

Tällä kurssilla lähinnä aiheet 1. - 3.

Käytännön sovelluksissa tehtävä 5. hyvin tärkeä !

Esim. "Pallot kulkossa" -esimerkissä edellä tilast. malli  
sm:lle  $T$  on

$$T \sim \text{Bin}(n, \theta)$$

( $\Rightarrow$ )

$$f(t; \theta) = P(T=t) = \binom{n}{t} \theta^t (1-\theta)^{n-t}, \quad t=0, 1, \dots, n$$

Parametri on  $\theta$ ,  $0 \leq \theta \leq 1$  (keltaisten suht. osuus)  
( $n$  = nostojen lkm ajatellaan tunnetuksi luvuksi)

Päätelyn kaksi tärkeää paradigmaa/koulukuntaa:

Frekventistinen päätely (kurssin alkuosassa)

- Aineisto  $\underline{x}$  on satunnaisvektorin  $\underline{X}$  toteutunut arvo.
- Satunnaisuus viittaa "torstetun aineistonkeruun" ideaan: frekventistinen  $n:n$  tulkinna
- $\theta$  on kiinteä mutta tuntematon luku tai prste
- $\theta$ :lla ei ole todennäköisyysjakaumaa !

Bayesläinen päätely (kurssin loppupuolella)

- Myös parametri tulkitaan satunnaismuuttujaksi
- Todennäköisyys kuvaa siihen liittyvää epävarmuutta, subjektiivinen  $n:n$  tulkinna
- Tyylikäs tapa yhdistää parametria koskeva ennakkotieto ja aineiston antama lisäinformaatio.
- Perustuu Bayesin kaavan käyttöön

### 3. TILASTOLLISEN MALLIN MUODOSTAMISESTA

Yleistä:

- Ei helppo tehtävä (ainan perusesimerkkejä lukuunottamatta)!
- Edellyttää ilmiöön liittyvän taustateorian tunteusta.
- Usein jatkuva ja iteratiivinen prosessi (vrt. (5) sivulla 5)
- Meillä mallit joko helposti muodostettavina tai "valmiiksi annettuja".

Mursta: Malli  $f(\underline{x}; \theta)$  (tai selyyden vuoksi  $f_{\underline{x}}(\underline{x}; \theta)$ ) on

\* diskreetissä tapauksessa (yhters) pistetnf

$$f(\underline{x}; \theta) = P_{\theta}(\underline{X} = \underline{x}) = P_{\theta}(X_1 = x_1, \dots, X_n = x_n)$$

\* jatkuvassa tapauksessa (yhters) tiheysfunktio, ts.

$$P_{\theta}(\underline{X} \in A) = \int_A f(\underline{x}; \theta) d\underline{x}, \quad \text{kun } A \subset \mathbb{R}^n,$$

(n-ulott. integraali)       $\underline{x} = (x_1, \dots, x_n)$

(Tarkempi käsittely aineopintojen tu-laskennan ja päätelyn kursseilla.)

Erikoistapaus: Riippumattomat samoin jakautuneet havainnot:

$$\begin{cases} X_1, \dots, X_n \perp \\ X_i \text{ :lla } \text{ptnf/} \text{tf } g(\cdot; \theta) \quad (\text{sama jakaisella } i) \end{cases}$$

Tällöin tu-laskennan perusteella

$$f(\underline{x}; \theta) = g(x_1; \theta) \cdots g(x_n; \theta)$$

$$= \prod_{i=1}^n g(x_i; \theta), \quad \underline{x} = (x_1, \dots, x_n)$$

Huom. Riippumattomuusoletus ei aina toteudu!

Esimerkkinä aikasarja-tyyppiset mallit, joissa havainnot  $X_1, \dots, X_n$  ovat saman muuttujan arvoja peräkkäisinä ajanhetkinä. Ajattele esim.

$X_i$  = lämpötila Karsaniemessä vuoden  $i$ :ntenä päivänä klo 12.00

Esim. Riippumaton otos normaalijakaumasta  $N(\mu, \sigma^2)$ :

$$X_1, \dots, X_n \sim N(\mu, \sigma^2) \quad \parallel$$

"kaikkien tilast. mallien äiti"

Parametri 2-ulotteinen  $(\mu, \sigma^2)$ ,  $\mu \in \mathbb{R}$ ,  $\sigma^2 > 0$ .  
(Joskus  $\sigma^2$  tunnettu, jolloin parametrina vain  $\mu$ .)

