

1. JOHDATTELEVIA ESIMERKKEJÄ

1.) Pallot kulhossa; vrt. [AS, jakso 1.1]

AS =
Arjas - Sirén
moniste

Kulhossa N palloa: K keltaista ja
 $N-K$ valkoista. Tark. satunnaiskoe

(*) $\left\{ \begin{array}{l} \text{nostetaan uupimähkään ja palauttaen (sekä} \\ \text{välillä hyvin sekoittaen) } n \text{ palloa} \end{array} \right.$

Pal. tn-laskennasta mieleen, miten tätä voidaan
mallintaa tn-laskennan avulla: Määr. sm:t

$$X_i = \begin{cases} 1, & \text{jos } i\text{:s pallo keltainen} \\ 0, & \text{jos } i\text{:s pallo valkoinen} \end{cases} \quad (i=1, \dots, n)$$

$\Rightarrow X_1, \dots, X_n \perp$ (riippumattomia) ja kullakin
 X_i :llä on pistetodennäköisyydet (lyh. ptn)

$$\begin{cases} P(X_i = 1) = K/N = \theta \\ P(X_i = 0) = 1 - K/N = 1 - \theta \end{cases} \quad (1)$$

jossa $\theta = K/N$ keltaisten suht. osuus kulhossa.

Merh. $\underline{X} = (X_1, \dots, X_n)$ ("n-ulott. satunnaisvektori")

Kun $\underline{x} = (x_1, \dots, x_n)$ jossa kukin $x_i = 0$ tai 1 ,
saadaan riippumattomuuden nojalla (yhters) ptn: t

$$\begin{aligned} P(\underline{X} = \underline{x}) &= P(X_1 = x_1, \dots, X_n = x_n) \\ &= \prod_{i=1}^n P(X_i = x_i) \stackrel{(1)}{=} \theta^t (1-\theta)^{n-t} \end{aligned} \quad (2)$$

jossa

$$\begin{cases} t = x_1 + \dots + x_n = \text{keltaisten lkm "otoksessa"} \\ n-t = \text{valkoisten lkm "otoksessa"} \end{cases}$$

Jos kiinnitetään huomiota vain keltaisten lkm:ään otoksessa eli sm:aan $T = X_1 + \dots + X_n$, tiedetään tulolaskennasta, että $T \sim \text{Bin}(n, \theta)$ eli T :llä on ptn:t

$$P(T=t) = \binom{n}{t} \theta^t (1-\theta)^{n-t}, \quad t=0, 1, \dots, n \quad (3)$$

Siispiä: Edellyttäen, että tunnemme keltaisten suht. osuuden θ kulhossa, kaavat (2) ja (3) kertovat täydellisesti kokeen (*) eri tulosvaihtoehtojen tn:t (riippuen siitä haluammeko pitää kokeen "tuloksena" täydellisistä tulospöytäkirjaa $\underline{x} = (x_1, \dots, x_n)$ vai arvoastaan keltaisten lkm:ää t).

Tilastollinen päättely tutkii käänterstä ongelmaa: Oletetaan, että θ on tuntematon ja suoritetaan koe (*). Havaitaan

- joko $\underline{x} = (x_1, \dots, x_n)$ eli nostettujen pallojen värit järjestyksessä
- tai vain keltaisten lkm $t = x_1 + \dots + x_n$.

Mitä voidaan päätellä θ -sta tämän perusteella?

Esim. Suoritetaan $n=100$ nostoa ja havaitaan $t=35$ keltaista. Mitä johtopäätöksiä θ -sta?

(2) Kliininen koe. Tautiin on kehitetty lupaava lääke, jota halutaan testata. Oletetaan:

- ilman lääkettä taudista paranee tn:llä θ_0
- lääkkeellä hoidettuna siitä paranee tn:llä θ_1

Jaetaan 100 potilasta satunnaisesti kahteen 50 hengen ryhmään. Torselle annetaan plaseboa, torselle lääkettä.

Seurantajakson kuluttua havaitaan:

- plaseboryhmässä parani x_0 potilasta
- lääkeeryhmässä parani x_1 potilasta

Mitä voidaan päätellä luvuista θ_0 ja θ_1 ?

Eristyisesti: voidaanko päätellä, että $\theta_1 > \theta_0$ (kuten tutkija luultavasti toivoi!)?

3. Metron odotusajat. Oletetaan, että metro kulkee säännöllisesti θ minuutin väliajoin. Umpimähkään laiturille saapuvan henkilön odotusajan voidaan siis ajatella olevan $Tas(0, \theta)$ -jakaunut.

