

# 11 Raja-arvolauseita ja approksimaatioita

- Tässä luvussa esitellään sellaisia kuuluisia todennäköisyysteorian raja-arvolauseita, joita sovelletaan usein tilastollisessa päättelyssä.
- Näiden raja-arvolauseiden tunteminen kuuluu jokaisen tilastotieteilijän yleissivistykseen.
- Päätulosten todistuksia ei voida käydä läpi tämän kurssin puitteissa.

# 11.1 Suurten lukujen laki

- Suurten lukujen laki (engl. *law of large numbers*) kertoo, että otoskeskiarvo suppenee kohti vastaavaa odotusarvoa, kun otoskoko kasvaa rajatta.
- Suurten lukujen lakeja on olemassa lukuisia sen mukaan, minkätyyppisiä objekteja tarkastellaan ja minkälaisia oletuksia niiden yhteisjakaumasta tehdään.
- Suuren lukujen laeista on olemassa sekä heikkoja versioita (sana heikko viittaa stokastiseen suppenemiseen) että vahvoja versioita (sana vahva viittaa melkein varmaan suppenemiseen).

# Stokastinen suppeneminen

## Määritelmä 11.1

Jono satunnaismuuttujia  $X_1, X_2, \dots$  **suppenee stokastisesti** eli **konvergoi stokastisesti** (engl. *converges in probability*) kohti satunnaismuuttujaa  $Y$ , jos

$$\mathbf{P}(|X_n - Y| \geq \epsilon) \xrightarrow{n \rightarrow \infty} 0,$$

kaikilla  $\epsilon > 0$ .

- Voitaisiin yhtäpitävästi vaatia, että  $\mathbf{P}(|X_n - Y| > \epsilon) \rightarrow 0$  kaikilla  $\epsilon > 0$ , tai  $\mathbf{P}(|X_n - Y| \leq \epsilon) \rightarrow 1$  kaikilla  $\epsilon > 0$ .
- Stokastista suppenemistä merkitään usein seuraavaan tapaan,

$$X_n \xrightarrow{\mathbf{P}} Y, \quad \text{tai} \quad \text{plim}_{n \rightarrow \infty} X_n = Y.$$

# Melkein varma suppeneminen

## Määritelmä 11.2

Jono satunnaismuuttujia  $X_1, X_2, \dots$  **suppenee (eli konvergoi) melkein varmasti** (engl. *converges almost surely*) kohti satunnaismuuttujaa  $Y$ , jos  $X_n(\omega) \rightarrow Y(\omega)$  perusjoukon osajoukossa, jonka tn on yksi, eli jos

$$\mathbf{P}(\lim_{n \rightarrow \infty} X_n = Y) = 1.$$

- Melkein varmaa suppenemista merkitään esim. seuraavilla tavoilla,

$$X_n \xrightarrow{\text{m.v.}} Y, \quad \text{tai} \quad X_n \xrightarrow{\text{a.s.}} Y.$$

- On mahdollista osoittaa, että melkein varmasta suppenemisestä seuraa stokastinen suppeneminen, eli että

$$X_n \xrightarrow{\text{a.s.}} Y \quad \Rightarrow \quad X_n \xrightarrow{\text{P}} Y \quad (11.1)$$

# Heikko suurten lukujen laki *i.i.d.*-jonolle

- Lyhenne **i.i.d.** on todennäköisyysteoriassa hyvin yleinen. Se tulee sanoista **independent, identically distributed** eli riippumattomat ja samoin jakautuneet.
- Todistimme jaksossa 6.1 Tšebyševin epäyhtälön avulla ns. heikon suurten lukujen lain (engl. *weak law of large numbers, WLLN*).
- Jos tätä heikkoa suurten lukujen lakia sovelletaan *i.i.d.*-jonoon satunnaismuuttujia  $X_1, X_2, \dots$ , niin se toteaa, että mikäli  $\sigma^2 = \text{var}(X_i)$  on äärellinen, niin  $n$  ensimmäisen satunnaismuuttujan aritmeettinen keskiarvo suppenee stokastisesti kohti vastaavaa odotusarvoa  $\mu = EX_i$ , eli

$$\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i \xrightarrow{P} \mu$$

# Vahva suurten lukujen laki *i.i.d.*-jonolle

## Lause 11.1 (Vahva suurten lukujen laki)

(Engl. *strong law of large numbers, SLLN.*) Olkoon  $X_1, X_2, \dots$  jono riippumattomia ja samoin jakautuneita satunnaismuuttujia, joiden odotusarvo  $\mu = \mathbf{E}X_1 \in \mathbb{R}$ . Tällöin keskiarvojen