Pal. mieleen tn-laskennasta:  $N(\mu, \sigma^2)$ -jakauman tf

$$g(x; \mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left\{-\frac{1}{2\sigma^2}(x-\mu)^2\right\}$$

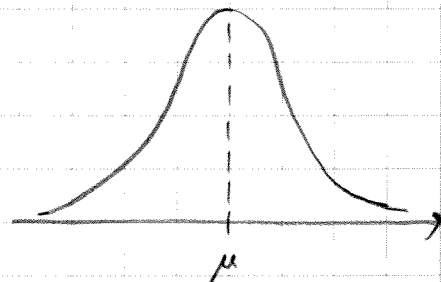
Siten tilast. mallin lauseke (yhteistf) on

$$\begin{aligned} f(\underline{x}; \mu, \sigma^2) &= \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left\{-\frac{1}{2\sigma^2}(x_i-\mu)^2\right\} \\ &= \frac{1}{(2\pi\sigma^2)^{n/2}} \exp\left\{-\frac{1}{2\sigma^2} \sum_{i=1}^n (x_i-\mu)^2\right\} \end{aligned}$$

Palaamme tähän myöhemmin.

Mursta: Jos  $X \sim N(\mu, \sigma^2)$ , niin

$\mu = E(X)$  odotusarvo ja  
 $\sigma^2 = \text{Var}(X) = D^2(X)$  varianssi



Muita paljon käytettyjä jakaumia: (tn-laskennan kurssi!)

\*  $X \sim \text{Poisson}(\lambda)$  Poisson-jakauma parametrina  $\lambda > 0$ ,  
ptstetnf

$$g(x; \lambda) = P(X=x) = e^{-\lambda} \frac{\lambda^x}{x!}, \quad x=0, 1, 2, \dots$$

\*  $X \sim \text{Exp}(\lambda)$  eksponenttijakauma parametrina  $\lambda > 0$ ,  
tiheysf

$$g(x; \lambda) = \lambda e^{-\lambda x}, \quad x > 0$$

\* tasajakauma välillä  $[\alpha, \beta]$ ,  $Tas(\alpha, \beta)$

#### 4. USKOTTAVUUSFUNKTIO JA SUURIMMAN USKOTTAVUUDEN ESTIMAATTI

vrt.  
Arjas-Sirén  
jakso 1.2.

Motivaatio: Tark. mallia  $T \sim \text{Bin}(n, \theta)$ , parametri  $0 \leq \theta \leq 1$ .

esim. "nostetaan  $n$  palloa kulhosta palauttaen",  $\theta = \begin{cases} \text{keltaisten} \\ \text{osuus} \end{cases}$   
 $T =$  keltaisten lkm otoksessa

tai yleisemmin:  $n$ -kertainen riippumaton torstokoe,  
jossa kunkin torston "onnistuminen"  $= \theta$   
 $T =$  "onnistumisten" lkm

(A) olet. että  $n=10$  ja havaittu  $T=6$

Ko. havainnon todennäköisyys on (ks. s. 6)

$$P(T=6) = f(6; \theta) = \binom{10}{6} \theta^6 (1-\theta)^4$$

Tark. tätä  $\theta$ :n funktiona: (ks. kuvaaja sivulla 11)

$$L(\theta) = 210 \cdot \theta^6 (1-\theta)^4, \quad 0 \leq \theta \leq 1$$

Siiis:  $L(\theta) = P_{\theta}(T=6) = t_n$  saada havainto " $T=6$ "  
silloin kun parametrilla arvo  $\theta$

suurimmillaan  $L$  on pist.  $\theta = 0.6$ :  $L(0.6) \approx 0.25$

sanomme:  $\theta = 0.6$  on (havainnon  $T=6$  valossa)  
uskottavin parametrinarvo, koska se maksimoi  
ko. havainnon saamisen  $t_n = n$ !

torsaalta esim.  $L(0.4) \approx 0.1$ , joten  
mihäli parametrilla on arvo  $\theta = 0.4$ , on havainnon  
"T=6" saaminen  $\sim 2.5$  kertaa epätodennäköisempää  
kuin siinä tapauksessa että  $\theta = 0.6$

sanomme:  $\theta = 0.6$  on  $\sim 2.5$  kertaa uskottavampi  
parametrin arvo kuin  $\theta = 0.4$

jne...

Huom. L:n kuvaaja on varsin "loakea":

"melko uskottavia" parametrin arvoja on paljon

$\Rightarrow$  tarkkoja ja luotettavia päätelmiä  $\theta$ :sta ei  
mahdollista tehdä! Ei ihme, sillä otoskoko  $n=10$   
on hyvin pieni!

(B) olet.  $n=300$  ja havaittu  $T=159$  (luennolla  
tehty koe!)

