

# MATH.APP.210 Johdatus todennäköisyyslaskentaan ja tilastolliseen päättelyyn

Jussi Kangas

# Sisältö (osa 5)

- 1 Otosjakaumien piste-estimointi
  - Otoskeskiarvo
  - Approksimaatio normaalijakaumalla
  - Otosvarianssi ja  $\chi^2$ -jakauma
  
- 2 Luottamusväliestimointi
  - Luottamusvälin käsite
  - Odotusarvon luottamusväli ja  $t$ -jakauma
  - Muita luottamusvälejä

# Satunnaismuuttujan otos

- tarkastellaan  $n$ -toistokoetta, jossa satunnaismuuttujalle  $X$  realisoituu jokin arvo
- jos  $X_1, X_2, \dots, X_n$  ovat satunnaismuuttujia, jotka liittyvät erillisiin toistoihin, ne muodostavat  $n$ :n kappaleen **otoksen** satunnaismuuttujasta  $X$
- muuttujat  $X_i$  ovat riippumattomia ja noudattavat samaa jakaumaa kuin  $X$ , jolloin  $E(X_i) = E(X)$  ja  $\text{Var}(X_i) = \text{Var}(X)$
- myös realisoituneiden arvojen joukkoa  $x_1, x_2, \dots, x_n$  kutsutaan otokseksi
- jatkossa otos viittaa satunnaismuuttujan otokseen

# Otossuureet

- satunnaismuuttujan  $X$  otoksesta  $X_1, X_2, \dots, X_n$  riippuva **otossuure**  $\Theta$  (iso theta) on otokseen liittyvä tunnusluku, kuten otoskeskiarvo  $\bar{X}$ , otosvarianssi  $S^2$  tai otoskeskihajonta  $S$ .
- otossuureet muodostetaan otosmuuttujien  $X_i$  funktioina, joten ne ovat satunnaismuuttujia
- jos ei tunneta  $X$ :n jakauman jotain parametria  $\theta$  (pieni theta), voidaan sitä arvioida eli **estimoida** otossuureilla tai käyttää otossuureita parametreihin liittyvien väitteiden testaamiseen
- usein estimoidaan mm.  $X$ :n odotusarvoa ja varianssia
- parametriin  $\theta$  liittyvä otossuure  $\Theta$  on parametrin **estimaattori** ja sille realisoituva arvo on **(piste-)estimaatti**

# Otoskeskiarvo

- jos  $X_1, X_2, \dots, X_n$  on otos satunnaismuuttujasta  $X$  ja  $x_1, x_2, \dots, x_n$  otosmuuttujien realisoituneet arvot, niin satunnaismuuttujan  $X$  **otoskeskiarvo** on satunnaismuuttuja

$$\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$$

ja sen realisoitunut arvo on

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

- koska  $\bar{X}$  on satunnaismuuttuja, voimme tarkastella sen jakaumaa

# Otoskeskiarvon jakauma

- jos satunnaismuuttujan  $X$  odotusarvo on  $\mu$  ja varianssi  $\sigma^2$ , niin otoskeskiarvon  $\bar{X}$  odotusarvo ja varianssi ovat

$$E(\bar{X}) = E(X) = \mu \quad \text{ja} \quad \text{Var}(\bar{X}) = \frac{\text{Var}(X)}{n} = \frac{\sigma^2}{n}$$

- otossuure  $\Theta$  on **harhaton estimaattori** parametrille  $\theta$  ja realisoitunut arvo on **harhaton estimaatti**, jos

$$E(\Theta) = \theta$$

- $E(\bar{X})$  on harhaton estimaattori  $E(X)$ :lle

# Keskiarvon keskivirhe

- otoskeskiarvon keskihajontaa kutsutaan **keskiarvon keskivirheeksi**, merkitään

$$D(\bar{X}) = \sqrt{\text{Var}(\bar{X})} = \frac{\sigma}{\sqrt{n}}$$

- otoksesta arvioitu  $X$ :n odotusarvo ja sen virhearvio esitetään usein karkeasti arvioiden muodossa  $\mu \pm \sigma/\sqrt{n}$
- luotettavampi tapa on muodostaa odotusarvolle  $\mu$  ns. **luottamusväli**, joka suurella todennäköisyydellä sisältää varsinaisen odotusarvon