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

muodostama jono suppenee melkein varmasti kohti arvoa  $\mu$ , eli  $\bar{X}_n \xrightarrow{\text{a.s.}} \mu$ .

- Oletuksena tarvitaan ainostaan se, että odotusarvon täytyy olla olemassa (reaalilukuna).
- Jos jono  $(\bar{X}_n)$  toteuttaa vahvan suurten lukujen lain, niin se tietenkin toteuttaa myös heikon suurten lukujen lain.

## 11.2 Jakaumasuppeneminen

### Määritelmä 11.3

Olkoon  $X_1, X_2, \dots$  jono satunnaismuuttujia, joiden kertymäfunktiot ovat  $F_1, F_2, \dots$ . Olkoon  $Y$  satunnaismuuttuja, jonka kertymäfunktio on  $G$ , ts. kaikilla  $x$

$$F_n(x) = \mathbf{P}(X_n \leq x), \quad G(x) = \mathbf{P}(Y \leq x).$$

Jono  $(X_n)$  **suppenee jakaumaltaan** (engl. *converges in distribution* tai *converges in law*) kohti  $Y$ :tä, mikäli

$$\lim_{n \rightarrow \infty} F_n(x) = G(x)$$

kaikissa rajajakauman kertymäfunktion  $G$  jatkuvuusasteissa  $x$ .

# Kommentteja jakaumasuppenemisestä

- Usein jakaumasuppenemisessä rajajakauma on jatkuva jakauma (esim.  $N(0, 1)$ ).
- Jatkuvan jakauman kertymäfunktio on jatkuva funktio koko reaaliakselilla. Jos rajajakauma on jatkuva, niin kertymäfunktioiden jonon pitää supeta jokaisessa reaaliakselin pisteessä,

$$\lim_{n \rightarrow \infty} F_n(x) = G(x), \quad \text{kaikilla } x \in \mathbb{R}.$$

- Merkitsemme jakaumasuppenemistä  $X_n \xrightarrow{d} Y$ . Nuolen yläindeksi  $d$  tulee sanasta *distribution*.
- Jos rajajakauma on jokin tuttu jakauma, kuten  $N(0, 1)$ , niin jakaumasuppenemistä voidaan merkitä siten, että satunnaismuuttujan sijasta käytetään rajajakauman tunnusta, esim.

$$X_n \xrightarrow{d} N(0, 1).$$



# Jakaumasuppenemisen hyödyntäminen

- Tilastotieteessä jakaumasuppenemista käytetään usein **jakaumien approksimointiin**.
- Jos tiedetään, että

$$X_n \xrightarrow{d} Y, \quad \text{kun } n \rightarrow \infty$$

niin jollakin äärellisellä indeksin  $n$  arvolla voidaan satunnaismuuttujan  $X_n$  jakaumaa approksimoida rajamuuttujan  $Y$  jakaumalla, mitä voidaan merkitä symbolisesti

$$X_n \stackrel{d}{\approx} Y.$$

- Edellisessä merkinnässä  $Y$ :n tilalla voidaan käyttää sen jakauman tunnusta.

# Esimerkki: asymptoottinen luottamusväli

- Jos rajamuuttujalla  $Y$  on jatkuva jakauma, niin suppenemisesta  $X_n \xrightarrow{d} Y$  seuraa esimerkiksi, että

$$\mathbf{P}(X_n \in I) \xrightarrow[n \rightarrow \infty]{} \mathbf{P}(Y \in I),$$

kun  $I \subset \mathbb{R}$  on mikä tahansa väli.

