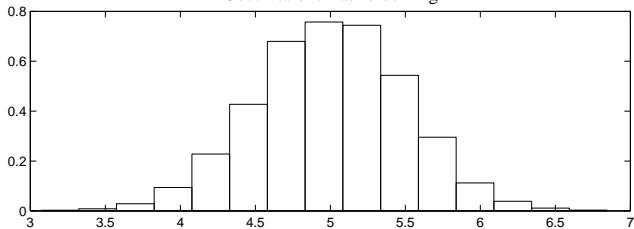
$$V\left(\sum_i a_i X_i + b\right) = \sum_i a_i^2 V(X_i) + 2 \underbrace{\sum_{i < j} a_i a_j C(X_i, X_j)}_{=0 \text{ om okorrelerade}}$$

# Variation i observationer ger variation i skattningen

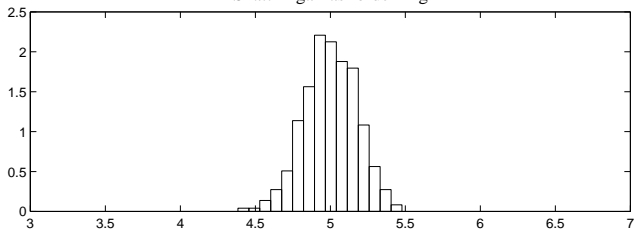
$$\mu_n^* = \frac{1}{n} \sum_{i=1}^n X_i \quad E(\mu_n^*) = \mu \quad V(\mu_n^*) = \frac{\sigma^2}{n}$$

Stickprov	Observationer $x_{jk}$								$\mu^* = \bar{x}_j$
1	4.83	4.93	5.24	5.12	5.10	4.69	5.62	4.73	5.03
2	5.09	5.13	4.53	4.59	4.70	4.10	4.96	5.26	4.79
3	5.53	5.10	4.34	5.05	5.21	4.43	4.30	4.56	4.82
4	4.48	5.10	4.75	5.17	4.98	5.01	5.82	5.12	5.05
5	5.14	5.10	4.79	5.48	4.70	5.89	5.22	5.91	5.28
6	4.80	5.33	5.22	5.26	4.45	4.12	5.29	5.09	4.95
7	5.20	5.26	5.49	5.60	4.83	5.28	4.38	5.18	5.15
8	4.48	4.81	4.62	4.61	5.04	4.81	4.32	4.41	4.64
⋮									

Observationernas fördelning



Skattningarnas fördelning



## Minsta kvadrat-metoden, MK

Om  $E(X_i) = \mu_i(\vartheta)$  så fås **MK-skattningen** av  $\vartheta$  genom att **minimera förlustfunktionen**

$$Q(\vartheta) = \sum_{i=1}^n (x_i - \mu_i(\vartheta))^2$$

med avseende på  $\vartheta$ .

- ▶ Bestäm **hur** väntevärdet beror av  $\vartheta$ ,  $E(X_i) = \mu_i(\vartheta)$ .
- ▶ Sätt upp  $Q(\vartheta)$
- ▶ Derivera  $Q(\vartheta)$ , sätt lika med noll och lös m.a.p.  $\vartheta$ .
- ▶ Det  $\vartheta$  som minimerar  $Q(\vartheta)$  är MK-skattningen,  $\vartheta_{MK}^*$ .

## Maximum likelihood-metoden, ML

**ML-skattningen** av  $\vartheta$  fås genom att **maximera likelihood-funktionen**  $L(\vartheta; \mathbf{x}_1, \dots, \mathbf{x}_n)$  m.a.p.  $\vartheta$ .

$$L(\vartheta) = p_X(\mathbf{x}_1) \cdot \dots \cdot p_X(\mathbf{x}_n) \quad (\text{diskret})$$

$$L(\vartheta) = f_X(\mathbf{x}_1) \cdot \dots \cdot f_X(\mathbf{x}_n) \quad (\text{kontinuerlig})$$

I det diskreta fallet anger L-funktionen:

”Sannolikheten att få det stickprov som vi fått”.

- ▶ Sätt upp  $L(\vartheta)$
- ▶ Logaritmera —  $\ln L(\vartheta)$  maximeras av samma  $\vartheta$  som  $L(\vartheta)$ .
- ▶ Derivera  $\ln L(\vartheta)$  med avs. på  $\vartheta$ , sätt lika med noll och lös för  $\vartheta$ .
- ▶ Det  $\vartheta$  som maximerar  $L(\vartheta)$  är ML-skattningen  $\vartheta_{ML}^*$ .

## Exempel: Radon

Radonkoncentrationen i inomhusluft kan mätas genom att hänga upp en  $\alpha$ -känslig film. Antalet hål i filmen beskrivs av en Poisson-process med

$$X_j \in \text{Po}(\mu K_j)$$

där  $\mu$  är den okända radonkoncentrationen och  $K_j$  är kända konstanter som beror på bl.a. filmens känslighet, storlek och exponeringstiden.

Radon-data återkommer i lab 4.