# Normaalijakautuneen satunnaismuuttujan otoskeskiarvo

- jos  $X_1, X_2, \dots, X_n$  on otos normaalijakautuneesta satunnaismuuttujasta  $X \sim N(\mu, \sigma^2)$ , niin otoskeskiarvo  $\bar{X}$  on normaalijakautunut,

$$\bar{X} \sim N\left(\mu, \frac{\sigma^2}{n}\right)$$

- aina ei voida olettaa, että  $X$ :llä on normaalijakauma tai että sen jakauma edes tunnetaan
- kuitenkin suurilla otoksilla otoskeskiarvon  $\bar{X}$  jakauma on lähellä normaalijakaumaa (merkissä  $\sim$  piste päällä = "noudattaa likimain") riippumatta satunnaismuuttujan  $X$  jakaumasta tai siitä onko jakauma diskreetti tai jatkuva



# Normaalijakautuneen satunnaismuuttujan otoskeskiarvo

## Esimerkki

PVC putkea valmistetaan asetuksilla joiden tiedetään tuottavan keskimäärin 2,565 senttiä paksua putkea hajonnan ollessa 0,0076 senttiä. Lisäksi tiedetään, että putken paksuus on normaalijakautunut. Millä todennäköisyydellä 9 kappaleen otoksessa otoskeskiarvo olisi suurempi kuin 2,56 senttiä ja pienempi kuin 2,57 senttiä?

Nyt siis  $\bar{X} \sim N(2.565, 0.0076^2/9)$ . Lasketaan todennäköisyys tapahtumalle  $P(2.56 \leq \bar{X} \leq 2.57)$ .

Matlabista komennolla

`normcdf(2.57,2.565,0.0076/3)-normcdf(2.56,2.565,0.0076/3)`

saadaan todennäköisyydeksi likimain 0,9516.

# Keskeinen raja-arvolause

## Keskeinen raja-arvolause

Olkoon  $X_1, X_2, \dots, X_n$  otos satunnaismuuttujasta  $X$ , jonka odotusarvo on  $\mu$  ja varianssi  $\sigma^2$ . Tällöin standardoidun otoskeskiarvon

$$\bar{X}^* = \frac{\bar{X} - \mu_{\bar{X}}}{\sigma_{\bar{X}}} = \frac{\bar{X} - \mu}{\sigma/\sqrt{n}}$$

kertymäfunktio  $F(t)$  lähenee  $N(0,1)$ -jakauman kertymäfunktia  $\Phi(t)$ , kun  $n \rightarrow \infty$  eli

$$F(t) = P(\bar{X}^* \leq t) \rightarrow \Phi(t), \quad \text{kun } n \rightarrow \infty.$$

# Normaaliapproksimaatio

- otoskeskiarvon **normaaliapproksimaatio**: suurilla  $n$ :n arvoilla  $\bar{X}^*$  on likimain standardinormaalijakautunut,  $\bar{X}^* \sim N(0, 1)$
- yleensä normaaliapproksimaatio antaa hyviä arvioita jo arvoilla  $n \geq 30$  riippumatta  $X$ :n jakaumasta
- “lähellä” normaalijakaumaa olevilla jakaumilla voidaan soveltaa myös pienemmillä  $n$ :n arvoilla

## Seuraus

Olkoon  $X_1, X_2, \dots, X_n$  otos satunnaismuuttujasta  $X$ , jonka odotusarvo on  $\mu$  ja varianssi  $\sigma^2$ . Tällöin otoskeskiarvolle ja otoksen summalle pätevät

$$\bar{X} \sim N\left(\mu, \frac{\sigma^2}{n}\right) \qquad X_1 + X_2 + \dots + X_n \sim N(n\mu, n\sigma^2)$$

# Esimerkki normaaliapproksimaatiosta

## Esimerkki

Hissi voi kuljettaa korkeintaan 25 henkilöä tai 2000 kg. Henkilöpaino (kg) on satunnaismuuttuja, jonka odotusarvo on  $\mu = 74$  ja varianssi  $\sigma^2 = 100$ . Millä todennäköisyydellä satunnaisesti valitun 25 henkilön kokonaispaino ylittää 2000 kg? Jos  $X =$  "25 hlön kokonaispaino", niin

$$X \sim N(25 \cdot 74, 25 \cdot 100) = N(1850, 2500).$$

Nyt siis kysytään, että mitä on  $P(X > 2000)$ . Tämä on jälleen helpointa selvittää ohjelmistolla: Matlab-komento `normcdf(2000,1850,sqrt(2500),'upper')` antaa todennäköisyydeksi likimain 0,0013.