- Tämän ansiosta suurilla  $n$

$$\mathbf{P}(X_n \in I) \approx \mathbf{P}(Y \in I).$$

- Usein tilastollisessa päättelyssä joudutaan tarkkojen luottamusvälien sijasta soveltamaan tähän ajatukseen perustuvia asymptoottisia luottamusvälejä.
- Tämän takia jakaumasuppeneminen on tärkeä käsite (frekventistisessä) tilastotieteessä.

## 11.3 Keskeinen raja-arvolause

- Olkoon  $X_1, X_2, \dots$  *i.i.d.*-jono, ja  $\bar{X}_n$  olkoon  $n$  ensimmäisen jonon muuttujan keskiarvo. Merkitään  $\mu = \mathbf{E}X_1$  ja  $\sigma^2 = \text{var } X_1$ .
- Tiedämme suurten lukujen lain nojalla, että  $\bar{X}_n$  suppenee kohti odotusarvoa  $\mu$ .
- Lisäksi tiedämme, että

$$\mathbf{E}\bar{X}_n = \mu, \quad \text{var } \bar{X}_n = \frac{\sigma^2}{n}.$$

- Keskiarvon  $\bar{X}_n$  jakauma keskittyy yhä tiiviimmin arvon  $\mu$  ympärille, kun  $n$  kasvaa. Keskittymisvauhtia kuvaa tietyllä tavalla keskiarvon  $\bar{X}_n$  keskihajonta  $\sigma/\sqrt{n}$ , joka suppenee nolaa kohti vauhdilla  $1/\sqrt{n}$ .

# Minkä suureen jakauma suppenee keskeisessä raja-arvolauseessa

- Jos keskiarvo  $\bar{X}_n$  **standardoidaan** vähentämällä siitä keskiarvon odotusarvo ja jakamalla keskiarvon keskihajonnalla, saadaan lauseke

$$\frac{\bar{X}_n - \mathbf{E}\bar{X}_n}{\sqrt{\text{var } \bar{X}_n}} = \sqrt{n} \frac{\bar{X}_n - \mu}{\sigma} = \frac{1}{\sqrt{n}} \sum_{i=1}^n \frac{X_i - \mu}{\sigma}.$$

- Standardoidun keskiarvon odotusarvo on tietenkin nolla, ja sen varianssi on yksi kaikilla  $n$ .
- Keskeisen raja-arvolauseen mukaan standardoitu keskiarvo suppenee jakaumaltaan kohti standardinormaalijakaumaa.

# Keskeinen raja-arvolause

## Lause 11.2 (Keskeinen raja-arvolause)

(Engl. *central limit theorem, CLT*.) Olkoon  $X_1, X_2, \dots$  jono riippumattomia ja samoin jakautuneita satunnaismuuttujia siten, että  $0 < \sigma^2 < \infty$ , jossa  $\sigma^2 = \text{var } X_1$ . Merkitään

$$\mu = \mathbf{E}X_1, \quad \bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i.$$

Tällöin

$$\sqrt{n} \frac{\bar{X}_n - \mu}{\sigma} \xrightarrow{d} N(0, 1).$$

- Keskeisen raja-arvolauseen väite voitaisiin yhtä hyvin muotoilla siten, että

$$\sqrt{n}(\bar{X}_n - \mu) \xrightarrow{d} N(0, \sigma^2).$$

- Standardoinnissa on oleellista vain se, että keskiarvon varianssin  $\text{var}(\bar{X}_n)$  riippuvuus otoskoosta  $n$  jaetaan pois.
- Yo. versio autta ymmärtämään moniulotteista versiota keskeisestä raja-arvolauseesta.

# Moniulotteinen versio keskeisestä raja-arvolauseesta

## Lause 11.3 (Keskeinen raja-arvolause, moniulotteinen versio)

Jos  $X_1, X_2, \dots$  on *i.i.d.*-jono satunnaisvektoreita, joiden yhteinen odotusarvovektori on  $\mu$  ja kovarianssimatriisi on  $\Sigma$ , niin keskeinen raja-arvolause on voimassa muodossa

$$\sqrt{n}(\bar{X}_n - \mu) \xrightarrow{d} N(0, \Sigma), \quad (11.2)$$

jossa  $\bar{X}_n$  on  $n$  ensimmäisen satunnaisvektorin keskiarvo  $\frac{1}{n} \sum_{i=1}^n X_i$ .