## Ex: Normalfördelning

Om  $x_1, \dots, x_n$  är observationer av  $X_j \in \mathbf{N}(\mu, \sigma)$  blir ML- och MK-skattningen av  $\mu$  och en **korrigerad** ML-skattning av  $\sigma^2$

$$\mu^* = \bar{x}$$

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

Dessa används även för att skatta väntevärde och varians vid okänd fördelning.

Värför korrigerad?

$$E(s^2(X)) = \sigma^2$$

# Medelfel

Om  $D(\vartheta^*(X)) = \sqrt{V(\vartheta^*(X))}$ , dvs. standardavvikelsen av en skattning för  $\vartheta^*$ , , innehåller okända parametrar  $\vartheta$  då kan man inte räkna ut ett numeriskt värde på den. Istället **stoppar man in skattningar på de okända parametrarna** och får **medelfelet**  $d(\vartheta^*)$ .

**Exempel:**  $p^* = \frac{X}{n}$  där  $X \in \text{Bin}(n, p)$  så  $V(X) = np(1 - p)$ .

$$V(p^*) = V\left(\frac{X}{n}\right) = \frac{1}{n^2}V(X) = \frac{1}{n^2}np(1 - p) = \frac{p(1 - p)}{n},$$

$$\text{så } D(p^*) = \sqrt{\frac{p(1 - p)}{n}} \quad \Rightarrow \quad d(p^*) = \sqrt{\frac{p^*(1 - p^*)}{n}}$$

**Exempel.**  $\mu^* = \bar{X}$ , där  $X \in N(\mu, \sigma)$ ,  $\sigma$  okänd

$$V(\mu^*) = \frac{\sigma^2}{n}, \quad d(\mu^*) = \frac{s}{\sqrt{n}}, \quad \text{där } s = \sqrt{\frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2}$$



## Exempel: två binomialfördelningar

Vi har två oberoende observationer:  $x_1$  från  $X_1 \in \text{Bin}(n_1, p)$  och  $x_2$  från  $X_2 \in \text{Bin}(n_2, 2p)$  där  $n_1$  och  $n_2$  är kända medan  $p$  är en okänd parameter (OBS!  $0 < p < 1/2$ ).

- ▶ Bestäm MK-skattningen av  $p$ .
- ▶ Bestäm ML-skattningen av  $p$ .
- ▶ Beräkna skattningarnas värde när  $n_1 = 5$  och  $x_1 = 2$ , samt  $n_2 = 6$  och  $x_2 = 3$ .
- ▶ Är skattningarna väntesvärderiktiga?
- ▶ Vilken av skattningarna har lägst varians?

## Exempel: två binomialfördelningar

- ▶ Enkelt: skatta  $p$  och  $2p$  var för sig:

$$p_1^*(x_1) = p^* = \frac{x_1}{n_1} = \frac{2}{5} = 0.40$$

$$2p_2^*(x_2) = (2p)^* = \frac{x_2}{n_2} = \frac{3}{6} = 0.50 \quad \Rightarrow \quad p^* = \frac{x_2}{2n_2} = 0.25$$

- ▶ Vilken är rätt? (Ingen av dem...)
- ▶ Vilken är bäst? ( $p_1^*$  kan bli  $> 1/2$  som är omöjligt ☹☹☹)

$$V(p_1^*) = \frac{p(1-p)}{n_1} = \frac{p(1-p)}{5};$$

$$V(p_2^*) = \frac{2p(1-2p)}{2^2 n_1} = \frac{p(1-2p)}{12}.$$

## Vikta ihop skattningarna, på olika sätt...

$$p_a^* = p_a^*(x_1, x_2) = \frac{p_1^* + p_2^*}{2} = \frac{\frac{x_1}{n_1} + \frac{x_2}{2n_2}}{2} = \frac{2n_2x_1 + n_1x_2}{4n_1n_2} = 0.325$$

$$p_b^* = p_b^*(x_1, x_2) = \frac{n_1p_1^* + n_2p_2^*}{n_1 + n_2} = \frac{x_1 + \frac{x_2}{2}}{n_1 + n_2} = \frac{2x_1 + x_2}{2(n_1 + n_2)} = 0.318$$

$$p_c^* = p_c^*(x_1, x_2) = \frac{p_1^* + 2p_2^*}{1 + 2} = \frac{\frac{x_1}{n_1} + \frac{x_2}{n_2}}{3} = \frac{n_2x_1 + n_1x_2}{3n_1n_2} = 0.300$$

$$p_d^* = p_d^*(x_1, x_2) = \frac{n_1p_1^* + 2n_2p_2^*}{n_1 + 2n_2} = \frac{x_1 + x_2}{n_1 + 2n_2} = 0.294$$

## MK- och ML-skattningar

$$p_{MK}^* = \frac{n_1 x_2 + 2n_2 x_1}{n_1^2 + 4n_2^2} = 0.272$$

$$p_{ML}^* = \frac{n_1 + 2n_2 + 2x_1 + x_2}{4(n_1 + n_2)} - \frac{\sqrt{(n_1 + 2n_2)^2 + (2x_1 + x_2)^2 - (4n_1 x_1 + 6n_1 x_2 + 4n_2 x_2)}}{4(n_1 + n_2)} = 0.280$$