

3. TILASTOLLISEN MALLIN MUODOSTAMISESTA

Yleistä:

- Ei helppo tehtävä (ainoa perusesimerkkejä lukuunottamatta)!
- Edellyttää ilmiöön liittyvän taustateorian tuntemusta.
- Usein jätuva ja iteratiivinen prosessi (vrt. (5) sivulla 5)
- Meillä mallit joko helposti muodostettavia tai "valmiiksi annettuja".

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

* diskreetissä tapauksessa (yhters) pistetnf

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

* jatluvassa tapauksessa (yhters) tiheysfunktio, ts.

$$P_{\theta}(\underline{X} \in A) = \int_A f(\underline{x}; \theta) d\underline{x}, \quad \text{kun } A \subset \mathbb{R}^n, \\ \text{(n-ulott. integraali)} \quad \underline{x} = (x_1, \dots, x_n)$$

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

Erikoistapaus: Riippumattomat samoin jakautuneet havainnot:

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

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 Kaisaniemessä 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 (μ, σ^2) , $\mu \in \mathbb{R}$, $\sigma^2 > 0$.
(Joskus σ^2 tunnettu, jolloin parametrina vain μ .)

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

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

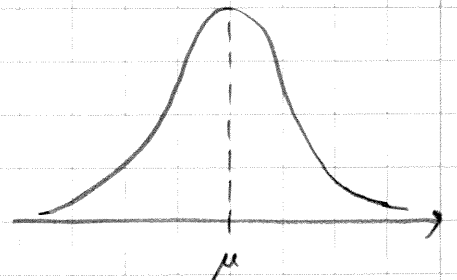
Siten tilast. mallin lauseke (yhdistettf) on

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

Palautamme 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-Sivén
jakso 1.2.

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

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

tai yleisemmin: n -kertainen riippumaton toistokoe,
jossa kunkin toiston "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 θ :n funktiona: (ks. kuvaaja sivulla 11)

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

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

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

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

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

sanomme: $\theta = 0.6$ on ~ 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 θ :sta ei
mahdollista tehdä! Ei ihme, sillä otoskoko $n=10$
on hyvin pieni!

(B) olet. $n=300$ ja havaittu $T=159$ (v. 2010
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)

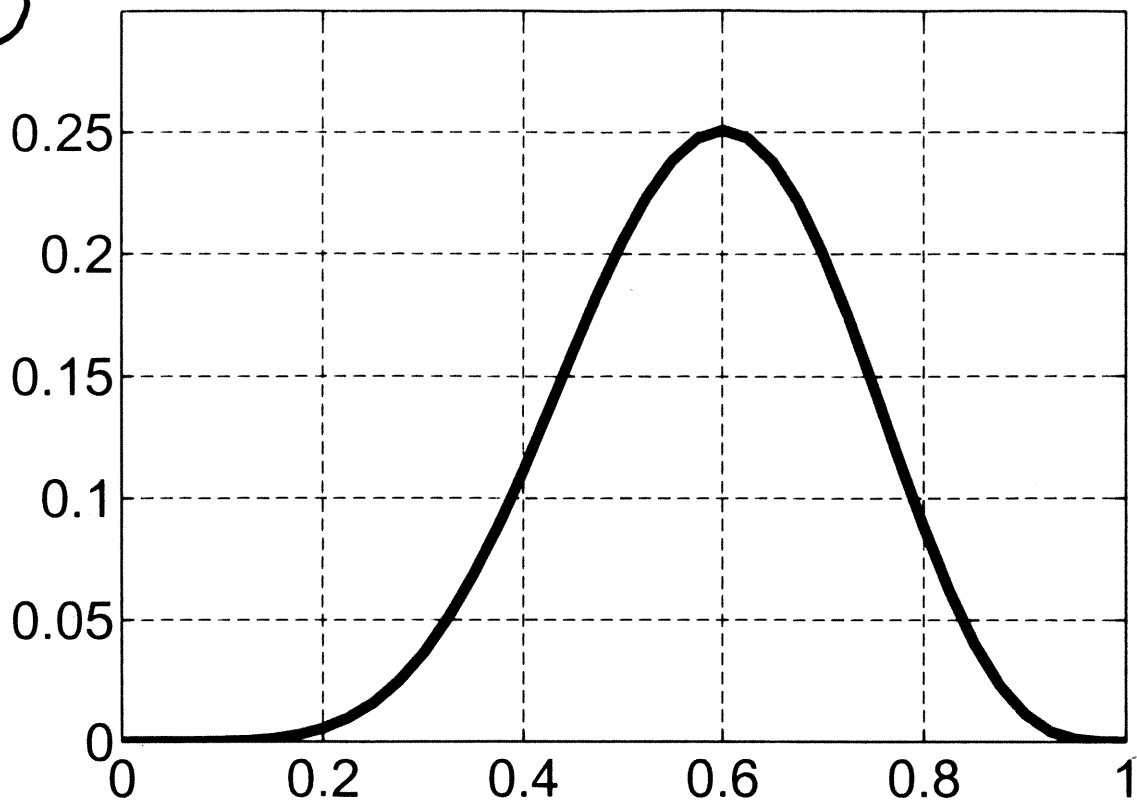
Nyt L:n globaali maksimikohta on $\theta = \frac{159}{300} = 0.53$,

sitten tämä on uskottavin parametrin arvo.