Maalta muuttanut opiskelija ei tunne θ :oa, mutta joka kerta laiturille saavuttuaan hän mittaa odotusajansa, saaden siten mittauksia t_1, \dots, t_n (minuuttia). Mitä hän voi päätellä θ :sta?

Esim. Kolmen kerran tulokset ovat $t_1 = 3.4$, $t_2 = 1.1$ ja $t_3 = 4.5$.

- Täysin varmasti voi päätellä, että $\theta \geq 4.5$.
- Kuinka luotettavasti voisi päätellä, että $\theta \leq 10$?

4. Appelssinikauppa. Hedelmäkauppias väittää, että hänen myymiensä appelsiinien keskipaino μ on ainakin 250 g. Asiakas ostaa umpimähkään 10 appelsinia ja havaitsee kotona, että niiden keskipaino on 236 g.

Onko tämän perusteella syytä asettaa kyseenalaiseksi kauppiaan väite " $\mu \geq 250$ " vai käviikö voim "huono tuuri"?

Voiko asiaan edes ottaa kantaa ilman lisätietoja/-oletuksia?

Pohdittavaksi yo. "tutkimusasetelma" liittyen:

- * Mitkä olivat kokeellisia, mitkä puhtaasti havainto-tutkimuksia? Mitkä olivat otanta tutkimuksia?
- * Olivatko havaintoja vastaavat muuttujat diskreettejä vai jatkuvia?
- * Miten todennäköisyyden käsite on tutkittava? Klassisesti (symmetriset alkeistapaukset)? Frekventistisesti? Jotain muuten?

2. TILASTOLLISEN PÄÄTTELYN YLEINEN ASETELMA JA TAVOITTEET

Tarkasteltavana aineisto $\underline{x} = (x_1, \dots, x_n)$, joka koostuu havainnoista x_1, \dots, x_n . Nihän liittyy satunnaisuutta - johtuen esim. satunnaisotannasta tai tutkittavan ilmiön itsensä satunnaisesta luonteesta

⇒ ajattelemme että ne ovat eräiden satunnaismuuttujien

x_1, \dots, x_n ; vektorina $\underline{x} = (x_1, \dots, x_n)$

toteutuneita arvoja eli realisatioita.

Näiden sm:ien (yhters)jakauma $f(\underline{x}) = f_{\underline{x}}(\underline{x})$ (käytännössä: yhteispistetyn- tai yhteistiheysfunktio) on (ainakin osittain) tuntematon ja tavoitteena on tehdä siitä päätelmiä aineiston \underline{x} perusteella.

Useimmiten (ja erit. tällä kurssilla) oletetaan, että jakauma on tunnettu jotakin parametria θ vaille, ts. \underline{x} :n yptnf/ytf on muotoa

$$f(\underline{x}; \theta) = f_{\underline{x}}(\underline{x}; \theta),$$

jossa θ on tuntematon.

Määr. $f(\underline{x}; \theta)$ on (parametrinen) tilastollinen malli.

Tässä parametri θ voi olla reaaliarvoinen (1-ulott.) tai koostua useammasta komponentista: $\theta = (\theta_1, \dots, \theta_k)$

Yleinen tavoite: tehdä johtopäätöksiä θ :sta havaitun aineiston \underline{x} perusteella.

Erityiskysymyksiä:

1. Piste-estimointi: määritettävä piste (luku tai vektori) kaikkien mahdollisten parametrin arvojen joukosta, joka on "hyvä arvaus" θ :n todelliseksi "oikealle" arvolle.
2. Väliestimointi: rajattava sellainen parametrin arvojen joukko (1-ulott. tapauksessa yleensä väli), joka suurella varmuudella sisältää "oikean" θ :n arvon.
3. Hypoteesintestaus: onko aineisto sopuisuudessa annetun hypoteesin " $\theta \in H$ " kanssa vai ei?

Muuta til. päättelyn piiriin kuuluvaa:

4. Ennustaminen: Aineiston (x_1, \dots, x_n) perusteella ennustettava samasta satunnaisilmiöstä/kokeesta tulevaa uutta/seuraavaa havaintoa x_{n+1} .
5. Malliarviointi (diagnostiikka): Onko malli $f(\underline{x}; \theta)$ sopiva ja oikea selittämään saadut havainnot? Totentuvatho mallin oletukset?

Huom. Johainen malli on aina todellisuuden yksinkertaistus ja approksimaatio!

Tällä kurssilla lähinnä aiheet 1.-3.

Käytännön sovelluksissa tehtävä 5. hyvin tärkeä!

Esim. "Pallot kulkossa" -esimerkissä edellä tilast. malli
smille 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 θ , $0 \leq \theta \leq 1$ (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
- θ on kiinteä mutta tuntematon luku tai p:ste
- θ :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