## 11.4 Normaaliapproksimaatio

- Kun  $(X_i)$  on *i.i.d.*-jono, niin keskeiseen raja-arvolauseeseen avulla usein approksimoidaan **äärellisellä, kiinteällä**  $n$

$$\sqrt{n} \frac{\bar{X}_n - \mu}{\sigma} \stackrel{d}{\approx} N(0, 1).$$

- Tämä on sama asia kuin approksimaatio

$$\bar{X}_n \stackrel{d}{\approx} N\left(\mu, \frac{\sigma^2}{n}\right).$$

(Teeskentele, että edellä saatiin täsmällinen jakaumatulos, kerro vakiolla  $\sigma/\sqrt{n}$  ja lopuksi lisää  $\mu$ .)

- Tämä on sama asia kuin approksimaatio

$$\sum_{i=1}^n X_i \stackrel{d}{\approx} N(n\mu, n\sigma^2).$$



# Jakauma-approksimaatiot otoskeskiarvon ja summan jakaumille

Huomaa, miten approksimaatioiden

$$\bar{X}_n \stackrel{d}{\approx} N\left(\mu, \frac{\sigma^2}{n}\right), \quad \sum_{i=1}^n X_i \stackrel{d}{\approx} N(n\mu, n\sigma^2).$$

normaalijakauman parametrit saadaan approksimoitavan suureen odotusarvosta ja varianssista:

$$\begin{aligned} \mathbf{E}\bar{X}_n &= \mu, & \text{var } \bar{X}_n &= \frac{1}{n}\sigma^2 \\ \mathbf{E}\sum_1^n X_i &= n\mu, & \text{var } \sum_1^n X_i &= n\sigma^2. \end{aligned}$$

# Suuren otoskoon jakauma-approksimaatioita

- Keskeinen raja-arvolause takaa, että nämä normaaliapproksimaatiot (eli normaaliset approksimaatiot tai normaalijakauma-approksimaatiot) saadaan mielivaltaisen tarkoiksi, kun otoskoko  $n$  valitaan riittävän suureksi.
- **Milloin otoskoko  $n$  on riittävän suuri?**
- Tämä asia riippuu toisaalta halutusta tarkkuudesta ja toisaalta approksimaation käyttötarkoituksesta ja toisaalta satunnaismuuttujien  $X_i$  yhteisen jakauman luonteesta.
- Symmetrisille ja yksihuippuisille jatkuville jakaumille saavutetaan pienehköllä (muutaman kymmenen) otoskoolla useisiin tarpeisiin riittävä tarkkuus, mutta vinojen jakaumien kohdalla vastaavaan tarkkuuteen tarvittava otoskoko voi olla monta kertaluokkaa suurempi.

# Esimerkki: binomijakauman normaalin approksimaatio

- Olkoot  $X_1, X_2, \dots, X_n$  riippumattomia Bernoulli( $p$ )-jakaumaa noudattavia sm:ia, jossa  $0 < p < 1$ . Tällöin

$$\mathbf{E}\bar{X}_n = p, \quad \text{var } \bar{X}_n = \frac{1}{n} p(1 - p).$$

- Normaaliapproksimaatiolla saadaan

$$\bar{X}_n \stackrel{d}{\approx} N\left(p, \frac{1}{n} p(1 - p)\right), \quad \text{ja} \quad \sum_{i=1}^n X_i \stackrel{d}{\approx} N(np, np(1 - p)).$$

- Tässä tapauksessa pätee tarkka tulos

$$\sum_{i=1}^n X_i \sim \text{Bin}(n, p).$$

# Numeerinen esimerkki binomijakauman approksimoinnista

- Esimerkiksi, kun  $n = 25$  ja  $p = 0.6$ , niin  $np = 15$  ja  $np(1 - p) = 6$ .
- Kun  $Z \sim N(0, 1)$  ja  $\Phi$  on standardinormaalijakauman kertymäfunktio, niin saadaan approksimaatio

$$\begin{aligned} \mathbf{P}\left(\sum_{1}^n X_i \leq 13\right) &\approx \mathbf{P}(15 + \sqrt{6} Z \leq 13) = P\left(Z \leq \frac{13 - 15}{\sqrt{6}}\right) \\ &= \Phi\left(\frac{13 - 15}{\sqrt{6}}\right) \approx 0.207. \end{aligned}$$

