

Biostatisztika és informatika

5. előadás: Becslés és megbízhatóság

2018. október 11.

Agócs Gergely

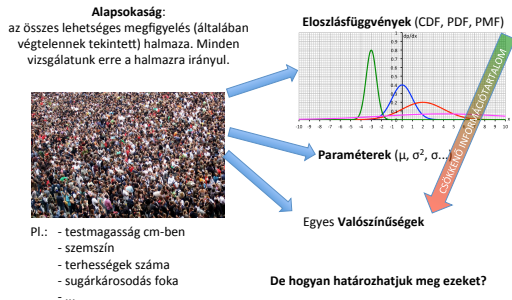
Források: – Herényi L (2016): **Statisztika és Informatika: 14. fejezet**
– Reiczigél J, Harnos A, Solymosi N (2014): Biostatisztika nem statisztikusoknak: 5. fejezet
– WolframMathWorld: Probability and Statistics:
<http://mathworld.wolfram.com/topics/ProbabilityandStatistics.html>

Az előadás céljai

- A **becslés** célja, típusai és folyamata
 - minta és alapsokaság más néven populáció
 - **pontbecslés** és **intervallumbecslés**
 - minta, becslés, becslült és becslő érték
- Mik a „jó becslés” tulajdonságai?
 - torzítatlan, hatásos, konzisztens, elégséges
- **Standard hiba (SE)** megértése
- **Konfidenciaintervallum (CI)** megértése
 - konfidenciaintervallumok helyes értelmezése
 - konfidenciaintervallumok helyes feltüntetése
- Megtanulni egyes paraméterek **becslését**:
 - valószínűség (rövid.: **valség**) más néven alapsokaságbeli arány
 - **várható érték** más néven alapsokaságbeli átlag
 - elméleti **variancia** (más néven szórásnégyzet) és elméleti **szórás**
- Adott SE-hoz **szükséges mintaelemszám** számítása
- **Excel**-függvények használata becsléshez

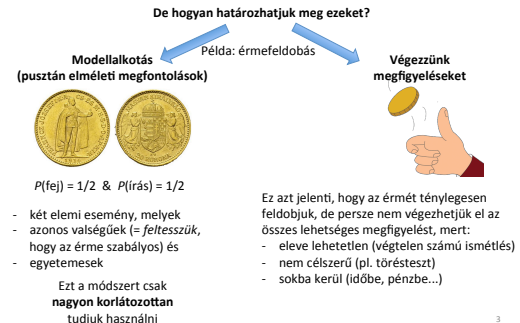
1

A becslés célja



2

A becslés célja



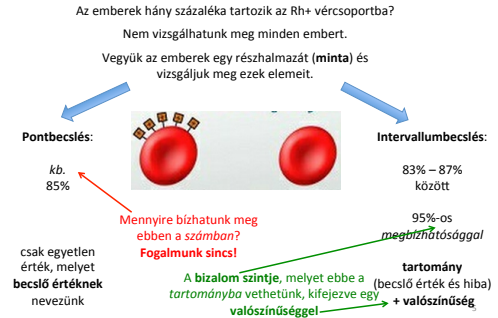
3

A becslés folyamata



4

Becslések típusai

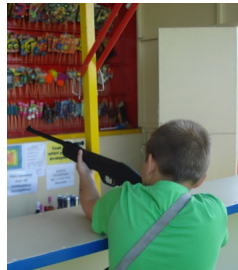


Pontbecslés

Elméleti (alapsokaság) értékek:
„**CÉL**”



becslő é.:
„**LÖVÉS**”



6

Pontbecslés

Elméleti (alapsokaság) értékek:
„**CÉL**”

- valószínűség vagy arány (p_i)
- várható érték („alapsokaság átlaga”) ($E(\xi)$ vagy μ)
- elméleti variancia ($Var(\xi)$ vagy σ^2)
- elméleti szórás ($SD(\xi)$ vagy σ)

$$SD(\xi) = \sigma = \sqrt{Var(\xi)} = \sqrt{\sum_{i=1}^n p_i \cdot (x_i - \mu)^2}$$

becslő é.:
„**LÖVÉS**”



7

Pontbecslés

Elméleti (alapsokaság) értékek:

(diszkrét valószínűségi változók)

„CEL”

- valószínűség vagy arány (p_i)

- várható érték („alapsokaság átlaga”) ($E(\xi)$ vagy μ)

$$E(\xi) = \mu = \sum_{i=1}^n p_i \cdot x_i$$

- elméleti variancia ($Var(\xi)$ vagy σ^2)

$$Var(\xi) = \sigma^2 = E[(\xi - E(\xi))^2] = \sum_{i=1}^n p_i \cdot (x_i - \mu)^2$$

- elméleti szórás ($SD(\xi)$ vagy σ)

$$SD(\xi) = \sigma = \sqrt{Var(\xi)} = \sqrt{\sum_{i=1}^n p_i \cdot (x_i - \mu)^2}$$

Behelyettesítéses („plug-in”) becslő é.:

„LÖVÉS”

- relatív gyakoriság $\hat{p}_i = (k_i/n)$ ✓

Excel: =DARABHATÓBB(adat)/DARAB2(adat)

- minta átlaga

Excel: =ÁTLAG(adat)

$$\bar{x} = \sum_{i=1}^n \frac{k_i}{n} \cdot x_i = \frac{1}{n} \sum_{i=1}^n k_i \cdot x_i \quad \checkmark$$

- „plug-in” variancia (s^{**2})

$$s^{**2} = \sum_{i=1}^n \frac{k_i}{n} \cdot (x_i - \mu)^2 = \frac{1}{n} \sum_{i=1}^n k_i \cdot (x_i - \mu)^2 \quad \times$$

- „plug-in” szórás (s^{**} , SD)

$$s^{**} = \sqrt{\frac{1}{n} \sum_{i=1}^n k_i \cdot (x_i - \mu)^2} \quad \times$$

8

A „jó becslés” ...

... torzítatlan

A valószínűség (p_i) pontbecslése relatív gyakorisággal (\hat{p}_i)

Tekintsük az imént tanult „plug in” becslő képletét:

$\hat{p}_i = (k_i/n)$ ✓ E becslés végtelen számú megismétléssel kapott becslő értékek átlaga maga a becslendő érték (valóság).