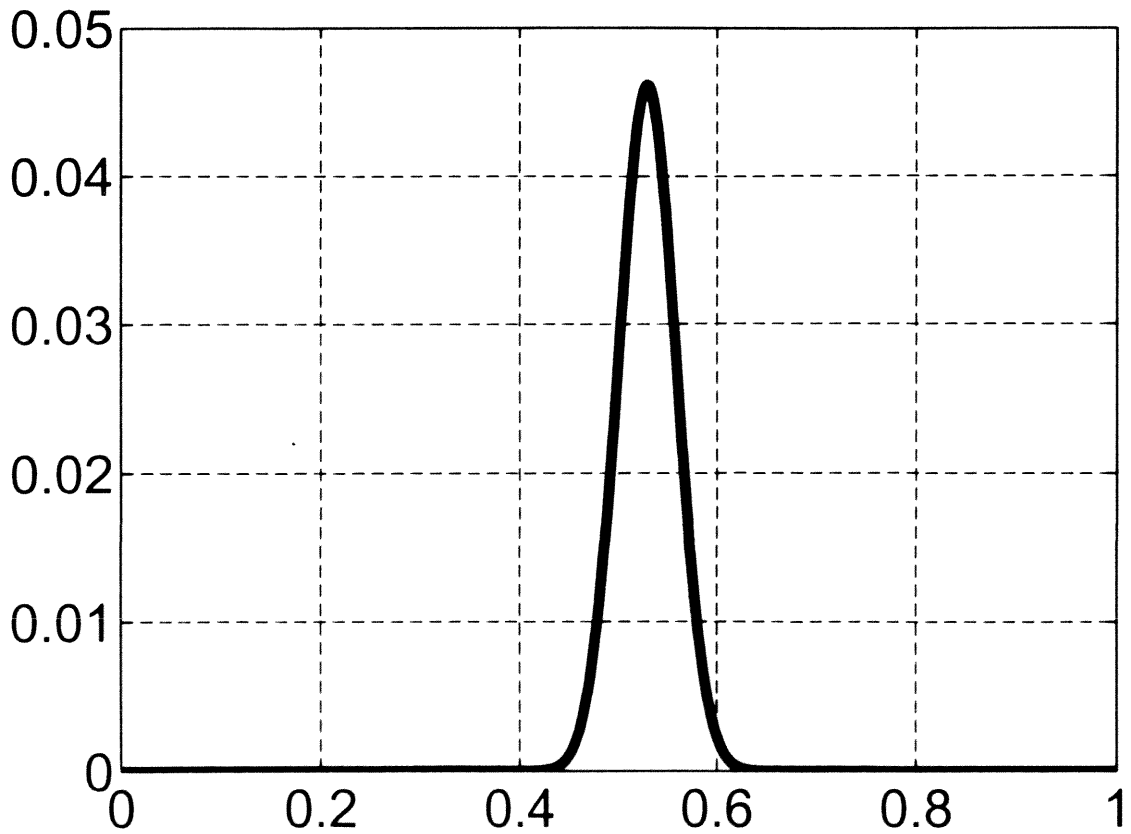
[Tarkistamme tämän kohta differentiaalilaskennan avulla!]

Huom. L:n kuvaaja "huipukkaampi" kuin (A)-tilanteessa
ts. "uskottavat" parametrin arvot ovat melko kapealla
välillä 0.53:n ympärillä \Rightarrow θ :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 θ kaikkien mahdollisten arvojen joukko) on Ω .

* 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 ehtö $L(\theta; \underline{x}) > L(\theta'; \underline{x})$, sanomme, että θ on (aineiston \underline{x} valossa) uskottavampi parametrinarvo kuin θ' .

* 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 θ 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 (θ ei ole sat.muuttuja)!

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

$$l(\theta) = \log L(\theta)$$

$$= \log \binom{n}{t} + t \log \theta + (n-t) \log (1-\theta)$$

$$l'(\theta) = \frac{t}{\theta} - \frac{n-t}{1-\theta} = \frac{t - \theta n}{\theta(1-\theta)}$$

(etumerkin määrää $t - \theta n$)

	t/n	
l'	+	-
l	↗	↘

Johtopäätös: l ja siten myös L saa suurimman arvonsa pisteessä $\theta = t/n$

Siten θ :n su-estimaatti on $\hat{\theta} = \frac{t}{n}$

kuten jo esimerkiksi s. 9-10 saatiin.

Huom. 1. Jos havaintoja vast. sm:t X_1, \dots, X_n ovat ii, on uskottavuusfunktio muotoa

$$L(\theta; \underline{x}) = g_{X_1}(x_1; \theta) \cdots g_{X_n}(x_n; \theta)$$

jossa $g_{X_i} = X_i$:n ptnf/ff

Jos X_1, \dots, X_n lisäksi samom jakautuneita (s. 8 tilanne), on tässä tetysti $g_{X_i} = g$ sama kaikille i .

Huom. 2. Yleensä su-estimaatin määrittäminen on laskuteknisesti mukavampaa uskottavuusfunktion logaritmin avulla! (Syy: logaritmointi muuntaa tulot summiksi.)

Huom. 3. Todellisissa sovellustilanteissa, jorssa mallit ovat monimutkaisia ja parametri θ yleensä moniulotteinen (ts. koostuu monesta reaalisesta komponentista), su-estimaatin määrittäminen ei onnistu analyttisesti "suljetussa muodossa" vaan joudutaan turvautumaan numeerisiin (optimointi) menetelmiin.