- Kun käytetään binomijakaumaa, eikä sen approksimaatiota, saadaan tulos

$$\mathbf{P}\left(\sum_{1}^n X_i \leq 13\right) \approx 0.268.$$

# Approksimaation parantaminen jatkuvuuskorjauksella

- Näin pienellä otoskoolla normaalin approksimaatio ei ole kovin tarkka, mutta sitä voidaan parantaa käyttämällä ns. **jatkuvuuskorjausta**.
- Koska summa on kokonaislukuarvoinen, niin seuraavat tapahtumat ovat samoja,

$$\left\{ \sum_1^n X_i \leq 13 \right\} = \left\{ \sum_1^n X_i \leq 13 + 0.5 \right\}.$$

- Kun normaaliapproksimaatiossa yläraja 13 korvataan luvulla 13.5, niin saadaan jo alkuperäistä approksimaatiota paljon tarkempi approksimaatio

$$\mathbf{P}\left(\sum_1^n X_i \leq 13\right) \approx \mathbf{P}\left(Z \leq \frac{13.5 - 15}{\sqrt{6}}\right) \approx 0.270.$$

# Kokonaislukuarvoisten satunnaismuuttujien summa

- Jos satunnaismuuttujat  $X_i$  ovat kokonaislukuarvoisia, niin myös niiden summa  $S = \sum_{i=1}^n X_i$  on kokonaislukuarvoinen.
- Sen yhteydessä on mielekästä tarkastella häntätodennäköisyyksiä vain kokonaislukuarvoisissa pisteissä.
- Jos  $x$  ja  $y$  ovat kokonaislukuja, niin

$$\{S \geq x\} = \{S \geq x - \frac{1}{2}\}$$

$$\{S \leq y\} = \{S \leq y + \frac{1}{2}\}$$

$$\{x \leq S \leq y\} = \{x - \frac{1}{2} \leq S \leq y + \frac{1}{2}\}$$

- Kun kokonaislukuarvoisten satunnaismuuttujien summan **diskreettiä** jakaumaa approksimoidaan **normaalijakaumalla**, niin saavutetaan parempi tarkkuus, mikäli tällaiset kokonaislukuarvoisen satunnaismuuttujan kokonaislukuarvoiset rajat korvataan tähän tapaan kokonaislukujen puoliväleissä olevilla arvoilla.
- Tämä idea on nimeltään jatkuvuuskorjaus (engl. *continuity correction*).
- Jos approksimoidaan otoskeskiarvon jakaumaa, niin diskreetin jakauman arvojoukkoon (hilapisteet) sattuvat pisteet vastaavasti korvataan hilapisteiden puolivälissä olevilla arvoilla.

# Jatkuvuuskorjaus *i.i.d.*-muuttujien summalle

- $S = \sum_1^n X_i$  ja  $X_i$ :t ovat kokonaislukuarvoisia.
- Jatkuvuuskorjaus johtaa summalle  $S$  approksimaatioihin

$$\mathbf{P}(S \geq x) = P(S \geq x - \frac{1}{2}) \approx 1 - \Phi\left(\frac{x - \frac{1}{2} - n\mu}{\sqrt{n}\sigma}\right)$$

$$\mathbf{P}(S \leq y) = P(S \leq y + \frac{1}{2}) \approx \Phi\left(\frac{y + \frac{1}{2} - n\mu}{\sqrt{n}\sigma}\right)$$

$$\mathbf{P}(x \leq S \leq y) = P(x - \frac{1}{2} \leq S \leq y + \frac{1}{2}) \approx \Phi\left(\frac{y + \frac{1}{2} - n\mu}{\sqrt{n}\sigma}\right) - \Phi\left(\frac{x - \frac{1}{2} - n\mu}{\sqrt{n}\sigma}\right)$$



## 11.5 Deltamenetelmä

Deltamenetelmän (engl. *delta method*) avulla voidaan normaaliapproksimaatiota käyttää myös silleille funktioille otoskeskiarvoista.