Excel: =DARABHATÓBB(adat)/DARAB2(adat)

A várható érték (μ) pontbecslése a mintaátlaggal (\bar{x})

$$\bar{x} = \sum_{i=1}^n \frac{k_i}{n} \cdot x_i = \frac{1}{n} \sum_{i=1}^n k_i \cdot x_i \quad \checkmark$$

Excel: =ÁTLAG(adat)

Egy becslés torzítatlan, ha a végtelenszer megismételt becslésekkel kapott becslő értékek várható értéke megegyezik a becslendő értékkel.

9

A „jó becslés” ...

... torzítatlan

Az elméleti variancia (σ^2) pontbecslése

Tekintsük az imént tanult „plug-in” becslő (s^{**2}) képletét:

$$s^{**2} = \frac{1}{n} \sum_{i=1}^n k_i \cdot (x_i - \mu)^2 \quad \times \quad \mu \text{ egy elméleti érték, melyet nem ismerünk. Helyette a mintaátlagot kell használnunk.}$$

Helyettesítsük be a várható érték helyére a mintaátlagot:

$$s^2 = \frac{1}{n} \sum_{i=1}^n k_i \cdot (x_i - \bar{x})^2 \quad \times \quad \text{Ez a képlet a mintaátlagra minimális, nem a várható értékre. Ez torzítást okoz, a képlet alábecsli az elméleti szórást.}$$

Az $n/(n-1)$ korrekciós szorzóval (Bessel's korrekciónak hívják)

megszüntethető a torzítás:

(bizonyítás a jegyzetben, nem kell tudni)

$$s^2 = \frac{n}{n-1} \cdot \frac{1}{n} \sum_{i=1}^n k_i \cdot (x_i - \bar{x})^2 = \frac{1}{n-1} \sum_{i=1}^n k_i \cdot (x_i - \bar{x})^2 \quad \checkmark$$

Excel: =VAR.M(adat)



Friedrich Bessel
1784–1846



10

A „jó becslés” ...

... kevésbé torzít

Az elméleti szórás (σ) pontbecslése

A variancia becslő értékéhez hasonlóan a szórás „plug in”-képlete is torzít:

$$s^* = \sqrt{\frac{1}{n} \sum_{i=1}^n k_i \cdot (x_i - \bar{x})^2} \quad \times$$

Itt a Bessel-korrekciós faktor $n/(n-1)$ használata csak csökkenti de teljesen nem szünteti meg a torzítást:

(csak aszimptotikusan, vagyis végtelen nagy minta esetén; ok: a gyökfüggvény aszimmetrikussága)

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

Excel: =SZÖR.M(adat)

Ezt a kevésbé torzító (de nem teljesen torzítatlan) becslést használjuk.

11

Pontbecslés

Elméleti (alapsokaság) értékek:

„CEL”

- valószínűség vagy arány (p_i)
- várható érték („alapsokaság átlaga”) ($E(\xi)$ vagy μ)

$$E(\xi) = \mu = \sum_{i=1}^n p_i \cdot x_i$$

- elméleti variancia ($Var(\xi)$ vagy σ^2)

$$Var(\xi) = \sigma^2 = E\left[(\xi - E(\xi))^2\right] = \sum_{i=1}^n p_i \cdot (x_i - \mu)^2$$

- elméleti szórás ($SD(\xi)$ vagy σ)

$$SD(\xi) = \sigma = \sqrt{Var(\xi)} = \sqrt{\sum_{i=1}^n p_i \cdot (x_i - \mu)^2}$$

Behelyettesítés („plug-in”) becslő é.: „LÖVÉS”

- relatív gyakoriság $\hat{p}_i = (k_i/n)$ ✓

Excel: =DARABHATOBBI(adat)/DARAB2(adat)

- minta átlaga Excel: =ÁTLAG(adat)

$$\bar{x} = \sum_{i=1}^n \frac{k_i}{n} \cdot x_i = \frac{1}{n} \sum_{i=1}^n k_i \cdot x_i \quad \checkmark$$

- tapasztalati variancia (s^2)

Excel: =VAR.M(adat)

$$s^2 = \frac{1}{n-1} \sum_{i=1}^n k_i \cdot (x_i - \bar{x})^2 \quad \checkmark$$

- tapasztalati szórás (s , SD)

Excel: =SZÖR.M(adat)

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

12

Pontbecslés

1. feladat: Meg szeretnénk becsülni a kékszeműek arányát, valamint a testmagasság várható értékét, elméleti varianciáját és elméleti szórását a SOTÉ-s golyók (elsőéves hallgatók) körében.

Egy 15-elemű mintát vettünk (lásd a táblázatot). Használjuk az Excel-t az elméleti paraméterek pontbecslésére!

$$p(\text{kék szem}) = \hat{p} = 4/15 = 0.2667 = 26,67\%$$

$$\mu(\text{testmagasság}) = \bar{x} = \text{ÁTLAG}(\text{adat}) = 170 \text{ cm}$$

$$\sigma^2(\text{testmagasság}) = s^2 = \text{VAR.M}(\text{adat}) = 107,5 \text{ cm}^2$$

$$\sigma(\text{testmagasság}) = s = \text{SZÖR.M}(\text{adat}) = 10,4 \text{ cm}$$

Megfigyelt szorzat	Személy szám	Magasság (cm)
1	HAMS	163
2	HAMS	163
3	HAMS	152
4	HAMS	158
5	HAMS	167
6	GAZ	184
7	GAZ	165
8	HAMS	184
9	GAZ	175
10	HAMS	167
11	HAMS	178
12	HAMS	168
13	HAMS	172
14	GAZ	178
15	HAMS	180

13

A „jó becslés” ...

... hatások

Egy becslés hatása, ha szórása (az ún. **standard hiba**, **SE**) minimális.

- A megismertelt becslések egy sor becslő értéket eredményeznek, amelyek különböznek egymástól a **mintavétel véletlenszerűsége** miatt (**véletlen hiba**)
- Tehát a **becslő érték maga is egy véletlen változó**, melynek van elméleti eloszlása, várható értéke, elméleti szórása stb.
- A becslő érték elméleti szórását standard hibának (**standard error**, **SE**) nevezzük

A következő diákon megtanuljuk hogyan kell kiszámolni az alábbi két becslő érték standard hibáját.