# Otosvarianssi

- jos  $X_1, X_2, \dots, X_n$  on otos satunnaismuuttujasta  $X$  ja  $x_1, x_2, \dots, x_n$  otosmuuttujien realisoituneet arvot, niin satunnaismuuttujan  $X$  **otosvarianssi** on satunnaismuuttuja

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

ja **otoshajonta**  $S = \sqrt{S^2}$ , joiden realisoituneet arvot ovat

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

- otosvarianssi on varianssin harhaton estimaattori eli

$$E(S^2) = \text{Var}(X) = \sigma^2$$

# $\chi^2$ -jakauma

- normaalijakautuneiden satunnaismuuttujien muunnoksina saadaan uusia jakaumia, mm.  $\chi^2$ -jakauma ( $\chi^2$  = khii toiseen).
- jatkuva satunnaismuuttuja  $W$  noudattaa  **$\chi^2$ -jakaumaa vapausastein  $n$** ,  $W \sim \chi^2(n)$ , jos sen tiheysfunktio on

$$f(x) = \frac{1}{2^{\frac{n}{2}} \Gamma\left(\frac{n}{2}\right)} x^{\frac{n}{2}-1} e^{-\frac{x}{2}}, \quad \text{kun } x \in \Omega = [0, \infty),$$

missä  $\Gamma(t)$  on Eulerin gammafunktio  $\Gamma(t) = \int_0^\infty e^{-x} x^{t-1} dx$

- satunnaismuuttujan  $W \sim \chi^2(n)$  odotusarvo ja varianssi ovat

$$E(W) = n \qquad \text{Var}(W) = 2n$$

## $\chi^2$ -jakautuneiden satunnaismuuttujien summa

- jos satunnaismuuttujat  $Z_i \sim N(0,1)$ ,  $i = 1, 2, 3, \dots, n$  ovat riippumattomia, niin niiden neliösummalle pätee

$$W = \sum_{i=1}^n Z_i^2 = Z_1^2 + Z_2^2 + \dots + Z_n^2 \sim \chi^2(n)$$

- jos alla esitetyt satunnaismuuttujat  $W_1$  ja  $W_2$  ovat riippumattomia ja jakautuneet seuraavasti

$$W_1 = \sum_{i=1}^n Z_i^2 \sim \chi^2(n) \quad \text{ja} \quad W_2 = \sum_{i=1}^m U_i^2 \sim \chi^2(m),$$

niin muuttujat  $Z_1, Z_2, \dots, Z_n, U_1, U_2, \dots, U_m$  ovat riippumattomia ja  $W_1 + W_2 \sim \chi^2(n+m)$

# Otosvarianssin jakauma

- sivun otsikosta huolimatta tutkitaan pikemminkin otosvarianssiin  $S^2$  liittyvän otossuureen  $W$  jakaumaa

## Lause

Jos  $X_1, X_2, \dots, X_n$  on otos muuttujasta  $X \sim N(\mu, \sigma^2)$ , niin

- 1  $\bar{X}$  ja  $S^2$  ovat riippumattomia
- 2 otossuurelle

$$W = \frac{(n-1)S^2}{\sigma^2} = \frac{1}{\sigma^2} \sum_{i=1}^n (X_i - \bar{X})^2 = \sum_{i=1}^n \left( \frac{X_i - \bar{X}}{\sigma} \right)^2$$

pätee  $W \sim \chi^2(n-1)$



## Esimerkki

Lasilevyn paksuuden hajonta ei saisi ylittää arvoa 0,2 mm. Kun  $n = 40$  satunnaisesti valitun lasilevyn paksuudet mitattiin saatiin otoshajonnaksi 0,214 mm. Paksuuden jakaumaa voidaan pitää normaalina. Mikä on todennäköisyys näinkin suurelle otoshajonnalle, jos tehtaan laitteet toimivat kuten pitää?

Tutkitaan siis kuinka suurella todennäköisyydellä otoshajonnaksi saadaan arvo 0,214 mm tai vielä suurempi arvo, jos todellisuudessa  $\sigma = 0,2$  mm. Nyt

$$P(S^2 \geq 0,214^2) = P\left(\frac{(n-1)S^2}{\sigma^2} \geq \frac{39 \cdot 0,214^2}{0,2^2}\right) \approx P(W \geq 44,65)$$

$\chi^2$ -jakaumasta vapausasteella  $n - 1 = 39$  saadaan Matlabilla komennolla `chi2cdf(44.65,39,'upper')` todennäköisyydeksi likimain 0.2464.

# Luottamusväliestimointi

- tilastotieteessä puhutaan usein **populaatiosta**, josta valitaan satunnaisotos ja jonka populaatiojakauma tunnetaan tai ei tunneta
- usein estimoidaan jotain populaation liittyvää parametria  $\theta$  otokseen liittyvillä harhaatomilla estimaattoreilla
- tässä piste-estimoinnilla saatavaa parametrin arvon tarkkuutta ei oikein voi päätellä
- tyypillisempi estimoinnin muoto on siksi **luottamusväliestimointi**
- otoksen perusteella reaalityökalu, jolla parametrin  $\theta$  **luotetaan** olevan

# Luottamusväli

- satunnaismuuttujan  $X$  parametrin  $\theta$   $100(1 - \alpha)\%:n$  **väliestimaattori** ( $\alpha = \text{alfa}$ ,  $\alpha \in [0, 1]$ ) on satunnaisväli  $[\hat{\theta}_1, \hat{\theta}_2]$ , missä  $\hat{\theta}_1$  ja  $\hat{\theta}_2$  ovat sellaisia satunnaismuuttujan  $X$  otoksesta  $X_1, X_2, \dots, X_n$  riippuvia välin päätepisteitä, joille

$$P(\hat{\theta}_1 \leq \theta \leq \hat{\theta}_2) = 1 - \alpha$$

- realisoitunutta väliä  $[\theta_1, \theta_2]$  sanotaan silloin parametrin  $\theta$   $100(1 - \alpha)\%:n$  **luottamusväliksi**
- luku  $1 - \alpha$  on välin **luottamusaste** tai **luottamustaso**
- päätepiseet ovat välin **alempi** ja **ylempi luottamusraja**

# Luottamusvälin tulkinta

- monet tekevät intuitiivisen, mutta hiukan virheellisen tulkinnan, että luottamusväli olisi reaalitykuväli, jolle parametrin arvo  $\theta$  kuuluu todennäköisyydellä  $1 - \alpha$
- parametriin  $\theta$  liittyvä väliestimaattori  $[\hat{\theta}_1, \hat{\theta}_2]$  muodostuu satunnaismuuttujista  $\hat{\theta}_1$  ja  $\hat{\theta}_2$ , joille realisoituu erilaisia arvoja otoksesta riippuen
- jos otoksia kerättäisiin hyvin monta erilaista, niin niiden perusteella lasketuista luottamusväleistä  $100(1 - \alpha)\%$  sisältää parametrin  $\theta$  todellisen arvon.
- luottamusväli on siis olennaisesti satunnaisväli

# Luottamusvälin valinnasta

- mitä suurempi luottamusaste vaaditaan, sitä leveämpi luottamusväli
- väliestimaattorin määritelmän ehdon  $P(\hat{\theta}_1 \leq \theta \leq \hat{\theta}_2) = 1 - \alpha$  täyttäviä välejä on useita samalle luottamusasteelle  $1 - \alpha$ , eli määritelmä ei kerro tarkempaa perustetta välin valitaan
- useimmiten perusteena on seuraavanlainen symmetria

$$P(\theta \leq \hat{\theta}_1) = P(\theta \geq \hat{\theta}_2) = \frac{\alpha}{2}$$

- usein käytettyjä luottamusasteita ovat 90%, 95% ja 99%, jotka vastaavat **riskitasoja**  $\alpha = 0.1$ ,  $\alpha = 0.05$  ja  $\alpha = 0.01$

# Odotusarvon luottamusväli, kun varianssi tunnetaan

## Lause

Jos  $X_1, X_2, \dots, X_n$  on otos muuttujasta, jonka **varianssi**  $\sigma^2$  **tunnetaan**, niin odotusarvon  $100(1 - \alpha)\%$ :n väliestimaattori on

$$\left[ \bar{X} - z_{\alpha/2} \frac{\sigma}{\sqrt{n}}, \bar{X} + z_{\alpha/2} \frac{\sigma}{\sqrt{n}} \right],$$

missä  $\Phi(z_{\alpha/2}) = 1 - \frac{\alpha}{2}$ . Jos otoskeskiarvolle realisoituu arvo  $\bar{x}$ , niin odotusarvon  $100(1 - \alpha)\%$ :n luottamusväli on

$$\left[ \bar{x} - z_{\alpha/2} \frac{\sigma}{\sqrt{n}}, \bar{x} + z_{\alpha/2} \frac{\sigma}{\sqrt{n}} \right].$$

# Odotusarvon luottamusväli, kun varianssi tunnetaan

## Esimerkki

Tietyn muovisen liittimen paksuutta tutkittaessa otettiin  $n = 75$  liittimen otos ja mitattiin paksuudet. Valmistusprosessista johtuen tiedetään paksuuden hajonnaksi  $\sigma = 0,04$  mm. Otoskeskiarvo on  $\bar{x} = 7,874$  mm. Etsitään paksuudelle 95 % ja 90 % luottamusvälit.

95 prosentin väliä määritettäessä  $\alpha = 0,05$  eli  $\alpha/2 = 0,025$ .

Matlabilla komennolla `norminv(1-0.025,0,1)` saadaan  $z_{\alpha/2} \approx 1,96$ .

Luottamusväliksi kaavalla  $\bar{X} \pm z_{\alpha/2} \frac{\sigma}{\sqrt{n}}$  saadaan (7,8649; 7,8831).

90 prosentin väliä määritettäessä  $\alpha = 0,1$  eli  $\alpha/2 = 0,05$ .

Matlabilla komennolla `norminv(1-0.05,0,1)` saadaan  $z_{\alpha/2} \approx 1,645$ .

Luottamusväliksi kaavalla  $\bar{X} \pm z_{\alpha/2} \frac{\sigma}{\sqrt{n}}$  saadaan (7,8664; 7,8816).

# Huomioita odotusarvon luottamusvälistä

- Edellisen lauseen soveltaminen vaatii sen, että 1) jakauman varianssi  $\sigma^2$  tunnetaan ja 2) otoskoko on kohtuullisen suuri (keskeisen raja-arvolauseen vaatimus)
- tavallisempi tilanne kuitenkin on, että varianssia ei tunneta, vaikka otosjakauma voidaankin olettaa normaalijakautuneeksi
- varienssin harhaton piste-estimaatti löydettiin otosvarienssin  $S^2$  avulla, mutta koska se on satunnaismuuttuja, väliestimointiin vaikuttavien satunnaistekijöiden lisääntymisen myötä otoskeskiarvo ei ole enää normaalijakautunut
- tuntemattoman varienssin tapauksessa käytetään (Studentin)  $t$ -jakaumaa



# Studentin $t$ -jakauma

- jatkuva satunnaismuuttuja  $T$  noudattaa Studentin  $t$ -jakaumaa vapausastein  $n$ ,  $T \sim t(n)$ , jos tiheysfunktio on

$$f(t) = \frac{1}{\sqrt{n\pi}} \frac{\Gamma\left(\frac{n+1}{2}\right)}{\Gamma\left(\frac{n}{2}\right)} \left(1 + \frac{t^2}{n}\right)^{-\frac{n+1}{2}}, \quad \text{kun } x \in \Omega = \mathbb{R},$$

missä  $\Gamma(t)$  on Eulerin gammafunktio  $\Gamma(t) = \int_0^\infty e^{-x} x^{t-1} dx$

- tiheysfunktion kuvaaja on yksihiippuinen ja symmetrinen keskikohdan 0 suhteen
- muistuttaa normaali-jakaumaa  $N(0, 1)$  ja suurilla vapausasteilla  $n$  lähestyy ko. jakaumaa

# $t$ -jakaumaa noudattavia satunnaismuuttujia

- jos  $Z \sim N(0,1)$  ja  $W \sim \chi^2(n)$  ovat riippumattomia, niin

$$T = \frac{Z}{\sqrt{\frac{W}{n}}} \sim t(n)$$

## Lause

Olkoon  $X_1, X_2, \dots, X_n$  otos muuttujasta  $X \sim N(\mu, \sigma^2)$ . Jos  $S^2$  on muuttujan  $X$  otosvarianssi, niin

$$T = \frac{\bar{X} - \mu}{\frac{S}{\sqrt{n}}} \sim t(n-1).$$

# Odotusarvon luottamusväli, kun varianssi on tuntematon

## Lause

Jos  $X_1, X_2, \dots, X_n$  on otos muuttujasta  $X \sim N(\mu, \sigma^2)$ , jonka **varianssi on tuntematon**, niin odotusarvon  $100(1 - \alpha)\%$ :n väliestimaattori on

$$\left[ \bar{X} - t_{\alpha/2} \frac{S}{\sqrt{n}}, \bar{X} + t_{\alpha/2} \frac{S}{\sqrt{n}} \right],$$

missä  $t_{\alpha/2}$  toteuttaa ehdon  $P(T \leq t_{\alpha/2}) = 1 - \frac{\alpha}{2}$ , kun  $T \sim t(n-1)$ . Jos otoskeskiarvolle realisoituu arvo  $\bar{x}$  ja otoshajonnalle arvo  $s$ , niin odotusarvon  $100(1 - \alpha)\%$ :n luottamusväli on

$$\left[ \bar{x} - t_{\alpha/2} \frac{s}{\sqrt{n}}, \bar{x} + t_{\alpha/2} \frac{s}{\sqrt{n}} \right].$$

## Esimerkki

Tietyllä seudulla aikuisten miesten pituuden tutkimiseksi valittiin satunnaisesti  $n = 50$  miestä. Saatu otoskeskipituus oli  $\bar{x} = 178,15$  cm ja otoshajonta oli  $s = 5,92$  cm. Pituus oletetaan kyllin tarkasti normaalijakautuneeksi. Etsitään 98 % luottamusväli seudun miesten keskipituudelle.

Nyt siis  $\alpha = 0,02$  eli  $\alpha/2 = 0,01$ . Matlabilla komennolla `tinvt(1-0.01,50-1)` saadaan  $t_{\alpha/2} \approx 2,405$ . Nyt kaavalla  $\bar{X} \pm t_{\alpha/2} \frac{s}{\sqrt{n}}$  saadaan luottamusväliksi (176,1365; 180,1635).

# Huomioita odotusarvon luottamusvälistä

- kun otos hyvin suuri, on vapausasteluku  $n - 1$  myös suuri, joten jakauma lähellä standardinormaalijakaumaa
- tällöin otosvarianssi, joka on varianssille harhaton piste-estimaatti, tarkentuu kohti tuntematonta varianssia
- siksi kahdessa odotusarvon luottamusvälin kaavassa esiintyvät luvut  $z_{\alpha/2}$  ja  $t_{\alpha/2}$  lähellä toisiaan, kun otoskoko hyvin suuri
- esim. 95%:n luottamusväleille  $z_{\alpha/2} \approx 1.9600$  ja otoskoolla 100 laskettu  $t_{\alpha/2} \approx 1.9840$

# Normaalijakauman varianssin luottamusväli

## Lause

Jos  $X_1, X_2, \dots, X_n$  on otos muuttujasta  $X \sim N(\mu, \sigma^2)$ , niin varianssin  $100(1 - \alpha)\%$ :n väliestimaattori on

$$\left[ \frac{(n-1)S^2}{w_2}, \frac{(n-1)S^2}{w_1} \right],$$

missä luvut  $w_1$  ja  $w_2$  on valittu siten, että  $P(W < w_1) = \frac{\alpha}{2}$  ja  $P(W > w_2) = \frac{\alpha}{2}$ , kun  $W \sim \chi^2(n-1)$ . Jos varianssille realisoituu arvo  $s^2$ , niin varianssin  $100(1 - \alpha)\%$ :n luottamusväli on

$$\left[ \frac{(n-1)s^2}{w_2}, \frac{(n-1)s^2}{w_1} \right].$$

## Esimerkki

Pakkauskone pakkaa nestettä purkkeihin. Koneen toiminnan kannalta oleellista on purkkiin pakatun nestemäärän varianssi, jonka pitäisi olla noin  $100 \text{ g}^2$ . Nestemäärä per purkki voidaan olettaa normaalijakautuneeksi. 40 purkin otoksessa saatiin otosvarianssiksi  $s^2 = 151,09 \text{ g}^2$ . Lasketaan 90 prosentin luottamusväli varianssille.

Nyt siis  $\alpha = 0,1$  ja  $\alpha/2 = 0.05$ . Matlabilla saadaan  $\chi^2$ -jakaumasta vapausasteella  $40 - 1$  komennoilla  $\text{chi2inv}(0.05,39)$  ja  $\text{chi2inv}(0.95,39)$  luvuiksi  $w_1$  ja  $w_2$  likimain arvot 25,6954 ja 54,5722. Täten luottamusväliksi saadaan

$$\left[ \frac{(n-1)s^2}{w_2}, \frac{(n-1)s^2}{w_1} \right] \approx \left[ \frac{39 \cdot 151,09}{54,5722}, \frac{39 \cdot 151,09}{25,6954} \right]$$

eli likimain (107,98; 229,32). Pakkauskone voisi siis olla säätämisen tarpeessa.

# Suhteellinen osuus

- oletetaan, että  $X \sim \text{Bin}(n, p)$ , missä  $p$  on onnistumisen todennäköisyys
- todennäköisyyden frekvenssitulkinnassa  $p$  on myös onnistumisten suhteellinen frekvenssi eli niiden **suhteellinen osuus** kaikista toistoista
- tulkitaan  $X$  jakaumaa  $\text{Ber}(p) = \text{Bin}(1, p)$  noudattavien satunnaismuuttujien summana  $X = Y_1 + Y_2 + \dots + Y_n$
- tällöin satunnaismuuttujat  $Y_1, Y_2, \dots, Y_n$  muodostavat otoksen satunnaismuuttujasta  $Y \sim \text{Bin}(1, p)$
- monesti satunnaismuuttujan  $Y$  parametri  $p$  on tuntematon, joten sitä on estimoitava



# Suhteellisen osuuden estimointi

- edellisen sivun satunnaismuuttujan  $Y$  otoskeskiarvo on

$$\hat{P} \stackrel{\text{merk.}}{=} \bar{Y} = \frac{Y_1 + Y_2 + \dots + Y_n}{n} = \frac{X}{n}$$

- binomijakautunut  $X$  kuvaa onnistumisten frekvenssiä  $n$ -toistokokeessa, joten  $\hat{P} = \frac{X}{n}$  on suhteellinen osuus
- satunnaismuuttuja  $\hat{P}$  on harhaton estimaattori parametrille  $p$ , sillä otoskeskiarvon odotusarvo on

$$E(\hat{P}) = E(\bar{Y}) = E(Y) \stackrel{Y \sim \text{Bin}(1,p)}{=} 1 \cdot p = p$$

- seuraavan sivun luottamusväli on parametrille  $p$  hyvä estimointi, jos otoskoko on riittävän suuri

# Suhteellisen osuuden luottamusväli

## Lause

Oletetaan, että  $X \sim \text{Bin}(n, p)$ . Parametrin  $p$   $100(1 - \alpha)\%$ :n väliestimaattori on

$$\left[ \hat{p} - z_{\alpha/2} \sqrt{\frac{\hat{p}(1 - \hat{p})}{n}}, \hat{p} + z_{\alpha/2} \sqrt{\frac{\hat{p}(1 - \hat{p})}{n}} \right],$$

missä  $\Phi(z_{\alpha/2}) = 1 - \frac{\alpha}{2}$ . Jos suhteelliselle osuudelle  $\hat{P}$  realisoituu arvo  $\hat{p}$ , niin parametrin  $p$ , niin  $100(1 - \alpha)\%$ :n luottamusväli on

$$\left[ \hat{p} - z_{\alpha/2} \sqrt{\frac{\hat{p}(1 - \hat{p})}{n}}, \hat{p} + z_{\alpha/2} \sqrt{\frac{\hat{p}(1 - \hat{p})}{n}} \right].$$

## Esimerkki

Vuoden 2024 helmikuun alussa Helsingin Sanomat teetti gallupin koskien ihmisten äänestyskäyttäytymistä presidentinvaalien tulevalla toisella kierroksella. Gallupissa haasteltiin 1154 henkilöä ja tulosten mukaan noin 57 prosenttia vastaajista äänestäisi tulevalla vaalien toisella kierroksella Alexander Stubbia. Lasketaan 95 prosentin luottamusväli Stubbin äänestäjien osuudelle äänioikeutettujen joukossa.

Nyt siis  $\alpha/2 = 0,025$  ja arvoksi  $z_{\alpha/2}$  saadaan likimain 1,96. Luottamusväliksi saadaan kaavalla

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

likimain (0,5414; 0,5986). (Huom. Lopullinen Stubbin ääniosuus 51,6 % ei siis mahtunut tälle luottamusvälille.)