### Lause 11.4 (Deltamenetelmä)

Jos

$$\sqrt{n}(X_n - a) \xrightarrow{d} N(0, \sigma^2),$$

jossa  $a$  on vakio, ja funktio  $g$  on derivoituva pisteessä  $a$ , niin

$$\sqrt{n}(g(X_n) - g(a)) \xrightarrow{d} N(0, \sigma^2(g'(a))^2). \quad (11.3)$$

# Deltamenetelmän heuristinen perustelu

- Oletuksesta  $\sqrt{n}(X_n - a) \xrightarrow{d} N(0, \sigma^2)$  seuraa, että suurella  $n$  sm:n  $X_n$  jakauma on voimakkaasti keskittynyt arvon  $a$  lähelle, sillä

$$X_n - a = \frac{1}{\sqrt{n}} \sqrt{n}(X_n - a).$$

- Tässä jono  $\frac{1}{\sqrt{n}}$  suppenee kohti nollaa ja jono  $\sqrt{n}(X_n - a)$  on siinä mielessä stabiili, että sillä on rajajakauma.
- Tämän takia suurilla  $n$  tuntuu järkevältä approksimoida satunnaismuuttujaa  $g(X_n)$  pisteessä  $a$  kehitetyllä ensimmäisen asteen Taylorin polynomilla, josta jäännöstermi voidaan (toivon mukaan) unohtaa.

# Deltamenetelmän heuristinen perustelu päättyy

- Tehdään ensimmäisen asteen Taylorin approksimaatio

$$g(X_n) \approx g(a) + g'(a)(X_n - a),$$

jossa jäännöstermi jätettiin pois.

- Tämän jälkeen

$$\begin{aligned}\sqrt{n}(g(X_n) - g(a)) &\approx g'(a) \sqrt{n}(X_n - a) \\ &\stackrel{d}{\approx} g'(a) N(0, \sigma^2) \\ &\stackrel{d}{=} N(0, \sigma^2 (g'(a))^2).\end{aligned}$$

- Deltamenetelmää sovelletaan tavallisesti siten, että äärellisellä, kiinteällä otoskoolla  $n$  approksimoidaan

$$\sqrt{n}(g(X_n) - g(a)) \stackrel{d}{\approx} N(0, \sigma^2(g'(a))^2).$$

## Esimerkki: jakauma-approksimaatio *log-odds*-suurelle

- Tarkastellaan riippumattomia muuttujia  $X_i \sim \text{Bernoulli}(p)$ ,  $i = 1, \dots, n$ , jossa  $0 < p < 1$ .
- Parametrin  $p$  SU-estimaatti on  $\bar{X}_n = s/n$ , jossa  $s$  on onnistumisten lukumäärä.
- Oletetaan, että tahdotaan estimoida ns. *log-odds* -suuretta (vedonlyöntisuhteen logaritmia),

$$\theta = \ln \frac{p}{1-p}.$$

Sen luonteva estimaatti on

$$\hat{\theta}_n = \ln \frac{\bar{X}_n}{1 - \bar{X}_n} = \ln \frac{s}{n - s},$$

jonka jakauma ei ole mikään tunnettu jakauma.

# Jakauma-approksimaatio *log-odds*-suureelle – esimerkki päätty

- Valitaan deltamenetelmässä

$$g(u) = \ln \frac{u}{1-u}, \quad 0 < u < 1,$$

- Koska  $\sqrt{n}(\bar{X}_n - p) \xrightarrow{d} N(0, p(1-p))$ , niin helppojen laskujen avulla saadaan approksimaatio

$$\sqrt{n}(\hat{\theta}_n - \theta) \approx N\left(0, \frac{1}{p(1-p)}\right).$$

- Tällä perusteella olisi esim. mahdollista johtaa asymptoottinen luottamusväli parametrille  $\theta$ .