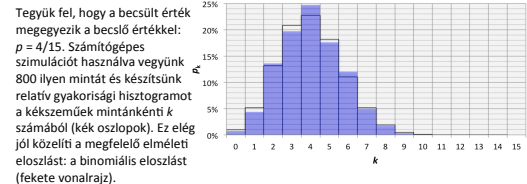
Arányok elméleti eloszlása:
Binomiális eloszlás

Átlagok elméleti eloszlása:
Student-féle t-eloszlás ($df = n - 1$)

14

Arány standard hibája

1. példa: Meg szeretnénk becsülni a kékszemű hallgatók arányát a golyók mint alapsokaság között. Egy $n = 15$ elemű mintát veszünk, melyben a kékszemek száma $k = 4$ -nek adódik. Már megtanultuk, hogy az alapsokaságon belüli arány (azaz valószínűség, p) becslő értéke a relatív gyakoriság ($\hat{p} = k/n$), mely esetünkben $4/15 = 0,2667$. Mekkora a becslés hibája?



15

Arány standard hibája

A korábbi előadásokban már megtanult képletekkel számoljuk ki ezen binomiális eloszlás paramétereit (μ , σ^2 és σ).

Várható érték:

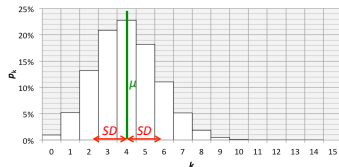
$$\mu = np = 4$$

Elméleti variancia:

$$\sigma^2 = np(1-p) = 44/15 = 2,933$$

Elméleti szórás:

$$SD = \sqrt{np(1-p)} = 1,713$$



16

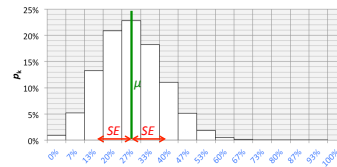
Arány standard hibája

Mivel arányokat becslünk, ezért a k változót k/n aránnyá kell alakítanunk, vagyis a vízszintes tengelyt át kell skálázunk n -nel való osztással.

Ugyanezt kell tennünk a binomiális eloszlás μ és σ paramétereivel is.

Várható érték

$$\mu = p = 4/15$$

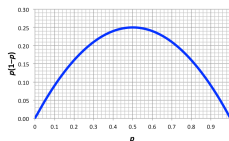


Elméleti szórás = az arány standard hibája

$$SE = \frac{\sqrt{np(1-p)}}{n} = \sqrt{\frac{p(1-p)}{n}} = 0,1142$$

$$SE_{arany} = \sqrt{\frac{p(1-p)}{n}} = \sqrt{\frac{k \left(1 - \frac{k}{n}\right)}{n}} \quad 17$$

Arány standard hibája



Theoretical Distribution of Proportions:
Binomial Distribution
(normalized to sample size)

$$\max(SE_{arany}) = \frac{1}{\sqrt{4n}}$$

Mekkora az n -elemű mintából számolható arányok lehető legmagasabb standard hibája?

A $p(1-p)$ szorzat akkor maximális, ha $p = 1-p = 0,5$. Ez esetben $p(1-p) = 0,25$ így a SE :

$$\max(SE_{arany}) = \sqrt{\frac{0,5(1-0,5)}{n}} = \sqrt{\frac{0,25}{n}} = \frac{1}{\sqrt{4n}}$$

Ha egy tanulmány ugyanazon mintából (vagy azonos elemszámú mintákból) becsült több valóságot is tartalmaz, akkor gyakran csak a maximális standard hibát adják meg a fenti képlet segítségével.

Pl. egy $n = 100$ -elemű mintából számolt arány esetén a SE legfeljebb 0,05 lehet.

18

Arány standard hibája

2. feladat: Mi a sárlósejtes anémia (SC) génhordozóinak prevalenciája (népességen belüli aránya) és a becslés standard hibája Nigériában, ha 172 vizsgált emberből 43 hordozta a betegség génjét?

Hordozók mintabeli aránya: $\hat{p}(SC) = k(SC)/n = 43/172 = 0,25$. A hiba az imént tanult képlettel

$$\text{számolható: } SE_{arany} = \sqrt{\frac{\hat{p}(1-\hat{p})}{n}} = \sqrt{\frac{\hat{p}(1-\hat{p})}{n}} = \sqrt{\frac{0,25 \cdot (1-0,25)}{172}} = 0,033 = 3,3\%$$

3. feladat: Szeretnénk egy vizsgálatban megbecsülni egyes krónikus betegségek budapesti prevalenciáját. Milyen mintaméretet javasolnál, ha nem akarjuk, hogy a becslés hibája az 1%-ot meghaladjon?

Mivel konkrét valóságot vagy arányt nem ismerünk, a „legrosszabb esetnek” tekinthető 0,5-ös valósággal kell számolnunk, melynek adott mintaméret esetén a legnagyobb a standard hiba:

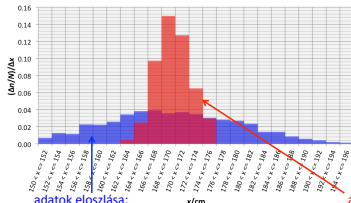
$$\text{A képletet átrendezve kapjuk } n\text{-t: } n \leq \frac{1}{4(SE_{arany})^2} = \frac{1}{4 \cdot 0,01^2} = \frac{2500}{4}$$

Vagyis egy 2500-as mintaelemszám garantálja, hogy egy esetleges 0,5-ös prevalencia becslési hibája sem haladja meg az 1%-ot. (A \leq jel azt jelenti, hogy 0,5 alatti vagy feletti prevalenciák esetén kisebb minták is elégségesek lennének az SE 1% alatt tartására.)

19

Átlag standard hibája

2. példa: Meg szeretnénk becsülni a golyák testmagasságát. Vettünk egy $n = 15$ hallgatóból álló mintát, melyből az átlag $\bar{x} = 170$ cm-nek, korrigált szórás $s = 10$ cm-nek adódott. Számítógépes simulációval állítsuk elő az **adateloszlást** és a **mintátlagok eloszlását** $n = 15$ -elemű mintákra.



adatok eloszlása:
 $\bar{x} = 170,01$ cm
 $SD = 9,947$ cm

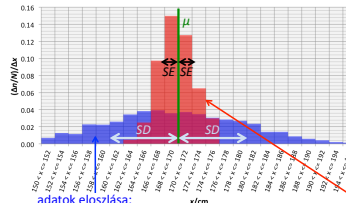
Kétszáz mintavétel (összesen 3000 elem) simulálása tisztán mutatja, hogy az átlagok eloszlása lényegesen „keskenyebb”, mint az adatoké. De mennyivel?