Menetellään kuten edellä: saadun havainnon tn on nyt

$$L(\theta) = P_{\theta}(T=159) = \binom{300}{159} \theta^{159} (1-\theta)^{141}$$

(ks. kuvaaja sivulla 11)

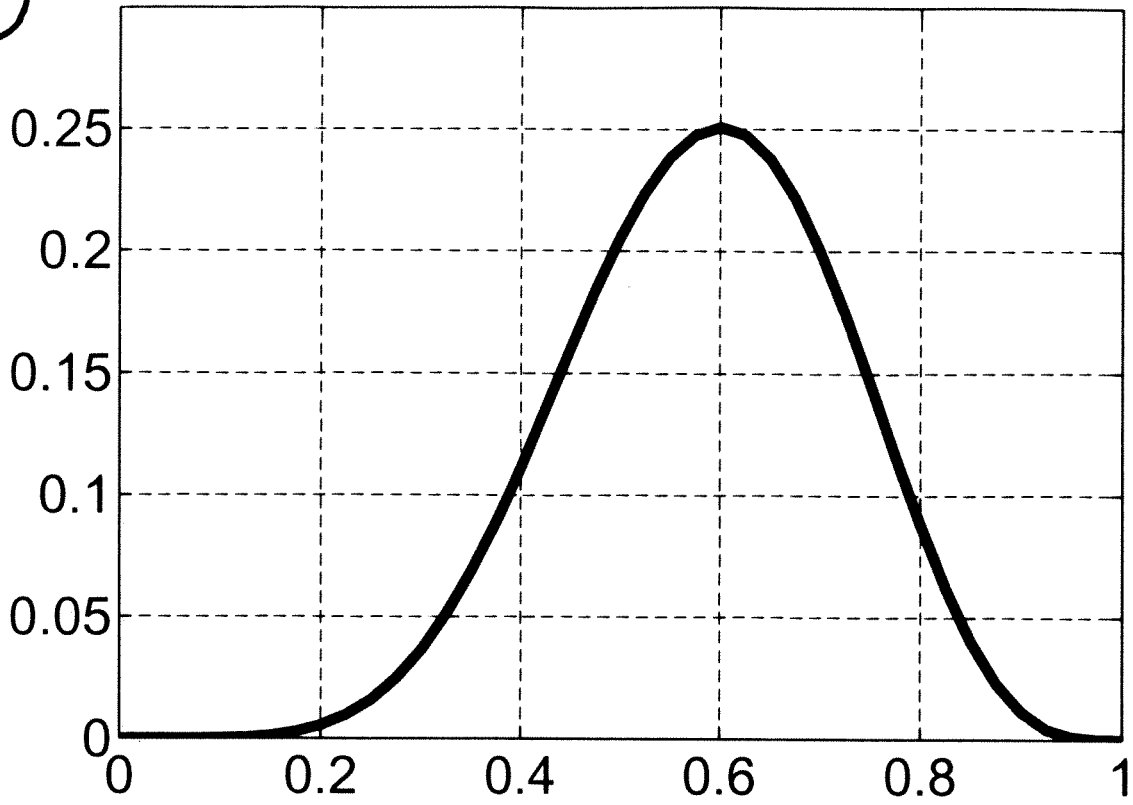
Nyt L:n globaali maksimikohta on  $\theta = \frac{159}{300} = 0.53$ ,

siis tämä on uskottavin parametrin arvo.

