

5. NORMAALIJAKAUMAMALLIN PARAMETRIEN ESTIMOINNISTA

[vrt. Arjas-Sirén, jaksot 2.6 (ja 2.4)]

Tark. mallia "riippumaton otos jakaumasta $N(\mu, \sigma^2)$ ":

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

Mallin yhteistihetyys on (ks. s. 8) (merk. $\underline{x} = (x_1, \dots, x_n)$)

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

Tutkitaan parametrien μ ja σ^2 estimointia su-menetelmällä.

(A) olet. aluksi: $\sigma^2 > 0$ on tunnettu luku ja parametriina vain μ , $\mu \in \mathbb{R}$. (harvoin realistinen tilanne)

Uskottavuusfunktio (vast. aineistoa $\underline{x} = (x_1, \dots, x_n)$)

$$(x) \quad L(\mu) = \frac{1}{(2\pi\sigma^2)^{n/2}} \exp \left\{ -\frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2 \right\}, \quad \mu \in \mathbb{R}$$

logaritmi:

$$l(\mu) = \log L(\mu) = -\frac{n}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2$$

$$l'(\mu) = -\frac{1}{2\sigma^2} \sum_{i=1}^n [-2(x_i - \mu)] = \frac{1}{\sigma^2} \left[\sum_{i=1}^n x_i - n\mu \right]$$

$$l'(\mu) = 0 \quad (\Leftrightarrow) \quad \mu = \frac{1}{n} \sum_{i=1}^n x_i$$

vaihtaa merkkinsä
+ :sta - :een

Saatu: μ :n su-estimaatti on

$$\hat{\mu} = \bar{x} = \frac{1}{n} \sum_{i=1}^n x_i \quad (\text{otokeskiarvo})$$

(B) Yleinen tilanne: parametri on 2-ulotteinen (μ, σ^2) , jossa $\mu \in \mathbb{R}$, $\sigma^2 > 0$.

Uskottavuusfunktio $L(\mu, \sigma^2)$ kuten edellä (lauseke (*)).
Pätee hajotelma (tarhista!)

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

$$\left(\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i \quad \text{otosheskiarvo} \right)$$

Logaritmoidaan ja käytetään tätä hajotelmaa:

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

Huom. tämä riippuu μ :stä vain viimeisen termin kautta!

Selvästi $-\frac{n}{2\sigma^2} (\bar{x} - \mu)^2 \leq 0$ aina ja "=" pätee

$$\Leftrightarrow \mu = \bar{x}.$$

Tark. sitten (vain σ^2 :sta riippuvaa!) funktiota

$$u(\sigma^2) = l(\bar{x}, \sigma^2) = -\frac{n}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \bar{x})^2,$$

$$u'(\sigma^2) = -\frac{n}{2\sigma^2} + \frac{1}{2\sigma^4} \sum_{i=1}^n (x_i - \bar{x})^2$$

Huom. Derivaatta σ^2 :n suhteen, ei σ :n!

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

$$= 0 \quad \Leftrightarrow \quad \sigma^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2$$

(tee merkkitarkastelu: tämä todella on maksimikohta!)

Saatu: parametrin (μ, σ^2) su-estimaatti on

$$\hat{\mu} = \bar{x} = \frac{1}{n} \sum_{i=1}^n x_i, \quad \hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2$$

Huomautuksia & lisätietoja:

* Tuntuu järkevältä estimoida jakauman odotusarvoa ("teoreettista keskiarvoa") otoskeskiarvolla $\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$

* Vastaavalle sm:lle $\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$ (eli estimaattorille) pätee

$$E(\bar{X}) = \frac{1}{n} \sum_{i=1}^n E(X_i) = \frac{1}{n} \sum_{i=1}^n \mu = \mu$$

Sanotaan: \bar{X} on harhaton estimaattori μ :lle.

Se siis tuottaa "odotusarvoisesti" oikean tuloksen!

Torstetun ainerstokerrun näkökulmasta: jos otanta torstetaan yhä uudelleen ja uudelleen ja jokaisesta ainerstosta \underline{x} lasketaan $\hat{\mu} = \bar{x}$, niin saadut estimaatit osuvat keskimäärin oikeaan. (Vrt. HT 2/5.)

* Sen sijaan $\hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2$ ei ole harhaton estimaattori σ^2 :lle, vaan sen odotusarvo on $E(\hat{\sigma}^2) = \frac{n-1}{n} \sigma^2 = (1 - \frac{1}{n}) \sigma^2$.

Tästä syystä σ^2 :n estimaattina useimmiten käytetään lukua

$$S^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2 \quad (\text{otosvarianssi})$$

Sitä vastaavalle sm:lle S^2 pätee $E(S^2) = \sigma^2$.

Ero $\hat{\sigma}^2$:een on merkityksetön, jollei n ole kovin pieni.

* Tarkemmin ottaen pätee seuraava lause:

Kun $X_1, \dots, X_n \sim N(\mu, \sigma^2) \perp$, niin

i) $\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i \sim N(\mu, \frac{\sigma^2}{n})$

ii) $\frac{(n-1) S^2}{\sigma^2} = \frac{1}{\sigma^2} \sum_{i=1}^n (X_i - \bar{X})^2 \sim \chi_{n-1}^2$

iii) $\bar{X} \perp S^2$

"khiin-toriseen" jakauma, $n-1$ vapausastetta

[Todistus aineopintojen tn-laskennan kurssilla.]

6. HIEMAN YLEISTÄ VÄLIESTIMOINNISTA

vrt.
Arjas-Sirén,
luku 3

Tark. yleisistä tilastollisista malleista $f(\underline{x}; \theta)$, jossa (tuntematon) parametri on θ ja aineisto $\underline{x} = (x_1, \dots, x_n)$.

Edellä tutustumme piste-estimointiin, jossa aineiston perusteella piti määrittää sellainen parametrin arvo, joka olisi (tavalla tai toisella) hyvä ja perusteltu "arvans" todelliselle parametrin arvolla. Su-menetelmä oli (eräs) ratkaisu tähän tehtävään.

Ongelma: piste-estimaatti harvoin osuu juuri oikeaan (emmehä tiedä milloin niin käy).

Haluaisimme määrittää sellaisen joukon parametrin arvoja (esim. $\hat{\theta}$:n ympäriltä), josta voitaisiin "suurella varmuudella" sanoa, että se sisältää oikean parametrin arvon.

1-ulotteisessa tapauksessa (ts. kun θ on 1-ulotteinen tai halutaan tehdä päätelmä vain sen yksittäisestä komponentista) tällainen joukko on yleensä väli ja puhutaan väli-estimoinnista.

KAKSI RATKAISUA tähän tehtävään:

1. Uskottavuusvälit tai -joukot (melko harvoin käytetty mutta helppo menetelmä uskottavuusfunktion pohjalta)

Olk. $L(\theta) = L(\theta; \underline{x}) = f(\underline{x}; \theta)$ aineistoa \underline{x} vast. uskottavuusfunktio.

Mursta: Su-estimaatti $\hat{\theta}$ on piste, jossa $L(\theta)$ saa suurimman arvonsa.

Määr. Kun $0 < c < 1$, niin joukko

$$U_c = U_c(x) = \{ \theta \in \Omega \mid L(\theta) \geq c L(\hat{\theta}) \}$$

Ω = kaikkien mahdollisten param. arvojen joukko

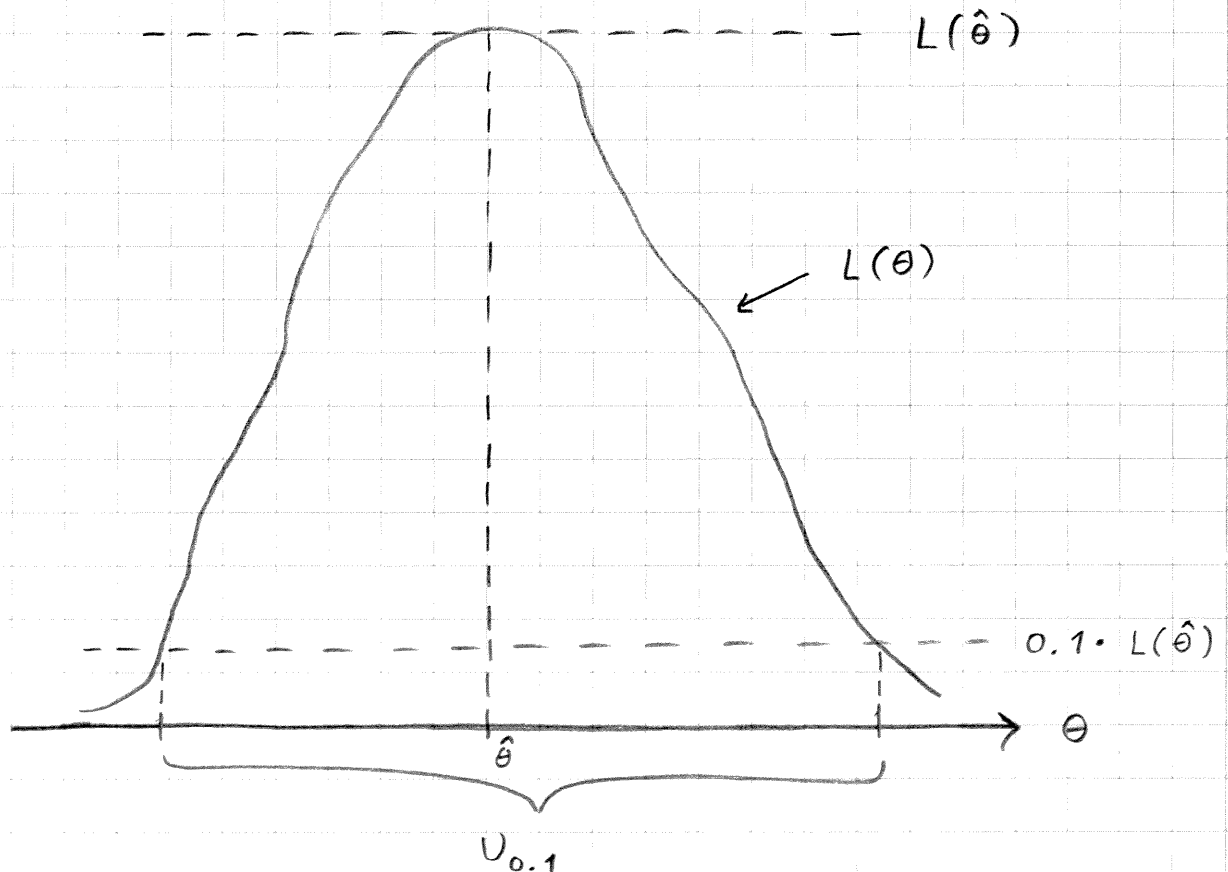
on $100 \cdot c$ %:n uskottavuusjoukko θ :lle.
(1-ulott. tapauksessa se on yleensä väli.)

Esim. $U_{0.1} = \{ \theta \in \Omega \mid L(\theta) \geq 0.1 \cdot L(\hat{\theta}) \}$

on 10 %:n uskottavuusjoukko.

Sen ulkopuolelle jäävät ne θ :n arvot, joiden uskottavuus on alle 10 % uskottavuuden suurimmasta arvosta.

(ks. kuva alla)



2. Luottamusvälit (laajassa käytössä tilastollisessa tutkimuksessa)

Kiinnitetään jokin (pieni) luku $0 < \alpha < 1$.

Pyritään aineiston \underline{x} perusteella määrittämään joukko $A(\underline{x})$ parametriarvossa Ω siten, että vastaavalle "satunnaiselle joukolle" $A(\underline{X})$ pätee

$$P(\theta \in A(\underline{X})) \geq 1 - \alpha.$$

Tällöin $A(\underline{x})$ on θ :n luottamusjoukko (luottamus)tasolla $1 - \alpha$.

Tyypillisesti esim. $\alpha = 0.05$, jolloin

$$P(\theta \in A(\underline{X})) \geq 0.95$$

ja kyse on 95 %:n (eli tason 0.95) luottamusjoukosta.

Satunnainen (toistetun aineistonkeruun mielessä!)

joukko $A(\underline{X})$ siis peittää todellisen parametrin arvon θ (ainakin) tn:llä 0.95.

1-ulotteisen parametrin tapauksessa $A(\underline{x})$ on yleensä väli reaalialueella: voimme siis kirjoittaa esim.

$A(\underline{x}) = (a(\underline{x}), b(\underline{x}))$, jossa päätepisteillä on ominaisuus

$$P(a(\underline{X}) < \theta < b(\underline{X})) \geq 1 - \alpha$$



Tutkimme seuraavassa jaksossa luottamusvälien muodostamista normaalijakaumamallin tapauksessa. Ylseremmien asioita käsitellään aineopintojen til. päättelyn kurssilla.

6. LUOTTAMUSVÄLIT NORMAALIJAKAUMALLE

[vrt. Arjas-Siren, luku 3]

Tark. mallia $X_1, \dots, X_n \sim N(\mu, \sigma^2) \perp$.

Tavoite: Muodostetaan luottamusväli odotusarvolle μ annetulla luottamustasolla $1-\alpha$ (esim. 0.95)

Pal. mieleen, että μ :n (piste-)estimaattina käytetään

$$\hat{\mu} = \bar{x} = \frac{1}{n} \sum_{i=1}^n x_i.$$

Vastaavasti σ^2 :n estimaattina (useimmiten) käytetään

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

(A) Olet. että $\sigma^2 > 0$ on tunnettu luku (ei siis parametri)

Sivulla 16 todettiin, että sm:lle $\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$ pätee

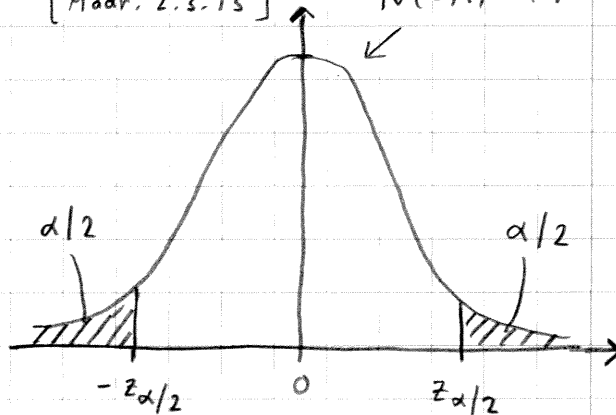
$$\bar{X} \sim N\left(\mu, \frac{\sigma^2}{n}\right) \quad [\text{vrt. Tuominen, TNI, Esim. 3.7.9}]$$

eli

$$(*) \quad \frac{\bar{X} - \mu}{\sigma/\sqrt{n}} \sim N(0,1) \quad \left[\begin{array}{l} \text{TNI,} \\ \text{Määr. 2.3.13} \end{array} \right] \quad N(0,1)\text{-tf}$$

Valitaan $N(0,1)$ -jakaumasta kohta $z_{\alpha/2}$, josta oikealle "häntätu" = $\alpha/2$

Symmetria \Rightarrow $-z_{\alpha/2}$:sta vasemmalle "häntätu" on samoin $\alpha/2$



Välille $(-z_{\alpha/2}, z_{\alpha/2})$ jäävä tn-massa = $1-\alpha$.

Siten

$$P\left(-z_{\alpha/2} < \frac{\bar{X} - \mu}{\sigma/\sqrt{n}} < z_{\alpha/2}\right) = 1 - \alpha$$

⇒

$$P\left(\bar{X} - z_{\alpha/2} \frac{\sigma}{\sqrt{n}} < \mu < \bar{X} + z_{\alpha/2} \frac{\sigma}{\sqrt{n}}\right) = 1 - \alpha \quad (**)$$

Siten (satunnaisten) pisteiden $\bar{X} - z_{\alpha/2} \frac{\sigma}{\sqrt{n}}$ ja $\bar{X} + z_{\alpha/2} \frac{\sigma}{\sqrt{n}}$ rajoittama väli sisältää μ :n todella $1 - \alpha$.

Tämä merkitsee:

Kun $\underline{x} = (x_1, \dots, x_n)$ on havaittu aineisto ja $\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$, niin väli

$$\left(\bar{x} - z_{\alpha/2} \frac{\sigma}{\sqrt{n}}, \bar{x} + z_{\alpha/2} \frac{\sigma}{\sqrt{n}}\right) \quad (\#)$$

on μ :n luottamusväli luottamustasolla $1 - \alpha$.

Esim. Kun $\alpha = 0.05$, niin

$$z_{\alpha/2} = z_{0.025} \approx 1.96 \quad (\approx 2)$$

joten 95 %:n luottamusväli μ :lle on

$$\left(\bar{x} - 1.96 \frac{\sigma}{\sqrt{n}}, \bar{x} + 1.96 \frac{\sigma}{\sqrt{n}}\right)$$

Huom. Lukujen $z_{\alpha/2}$ arvoja on normaalijakauman taulukoissa!

(B) Oletetaan, että myös σ^2 on (tuntematon) parametri

Ongelma: väli (#) yllä ei käytettävissä, koska σ ei tunnettu ⇒ joudumme estimoimaan sen aineistosta.

Estimaattina käytämme

$$s = \sqrt{s^2} = \sqrt{\frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2} \quad \left(\frac{\text{otos-keskihajonta}}{\quad}\right)$$

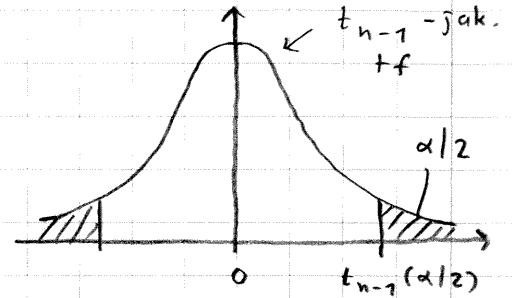
Olkoon S vastaava satunnaismuuttuja.

Tällöin pätee (vrt. A-tapauksen jakaumatulos (*) !)

$$\frac{\bar{X} - \mu}{S/\sqrt{n}} \sim t_{n-1} \quad \left(\begin{array}{l} \text{Studentin } t\text{-jakauma,} \\ n-1 \text{ vapausastetta} \end{array} \right)$$

Menetellään kuten A-kohdassa:

Valitaan piste $t_{n-1}(\alpha/2)$, josta oikealle t_{n-1} -jakauman häntäosuus $= \alpha/2$. Päädytään (**)-in sijasta tulokseen



$$P\left(\bar{X} - t_{n-1}(\alpha/2) \frac{S}{\sqrt{n}} < \mu < \bar{X} + t_{n-1}(\alpha/2) \frac{S}{\sqrt{n}}\right) = 1 - \alpha$$

Saatu:

Kun $\underline{x} = (x_1, \dots, x_n)$ on havaittu ainesisto ja \bar{x} = otoskeskiarvo sekä s = otoskeskihajonta (ks. edell. sivu), niin väli

$$\left(\bar{x} - t_{n-1}(\alpha/2) \frac{s}{\sqrt{n}}, \bar{x} + t_{n-1}(\alpha/2) \frac{s}{\sqrt{n}} \right) \quad (\#\#)$$

on μ :n luottamusväli luottamustasolla $1 - \alpha$.

Huom. * Verrattuna lv:iin (#) tässä lasketaan ainesistosta \bar{x} :n lisäksi myös keskihajonta s . lisäksi kerroin $t_{n-1}(\alpha/2)$ riippuu otoskoosta n .

* t_{n-1} -jakaumalla on "paksimmat hännät" kuin $N(0,1)$ -jakaumalla \Rightarrow kerroin $t_{n-1}(\alpha/2)$ on suurempi kuin $z_{\alpha/2}$. Käytännössä ero on merkityksellinen vain pienillä n :n arvoilla (\sim alle 50), sillä t_{n-1} lähestyy $N(0,1)$ -jakaumaa kun $n \rightarrow \infty$.

* Lukuja $t_{n-1}(\alpha/2)$ löytyy taulukoista ja tietokoneohjelmita, joten t -jakauman tiheysfunktioita ei tarvitse tietää.