átlagok eloszlása:
 $\bar{x} = 170,01$ cm
 $szórás = 2,507$ cm

20

Átlag standard hibája

Ha az adatok szórását osztjuk a mintaelemszám (n) gyökével, megkapjuk az átlagok szórását. Ez utóbbit nevezzük az **átlag standard hibájának** ($SE_{\text{átlag}}$). Az átlagok eloszlása (legalábbis közelítőleg) normális a **centrális határeloszlás tétele** miatt.



adatok eloszlása:
 $\bar{x} = 170,01$ cm
 $SD = 9,947$ cm

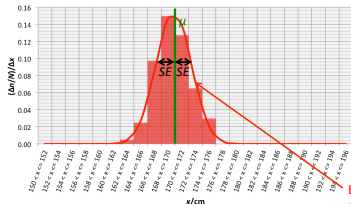
átlagok eloszlása:
 $\bar{x} = 170,01$ cm
 $SD/n^{0.5} = 9,947 \text{ cm}/3,872 = 2,568 \approx SE_{\text{átlag}} = 2,507 \text{ cm}$

$$SE_{\text{átlag}} = s_x = \frac{\sigma}{\sqrt{n}}$$

21

Átlag standard hibája

Azonban az átlag hibáját ($SE_{\text{átlag}}$) többnyire a minta szórásából számoljuk, így eloszlása $n-1$ szabadsági foku (degrees of freedom, df) **Student t-eloszlás** lesz. Ez az eloszlás nagy df esetén hasonlít a normális eloszláshoz, kis df esetén azonban jelentősen „nehézzebbek” a szélei (leptokurtotikus).



Elméleti modell: t-eloszlás
 $\mu \rightarrow \bar{x} = 170,01$ cm
 $\sigma \rightarrow SE = 2,507$ cm
 $df = n - 1 = 14$

$$SE_{\text{átlag}} = s_x = \frac{\sigma}{\sqrt{n}} = \frac{s}{\sqrt{n}}$$



William S. Gosset
1876–1937
"Student"

22

Átlag standard hibája

4. feladat: Meg szeretnénk becsülni a banán átlagos tömegét (várható értékét). Lemértünk öt banánt, az eredmények: 134 g, 152 g, 158 g, 141 g, 170 g. Add meg a becslés hibáját.

$n = 5$
 $\bar{x} = 151$ g
 $SD = 14,14$ g
 $SE_{\text{átlag}} = 6,32$ g

5. feladat: egy tudományos cikkben a következő mondat szerepelt: „... a kísérletben használt patkányok átlagtömege 420 g ($SE = 20$ g), míg átlagos életkora 5 hónap volt ...”. A patkányok számáról azonban nem tesznek említést. Mit gondolsz, hány patkányt használtak, ha máshonnan tudjuk, hogy az ilyen korú patkányok tömegének szórása kb. 40 g?

Az átlag szórása egyenlő a változó szórása osztva a mintaelemszám gyökével. Rendezzük át ezt a képletet n -re:
 $n = (SD/SE)^2 = (40 \text{ g}/20 \text{ g})^2 = 2^2 = 4$

Megjegyzés: Ez egy eléggé alacsony elemszám egy vizsgálathoz, mely megmagyarázhatja, hogy a szerzők miért nem említették meg. Ez viszont az egész cikk hitelét megkérdőjelezi...

23

A „jó becslés” ...

... hatásos



Valószínűség becslésének standard hibája:

$$SE_{\text{arány}} = \sqrt{\frac{p(1-p)}{n}} \approx \sqrt{\frac{k \left(1 - \frac{k}{n}\right)}{n}}$$

Általában kimondhatjuk, hogy a standard hiba négyzete egyenes arányban áll a statisztikai változó varianciájával, és fordított arányban a minta elemszámával.

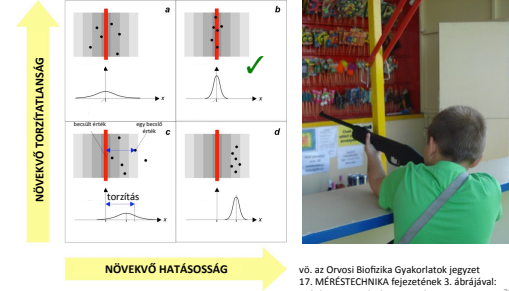
Vagyis a hatásosság megduplázásához négyzeres mintaméretre van szükség! ²⁴

Várható érték becslésének standard hibája:

$$SE_{\text{átlag}} = s_x = \frac{\sigma}{\sqrt{n}} \approx \frac{s}{\sqrt{n}}$$

A „jó becslés” ...

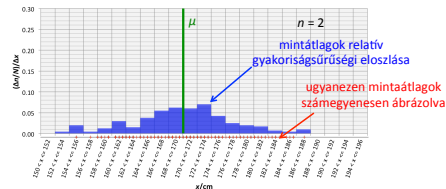
... torzítatlan és hatásos



A „jó becslés” ...

... konzisztens

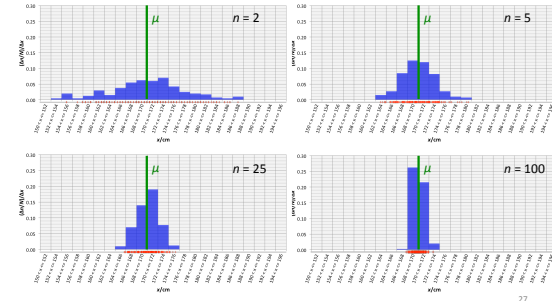
Képzeld el egyre nagyobb méretű (n) mintákon elvégzett becslések sorozatát. Egy becslés konzisztens, ha a becslő érték a mintaméret növelésével egyre kevésbé tér el a becsült értéktől. Más szóval nagyobb elemszám esetén kisebb lesz a torzítás és a standard hiba is.



A testmagasság várható értékének becslésére 200 mintából számolt mintátlagok eloszlása egyre növekvő mintaméretre: $n = 2, 5, 25$, and 100

A „jó becslés” ...

... konzisztens



A „jó becslés” ...

... konzisztens

A konzisztenciát szigorú matematikai módszerekkel kell vizsgálni (mi nem fogjuk). Ehelyett csak szemléltessük a képleteiken, hogy a valószínűség és a várható érték becslésére tanult módszereink konzisztensek.

Valószínűség
becslése

$$\hat{p}_i = (k_i/n) \checkmark$$

A pontbecslések
torzítatlanok

Várható érték
becslése

$$\bar{x} = \sum_{i=1}^k \frac{k_i}{n} \cdot x_i = \frac{1}{n} \sum_{i=1}^k k_i \cdot x_i \checkmark$$

$$SE_{arány} = \sqrt{\frac{p(1-p)}{n}} \approx \sqrt{\frac{\frac{k}{n}(1-\frac{k}{n})}{n}} \checkmark$$

$$SE_{átlag} = s_x = \frac{\sigma}{\sqrt{n}} \approx \frac{s}{\sqrt{n}} \checkmark$$

A hiba nullához tart, ha növeljük az
elemszámot (n a nevezőben szerepel)

28

A „jó becslés” ...

... elégséges

Egy becslés akkor elégséges, ha a becslő érték **minden**, a **mintából kinyerhető információt** tartalmaz, amely a becslőt érték szempontjából releváns.

Példa: Ha egy statisztikai változót legalább intervallumskálán mérünk, az átlag elégséges becslése a várható értéknek, mivel a minta minden megfigyelt elemének értékét felhasználja. Azaz a teljes minta ismerete nem ad több információt a várható értékre vonatkozóan annál, mintha csak az átlagot ismerjük.

Ellenpélda: Ugyanebben az esetben a medián nem elégséges, mivel csak a mintaelemek rangját használja.

Ellen-ellenpélda: Ha azonban a statisztikai változót csak ordinális szinten mérjük, a medián már elégséges becslő értéke a középértéknek (mivel az átlag nem használható).

29

A „jó becslés” ...

... torzítatlan

Végtelessok mintavételből számolt becslő értékek
átlaga egyenlő a becslőt értékkel.

... hatásos

A becslés standard hibája
(vagyis a becslő értékek szórása) kicsi.

... konzisztens

A mintaméret növelésével egyre kevésbé tér el
a becslőt értéktől a becslő érték.

... elégséges

Ugyannyi információt tartalmaz,
mintha a teljes mintát adnánk meg.

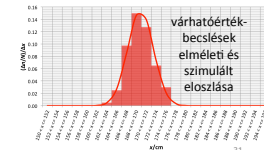
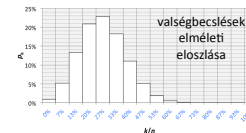
30

Intervallumbecslés

Az intervallum becslés során a becslőt értékhez tartományt rendelünk (a becslő érték mint valósági változó skáláján egy alsó és egy felső határt), melynek neve **konfidenciaintervallum** (angolul *confidence interval*, jele: CI) és ehhez egy valótséget, melynek neve **konfidenciaszint** (jele: $1-\alpha$), mely a becslési eljárás megbízhatóságát fejezi ki (konfidencia = megbízhatóság).

A CI létrehozásának tipikus lépései:

- mintavétel
- pontbecslés kiszámítása a mintaadatokból
- pontbecslés valószínűségi eloszlásának meghatározása
- az eloszlás alapján standard hiba számítása (opcionális)
- mindezek ismeretében egy tartomány megadása, ahol a becslőt érték „lehet”



31

Intervallumbecslés ...

... valségre

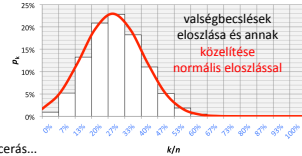
Mi a valsége, hogy egy véletlenül kiválasztott golyának kék szeme van?

A CI létrehozásának lépései:

- mintavétel: n elemű mintában k kék szemű
- pontbecslés: $p_k = k/n = \dots$
- a becslő érték elméleti eloszlása: binomiális eloszlás

$$SE_{\text{arány}} = \sqrt{\frac{k}{n} \left(1 - \frac{k}{n}\right)} = \dots$$

- a megfigyelt k értékek binomiális eloszlást követnek, melyet át kell skálázunk k/n -re. Lehet ezzel az is számolni (ún. egzakt eljárás) de macerás...
- könnyebb ehelyett a **normális eloszlással való közelítést** használni: Wald-intervallum (egyszerű és elterjedt, de eléggé megbízhatatlan módszer)



32

Intervallumbecslés ...

... valségre

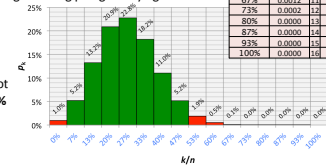
1. módszer: egzakt eljárás

- számold ki minden egyes kimenetel (k) valségét a becslő p felhasználásával, és rendezd őket csökkenő valség szerint
- kezd el ezeket összeadogatni kezdve a legnagyobb valségűvel, majd a második legvalószínűbb s.í.t.
- ha az összeg meghalad egy előre kiválasztott limitet, állj meg
- a kimenetek tartománya, mely bekerült a végső összegbe, lesz az egzakt CI, maga a végső összeg pedig a tényleges konfidenciaszint $(1-\alpha)$

k/n	p_k	Σ	szumma
0%	0.0005	9	0.9993
7%	0.0520	6	0.9154
13%	0.1324	4	0.7510
20%	0.2087	2	0.4364
27%	0.2277	1	0.2277
33%	0.1821	1	0.0456
40%	0.1104	5	0.0014
47%	0.0516	7	0.0000
53%	0.0188	8	0.0000
60%	0.0051	10	0.0000
67%	0.0012	11	0.0000
73%	0.0002	12	0.0000
80%	0.0000	13	0.0000
87%	0.0000	14	0.0000
93%	0.0000	15	0.0000
100%	0.0000	16	0.0000

3. példa:

- használjuk az 1. példa adatait
- lásd a kiszámolt táblázatot és a grafikon, mely a **95%-os (valójában 96,5%-os) CI-t mutatja: 7% – 47%**



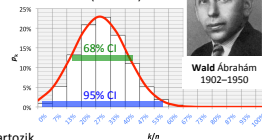
33

Intervallumbecslés ...

... valségre

A közelítés $n > 30$ elemzés esetén működik „jó”.

2. módszer: közelítés normális eloszlással (Wald A.)



Wald Abraham
1902-1950

A $\mu \pm \sigma$ közötti tartományhoz kb. 68% valség tartozik.

Az így definiált becslési tartomány neve **68%-os konfidencia-intervallum (CI)**, a valsége pedig **konfidenciaszint $(1-\alpha)$** .

A 3. példában: **68% CI = 15% – 38%**. (lásd a jobb oldali ábrát is)

A $\mu \pm 2\sigma$ közötti tartományhoz kb. 95% valség tartozik.

Ez a **95%-os konfidenciaintervallum**unk felül meg, és igen gyakran használják az élettudományokban. Kis mintáknál ez a CI akár ki is lóghat a $[0,1]$ tartományból! \rightarrow széleit le kell nyelni.

A 3. példában: **95% CI = 4% – 50%**. (lásd a jobb oldali ábrát és vö. az egzakt intervallummal.)

$$68\% CI \approx \frac{k}{n} \pm \sqrt{\frac{k}{n} \left(1 - \frac{k}{n}\right)}$$

$$95\% CI \approx \frac{k}{n} \pm 2 \cdot \sqrt{\frac{k}{n} \left(1 - \frac{k}{n}\right)}$$

34

Intervallumbecslés ...

... valségre

6. feladat: Meg szeretnénk becsülni a Rhesus faktor **prevalenciáját** (alapsokaságon belüli arányát) a budapesti lakosság körében. Véletlenszerűen kiválasztottunk 42 embert, s meghatároztuk a vércsoportjukat: 35 volt Rh+.

a) Adj pontbecslést a Rhesus faktor prevalenciájára.

$$p(Rh+) = \frac{k}{n} = \frac{35}{42} = 0.833$$

b) Add meg a 95%-os CI-t a Wald-intervallum eljárással.

Számoljuk ki a SE-t:

$$SE = \sqrt{\frac{k}{n} \left(1 - \frac{k}{n}\right)} = \sqrt{\frac{35}{42} \left(1 - \frac{35}{42}\right)} = 0.00331$$

A 95%-os CI kb. $\hat{p}(Rh+) \pm 2SE$ között van, azaz **0.833 \pm 0.006** vagy **0.826 – 0.839**.

7. feladat: add meg a 95%-os konfidenciaintervallumot a kék szemszín prevalenciájára a golyák közt, ha egy 10 hallgatóból álló mintában kettőnek volt kék a szeme.

Az ismert képlettel számolva: $\hat{p}(Rh+) = 0.200$ és $SE = 0.126$. Ebből a következő 95%-os CI adódna: **-0.052 – 0.452**. Azonban mivel valséget becsülünk, a CI nem lóghat ki a $[0,1]$ intervallumból; szélek nyesése után a CI: **0 – 0.452**. Itt a konfidenciaszint már biztos nem 95%-os. Ez csak egy példa, hogy mennyire megbízhatatlan a Wald-intervallum – de ennek ellenére szoktuk használni nem túl kis minták esetén, ha gyors becslésre van szükség.

35

Intervallumbecslés ...

... várható értékre

Student-féle t-eloszlást használva (W. S. Gosset)

Az arányokhoz hasonló intervallumokat lehet definiálni, ez esetben azonban az $n-1$ szabadsági fokú (df) Student-féle t-eloszlást fogjuk használni. Így nem használhatjuk a korábbi közelítő intervallumokat sem ($\mu \pm \sigma$ és $\mu \pm 2\sigma$).
Helyette függvénytáblában vagy Excel paranccsal határozzuk meg a kétszélű p-t: $=T.INV.2S(2S2)$ (valószínűség, szabadságfok)

4. példa:

- használjuk a 2. feladat adatait
- add meg a 95%-os CI t-értékét:
 - Excel-ben: $a = T.INV.2S(2S2)$ (5%)
 - (a Képlettárbán lévő függvénytáblából: 2.14)

T-ELOSZLÁS

szabadságfok	0,5	0,2	0,1	0,05	0,025	0,01	0,005	0,001
1	1,00	0,99	0,95	0,90	0,82	0,71	0,60	0,50
2	0,92	0,89	0,86	0,82	0,76	0,69	0,62	0,56
3	0,76	0,74	0,71	0,68	0,64	0,59	0,54	0,50
4	0,74	0,71	0,68	0,65	0,61	0,56	0,52	0,49
5	0,73	0,70	0,67	0,64	0,60	0,55	0,51	0,48
6	0,72	0,69	0,66	0,63	0,59	0,54	0,50	0,47
7	0,71	0,68	0,65	0,62	0,58	0,53	0,49	0,46
8	0,71	0,68	0,65	0,62	0,58	0,53	0,49	0,46
9	0,70	0,67	0,64	0,61	0,57	0,52	0,48	0,45
10	0,70	0,67	0,64	0,61	0,57	0,52	0,48	0,45
11	0,70	0,67	0,64	0,61	0,57	0,52	0,48	0,45
12	0,69	0,66	0,63	0,60	0,56	0,51	0,47	0,44
13	0,69	0,66	0,63	0,60	0,56	0,51	0,47	0,44
14	0,69	0,66	0,63	0,60	0,56	0,51	0,47	0,44
15	0,69	0,66	0,63	0,60	0,56	0,51	0,47	0,44
16	0,69	0,66	0,63	0,60	0,56	0,51	0,47	0,44

Ábra: A t-érték meghatározása a függvénytáblából, ha $\alpha = 5\%$ és $df = 14$.

36

Intervallumbecslés ...

... várható értékre

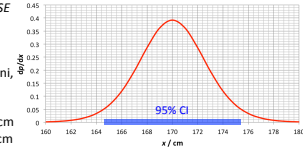
Student-féle t-eloszlást használva (W. S. Gosset)

A statisztika gyakorlaton Excel-t használunk a t kiszámolására. Azonban t eloszlása standard ($\mu = 0$ és $\sigma = 1$), így az $x = \bar{x} + t \cdot SE$ képletet kell használnunk a megfelelő x-értékek kiszámolására. (x a statisztikai változónk, jelen esetben a magasság).
-t-ből tudjuk a CI alsó határát (x_1) számolni, míg +t-ből a felsőt (x_2):

$$X_1 = \bar{x} - t \cdot SE = 170 \text{ cm} - 2,1445 \cdot 2,5 \text{ cm} = 164,6 \text{ cm}$$

$$X_2 = \bar{x} + t \cdot SE = 170 \text{ cm} + 2,1445 \cdot 2,5 \text{ cm} = 175,4 \text{ cm}$$

A 95%-os CI tehát 164,6 cm – 175,4 cm.



Ábra: Az átlagok elméleti eloszlásának grafikonja (egy df = n-1 szabadsági fokú nem-standard t-eloszlás, melynek paraméterei: $\mu = \bar{x}$ és $\sigma = SE$) és az imént számolt 95%-os CI.

37

Intervallumbecslés ...

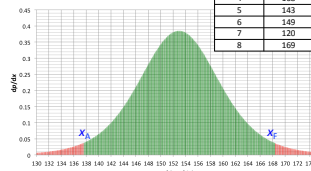
... várható értékre

8. feladat: Becslést szeretnénk adni a vérkoleszterinszint várható értékére és meg szeretnénk adni a becslés 95%-os CI-t. Nyolcelemű mintát vettünk, a megfigyelt értékek a táblázatban vannak.

Excel-ben fogjuk elvégezni a számításokat.
A várható érték pontbecslése a mintáztól:
 $\mu = \bar{x} = \text{ÁTLAG}(adat) = 152,9 \text{ mg/dL}$

Számoljuk ki a SE-t a már tanult módon:
 $n = \text{DARAB}(adat) = 8$
 $SD = \text{SZÖR.M}(adat) = 18,47 \text{ mg/dL}$
 $SE = SD / \text{GYÖK}(n) = 6,53 \text{ mg/dL}$

A mintaátlagok t-eloszlást követnek, a 95%-os CI adatai:
 $df = n - 1 = 7$
 $t = T.INV.2S(2S2) = 2,365$
 $x_1 = \bar{x} - t \cdot SE = 137,4 \text{ mg/dL}$
 $x_2 = \bar{x} + t \cdot SE = 168,3 \text{ mg/dL}$



Ábra: 95%-os konfidenciaintervallum (zölddel).

38

Intervallumbecslés ...

... várható értékre

9. feladat: Meg szeretnénk határozni a golyók átlagmagasságát 95%-os CI-mal, de azt szeretnénk, hogy a CI ne legyen 1 cm-nél szélesebb. Mekkora legyen a minta mérete? Az irodalomból tudjuk, hogy az emberi testmagasság szórása általában 5 cm, és ezt a golyókra is érvényesnek tekintjük.

Itt most „visszafelé” kell gondolkodnunk:

– Vann egy CI-ünk (most még) ismeretlen x_1 és x_2 között.

– A CI szélessége a határoló értékek különbsége: $\text{szélesség} = x_2 - x_1$

– Helyettesítsük be a képleteiket: $\text{szélesség} = (\bar{x} + t \cdot SE) - (\bar{x} - t \cdot SE) = 2 \cdot t \cdot SE$

– Most van egy kis problémánk: a t-érték maga is függ a mintamérettől (pontosabban a szabadsági fokok számától)! Egy feltevéssel kell eljárunk: vegyük úgy, hogy a minta elég nagy, és egyszerűen használjunk a $t = 2$ értéket; ekkor a képlet: $\text{szélesség} = 4 \cdot SE$

– Tudjuk, hogy $SE = SD / n^{0.5}$, helyettesítsük be a fenti képletbe: $\text{szélesség} = 4 \cdot SD / n^{0.5}$

– Rendezzük át a képletet a mintaméretre: $n = (4 \cdot SD / \text{szélesség})^2$

– Helyettesítsük be az adatainkat: $n = (4 \cdot (5 \text{ cm}) / (1 \text{ cm}))^2 = 400$

39

Intervallumbecslés

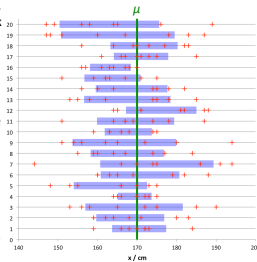
Konfidenciaszint $(1 - \alpha)$: annak a valószínűsége, hogy egy adott módszerrel meghatározott CI-ok tartalmazza-e a becsült értéket.

Ebből nem tudjuk meg, hogy egy konkrét CI ténylegesen tartalmazza-e a becsült értéket.

Magát a becsülési eljárást jellemzi általában, nem annak egy konkrét eredményét. Azt nem tudjuk megmondani, hogy a becsült érték benne van-e az éppen kiszámolt CI-ban, csak akkor, ha valahonnan ismerjük a becsült értéket – mely esetben az egész becslésre nincs szükség.

Szignifikanciaszint (α) : a komplementer valószínűség, pl. ha a konfidenciaszint 95%, akkor a szignifikanciaszint 5%. Hogy melyiket használjuk, az a kontextustól függ.

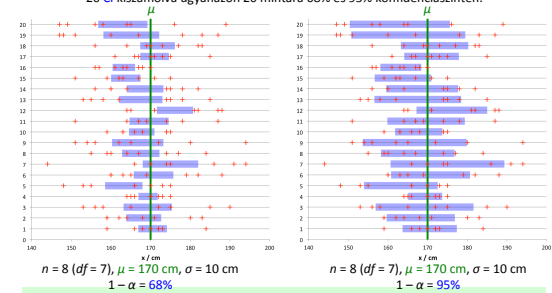
Ábra: A testmagasság várható értéke (μ) ugyanazon módszerrel végzett 20 becslés szimulációja: 8-elemű minta vételezése (piros +), átlag és standard hibájának számlálása, majd a CI megadása: $\text{átlag} \pm 2,36 \times \text{SE}$ (kék sáv). A becsült értékek: $\mu = 170$ cm, $\sigma = 10$ cm.



40

Intervallumbecslés

20 CI kiszámolva ugyanazon 20 mintára 68% és 95% konfidenciaszinten.



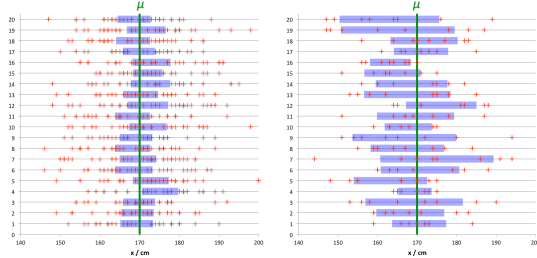
$n = 8$ ($df = 7$), $\mu = 170$ cm, $\sigma = 10$ cm
 $1 - \alpha = 68\%$

$n = 8$ ($df = 7$), $\mu = 170$ cm, $\sigma = 10$ cm
 $1 - \alpha = 95\%$

Magasabb konfidenciaszint = kisebb valószínűséggel esik a CI-on kívül a becsült érték, de kevesebb az információtartalma.

Intervallumbecslés

20–20 db 95%-os CI kiszámolva 20 db 34 elemű és 20 db 8-elemű mintára.



$n = 32$ ($df = 31$), $\mu = 170$ cm, $\sigma = 10$ cm
 $1 - \alpha = 95\%$

$n = 8$ ($df = 7$), $\mu = 170$ cm, $\sigma = 10$ cm
 $1 - \alpha = 95\%$

Nagyobb mintaelemszám ugyanazon konfidenciaszinten = szűkebb CI.
(4-szeres elemszám átlagosan fele szélességű CI-kat eredményez.)

42

Függelék: Normáltartomány

Normáltartomány, referenciatartomány vagy referenciaintervallum: a statisztikai változóra vonatkozó tartomány, mely a változó egy elemét 95%-os valószínűséggel tartalmazza.

A normáltartomány is egyfajta intervallumbecslés, de itt a statisztikai változó eloszlását, nem pedig eloszlásának egy paraméterét jellemizzük. Más szóval: a normáltartomány a magukra az adatokra vonatkozó 95%-os CI. Ezt a tartományt szokták feltüntetni az orvosi laborleleteken is.

Normális eloszlású változó esetén a normáltartomány kb.: $\bar{x} \pm 2SD$.
(pontosabban: $\bar{x} \pm 1,96SD$)

PARAMÉTER	EGYSÉG	ÁTLAG	STANDA	REFERENCIA	REFERENCIA
Glukózeszint	mg/dL	45-99	45-99	45-99	45-99
Glukózeszint, éhvért	mg/dL	70-100	70-100	70-100	70-100
HbA1c	%	5,7-6,4	5,7-6,4	5,7-6,4	5,7-6,4
Cholesterol, éhvért	mg/dL	125-200	125-200	125-200	125-200
Cholesterol, éhvért, HDL	mg/dL	40-120	40-120	40-120	40-120
LDL cholesterol	mg/dL	130-200	130-200	130-200	130-200
Triglycerides	mg/dL	50-150	50-150	50-150	50-150
Urea Nitrogen	mg/dL	7-20	7-20	7-20	7-20
Bilirubin	mg/dL	0,2-1,2	0,2-1,2	0,2-1,2	0,2-1,2
Alkaline Phosphatase	U/L	44-147	44-147	44-147	44-147
ALT (GPT)	U/L	7-56	7-56	7-56	7-56
AST (GOT)	U/L	10-40	10-40	10-40	10-40
Cr	mg/dL	0,6-1,3	0,6-1,3	0,6-1,3	0,6-1,3

$x_1 = \mu - 2\sigma = 65$ mg/dL

$x_2 = \mu + 2\sigma = 99$ mg/dL

$x_1 - x_2 = (\mu - 2\sigma) - (\mu - 2\sigma) = 4\sigma$

$\sigma = (x_2 - x_1)/4 = 8,5$ mg/dL

$x_1 + x_2 = (\mu + 2\sigma) + (\mu - 2\sigma) = 2\mu$

$\mu = (x_2 + x_1)/2 = 82$ mg/dL

$\sigma^2 = 72,25$ (mg/dL)²

43

10. feladat: Számold ki a szérum glükózeszint μ , σ és σ^2 paramétereit a laborlelet adatainak segítségével.

Ellenőrző kérdések

- Mi a populáció?
- Mi a minta?
- Hogyan szerezhetünk információt egy statisztikai változóról?
- Milyen becsléstípusokat ismersz?
- Mik a becslés lépései?
- Mi a „plug-in” becslés?
- Mi a pontbecslés?
- Mi a pontbecslés hátránya?
- Mi a becslő értéke a valségnek, várható értéknek, elméleti varianciának és elméleti szórásnak?
- Mik a „jó becslés” tulajdonságai?
- Mi a torzítatlanság? Szemléltess példával.
- Mi a hatásosság? Hogyan lehet matematikailag kifejezni?
- Milyen eloszlást követnek az arányok?
- Milyen eloszlást követ a várható érték becslő értéke?
- Mi a standard hiba?
- Hogyan számolható egy arány standard hibája?
- Hogyan számolható egy átlag standard hibája?
- Hogyan függ a standard hiba a változó szórásától?
- Hogyan függ a standard hiba a minta méretétől (elemszámától)?
- Meg szeretnénk háromszorozni egy becslés hatásosságát. Hogyan változtassunk a minta méretén?
- Mekkora egy 25 elemű mintából számolt arány legnagyobb lehetséges hibája?

44

Ellenőrző kérdések

- Mik egy konzisztens becslés ismérvei?
- Mit jelent, hogy egy becslés elégséges?
- Mi a Bessel-korrektió? Hol használjuk és mi a célja?
- Mi a CI jelentése?
- Mi a konfidenciaszint jele és jelentése?
- Hogyan adhatunk meg egzakt CI-t egy arányhoz?
- Mi a Wald-intervallum alapja, mi az előnye és a hátránya?
- Hogyan változik a CI, ha emelem a konfidenciaszintet?
- Hogyan változik a CI, ha növelem a mintaméretet?
- Hogyan változik a CI, ha a vizsgált változó szórása kisebb?
- Hogyan adunk meg CI-t várható érték becslése esetén?
- Mi a referenciatartomány?
- Hogyan adható meg egy normális eloszlású valségi változó referenciatartománya?
- Egy laborleletünkben különféle laborértékek mellett referenciatartományokat látunk. Hogyan lehet ezekből a kórház által használt várható értékeket és szórásokat kiszámolni?
- Mi a kapcsolat a konfidenciaszint és a szignifikanciaszint között?
- Miért nem lehetséges megmondani, hogy egy adott CI ténylegesen tartalmazza-e a becsült értéket?

45