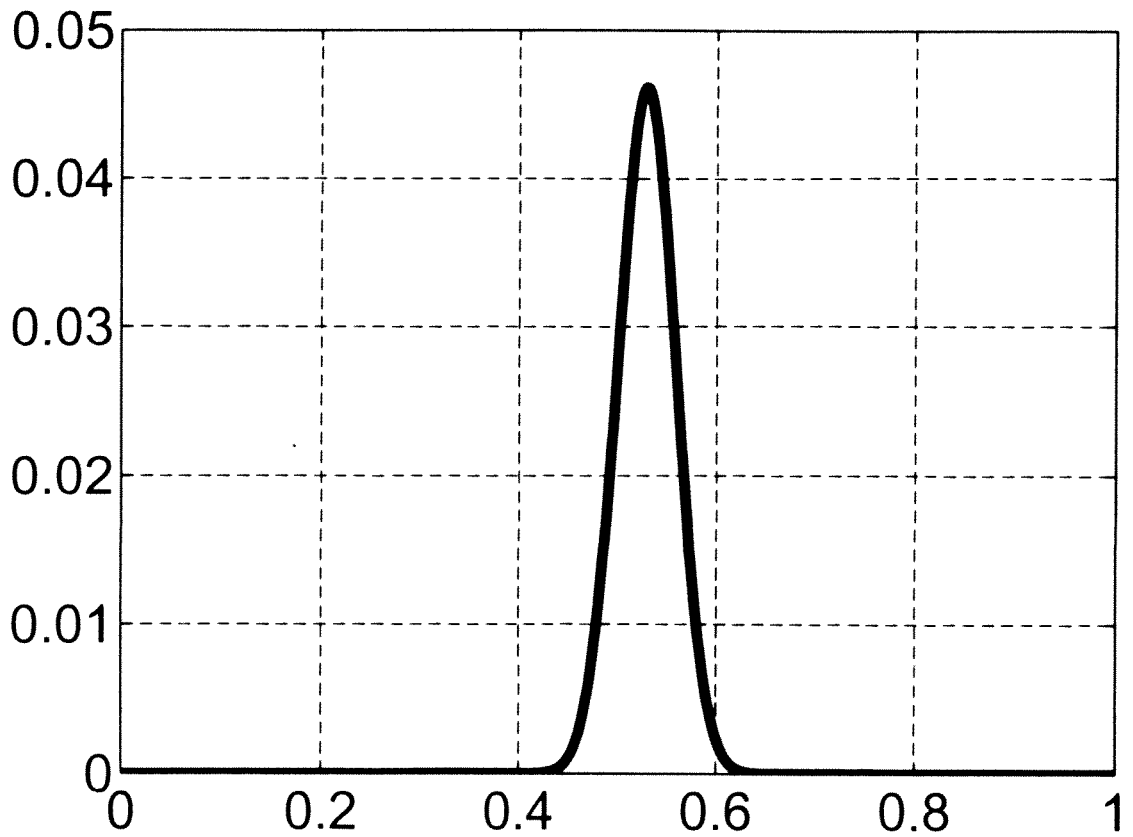
[Tarkistamme tämän kohta differentiaalilaskennan avulla!]

Huom. L:n kuvaaja "huipukkaampi" kuin (A)-tilanteesta  
ts. "uskottavat" parametrin arvot ovat melko kapealla  
välillä 0.53:n ympärillä  $\Rightarrow \theta$ :sta voi tehdä  
tarkempia päätelmiä kuin (A):ssa (suuremmasta otos-  
koosta johtuen).

(A)



(B)



Määritelmiä. Olkoon  $f(\underline{x}; \theta)$  tilastollinen malli, jonka parametriavaruus (s.o. parametrin  $\theta$  kaikkien mahdollisten arvojen joukko) on  $\Omega$ .

\* Aineistoon  $\underline{x} = (x_1, \dots, x_n)$  liittyvä uskottavuusfunktio on

$$L(\theta) = L(\theta; \underline{x}) = f(\underline{x}; \theta), \quad \theta \in \Omega$$

\* Jos  $\theta \in \Omega$  ja  $\theta' \in \Omega$  siten että  $L(\theta; \underline{x}) > L(\theta'; \underline{x})$ , sanomme, että  $\theta$  on (aineiston  $\underline{x}$  valossa) uskottavampi parametrinarvo kuin  $\theta'$ .

\* Sellainen piste  $\hat{\theta} = \hat{\theta}(\underline{x}) \in \Omega$ , jossa  $L$  saavuttaa suurimman arvonsa, ts. jossa

$$L(\hat{\theta}; \underline{x}) \geq L(\theta; \underline{x}) \quad \forall \theta \in \Omega,$$

on parametrin  $\theta$  suurimman uskottavuuden estimaatti (lyh. su-estimaatti).

Tulkinta: Su-estimaatti on sellainen parametrinarvo, jonka vallitessa "käsilläolevan" havainnon  $\underline{x}$  saamisen todennäköisyys (tai jatkuvan jak. tapauksessa "tn-tiheys") on suurin.

Huom.  $L$  ei ole tn-jakauma ( $\theta$  ei ole sat.muuttuja)!

Esim. Palataan torstokoemalliin  $T \sim \text{Bin}(n, \theta)$  (parametri  $\theta$ ) eli

$$f(t; \theta) = P_{\theta}(T=t) = \binom{n}{t} \theta^t (1-\theta)^{n-t}, \quad t=0, 1, \dots, n$$

Havaintoa  $t$  vastaava uskottavuusfunktio on

$$L(\theta) = L(\theta; t) = \binom{n}{t} \theta^t (1-\theta)^{n-t}, \quad 0 \leq \theta \leq 1.$$

Su-estimaatin laskemiseksi tutkitaan tämän logaritmsa: