

# Biostatisztika

Kunosné Nedényi Fanni, Szűcs Gábor

Szegedi Tudományegyetem, Bolyai Intézet

2018/19 őszi félév

# Mi is az a biostatisztika?

Ide majd még jön valami, valamikor...

# Események valószínűsége

A valószínűségszámítás a matematika egyik ága, melynek célja a véletlen jelenségekhez kapcsolódó valószínűségek meghatározása. Alapfogalmak:

- **Véletlen kísérlet:** Egy véletlen jelenség megfigyelése.
- **Kimenetek:** A véletlen kísérlet lehetséges eredményei.
- **Esemény:** A kísérlet aktuális kimenetelével kapcsolatos állítás. Egy esemény akkor **következik be**, ha a véletlen kísérlet olyan kimenetelt ad, melyre az állítás igaz.
- **Valószínűség:** Annak az esélye, hogy az esemény bekövetkezik.

## Példa:

- Véletlen kísérlet: feldobunk egy szabályos dobókockát.
- Kimenetek: 1, 2, 3, 4, 5, 6.
- Egy esemény:  $A =$  páros számot dobunk. Ez akkor következik be, ha a 2, 4, 6 értékek valamelyikét dobjuk, egyébként nem következik be.
- Az  $A$  esemény valószínűsége:  $P(A) = 3/6 = 50\%$ .

Ezen a kurzuson jellemzően az lesz majd a kísérlet, hogy véletlenszerűen kiválasztunk egy vagy több egyedet egy ember/állat/növény populációból. A „véletlenszerűen” szó itt azt jelenti, hogy mindegyik egyedet ugyanakkora eséllyel választjuk ki.

**Feladat:** Magyarországon az emberek 52 illetve 24 százalékának a vérében található meg az A illetve a B típusú antigén. Mindkét antigén az emberek 8 százalékánál található meg. Véletlenszerűen kiválasztunk egy magyar embert, és leteszteljük az antigénekre. Tekintsük a következő eseményeket:

$A$  = a kiválasztott ember rendelkezik az A típusú antigénnel

$B$  = a kiválasztott ember rendelkezik a B típusú antigénnel

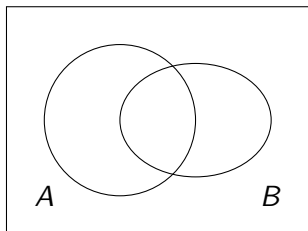
Most minden embert azonos eséllyel választunk ki, ezért a tulajdonságok bekövetkezési valószínűsége azonos lesz a tulajdonságok teljes populáción belül mért részarányával:

$P(A)$  = az A antigén aránya a teljes populáción belül = 52% = 0,52

$P(B)$  = a B antigén aránya a teljes populáción belül = 24% = 0,24

$P(A \text{ és } B)$  = a két antigén együttes megjelenésének aránya = 8% = 0,08

Az alábbi ábrán a magyar emberek populációját ábrázoljuk a két antigén szempontjából. A mellette lévő táblázat a vércsoportokat foglalja össze.



	van A	nincs A
van B	AB	B
nincs B	A	0

Amit tudunk:  $P(A) = 52\%$ ,  $P(B) = 24\%$ ,  $P(A \text{ és } B) = 8\%$ .

**Feladat:** Határozzuk meg a vércsoportok részarányát!

$P(\text{a kiválasztott ember az AB vércsoportba esik}) = P(A \text{ és } B) = 8\%$

$P(A \text{ vércsoport}) = P(A \text{ igen, de } B \text{ nem}) = P(A) - P(A \text{ és } B) = 44\%$

$P(B \text{ vércsoport}) = P(B \text{ igen, de } A \text{ nem}) = P(B) - P(A \text{ és } B) = 16\%$

$P(0 \text{ vércsoport}) = 100\% - \text{az előző három összege} = 32\%$

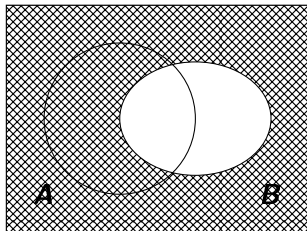
A kurzuson a valószínűség a teljes populáción belüli arányt jelenti. Időnként szükségünk lesz arra, hogy az arányokat egy részpopuláción belül vizsgáljuk. Az  $A$  eseménynek a  $B$  eseményre vett **feltételes valószínűsége**:

$$P(A|B) = \frac{P(A \text{ és } B)}{P(B)}.$$

A feltételes valószínűség jelentése:

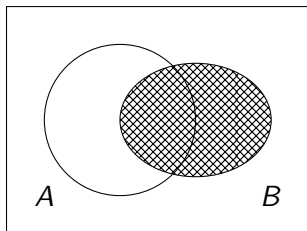
$P(A|B)$  = az  $A$  tulajdonság aránya a  $B$  részpopuláción belül  
 = az  $A$  esemény valószínűsége, ha tudjuk, hogy  $B$  bekövetkezik

**Feladat:** Mennyi  $P(A|B)$  az előző feladatban?



$$P(A|B) = \frac{P(A \text{ és } B)}{P(B)} = \frac{8\%}{24\%} = \frac{1}{3} = 33\%$$

**Feladat:** Mekkora az A típusú antigénnel rendelkező emberek aránya azon emberek között, akik nem rendelkeznek a B antigénnel?



$$P(A \mid \text{nem } B) = \frac{P(A \text{ és nem } B)}{P(\text{nem } B)} = \frac{44\%}{76\%} = 58\%,$$

$$P(\text{nem } B) = 100\% - P(B) = 76\%,$$

$$P(A \text{ és nem } B) = P(A \text{ vércsoport}) = 44\%.$$

Értelmezzük, hogy mit kaptunk:

- Ha véletlenszerűen kiválasztunk egy embert a teljes populációból, akkor 52% valószínűséggel található meg nála az A típusú antigén.
- Ha tudjuk, hogy a kiválasztott ember rendelkezik a B antigénnel, akkor 33% az esélye, hogy az A is antigén megtalálható nála.
- Ha viszont azt tudjuk, hogy nem rendelkezik a B típusú antigénnel, akkor 58% az esélye, hogy az A antigén megtalálható nála.
- Tehát a B antigén jelenléte csökkenti az A antigén megjelenési esélyét:

$$P(A|B) = 33\% < 58\% = P(A \mid \text{nem } B)$$

Legyenek  $A$  és  $B$  tetszőleges események. Bebizonyítható, hogy ekkor az alábbi három egyenlőség ekvivalens, tehát következnek egymásból:

$$1 \quad P(A \text{ és } B) = P(A)P(B)$$

$$2 \quad P(A|B) = P(A)$$

$$3 \quad P(B|A) = P(B)$$

Amennyiben ezen egyenlőségek közül bármelyik (és ezáltal mindegyik) teljesül, akkor azt mondjuk, hogy  $A$  és  $B$  **független események**.

A függetlenség szemléletesen azt jelenti, hogy a két esemény nem hat egymásra, nem akadályozzák, és nem is segítik elő egymás bekövetkezését.

Lássuk, hogyan következik az első egyenlőségből a második:

$$P(A|B) = \frac{P(A \text{ és } B)}{P(B)} = \frac{P(A)P(B)}{P(B)} = P(A).$$



**Feladat:** A vércsoportos feladatban  $A$  és  $B$  független események?

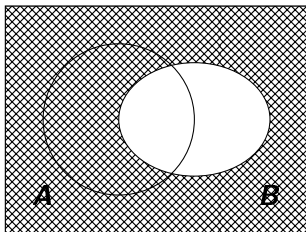
Nem, ugyanis  $P(A|B) = 33\% < 52\% = P(A)$ .

**Feladat:** A feladatban a két antigén aránya:  $P(A) = 52\%$  és  $P(B) = 24\%$ . Mikor lenne a két antigén megjelenése független egymástól?

A két antigén megjelenése akkor független, ha

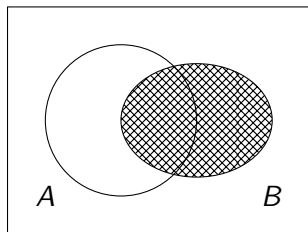
$$P(A \text{ és } B) = P(A)P(B) = 0,52 \cdot 0,24 = 0,125 = 12,5\%$$

**Feladat:** Mennyi lenne a feltételes valószínűségek értéke ebben az esetben?



Az  $A$  típusú antigénnel rendelkező emberek aránya a  $B$  csoporton belül:

$$P(A|B) = \frac{P(A \text{ és } B)}{P(B)} = \frac{12,5\%}{24\%} = 52\%$$



Az A típusú antigénnel rendelkező emberek aránya a B csoporton kívül:

$$P(A | \text{nem } B) = \frac{P(A \text{ és nem } B)}{P(\text{nem } B)} = \frac{39,5\%}{76\%} = 52\%,$$

$$P(A \text{ és nem } B) = P(A) - P(A \text{ és } B) = 39,5\%.$$

Tehát a független esetben az A antigénnel rendelkező emberek aránya (=kiválasztási valószínűsége) azonos az alábbi három populáción belül:

- a teljes populáción belül:  $P(A) = 52\%$ ,
- a B típusú antigénnel rendelkező emberek részpopulációján belül:  $P(A|B) = 52\%$ ,
- a B típusú antigénnel nem rendelkező emberek részpopulációján belül:  $P(A | \text{nem } B) = 52\%$ .

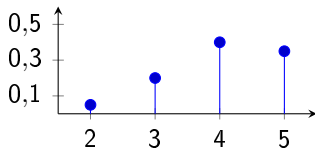
# Diszkrét valószínűségi változók

A biológiai vizsgálatok során gyakran felmerül az a kérdés, hogy mi az eloszlása egy mennyiségnek (életkor, testtömeg, utódok száma, stb.) egy populáción belül. Válasszunk ki véletlenszerűen egy egyedet a populációból, és legyen  $\xi$  a vizsgált mennyiség értéke a kiválasztott egyed esetében. Mivel az egyedet véletlenszerűen választjuk, a  $\xi$  érték egy véletlen szám lesz.

- **Valószínűségi változó:** Egy véletlen kísérletből származó véletlen szám (véletlen mennyiség). Jele:  $\xi$  (kszi),  $\eta$  (éta), stb.
- **Értékkészlet:** A változó lehetséges értékeinek a halmaza. Jele:  $R_\xi, R_\eta$
- **Diszkrét valószínűségi változó:** A változó értékkészlete egy véges vagy végtelen sorozat. Mi tipikusan két esettel fogunk találkozni:
  - az értékkészlet véges halmaz;
  - VAGY minden lehetséges érték egész szám.
- **Folytonos valószínűségi változó:** A változó értékkészlete egy véges vagy végtelen intervallum.

**Feladat:** Egy lengyel felmérés alapján a fehér gólyák 2-5 tojást raknak az alábbi táblázatban található megoszlásban. Véletlenszerűen kiválasztunk egy gólyafészket, és jelölje  $\xi$  a fészkekben található tojások számát.

$x$	2	3	4	5
$p_x$	5%	20%	40%	35%



A  $\xi$  egy valószínűségi változó, értékészlete  $R_\xi = \{2, 3, 4, 5\}$ . Ez egy véges halmaz, tehát a  $\xi$  diszkrét változó. A fészket véletlenszerűen választottuk, ezért a  $\xi$  pontosan akkora valószínűséggel veszi fel az egyes értékeket, amennyi ezen értékek aránya a teljes (fészkek-) populáción belül:

$$P(\xi = 2) = 0,05, \quad P(\xi = 3) = 0,2, \quad P(\xi = 4) = 0,4, \quad P(\xi = 5) = 0,35.$$

Legyen  $\xi$  diszkrét valószínűségi változó. A  $p_x = P(\xi = x)$  valószínűségeket a változó **valószínűségeloszlásának** nevezzük. Véletlenszerű kiválasztás esetén a valószínűségeloszlás azonos a populáción belül mért arányokkal.

**Feladat:** Mennyi a  $\xi$  változó lehetséges értékeinek összvalószínűsége?

$$P(\xi = 2) + P(\xi = 3) + P(\xi = 4) + P(\xi = 5) = 0,05 + 0,2 + 0,4 + 0,35 = 1$$

**Feladat:** A fészkek mekkora hányadában található legfeljebb 3 tojás?

$$P(\text{legfeljebb 3 tojás}) = P(\xi \leq 3) = P(\xi = 2) + P(\xi = 3) = 0,25.$$

**Feladat:** Melyik tojásszám a leggyakoribb a populációban?

A 4-es érték a leggyakoribb, a fészkek 40%-ában ennyi tojás található.

**Feladat:** Átlagosan hány tojás található a fészkekben?

A tojások átlagos száma:  $E(\xi) = 2 \cdot 0,05 + 3 \cdot 0,2 + 4 \cdot 0,4 + 5 \cdot 0,35 = 4,05.$

Legyen  $\xi$  diszkrét valószínűségi változó.

- **Módusz:** A  $\xi$  változó legnagyobb valószínűségű értéke.  
Jelentése: a  $\xi$  változó leggyakoribb értéke a teljes populáción belül.
- **Várható érték:**  $E(\xi) = \sum_{x \in R_\xi} xP(\xi = x).$   
Jelentése: a  $\xi$  változó átlagos értéke a teljes populáción belül.

Milyen módon számszerűsíthetjük egy  $\xi$  diszkrét változó szóródását?

- **Várható értéktől való átlagos eltérés:**  $\sum_{x \in R_\xi} |x - E(\xi)| P(\xi = x)$
- **Variancia:**

$$\text{Var}(\xi) = \sum_{x \in R_\xi} (x - E(\xi))^2 P(\xi = x)$$

- **Szórás:**  $D(\xi) = \sqrt{\text{Var}(\xi)}$

A szóródás mérésére a várható értéktől való átlagos eltérés egy egyszerű mutatószám lenne, de sajnos ennek rosszak a matematikai tulajdonságai. Emiatt inkább a szórást szoktuk alkalmazni a szóródás mérésére. A két érték jellemzően közel van egymáshoz:

$$\text{szórás} \approx \text{várható értéktől való átlagos eltérés}$$

Emiatt az alkalmazásokban a szórást magát is úgy értelmezzük, mint az átlagtól való átlagos eltérés. A varianciára csak azért van szükségünk, mert abból számoljuk ki a szórást.

**Feladat:** Mennyi a várható értéktől való átlagos eltérés és a szórás a jelen feladatban? (A várható érték  $E(\xi) = 4,05$ .)

$x$	2	3	4	5
$ x - E(\xi) $	2,05	1,05	0,05	0,95
$(x - E(\xi))^2$	$2,05^2$	$1,05^2$	$0,05^2$	$0,95^2$
$P(\xi = x)$	0,05	0,2	0,4	0,35

Várható értéktől való átlagos eltérés:

$$2,05 \cdot 0,05 + 1,05 \cdot 0,2 + 0,05 \cdot 0,4 + 0,95 \cdot 0,35 = 0,665$$

Variancia:

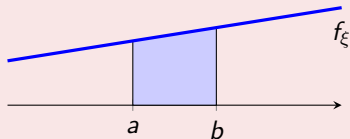
$$\text{Var}(\xi) = 2,05^2 \cdot 0,05 + 1,05^2 \cdot 0,2 + 0,05^2 \cdot 0,4 + 0,95^2 \cdot 0,35 \approx 0,75$$

Szórás:  $D(\xi) = \sqrt{0,75} \approx 0,87$ .

# Folytonos valószínűségi változók

Egy valószínűségi változó **folytonos**, ha értékészlete egy véges vagy végtelen intervallum. A  $\xi$  folytonos változó **sűrűségfüggvénye** egy olyan  $f_\xi : \mathbb{R} \rightarrow \mathbb{R}$  függvény, melyre tetszőleges  $a$  és  $b$  számok esetén:

$$P(a \leq \xi \leq b) = \int_a^b f_\xi(x) dx$$



Tekintünk egy mennyiséget (például a testtömeget) egy populáción belül. Véletlenszerűen kiválasztunk egy egyedet, és legyen  $\xi$  a mennyiség értéke ezen egyed esetében. Ekkor:

azon egyedek aránya, melyeknél a vizsgált mennyiség  $a$  és  $b$  közé esik  
 $= P(a \leq \xi \leq b) =$  görbe alatti terület  $a$  és  $b$  között



A folytonos változók és a sűrűségfüggvények néhány tulajdonsága:

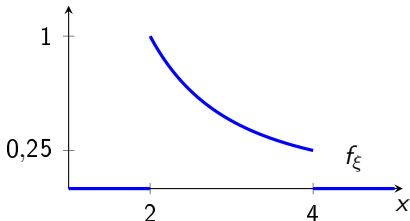
- 1  $\int_{-\infty}^{\infty} f_{\xi}(x) dx = 1$ .
- 2  $f_{\xi}(x) \geq 0$  minden  $x$  valós szám esetén.
- 3 A  $\xi$  változó értékkészlete azon  $x$  számok halmaza, melyekre  $f_{\xi}(x) > 0$ .
- 4 Tetszőleges  $a$  szám esetén  $P(\xi = a) = 0$ .

Rövid indoklás a fenti állításokhoz:

- 1  $\int_{-\infty}^{\infty} f_{\xi}(x) dx = P(-\infty \leq \xi \leq \infty) = 1$ .
- 2 Tegyük fel, hogy az  $f_{\xi}$  függvény negatív egy  $[a, b]$  intervallumon. Ekkor  $\int_a^b f_{\xi}(x) dx < 0$ , tehát  $\int_a^b f_{\xi}(x) dx \neq P(a \leq \xi \leq b)$ , ami ellentmondás.
- 3 Ha  $f_{\xi} = 0$  az  $[a, b]$  intervallumon, akkor  $P(a \leq \xi \leq b) = \int_a^b 0 dx = 0$ .  
Ha  $f_{\xi} > 0$  az  $[a, b]$  halmazon, akkor  $P(a \leq \xi \leq b) = \int_a^b f_{\xi}(x) dx > 0$ .  
Tehát a  $\xi$  változó oda eshet, ahol  $f_{\xi} > 0$ .
- 4  $P(\xi = a) = P(a \leq \xi \leq a) = \int_a^a f_{\xi}(x) dx = 0$ .

**Feladat:** Egy állatpopulációban legyen  $\xi$  egy véletlenszerűen kiválasztott egyed tömege. A változó az alábbi sűrűségfüggvénnyel írható le.

$$f_{\xi}(x) = \begin{cases} 4/x^2, & \text{ha } 2 \leq x \leq 4, \\ 0, & \text{különben.} \end{cases}$$



**Feladat:** A teljes görbe alatti terület valóban 1?

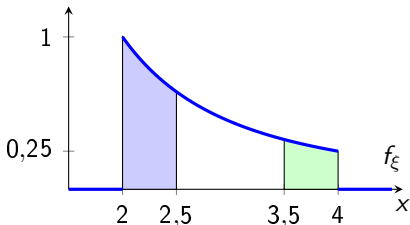
$$\begin{aligned} \int_{-\infty}^{\infty} f_{\xi}(x) dx &= \int_{-\infty}^2 0 dx + \int_2^4 \frac{4}{x^2} dx + \int_4^{\infty} 0 dx = 0 + 4 \int_2^4 x^{-2} dx + 0 \\ &= 4 \left[ \frac{x^{-1}}{-1} \right]_2^4 = 4 \left[ -\frac{1}{x} \right]_2^4 = 4 \left[ \left( -\frac{1}{4} \right) - \left( -\frac{1}{2} \right) \right] = 4 \cdot 0,25 = 1 \end{aligned}$$

**Feladat:** Milyen értékeket vehet fel a  $\xi$  változó?

A változó értékészlete:  $R_{\xi} = [2, 4]$ .

**Feladat:** Mennyi az esélye annak, hogy a  $\xi$  változó 2,5-nél kisebb értéket vesz fel? Mennyi a valószínűsége annak, hogy a  $\xi$  nagyobb, mint 3,5?

$$f_{\xi}(x) = \begin{cases} 4/x^2, & \text{ha } 2 \leq x \leq 4, \\ 0, & \text{különben.} \end{cases}$$



$$\begin{aligned} P(\xi < 2,5) &= P(2 \leq \xi \leq 2,5) = 4 \int_2^{2,5} x^{-2} dx = 4 \left[ -\frac{1}{x} \right]_2^{2,5} \\ &= 4 \left[ \left( -\frac{1}{2,5} \right) - \left( -\frac{1}{2} \right) \right] = 4 \cdot 0,1 = 0,4, \end{aligned}$$

$$P(\xi > 3,5) = P(3,5 \leq \xi \leq 4) = 4 \int_{3,5}^4 x^{-2} dx = \dots \approx 0,14.$$

Legyen  $\xi$  folytonos valószínűségi változó!

- **Móduszok:** Az  $f_\xi$  függvény lokális maximumhelyei.
- **Várható érték:**  $E(\xi) = \int_{-\infty}^{\infty} x f_\xi(x) dx$ .  
Jelentése:  $\xi$  átlagos értéke a teljes populációban.
- **Variancia:**  $\text{Var}(\xi) = \int_{-\infty}^{\infty} (x - E(\xi))^2 f_\xi(x) dx$ .
- **Szórás:**  $D(\xi) = \sqrt{\text{Var}(\xi)}$ .  
Jelentése: a várható értéktől való átlagos eltérés a populációban.

A varianciára adható egy könnyebben számolható formula is:

$$\begin{aligned}
 \text{Var}(\xi) &= \int_{-\infty}^{\infty} (x - E(\xi))^2 f_\xi(x) dx \\
 &= \int_{-\infty}^{\infty} x^2 f_\xi(x) dx - \int_{-\infty}^{\infty} 2E(\xi)x f_\xi(x) dx + \int_{-\infty}^{\infty} (E(\xi))^2 f_\xi(x) dx \\
 &= \int_{-\infty}^{\infty} x^2 f_\xi(x) dx - 2E(\xi) \int_{-\infty}^{\infty} x f_\xi(x) dx + (E(\xi))^2 \int_{-\infty}^{\infty} f_\xi(x) dx \\
 &= \int_{-\infty}^{\infty} x^2 f_\xi(x) dx - 2E(\xi)E(\xi) + (E(\xi))^2 \cdot 1 = \int_{-\infty}^{\infty} x^2 f_\xi(x) dx - (E(\xi))^2
 \end{aligned}$$

**Feladat:** Mennyi a  $\xi$  változó módusza, várható értéke és szórása a jelen feladatban?

A függvénynek csak egy maximumhelye van, az  $x = 2$  helyen, ez a módusz.

$$\begin{aligned} E(\xi) &= \int_{-\infty}^{\infty} x f_{\xi}(x) dx = \int_{-\infty}^2 x \cdot 0 dx + \int_2^4 x \cdot \frac{4}{x^2} dx + \int_4^{\infty} x \cdot 0 dx \\ &= 0 + 4 \int_2^4 x^{-1} dx + 0 = 4 \left[ \ln x \right]_2^4 = 4 [\ln 4 - \ln 2] = 2,77, \end{aligned}$$

$$\int_{-\infty}^{\infty} x^2 f_{\xi}(x) dx = \int_2^4 x^2 \frac{4}{x^2} dx = \int_2^4 4 dx = (4 - 2) \cdot 4 = 8,$$

$$\text{Var}(\xi) = \int_{-\infty}^{\infty} x^2 f_{\xi}(x) dx - (E(\xi))^2 = 8 - (2,77)^2 \approx 0,33,$$

$$D(\xi) = \sqrt{\text{Var}(\xi)} = \sqrt{0,33} = 0,57.$$

Tehát a  $\xi$  változó a 2 érték (módusz) közelébe esik a legnagyobb eséllyel. A változó átlagos értéke 2,77, a várható értéktől való átlagos eltérés 0,57.

Egy tetszőleges  $\xi$  valószínűségi változó **eloszlásfüggvénye** a következő módon van definiálva:  $F_\xi : \mathbb{R} \rightarrow [0, 1]$ ,  $F_\xi(t) = P(\xi < t)$ .

Ha a  $\xi$  értéket úgy kapjuk, hogy véletlenszerűen kiválasztunk egy egyedet egy populációból, és megmérünk egy kérdéses mennyiséget, akkor

$F_\xi(t)$  = azon egyedek aránya a populációban, melyeknél  $\xi$  kisebb, mint  $t$

Tetszőleges  $a$  és  $b$  valós számok esetén teljesülnek az alábbi egyenlőségek:

- 1  $P(\xi < a) = F_\xi(a)$ ,
- 2  $P(\xi \geq a) = 1 - F_\xi(a)$ ,
- 3  $P(a \leq \xi < b) = F_\xi(b) - F_\xi(a)$ .

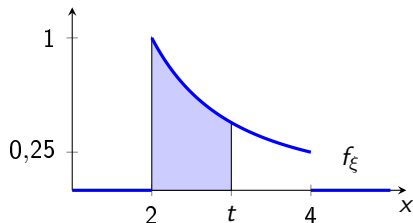
Hogyan kapjuk meg ezeket az azonosságokat?

- 1 Ez csak az eloszlásfüggvény definíciója.
- 2  $P(\xi \geq a) = 100\% - P(\xi < a) = 1 - F_\xi(a)$ ,
- 3  $P(a \leq \xi < b) = P(\xi < b) - P(\xi < a) = F_\xi(b) - F_\xi(a)$ .

**Feladat:** Hogyan írható fel az eloszlásfüggvény a jelen feladatban?

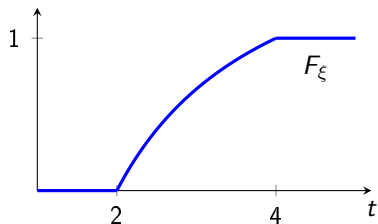
Sűrűségfüggvény:

$$f_{\xi}(x) = \begin{cases} 4/x^2, & \text{ha } 2 \leq x \leq 4, \\ 0, & \text{különben.} \end{cases}$$



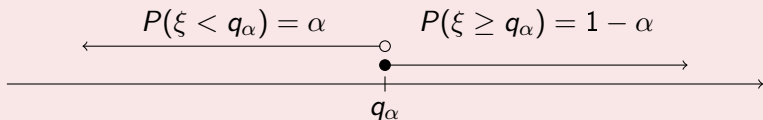
Eloszlásfüggvény:  $F_{\xi}(t) = P(\xi < t)$

- Ha  $t < 2$ :  $F_{\xi}(t) = 0$ .
- Ha  $t > 4$ :  $F_{\xi}(t) = 1$ .
- Ha  $2 \leq t \leq 4$ :



$$F_{\xi}(t) = P(2 \leq \xi \leq t) = \int_2^t \frac{4}{x^2} dx = 4 \int_2^t x^{-2} dx = 4 \left[ -\frac{1}{x} \right]_2^t = 2 - \frac{4}{t}.$$

Legyen  $\xi$  tetszőleges valószínűségi változó, és legyen  $\alpha \in (0, 1)$ . A  $\xi$  változó  **$\alpha$ -kvantilise** egy olyan  $q_\alpha$  valós szám, melyre  $P(\xi < q_\alpha) = \alpha$ .



A kvantilis jelentése: a vizsgált  $\xi$  mennyiség a teljes populáción belül

- az egyedek  $\alpha$  hányadánál kisebb, mint  $q_\alpha$ ,
- az egyedek  $1 - \alpha$  hányadánál nagyobb vagy egyenlő, mint  $q_\alpha$ .

Megjegyzés: Az  $\alpha$ -kvantilis nem mindig létezik, és ha létezik, akkor nem feltétlenül egyértelmű.

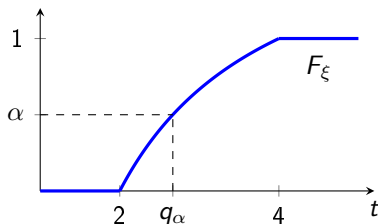
Nevezetes kvantilisek:

- **Medián:**  $q_{50\%}$
- **Alsó és felső kvartilisek:**  $q_{25\%}$  és  $q_{75\%}$
- **Decilisek:**  $q_{10\%}, q_{20\%}, \dots, q_{90\%}$



**Feladat:** Adjuk meg a mediánt valamint az alsó és a felső kvartilist a jelen feladatban.

$$F_{\xi}(t) = \begin{cases} 0, & t < 2, \\ 2 - 4/t, & 2 \leq t \leq 4, \\ 1, & t > 4. \end{cases}$$

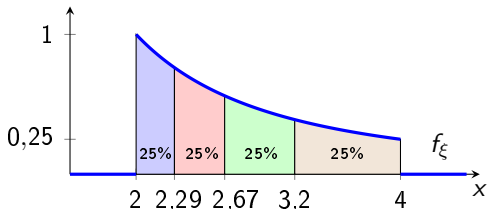


Tetszőleges  $0 < \alpha < 1$  szám esetén

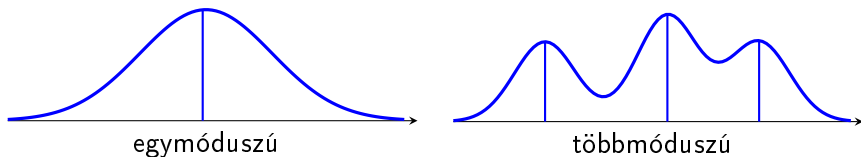
$$\alpha = P(\xi < q_\alpha) = F_{\xi}(q_\alpha) = 2 - 4/q_\alpha.$$

Ebből következik, hogy  $q_\alpha = 4/(2 - \alpha)$ . A kérdéses értékeket az alábbi táblázat tartalmazza. A medián és a két kvartilis négy részre bontja fel a változó értékkészletét, és a változó mindegyik részbe 25% eséllyel esik bele.

$\alpha$	25%	50%	75%
$q_\alpha$	2,29	2,67	3,2



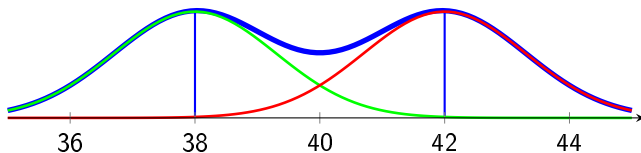
A móduszok száma alapján kétfajta sűrűségfüggvényt különböztetünk meg: **egymódusú** és **többmódusú** sűrűségfüggvényt.



A több módusz gyakran arra utal, hogy a populációt több részpopulációra lehet felbontani, melyeken belül a vizsgált  $\xi$  mennyiség már egymódusú.

Példa: lábméret eloszlása a felnőtt népességben belül.

- kék görbe: a lábméret sűrűségfüggvénye a felnőtt népességben belül,
- zöld görbe: a lábméret sűrűségfüggvénye a nők körében,
- piros görbe: a lábméret sűrűségfüggvénye a férfiak körében.

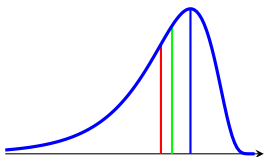


Tegyük fel, hogy a sűrűségfüggvénynek csak egyetlen módusza van. A módusz, a medián és a várható érték jelentése:

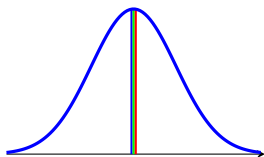
- **Módusz:** A változó ezen érték közelébe esik a legnagyobb eséllyel.
- **Medián:** A változó „középső” értéke.
- **Várható érték:** A változó átlagos értéke.

Ha a sűrűségfüggvény szimmetrikus, akkor a három mennyiség megegyezik.  
Ha a sűrűségfüggvény nem szimmetrikus, akkor jellemzően(!):

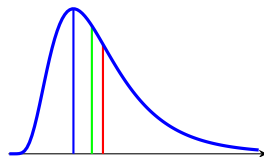
- Balra ferde sűrűségfüggvény esetén: **várható érték** < **medián** < **módusz**
- Jobbra ferde sűrűségfüggvény esetén: **módusz** < **medián** < **várható érték**



balra ferde függvény



szimmetrikus függvény

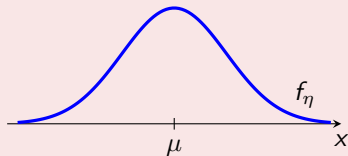


jobbra ferde függvény

# A normális eloszlás

Az  $\eta$  valószínűségi változó **normális** (másnéven **normál** vagy **Gauss-**) eloszlást követ  $\mu \in \mathbb{R}$  (mű) és  $\sigma > 0$  (szigma) paraméterekkel, ha a sűrűségfüggvénye:

$$f_{\eta}(x) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$



A sűrűségfüggvény neve: **Gauss-görbe**, **haranggörbe**.

A normális eloszlás fontosabb alkalmazásai:

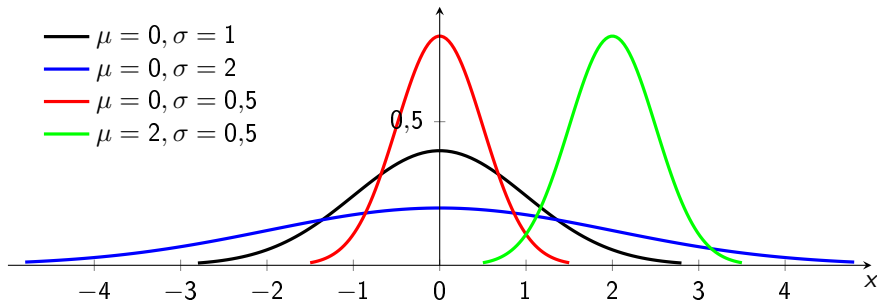
- Mérési hibák modellezése: mért érték = igazi érték + mérési hiba, ahol a mérési hiba normális eloszlást követ.
- Élettudományok: számos mennyiség (testmagasság, vérnyomás, IQ) normális vagy a normálisból származtatott eloszlást követ.

A normális eloszlás tulajdonságai:

- $f_{\eta}(x) > 0$  minden  $x$  valós számra, ezért  $R_{\eta} = \mathbb{R}$ .
- $E(\eta) = \mu$  és  $D(\eta) = \sigma$ .
- A sűrűségfüggvény szimmetrikus, ezért módusz = medián =  $E(\xi) = \mu$ .

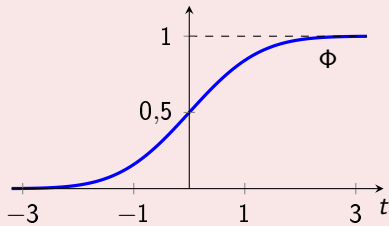
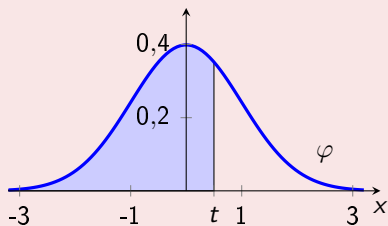
Hogyan hat a két paraméter a sűrűségfüggvényre:

- $\sigma$ : a sűrűségfüggvény alakját határozza meg,
- $\mu$ : eltolás, a sűrűségfüggvény szimmetriatengelye.



A  $\mu = 0$  és  $\sigma = 1$  paraméteres normális eloszlást **standard normális eloszlásnak** nevezzük. Jelölésben:  $\eta_{0,1}$ . Sűrűség- és eloszlásfüggvénye:

$$\varphi(x) = \frac{1}{\sqrt{2\pi}} e^{-x^2/2}, \quad \Phi(t) = P(\eta_{0,1} < t) = \int_{-\infty}^t \varphi(x) dx.$$



A  $\Phi$  függvény tulajdonságai:

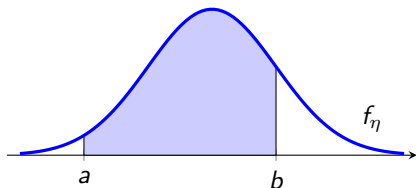
$$\Phi(t) \begin{cases} < 0,5, & \text{ha } t < 0, \\ = 0,5, & \text{ha } t = 0, \\ > 0,5, & \text{ha } t > 0, \end{cases}$$

$$\Phi(-t) = 1 - \Phi(t).$$

Ha  $\eta$  normális eloszlású, akkor tetszőleges  $a$  és  $b$  valós számokra:

$$P(a \leq \eta \leq b) = \int_a^b f_\eta(x) dx = \int_a^b \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}} dx$$

Probléma: ezt az integrált nem tudjuk papíron kiszámolni.



Legyen  $\eta$  normális eloszlású  $\mu$  várható értékkel és  $\sigma$  szórással. Ekkor az  $(\eta - \mu)/\sigma$  valószínűségi változót  $\eta$  **standardizáltjának** nevezzük. Megmutatható, hogy ez az új változó standard normális eloszlás követ.

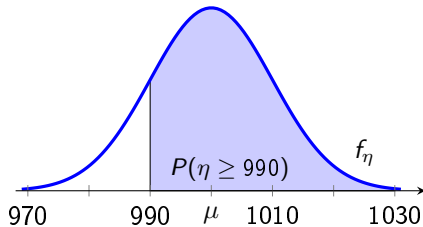
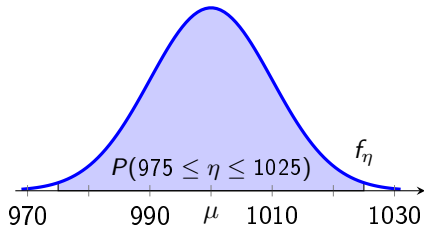
Ha  $\eta$  normális eloszlású változó, akkor standardizálással:

$$P(a \leq \eta < b) = P\left(\frac{a - \mu}{\sigma} \leq \frac{\eta - \mu}{\sigma} < \frac{b - \mu}{\sigma}\right) = \Phi\left(\frac{b - \mu}{\sigma}\right) - \Phi\left(\frac{a - \mu}{\sigma}\right).$$

**Feladat.** Egy tejgyárban az 1 literes dobozos tej csomagolását automata töltőberendezés végzi, és a dobozokba töltött mennyiség egy normális eloszlású valószínűségi változó, melynek várható értéke a névleges tartalom és szórása  $\sigma = 10$  ml. Véletlenszerűen kiválasztunk egy dobozt.

- Mennyi annak a valószínűsége, hogy a doboz legfeljebb 2,5%-kal tér el a névleges tartalomtól?
- Mennyi annak az esélye, hogy a doboz legalább 990 ml tejet tartamaz?

Legyen  $\eta$  a kiválasztott dobozban található mennyiség. Az  $\eta$  változó normális eloszlású  $\mu = 1000$  ml várható értékkel és  $\sigma = 10$  ml szórással. A következő valószínűségekre (=területekre) vagyunk kíváncsiak, de ezek most nem számolhatóak ki integrálással:





Az első valószínűség standardizálással határozható meg:

$$\begin{aligned} P(975 \leq \eta \leq 1025) &= P\left(\frac{975 - 1000}{10} \leq \frac{\eta - \mu}{\sigma} \leq \frac{1025 - 1000}{10}\right) \\ &= P(-2,5 \leq \eta_{0,1} \leq 2,5) = \Phi(2,5) - \Phi(-2,5) = 0,9938 - 0,0062 = 0,9876, \end{aligned}$$

Ez azt jelenti, hogy a tejesdobozok 98,76%-a tartalmaz 975 ml és 1025 ml közötti tejet. Itt felhasználtuk azt, hogy

$$\Phi(-2,5) = 1 - \Phi(2,5) = 1 - 0,9938 = 0,0062.$$

A második valószínűség az első mintájára:

$$\begin{aligned} P(\eta \geq 990) &= P\left(\frac{\eta - \mu}{\sigma} \geq \frac{990 - 1000}{10}\right) = P(\eta_{0,1} \geq -1) \\ &= 1 - P(\eta_{0,1} < -1) = 1 - \Phi(-1) = 1 - [1 - \Phi(1)] \\ &= 1 - [1 - 0,84] = 0,84. \end{aligned}$$

**Feladat.** Adjunk meg egy olyan  $[a, b]$  intervallumot, amire teljesül, hogy a tejesdobozok 95%-a ebbe az intervallumba esik:  $P(a \leq \eta \leq b) = 0,95$ .

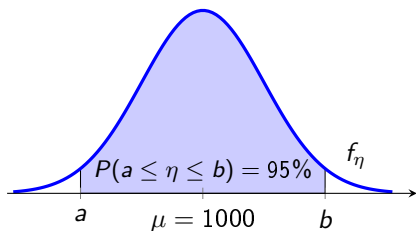
Az intervallumot  $[\mu - c\sigma, \mu + c\sigma]$  alakban fogjuk keresni. Ismét csak standardizálással:

$$0,95 = P(\mu - c\sigma \leq \eta \leq \mu + c\sigma) = P\left(-c \leq \frac{\eta - \mu}{\sigma} \leq c\right)$$

$$= P(-c \leq \eta_{0,1} \leq c) = \Phi(c) - \Phi(-c) = \Phi(c) - [1 - \Phi(c)] = 2\Phi(c) - 1.$$

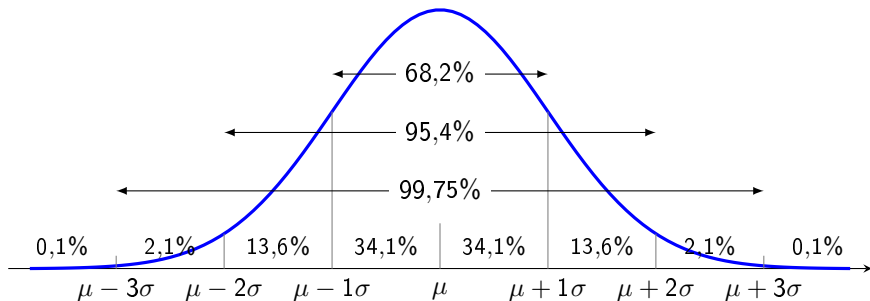
Ebből azt kapjuk, hogy  $\Phi(c) = 0,975 = \Phi(1,96)$ , tehát  $c = 1,96$ .

Tehát a kérdéses intervallum:  $[\mu - 1,96\sigma, \mu + 1,96\sigma] = [980,4, 1019,6]$ .



Közelítő intervallum a  $2\sigma$ -szabállyal:  $[\mu - 2\sigma, \mu + 2\sigma] = [980, 1020]$ .

Az alábbi ábra azt mutatja meg, hogy egy  $\eta$  normális eloszlású változó mekkora eséllyel esik a várható érték két oldalára felmért intervallumokba:



Legyen  $\eta$  normális eloszlású változó. Ekkor:

- 1 $\sigma$ -szabály:  $P(\mu - \sigma \leq \eta \leq \mu + \sigma) \approx 68\%$ ,
- 2 $\sigma$ -szabály:  $P(\mu - 2\sigma \leq \eta \leq \mu + 2\sigma) \approx 95\%$ ,
- 3 $\sigma$ -szabály:  $P(\mu - 3\sigma \leq \eta \leq \mu + 3\sigma) \approx 99,75\%$ .

# Statisztikai alapfogalmak

Legyen adva egy populáció, és tekintsünk egy mennyiséget az egyedeken (életkor, testtömeg, utódok száma, stb.). Véletlenszerűen kiválasztunk egy egyedet, és  $\xi$  jelöli a vizsgált mennyiséget a kiválasztott egyed esetében.

Valószínűségszámítás: Ha ismerjük a  $\xi$  változó valószínűségeloszlását vagy sűrűségfüggvényét, akkor ki tudjuk számolni a következő értékeket:

- $E(\xi)$  = a vizsgált mennyiség átlagos értéke a populáción belül,
- $D(\xi)$  = a vizsgált mennyiség szórása a populáción belül,
- $P(a \leq \xi \leq b)$  = arány a teljes populáción belül.

Matematikai statisztika: Nem ismerjük a  $\xi$  változó valószínűségeloszlását vagy sűrűségfüggvényét, ezért nem tudjuk kiszámolni ezeket az értékeket. Ehelyett megfigyeléseket végzünk a  $\xi$  változóra, és a kapott minta alapján vonunk le következtetéseket. Célok:

- Becslélmélet: Adjunk becslést a várható értékre, szórásra, stb.
- Hipotézisvizsgálat: Adott egy állítás a  $\xi$  mennyiséggel kapcsolatban. (Pl:  $E(\xi) = 2.$ ) Döntsük el, hogy ez az állítás igaz vagy hamis.

## Statisztikai alapfogalmak:

- **Háttérváltozó:** Az a  $\xi$  valószínűségi változó, melyet vizsgálunk.
- **Statisztikai minta (statistical sample):**  $\xi_1, \dots, \xi_n$  valószínűségi változók, független megfigyelések a  $\xi$  változóra. Jellemzően: véletlenszerűen kiválasztunk  $n$  egyedet a teljes populációból.
- **Mintarealizáció (realization, observations):** a  $\xi_1, \dots, \xi_n$  változók megfigyelés során kapott konkrét értékei.
- **Mintaméret (sample size):** a megfigyelések száma ( $n$ ).

## Hogyan is történik ez a gyakorlatban:

- Kíváncsiak vagyunk egy  $\xi$  mennyiség eloszlására egy populációban.
- Megtervezzük a mintavételezést és a statisztikai kiértékelést. Ezen a ponton a mintaelemek valószínűségi változók: még nem tudjuk, hogy mik lesznek a megfigyelt értékek.
- Elvégezzük a mintavételezést, ezzel megkapjuk a realizációt, tehát a mintaelemek konkrét értékeit.
- Elvégezzük a statisztikai elemzést a realizáción. (Mi a továbbiakban nagyrészt ezzel a lépéssel foglalkozunk.)

# Leíró statisztikák (descriptive statistics)

Egy  $\xi$  háttérváltozó várható értékét, varianciáját és szórását a következő módon becsülhetjük meg egy  $\xi_1, \dots, \xi_n$  minta alapján:

- **Empirikus várható érték, mintaátlag (sample mean):**

$$\bar{\xi} = E_n(\xi) = \frac{\xi_1 + \dots + \xi_n}{n} \approx E(\xi)$$

- **Empirikus variancia (sample variance):**

$$\text{Var}_n(\xi) = \frac{(\xi_1 - \bar{\xi})^2 + \dots + (\xi_n - \bar{\xi})^2}{n} \approx \text{Var}(\xi)$$

- **Empirikus szórás (standard deviation):**  $D_n(\xi) = \sqrt{\text{Var}_n(\xi)} \approx D(\xi)$

Miért így van definiálva az empirikus variancia?

$$\text{Var}(\xi) = E\left([\xi - E(\xi)]^2\right) \approx \frac{[\xi_1 - E(\xi)]^2 + \dots + [\xi_n - E(\xi)]^2}{n} \approx \text{Var}_n(\xi)$$

Az előző oldalon felsorolt becslések **erősen konzisztensek**, tehát

$$E_n(\xi) \rightarrow E(\xi), \quad \text{Var}_n(\xi) \rightarrow \text{Var}(\xi), \quad D_n(\xi) \rightarrow D(\xi), \quad n \rightarrow \infty.$$

Ez azt jelenti, hogy ezek a becslések nagy  $n$  esetén pontosak lesznek.

Probléma: kis  $n$  esetén  $\text{Var}_n(\xi)$  és  $D_n(\xi)$  tipikusan alábecsli az igazi varianciát és szórást. Megoldás: kicsit megnöveljük ezeket az értékeket.

**Korrigált empirikus variancia és korrigált empirikus szórás:**

$$\text{Var}_n^*(\xi) = \frac{n}{n-1} \text{Var}_n(\xi) \approx \text{Var}(\xi), \quad D_n^*(\xi) = \sqrt{\text{Var}_n^*(\xi)} \approx D(\xi).$$

Nagy mintaméret esetén a korrigálás csak kis mértékben változtat a becsléseken. Kis mintaméret esetén viszont jelentős a növekedés.

A korrigálás során kapott becslések kis  $n$  esetén pontosabban, mint az eredeti becslések, de az erős konzisztencia is megmarad:

$$\text{Var}_n^*(\xi) \rightarrow \text{Var}(\xi), \quad D_n^*(\xi) \rightarrow D(\xi), \quad n \rightarrow \infty.$$

**Feladat:** A kar férfi hallgatóinak testmagasságát vizsgáljuk, jelölje  $\xi$  egy véletlenszerűen kiválasztott férfi hallgató magasságát. Megfigyeléseket végzünk a változóra, a következő realizációt kapjuk: 180, 175, 188, 168, 173, 183. Adjunk becslést a testmagasság átlagára és szórására.

$$\bar{\xi} = E_6(\xi) = \frac{180 + 175 + 188 + 168 + 173 + 183}{6} = 177,8 \approx E(\xi),$$

$$\text{Var}_6(\xi) = \frac{(180 - 177,8)^2 + \dots + (183 - 177,8)^2}{6} = 43,81 \approx \text{Var}(\xi),$$

$$D_6(\xi) = \sqrt{43,81} = 6,62 \approx D(\xi).$$

A kis mintaméret miatt ( $n = 6$ ) a szórást jobb a korrigált szórással becsülni:

$$\text{Var}_6^*(\xi) = \frac{6}{5} 43,81 = 52,57, \quad D_6^*(\xi) = \sqrt{52,57} = 7,25 \approx D(\xi).$$

Foglaljuk össze, hogy mit kaptunk:

- átlagos testmagasság a populációban =  $E(\xi) \approx 177,8$ ,
- a testmagasság szórása a populációban =  $D(\xi) \approx 7,25$ .

Ezt a két értéket publikációkban így szokták közölni:  $177,8 \pm 7,25$  cm.



Ha van egy mintarealizációnk, akkor a mintaátlag egy becslés az ismeretlen várható értékre. Ha egy másik mintavételből származó másik realizációval dolgozunk, akkor egy másik becslést kapunk ugyanarra a várható értékre. A mintaátlag egy valószínűségi változó, ami a realizációtól függ.

**Tétel.** A mintaátlag várható értéke és szórása:

$$E(\bar{\xi}) = E(\xi) \quad \text{és} \quad D(\bar{\xi}) = D(\xi)/\sqrt{n}.$$

Értelmezzük a kapott eredményeket:

- Ha minden lehetséges realizációból kiszámolnánk a mintaátlagot, akkor átlagban a várható értéket kapnánk. Ez egy jó tulajdonság, amit **tozítatlanságnak** nevezünk.
- Ha minden lehetséges realizációból kiszámolnánk a mintaátlagot, akkor ezek az értékek átlagosan  $D(\xi)/\sqrt{n}$  mértékben térnek el a becsülni kívánt  $E(\xi)$  várható értéktől. Tehát átlagosan ennyit tévedünk a becslés során.

Vegyük észre:  $D(\xi)/\sqrt{n} \rightarrow 0$ , amint  $n \rightarrow \infty$ . Ez azt jelenti, hogy egyre nagyobb minta alapján egyre kisebb hibával tudunk becsülni.

**Standard hiba (standard error of the mean, s.e.m.):**  $SE = D_n^*(\xi)/\sqrt{n}$ .

Jelentése: a  $D(\bar{\xi})$  szórás becslése a minta alapján.

- Ha a standard hiba kicsi, akkor a mintaátlag minden realizáció esetén pontos becslése lesz a várható értéknek.
- Ha a standard hiba nagy, akkor vannak olyan realizációk, melyekre a mintaátlag pontatlan becslést ad a várható értékre.

**Feladat:** Határozzuk meg a standard hibát a jelen feladatban.

Amit tudunk:  $n = 6$ ,  $E_6(\xi) = 177,8$ ,  $D_6^*(\xi) = 7,25$ .

Ekkor:  $SE = 7,25/\sqrt{6} = 2,96$ .

Foglaljuk össze, hogy mit kaptunk:

- Az ismeretlen várható értékre adott becslésünk: 177,8. Ez csak egy becslés, nem fogja pontosan telibe találni az igazi várható értéket.
- A standard hiba: 2,96. A mintaátlag várhatóan ennyivel tér el az igazi várható értéktől, várhatóan ennyi a becslés hibája.
- Ezt a két értéket így szokták közölni:  $177,8 \pm 2,96$  (SE).

A  $\xi$  valószínűségi változó  $\alpha$ -kvantilise egy olyan  $q_\alpha$  valós szám, melyre  $P(\xi < q_\alpha) = \alpha$ . Jelentése: a populáción belül a vizsgált  $\xi$  mennyiség az egyedek  $\alpha$  hányadánál kisebb, mint  $q_\alpha$ .

Az  $\alpha$ -kvantilis becslésére egy  $\xi_1, \dots, \xi_n$  statisztika minta alapján több módszer is létezik. Mi most nem adunk precíz matematikai formulát a becslésre, csak a becslés alapötletét ismertetjük.

**Empirikus kvantilis, percentilis (percentile):** Az a  $\hat{q}_\alpha$  szám, melyre teljesül, hogy a  $\xi_1, \dots, \xi_n$  értékek  $\alpha$  hányada kisebb, mint  $\hat{q}_\alpha$ .

Például: **empirikus medián:**

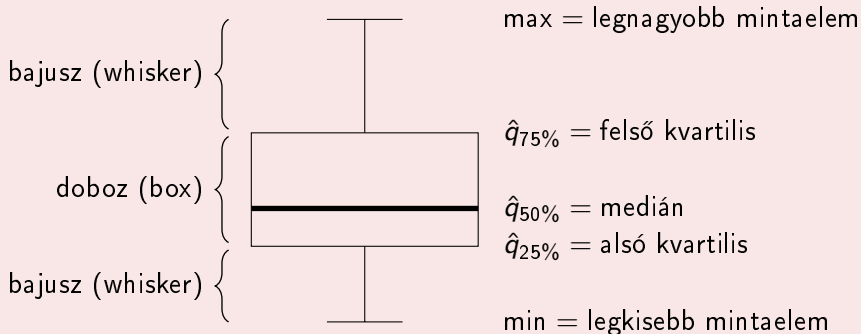
$$\hat{q}_{50\%} = \begin{cases} \text{a középső mintaelem,} & \text{ha } n \text{ páratlan,} \\ \text{a két középső átlaga,} & \text{ha } n \text{ páros.} \end{cases}$$

**Feladat:** Adjunk becslést a testmagasság elméleti mediánjára a kar férfi hallgatóinak populációjában.

A rendezett minta: 168, 173, 175, 180, 183, 188. A becslés:

$$q_{50\%} \approx \hat{q}_{50\%} = \text{két középső mintaelem átlaga} = 177,5.$$

A **boxplot** egy olyan grafikon, mely az alábbi statisztikai mutatószámokat ábrázolja egyszerű formában:



További mutatószámok:

- **Terjedelem (range)** =  $\max - \min$  = a boxplot magassága,
- **Interkvartilis távolság (interquartile range):**

$\text{IQR} = \text{felső kvartilis} - \text{alsó kvartilis} = \text{a doboz magassága.}$

# Konfidencia intervallumok (confidence intervals)

A statisztikában egy minta alapján kétféle formában becsülhetjük meg az ismeretlen mennyiségeket (várható érték, szórást, stb.):

- **Pontbecslés:** Az ismeretlen mennyiséget egyetlen számmal becsüljük meg, és reménykedünk benne, hogy nem tévedünk nagyot.
- **Intervallumbecslés:** Egy intervallumot adunk meg, mely nagy megbízhatósággal tartalmazza a kérdéses mennyiséget.

Legyen  $\xi_1, \dots, \xi_n$  statisztikai minta egy  $\xi$  valószínűségi változóra, és legyen  $\alpha \in (0, 1)$ . A minta alapján felírt  $[a, b]$  intervallum egy  $1 - \alpha$  megbízhatóságú **konfidencia intervallum a várható értékre**, ha

$$P(E(\xi) \in [a, b]) = 1 - \alpha.$$

- A megbízhatóság általában 90%, 95% vagy 99% szokott lenni, a biostatisztikában tipikusan a 95%-ot használják.
- A konfidencia intervallum hasonló módon definiálható tetszőleges más mutatószámra is (szórás, variancia, medián, stb.)

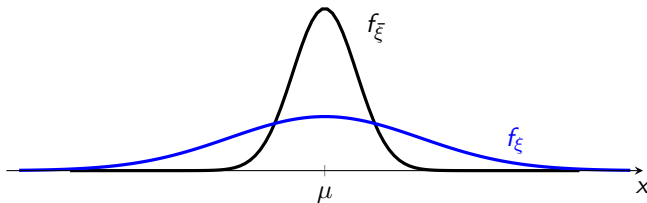
**Feladat:** Legyen  $\xi$  normális eloszlású valószínűségi változó ismeretlen  $\mu$  várható értékkel és ismert  $\sigma$  szórással. Egy  $\xi_1, \dots, \xi_n$  statisztikai minta alapján adjunk konfidencia intervallumot a várható értékre.

**Tétel.** Ha a  $\xi$  háttérváltozó normális eloszlású, akkor a  $\xi_1 + \dots + \xi_n$  összeg és a  $\bar{\xi} = (\xi_1 + \dots + \xi_n)/n$  mintaátlag is normális eloszlású változó.

Jelölje  $\mu_{\bar{\xi}}$  és  $\sigma_{\bar{\xi}}$  a mintaátlag várható értékét és szórását. Ekkor

- $\mu_{\bar{\xi}} = E(\bar{\xi}) = E(\xi) = \mu,$
- $\sigma_{\bar{\xi}} = D(\bar{\xi}) = D(\xi)/\sqrt{n} = \sigma/\sqrt{n}.$

Az alábbi ábrán a  $\xi$  háttérváltozó és a  $\bar{\xi}$  mintaátlag sűrűségfüggvénye látható:



Először megadunk egy olyan intervallumot, mely  $1 - \alpha$  valószínűséggel tartalmazza a  $\bar{\xi}$  változót. Az intervallumot most is  $[\mu_{\bar{\xi}} - c\sigma_{\bar{\xi}}, \mu_{\bar{\xi}} + c\sigma_{\bar{\xi}}]$  alakban keressük. Standardizálással:

$$\begin{aligned} 1 - \alpha &= P(\mu_{\bar{\xi}} - c\sigma_{\bar{\xi}} \leq \bar{\xi} \leq \mu_{\bar{\xi}} + c\sigma_{\bar{\xi}}) = P\left(-c \leq \frac{\bar{\xi} - \mu_{\bar{\xi}}}{\sigma_{\bar{\xi}}} \leq c\right) \\ &= P(-c \leq \eta_{0,1} \leq c) = \Phi(c) - \Phi(-c) = \Phi(c) - [1 - \Phi(c)] = 2\Phi(c) - 1 \end{aligned}$$

Tehát  $\Phi(c) = 1 - \alpha/2$ , amiből  $c = \Phi^{-1}(1 - \alpha/2)$ . Ezt az értéket ki tudjuk keresni a táblázatból tetszőleges  $\alpha \in (0, 1)$  esetén.

A fenti nagy formulát a következő módon tudjuk továbbalakítani:

$$\begin{aligned} 1 - \alpha &= P(\mu_{\bar{\xi}} - c\sigma_{\bar{\xi}} \leq \bar{\xi} \leq \mu_{\bar{\xi}} + c\sigma_{\bar{\xi}}) = P(-\bar{\xi} - c\sigma_{\bar{\xi}} \leq -\mu_{\bar{\xi}} \leq -\bar{\xi} + c\sigma_{\bar{\xi}}) \\ &= P(\bar{\xi} + c\sigma_{\bar{\xi}} \geq \mu_{\bar{\xi}} \geq \bar{\xi} - c\sigma_{\bar{\xi}}) = P\left(\bar{\xi} + c\frac{\sigma}{\sqrt{n}} \geq \mu \geq \bar{\xi} - c\frac{\sigma}{\sqrt{n}}\right) \end{aligned}$$

De hát ez éppen egy konfidencia intervallum az  $E(\xi) = \mu$  ismeretlen várható értékre:

$$1 - \alpha = P\left(E(\xi) \in \left[\bar{\xi} - c\frac{\sigma}{\sqrt{n}}, \bar{\xi} + c\frac{\sigma}{\sqrt{n}}\right]\right)$$

Legyen  $\xi$  normális eloszlású változó ismert  $\sigma$  szórással. Ekkor a változó várható értékére a következő formában adható  $1 - \alpha$  megbízhatóságú konfidencia intervallum:

$$\left[ \bar{\xi} - c \frac{\sigma}{\sqrt{n}}, \bar{\xi} + c \frac{\sigma}{\sqrt{n}} \right], \quad c = \Phi^{-1} \left( 1 - \frac{\alpha}{2} \right).$$

**Feladat:** Tegyük fel, hogy a kar férfi hallgatóinak testmagassága normális eloszlású  $\sigma = 7$  cm szórással. Adjunk 95% megbízhatóságú konfidencia intervallumot a testmagasság várható értékére (az átlagos testmagasságra).

A minta: 180, 175, 188, 168, 173, 183.

A mintaméret és a mintaátlag:  $n = 6$ ,  $\bar{\xi} = 177,8$ .

Most  $\alpha = 5\% = 0,05$ , tehát  $c = \Phi^{-1}(0,975) = 1,96$ .

Az intervallum:

$$\left[ 177,8 - 1,96 \frac{7}{\sqrt{6}}, 177,8 + 1,96 \frac{7}{\sqrt{6}} \right] = [172,2, 183,4].$$

De mi ennek az intervallumnak a jelentése?



Probléma: a  $\xi$  háttérváltozó igazi szórását sosem tudjuk.

Megoldás: helyettesítsük a szórást a becslésével:  $\sigma \approx D_n^*(\xi)$ . Ennek az az ára, hogy a  $c$  értéket a **Student-eloszlás** táblázatából kell kikeresni.

Legyen  $\xi$  normális eloszlású változó ismeretlen szórással. Egy  $1 - \alpha$  megbízhatóságú konfidencia intervallum a változó várható értékére:

$$\left[ \bar{\xi} - c \frac{D_n^*(\xi)}{\sqrt{n}}, \bar{\xi} + c \frac{D_n^*(\xi)}{\sqrt{n}} \right] = [\bar{\xi} - c SE, \bar{\xi} + c SE], \quad c = \Phi_{n-1}^{-1} \left( 1 - \frac{\alpha}{2} \right).$$

Itt  $\Phi_{n-1}$  az  $n - 1$  szabadsági fokú Student-eloszlás eloszlásfüggvénye.

**Feladat:** Adjunk 95% megbízhatóságú konfidencia intervallumot a kar férfi hallgatóinak átlagos testmagasságra ismeretlen szórás esetén!

Most:  $n = 6$ ,  $\bar{\xi} = 177,8$ ,  $D_6^*(\xi) = 7,25$ ,  $c = \Phi_5^{-1}(0,975) = 2,57$ .

Az intervallum:

$$\left[ 177,8 - 2,57 \frac{7,25}{\sqrt{6}}, 177,8 + 2,57 \frac{7,25}{\sqrt{6}} \right] = [170,2, 185,4].$$

**Kérdés:** Hogyan értelmezhető a kapott eredmény?

A mintavételezés során a véletlen sok különböző mintarealizációt sorsolhat ki nekünk. Ezek két csoportba sorolhatóak:

- „Jó” mintarealizációk: az ezekből számolt konfidencia intervallum tartalmazza az ismeretlen várható értéket. Ezek teszik ki az összes lehetséges mintarealizáció  $1 - \alpha = 0,95$  hányadát.
- „Rossz” mintarealizációk: ezek félrevezetőek, ugyanis a belőlük számolt konfidencia intervallum nem tartalmazza a várható értéket. Ezek alkotják az összes realizáció  $\alpha = 0,05$  hányadát.

**Kérdés:** Ebben a feladatban jó vagy rossz mintarealizációt kaptunk?

Ezt nem tudjuk eldönteni. Csak reménykedhetünk benne, hogy a jók közül kaptunk egyet, ugyanis ezek vannak többségben.

**Kérdés:** Ismeretlen szórás esetén miért kaptunk bővebb intervallumot?

Nem volt ismert a szórás, ami további bizonytalanságot jelentett. Emiatt egy kis „ráhagyással” kellett számolnunk: nagyobb lett a  $c$  érték, ami bővebb intervallumot eredményezett.

**Kérdés:** Hogyan értelmezhető az intervallum:  $[\bar{\xi} - cSE, \bar{\xi} + cSE]$ ?

A konfidencia intervallum felírásakor a  $\bar{\xi}$  mintaátlagból indulunk ki, ugyanis ez egy jó becslése a várható értéknek. Erre a becslésre mérjük fel a  $cSE$  szorzatot két oldalra. Ebben a szorzatban két dolog jelenik meg:

- A standard hiba számszerűsíti, hogy mennyire jól becsli a mintaátlag a várható értéket, mekkora „ráhagyással” kell számolni a konfidencia intervallum felírásakor.
- A  $c$  értékben a megbízhatóság jelenik meg:  
nagyobb megbízhatóság  $\Rightarrow$  magasabb  $c$  érték  $\Rightarrow$  bővebb intervallum.

**Kérdés:** Miért nem számolunk 99,99%-os megbízhatósággal?

A magasabb megbízhatóság szélesebb intervallumot jelent. A túl széles intervallum viszont nehezíti az eredmény alkalmazhatóságát.

A 95%-os választás jó egyensúlyt jelent a két cél (magas megbízhatóság és szűk konfidencia intervallum) között. A megbízhatóság további növelése drasztikusan szélesebb intervallumot eredményez. Csak akkor dolgozunk magasabb megbízhatósággal, ha a standard hiba alacsony.

**Kérdés:** Mi a helyzet akkor, ha a  $\xi$  nem normális eloszlású?

A levezetésnek a következő tétel volt az alapja: ha a  $\xi$  háttérváltozó normális eloszlású, akkor a  $\bar{\xi}$  mintaátlag is normális eloszlású változó.

**Tétel.** Ha a minta nem normális eloszlásból jön, de a mintaméret elég nagy, akkor a  $\bar{\xi}$  mintaátlag közel normális eloszlású.

A tételnek az a következménye, hogy a kapott intervallum egy közelítő konfidencia intervallum a várható értékre tetszőleges  $\xi$  háttérváltozó esetén:

$$P\left(E(\xi) \in [\bar{\xi} - c SE, \bar{\xi} + c SE]\right) \approx 1 - \alpha.$$

**Kérdés:** Mit jelent ebben az esetben az „elég nagy mintaméret”?

Erre a kérdésre nincs egyszerű válasz, a szükséges mintaméret attól függ, hogy a  $\xi$  változó eloszlása mennyire hasonlít a normális eloszláshoz:

- (közel) szimmetrikus eloszlás esetén 20–30 mintaelem tipikusan elég szokott lenni a pontos közelítéshez,
- ferde eloszlás esetén jellemzően kell legalább 50, vagy akár még annál is több mintaelem.

# Hipotézisvizsgálat

A hipotézisvizsgálat (hypothesis testing) alapfogalmai:

- Adott egy  $\xi$  háttérváltozó és egy  $\xi_1, \dots, \xi_n$  statisztikai minta.
- **Null-hipotézis ( $H_0$ , null hypothesis):** Egy állítás a  $\xi$  változóra.
- **Alternatív hipotézis ( $H_A$ , alternative hypothesis):** Egy másik állítás a  $\xi$  változóra.
- A hipotézisvizsgálat célja: A két hipotézis közül valamelyik igaz. Döntsük el a statisztikai minta alapján, hogy  $H_0$  vagy  $H_A$  igaz.

Például:  $H_0 : E(\xi) = 2$ ,  $H_A : E(\xi) = 4$ .

A továbbiakban a kurzuson az alternatív hipotézis mindig a nullhipotézis tagadása lesz. Azt kell eldönteni, hogy  $H_0$  igaz vagy nem. Például:

- $H_0 : P(\xi = 5) = 1/2$ ,  $H_A : P(\xi = 5) \neq 1/2$ .
- $H_0 : \xi$  normális eloszlású,  $H_A : \xi$  nem normális eloszlású.

A hipotézisvizsgálat menete:

- Eldöntjük, hogy milyen módszerrel tesztlünk.
- A statisztikai minta alapján kiszámoljuk a **próbat statisztika (test statistic)** értékét:  $s_n$ .
- Meghatározzuk a **kritikus értéket (critical value)**:  $c$ .
- Ha  $|s_n| \leq c$ , akkor **elfogadjuk (accept)** a nullhipotézist.  
Ha  $|s_n| > c$ , akkor **elvetjük (reject)** a nullhipotézist.

Az egész olyan, mint egy bírósági tárgyalás:

- A nullhipotézis a vádlott szava („ártatlan vagyok”).
- A statisztikai minta a bizonyítékok halmaza.
- A próbat statisztika ( $s_n$ ) azt fejezi ki, hogy a vádlott szava mennyire van ellentmondásban a bizonyítékokkal.
- A  $c$  kritikus érték egy küszöbérték. Ha  $|s_n| \leq c$ , akkor a bíró hisz a vádlottnak, és felmenti. Ha  $|s_n| > c$ , akkor nem hisz neki, és elítéli.

**Feladat:** A kar férfi hallgatóinak testmagasságát vizsgáljuk, jelölje  $\xi$  egy véletlenszerűen kiválasztott férfi hallgató magasságát. Mit állíthatunk  $\xi$  várható értékéről, az átlagos testmagasságról a teljes populáción belül?

- Megfigyelt értékek: 180, 175, 188, 168, 173, 183.
- Becslések:  $E(\xi) \approx \bar{\xi} = 177,8$ ,  $D(\xi) \approx D_n^*(\xi) = 7,25$ .

Teszteljük a következő nullhipotézist:  $H_0 : E(\xi) = 175$ .

Látni fogjuk, hogy a várható értéket a  $t$ -próba segítségével lehet tesztelni:

- Próbastatisztika:

$$s_n = \frac{\bar{\xi} - 175}{D_n^*(\xi)/\sqrt{n}} = \frac{177,8 - 175}{7,25/\sqrt{6}} = 0,946,$$

- A kritikus érték:  $c = 2,571$ . (Miért ennyi? Majd később kiderül.)
- Döntés:  $|s_n| \leq c$ , tehát a nullhipotézist elfogadjuk. A megfigyelt értékek nincsenek ellentmondásban a nullhipotézis állításával.

**Kérdés:** Biztosan jól döntöttünk? Biztos, hogy a nullhipotézis igaz?

Sajnos nem: ha félrevezető a minta, amivel dolgozunk, akkor helytelen következtetést vonhatunk le, és hibás döntést hozunk?

Milyen hibákat véthetünk a hipotézisvizsgálat során:

- **Elsőfajú hiba (type I error):** Elvetjük az igaz nullhipotézist, tehát börtönbe küldünk egy ártatlant. Valószínűsége:

$$\alpha = P(\text{elvetjük } H_0\text{-t} \mid H_0 \text{ igaz}).$$

- **Másodfajú hiba (type II error):** Elfogadjuk a hamis nullhipotézist, tehát felmentünk egy bűnöst. Valószínűsége:

$$\beta = P(\text{elfogadjuk } H_0\text{-t} \mid H_0 \text{ hamis}).$$

Még egy fogalom:

$$\text{erő (power)} = P(\text{elvetjük } H_0\text{-t} \mid H_0 \text{ hamis}) = 1 - \beta.$$

A lehetőségeket az alábbi táblázatban foglalhatjuk össze:

	elfogadjuk	elvetjük
$H_0$ igaz	helyes döntés	elsőfajú hiba
$H_0$ hamis	másodfajú hiba	helyes döntés



Mire hathatunk és mire nem a hipotézisvizsgálat során?

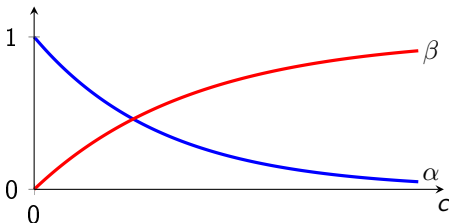
- Akkor vetjük el a nullhipotézist, ha  $|s_n| > c$ .
- A nullhipotézis, a tesztelési módszer és a statisztikai minta adott: az  $s_n$  próbastatisztika értékét nem tudjuk befolyásolni.
- A  $c$  kritikus értéket (=mennyire szigorú a bíró) mi választjuk.

Meg lehet választani úgy a kritikus értéket, hogy mindkét hiba alacsony maradjon? Erre sajnos nincs lehetőség:

alacsony elsőfajú hiba  $\Rightarrow$  magas kritikus érték  $\Rightarrow$  magas másodfajú hiba

alacsony másodfajú hiba  $\Rightarrow$  alacsony kritikus érték  $\Rightarrow$  magas elsőfajú hiba

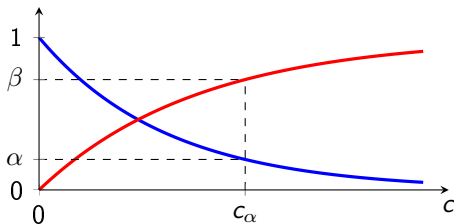
Adott  $n$  mintaméret esetén a kétfajta hiba nagysága egymással ellentétesen változik, ha módosítjuk a kritikus értéket:



A hipotézisvizsgálat során az  $\alpha$  elsőfajú hibát (szignifikancia szintet) előre meg szoktuk adni, és a kritikus értéket ennek megfelelően választjuk. A szignifikancia szint kicsi (tipikusan 1%, 5% vagy 10%) szokott lenni (ártatlanok védelme). A  $\beta$  másodfajú hibára nincsen ráhatásunk.

A kritikus érték meghatározása:

- A feladat megadja az  $\alpha$  szignifikancia szintet (=elsőfajú hiba).
- Meghatározzuk a hozzá tartozó kritikus értéket ( $c_\alpha$ ) és tesztelünk.
- A  $\beta$  másodfajú hiba lehet kicsi vagy nagy is, erre nincs ráhatásunk.



A hipotézisvizsgálat során megjelenő valószínűségek:

	elfogadjuk	elvetjük
$H_0$ igaz	$1 - \alpha$ (nagy)	$\alpha$ (kicsi)
$H_0$ hamis	$\beta$ (nem ismert)	$1 - \beta$ (nem ismert)

Hogyan lehet értelmezni a hipotézisvizsgálat eredményét?

- Ha elfogadjuk a nullhipotézist, az nem jelent semmit sem:
  - lehetséges, hogy a nullhipotézis igaz, tehát jól döntöttünk,
  - lehetséges, hogy hamis, és másodfajú hibát vétettünk.
- Ha elvetjük a nullhipotézist, az már jelent valamit:
  - lehetséges ugyan, hogy a nullhipotézis igaz, és elsőfajú hibát vétettünk, de ennek kicsi az esélye, ez ritkán történik meg,
  - a nullhipotézis elvetése tipikusan azt jelenti, hogy a nullhipotézis hamis.

Az általunk tanult tesztelési módszerek esetében  $\beta \rightarrow 0$ , ha  $n \rightarrow \infty$ .

Tehát ha növeljük a mintaméretet, akkor a másodfajú hiba is alacsony lesz.

Ez azt jelenti, hogy ezeknél módszereknél nagy mintaméret esetén a nullhipotézis elfogadása már tényleg arra utal, hogy a nullhipotézis igaz.

# Az egymintás $t$ -próba

## Egymintás $t$ -próba (One sample $t$ test)

Cél a  $\xi$  valószínűségi változó várható értékének tesztelése egy  $\xi_1, \dots, \xi_n$  statisztikai minta alapján.

- Feltevések:
  - $\xi$  normális eloszlású változó ismeretlen  $\mu$  várható értékkel,
  - $\mu_0$  egy tetszőleges hipotetikus érték.
- Nullhipotézis:  $H_0 : \mu = \mu_0$ .
- Próbastatisztika: ( $t$ -próba esetén hagyományosan  $t_n$  a jele)

$$t_n = \frac{\bar{\xi} - \mu_0}{D_n^*(\xi)/\sqrt{n}} = \frac{\bar{\xi} - \mu_0}{SE}.$$

- Kritikus érték:  $c_\alpha = \Phi_{n-1}^{-1}(1 - \alpha/2)$ .
- Döntés: akkor fogadjuk el a nullhipotézist, ha  $|t_n| \leq c_\alpha$ .

**Feladat:** A kar férfi hallgatóinak testmagasságát vizsgáljuk, jelölje  $\xi$  egy véletlenszerűen kiválasztott férfi hallgató magasságát. Mit állíthatunk  $\xi$  várható értékéről, az átlagos testmagasságról a teljes populáción belül?

- Megfigyelt értékek: 180, 175, 188, 168, 173, 183.
- Becslések:  $E(\xi) \approx \bar{\xi} = 177,8$ ,  $D(\xi) \approx D_6^*(\xi) = 7,25$ .

Teszteljük 5%-os szignifikancia szinten azt, hogy  $H_0 : E(\xi) = 175$ .

Tegyük fel, hogy a testmagasság normális eloszlást követ a populáción belül. Ekkor a  $t$ -próba alkalmazható.

- Hipotetikus érték, szignifikancia szint:  $\mu_0 = 175$ ,  $\alpha = 0,05$ .
- Próbastatisztika:

$$t_n = \frac{\bar{\xi} - \mu_0}{D_n^*(\xi)/\sqrt{n}} = \frac{177,8 - 175}{7,25/\sqrt{6}} = 0,946,$$

- A kritikus érték:  $c_\alpha = \Phi_{n-1}^{-1}(1 - \alpha/2) = \Phi_5^{-1}(0,975) = 2,571$ .
- Döntés:  $|t_n| \leq c$ , tehát a nullhipotézist elfogadjuk. A várható érték nem különbözik **szignifikáns (=statisztikailag kimutatható) mértékben** a 175-ös értéktől.

Mi a gondolat a  $t$ -próba mögött? A mintaátlag jó becslése a  $\mu$  igazi várható értéknek, tehát

$$t_n = \frac{\bar{\xi} - \mu_0}{SE} \approx \frac{\mu - \mu_0}{SE}.$$

A  $H_0 : \mu = \mu_0$  nullhipotézist teszteljük.

- Ha a nullhipotézis igaz, akkor

$$t_n \approx \frac{\mu - \mu_0}{SE} = 0.$$

- Ha a nullhipotézis nem igaz, akkor

$$t_n \approx \frac{\mu - \mu_0}{SE} \neq 0.$$

A nullhipotézist akkor fogadjuk el, ha  $|t_n| \leq c_{\alpha}$ , tehát ha  $t_n$  nullához közeli szám. Ez logikus ötlet, hiszen

- ha  $t_n \approx 0$ , akkor az arra utal, hogy  $H_0$  igaz,
- ha  $t_n \not\approx 0$ , akkor az arra utal, hogy  $H_0$  nem igaz.

Fejtsük ki egy kicsit jobban az előző oldalt! Mikor fogadjuk el  $H_0$ -t?

$$\begin{aligned} |t_n| \leq c_\alpha &\iff -c_\alpha \leq t_n \leq c_\alpha \iff -c_\alpha \leq \frac{\bar{\xi} - \mu_0}{SE} \leq c_\alpha \\ &\iff \bar{\xi} - c_\alpha SE \leq \mu_0 \leq \bar{\xi} + c_\alpha SE \iff \mu_0 \in [\bar{\xi} - c_\alpha SE, \bar{\xi} + c_\alpha SE] \end{aligned}$$

Amit kaptunk, az az  $1 - \alpha$  megbízhatóságú konfidencia intervallum a normális eloszlás várható értékére. Ekkor

$$\begin{aligned} P(\text{elfogadjuk } H_0\text{-t} \mid H_0 \text{ igaz}) &= P(\mu_0 \in [\bar{\xi} - c_\alpha SE, \bar{\xi} + c_\alpha SE] \mid \mu = \mu_0) \\ &= P(\mu \in [\bar{\xi} - c_\alpha SE, \bar{\xi} + c_\alpha SE]) = 1 - \alpha. \end{aligned}$$

Ebből következik, hogy

$$P(\text{elvetjük } H_0\text{-t} \mid H_0 \text{ igaz}) = 1 - P(\text{elfogadjuk } H_0\text{-t} \mid H_0 \text{ igaz}) = \alpha.$$

Az előző oldalon levezetett számolásnak több fontos következménye van:

- A próba pontosan akkor fogadja el a  $\mu_0$  hipotetikus várható értéket, ha  $\mu_0$  az  $1 - \alpha$  megbízhatóságú konfidencia intervallumba esik. A konfidencia intervallum értelmezhető olyan módon, mint a „hihető” várható értékek halmaza.
- Ha a minta normális eloszlásból jön, akkor a  $t$ -próba pontosan betartja az előírt elsőfajú hibát:

$$P(\text{elvetjük } H_0\text{-t} \mid H_0 \text{ igaz}) = \text{megadott szignifikancia szint.}$$

- Ha a minta nem normális eloszlásból származik, de a mintaméret elég nagy, akkor a  $t$ -próba használható a várható érték tesztelésére. Ebben az esetben a próba csak közelítőleg tartja be az előírt elsőfajú hibát:

$$P(\text{elvetjük } H_0\text{-t} \mid H_0 \text{ igaz}) \approx \text{megadott szignifikancia szint.}$$



Lefutattam a  $t$ -próbát 5%-os szignifikancia szinten a testmagasságokra az R programmal, az alábbi eredményt kaptam:

```
One Sample t-test
data:  magassag
t = 0.95723, df = 5, p-value = 0.3824
alternative hypothesis: true mean is not equal to 175
95 percent confidence interval:  170.2246 185.4420
sample estimates: mean of x 177.8333
```

Értelmezzük, hogy milyen információ van az outputban:

- Egymintás  $t$ -próba a „magassag” nevű adatsoron.
- Próbastatisztika:  $t = 0.95723$ , szabadsági fok (degrees of freedom):  $df = 5$ .
- Nullhipotézis és alternatív hipotézis:  $H_0 : \mu = 175$ ,  $H_A : \mu \neq 175$ .
- 95%-os konfidencia intervallum:  $[170.2246, 185.4420]$ .
- Mintaátlag: 177.8333

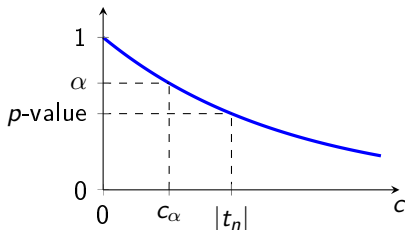
A program által adott értékek kissé eltérnek attól, amit mi kaptunk: nálunk sok volt a kerekítési hiba. Felmerülő kérdések:

- Hol a kritikus érték és a döntés? És mi az a „ $p$ -value”?

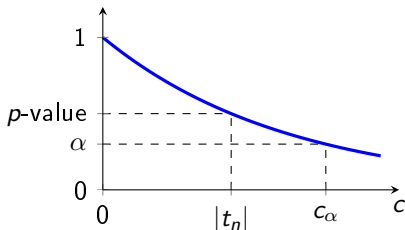
A  **$p$ -érték ( $p$ -value)** az a határ szignifikancia szint, amikor még éppen elfogadjuk a nullhipotézist, tehát  $c_{p\text{-value}} = t_n$ . Ekkor

$$\text{elvetjük } H_0\text{-t} \iff |t_n| > c_\alpha \iff p\text{-value} < \alpha.$$

A statisztikai programok tesztelés során gyakran nem a kritikus értéket, hanem  $p$ -értéket adják meg. A  $p$ -érték 0 és 1 közé esik, és értelmezhető olyan módon, hogy mennyire „hihető” a nullhipotézis az adott statisztikai minta mellett. A nullhipotézist akkor vetjük el, ha a  $p$ -érték alacsony.



elvetjük  $H_0$ -t



elfogadjuk  $H_0$ -t

Az előző feladatban:  $p\text{-value} = 0.3824 > \alpha = 0.05$ , tehát  $H_0$ -t elfogadjuk.

# Kétmintás hipotézisvizsgálat

Legyenek  $\xi$  és  $\eta$  valószínűségi változók. Két statisztikai minta:

- $\xi_1, \dots, \xi_n$  független megfigyelések  $\xi$ -re,
- $\eta_1, \dots, \eta_m$  független megfigyelések  $\eta$ -ra.

A minták segítségével becsléseket végezhetünk:

- $\bar{\xi} = E_n(\xi) \approx E(\xi)$ ,  $D_n^*(\xi) \approx D(\xi)$ ,
- $\bar{\eta} = E_m(\eta) \approx E(\eta)$ ,  $D_m^*(\eta) \approx D(\eta)$ .

A minták tipikusan kétfajta kapcsolatban állhat egymással:

- **Független minták (independent samples):** A minták között nincs kapcsolat, mert a két minta független mintavételezésből származik. Például: egymástól függetlenül veszünk mintát két részpopulációból.
- **Összetartozó minták (paired samples, related samples):** A  $\xi_i$  és az  $\eta_i$  megfigyelés minden  $i$  esetén a populáció ugyanazon egyedére vonatkozik, ezért ezek az értékek nem függetlenek egymástól. Ebben az esetben mindig  $n = m$ .

**Feladat:** Döntsük el, hogy az alábbi példákban független vagy összetartozó mintákról van szó.

- ①
- $\xi_1, \dots, \xi_n$ :  $n$  véletlenszerűen kiválasztott férfi hallgató testmagassága,
  - $\eta_1, \dots, \eta_m$ :  $m$  véletlenszerűen kiválasztott női hallgató testmagassága.

Független minták: a minták független megfigyelésekből jönnek.

- ②
- $\xi_1, \dots, \xi_n$ :  $n$  férfi hallgató testmagassága egy mai felmérésben,
  - $\eta_1, \dots, \eta_n$ : ugyanezen hallgatók édesanyjának testmagassága.

Összetartozó minták: a megfigyelések azonos egyedekre vonatkoznak.

- ③
- $\xi_1, \dots, \xi_n$ :  $n$  férfi hallgató testmagassága egy mai felmérésben,
  - $\eta_1, \dots, \eta_n$ :  $n$  férfi hallgató testmagassága egy 5 évvel ezelőtti független felmérésben.

Független minták: a minták független megfigyelésekből jönnek.

- ④
- $\xi_1, \dots, \xi_n$ :  $n$  férfi hallgató testmagassága egy mai felmérésben,
  - $\eta_1, \dots, \eta_n$ : ugyanezen hallgatók testmagassága egy 5 évvel ezelőtti felmérésben.

Összetartozó minták: a megfigyelések azonos egyedekre vonatkoznak.

A  $\xi$  és az  $\eta$  valószínűségi változó **együttesen normális eloszlást követ**, ha tetszőleges  $a$  és  $b$  valós számok esetén  $a\xi + b\eta$  normális eloszlású. Ez egy kicsivel több annál, hogy  $\xi$  és  $\eta$  normális eloszlású.

Tegyük fel, hogy

- $\xi$  és  $\eta$  együttesen normális eloszlásúak,
- a várható értékek ( $\mu_\xi$  és  $\mu_\eta$ ) ismeretlenek,
- $\xi_1, \dots, \xi_n$  és  $\eta_1, \dots, \eta_n$  összetartozó minták,

Célunk a következő nullhipotézist tesztelni:  $H_0 : \mu_\xi = \mu_\eta$ .

Gondolatmenet:

- $(+1)\xi + (-1)\eta = \xi - \eta$  normális eloszlású változó.
- $\xi_1 - \eta_1, \dots, \xi_n - \eta_n$  statisztikai minta a  $\xi - \eta$  változóra.
- $E(\xi - \eta) = E(\xi) - E(\eta) = \mu_\xi - \mu_\eta$ , ezért  $H_0 \Leftrightarrow E(\xi - \eta) = 0$ .
- Teszteltjük a  $H_0 : E(\xi - \eta) = 0$  nullhipotézist egymintás  $t$ -próbával.

Foglaljuk össze, hogyan lehet összetartozó minták várható értékét tesztelni:

## Páros $t$ -próba (paired samples $t$ test)

Cél a várható értékek tesztelése összetartozó minták esetén.

- Feltevések:
  - $\xi$  és  $\eta$  együttesen normális eloszlásúak,
  - a várható értékek ( $\mu_\xi$  és  $\mu_\eta$ ) ismeretlenek,
  - $\xi_1, \dots, \xi_n$  és  $\eta_1, \dots, \eta_n$  összetartozó minták.
- Nullhipotézis:  $H_0 : \mu_\xi = \mu_\eta$ .

- Próbastatisztika:

$$t_n = \frac{\overline{\xi - \eta} - 0}{D_n^*(\xi - \eta)/\sqrt{n}}.$$

- Kritikus érték:  $c_\alpha = \Phi_{n-1}^{-1}(1 - \alpha/2)$ .
- Döntés: akkor fogadjuk el a nullhipotézist, ha  $|t_n| \leq c_\alpha$ .

**Feladat:** Teszteljük azt, hogy a kar jelenlegi férfi hallgatóinak átlagos testmagassága nem változott ez elmúlt 5 év folyamán. ( $\alpha = 5\%$ )

Legyen

- $\xi$  = véletlenszerűen kiválasztott férfi hallgató testmagassága ma,
- $\eta$  = ugyanezen hallgató testmagassága 5 évvel ezelőtt,
- $\xi - \eta$  = testmagasság változása 5 év alatt.

Várható értékek:

- $\mu_\xi = E(\xi)$  = férfi hallgatók átlagos testmagassága ma,
- $\mu_\eta = E(\eta)$  = ugyanezen hallgatók átlagos magassága 5 évvel ezelőtt,
- $\mu_\xi - \mu_\eta = E(\xi - \eta)$  = átlagos magasságváltozás 5 év alatt.

Nullhipotézis:  $H_0 : \mu_\xi = \mu_\eta$ . Helyette:  $H_0 : E(\xi - \eta) = 0$ .

Statisztikai minták: ( $n = 6$ )

- Kiválasztott hallgatók magassága ma: 180, 175, 188, 168, 173, 183.
- Ugyanezen hallgatók magassága 1 éve: 175, 172, 184, 167, 170, 178.
- Minta a  $\xi - \eta$  változóra: 5, 3, 4, 1, 3, 5.

Minta a  $\xi - \eta$  változóra: 5, 3, 4, 1, 3, 5.

Becslések:

- Mintaátlag:  $\overline{\xi - \eta} = 3,5 \approx E(\xi - \eta)$ .
- Korrigált empirikus szórás:  $D_6^*(\xi - \eta) = 1,52 \approx D(\xi - \eta)$ .

Nullhipotézis:  $H_0 : E(\xi - \eta) = 0$ .

Egymintás  $t$ -próba:

- Próbastatisztika:

$$t_n = \frac{\overline{\xi - \eta}}{D_n^*(\xi - \eta)/\sqrt{n}} = \frac{3,5 - 0}{1,52/\sqrt{6}} = 5,64.$$

- A kritikus érték:  $c_\alpha = \Phi_{n-1}^{-1}(1 - \alpha/2) = \Phi_5^{-1}(0,975) = 2,571$ .
- Döntés:  $|t_n| > c_\alpha$ , tehát a nullhipotézist elvetjük. A populációban az átlagos testmagasság szignifikáns módon változott az elmúlt 5 év folyamán.



A következő módszerrel független minták várható értékét tesztelhetjük:

## Kétmintás $t$ -próba (independent samples $t$ test)

Cél a várható értékek tesztelése független minták esetén.

- Feltevések:
  - $\xi$  és  $\eta$  normális eloszlásúak,
  - a várható értékek ( $\mu_\xi$  és  $\mu_\eta$ ) és a szórások ( $\sigma_\xi$  és  $\sigma_\eta$ ) ismeretlenek,
  - a szórások megegyeznek:  $\sigma_\mu = \sigma_\eta$ ,
  - $\xi_1, \dots, \xi_n$  és  $\eta_1, \dots, \eta_m$  független minták.
- Nullhipotézis:  $H_0 : \mu_\xi = \mu_\eta$ .
- Próbastatisztika:  $t_{n,m} = (\bar{\xi} - \bar{\eta}) / D_{n,m}$ , ahol

$$D_{n,m} = \sqrt{\left[ (n-1) \text{Var}_n^*(\xi) + (m-1) \text{Var}_m^*(\eta) \right] \frac{n+m}{nm(n+m-2)}}$$

- Kritikus érték:  $c_\alpha = \Phi_{n+m-2}^{-1}(1 - \alpha/2)$ .
- Döntés: akkor fogadjuk el a nullhipotézist, ha  $|t_{n,m}| \leq c_\alpha$ .

Mi a kétmintás  $t$ -próba alapötlete? A mintaátlag jó becslés a várható értékre, ezért

$$t_{n,m} = \frac{\bar{\xi} - \bar{\eta}}{D_{n,m}} \approx \frac{\mu_{\xi} - \mu_{\eta}}{D_{n,m}}.$$

A  $H_0 : \mu_{\xi} = \mu_{\eta}$  nullhipotézist teszteljük.

- Ha a nullhipotézis igaz, akkor

$$t_{n,m} \approx \frac{\mu_{\xi} - \mu_{\eta}}{D_{n,m}} = 0.$$

- Ha a nullhipotézis nem igaz, akkor

$$t_{n,m} \approx \frac{\mu_{\xi} - \mu_{\eta}}{D_{n,m}} \neq 0.$$

A nullhipotézist akkor fogadjuk el, ha  $t_{n,m}$  nullához közeli szám. Ez logikus ötlet, hiszen

- ha  $t_{n,m} \approx 0$ , akkor az arra utal, hogy  $H_0$  igaz,
- ha  $t_{n,m} \not\approx 0$ , akkor az arra utal, hogy  $H_0$  nem igaz.

**Feladat:** Teszteljük azt a nullhipotézist, hogy a kar férfi hallgatóinak átlagos testmagassága nem változott ez elmúlt 5 év folyamán. ( $\alpha = 5\%$ )

Legyen

- $\xi$  = véletlenszerűen kiválasztott férfi hallgató testmagassága ma,
- $\eta$  = véletlenszerű hallgató magassága az 5 évvel ezelőtti populációban.

Várható értékek:

- $\mu_\xi = E(\xi)$  = férfi hallgatók átlagos testmagassága ma,
- $\mu_\eta = E(\eta)$  = férfi hallgatók átlagos testmagassága 5 évvel ezelőtt.

Statisztikai minták két független felmérésből:

- Kiválasztott hallgatók magassága ma: 180, 175, 188, 168, 173, 183.
- Kiválasztott hallgatók magassága 5 éve: 171, 178, 183, 168, 175.

Becslések:

- $\mu_\xi \approx \bar{\xi} = 177,8$ ,  $\sigma_\xi \approx D_6^*(\xi) = 7,25$ .
- $\mu_\eta \approx \bar{\eta} = 175$ ,  $\sigma_\eta \approx D_5^*(\eta) = 5,87$ .

Nullhipotézis:  $H_0 : \mu_\xi = \mu_\eta$ .

Feltehető, hogy a minták normális eloszlásból származnak és az elméleti szórások azonosak ( $\sigma_\xi = \sigma_\eta$ ).

Kétmintás  $t$ -próba:

- Próbastatisztika:  $t_{n,m} = (177,8 - 175)/3,845 = 0,73$ , ugyanis

$$D_{n,m} = \sqrt{\left[ (6-1)7,25^2 + (5-1)5,87^2 \right] \frac{6+5}{6 \cdot 5 \cdot (6+5-2)}} = 3,845.$$

- Kritikus érték:  $c_\alpha = \Phi_9^{-1}(0,975) = 2,262$ .
- Döntés:  $|t_{n,m}| \leq c_\alpha$ , ezért a nullhipotézist elfogadjuk. A kar férfi hallgatóinak átlagos testmagassága az elmúlt 5 évben nem változott szignifikáns mértékben.

## F-próba ( $F$ test)

Cél a szórások tesztelése független minták esetén.

- Feltevések:
  - $\xi$  és  $\eta$  normális eloszlásúak, a szórások ( $\sigma_\xi$  és  $\sigma_\eta$ ) ismeretlenek,
  - $\xi_1, \dots, \xi_n$  és  $\eta_1, \dots, \eta_m$  független minták.
- Nullhipotézis:  $H_0 : \sigma_\xi = \sigma_\eta$ .
- Próbastatisztika, kritikus érték: nem tanuljuk.

## Welch-próba (Welch test)

Ugyanaz, mint a kétmintás  $t$ -próba, de nem kell a szórások egyenlősége.

- Feltevések:
  - $\xi$  és  $\eta$  normális eloszlásúak,
  - a várható értékek ( $\mu_\xi$  és  $\mu_\eta$ ) ismeretlenek,
  - $\xi_1, \dots, \xi_n$  és  $\eta_1, \dots, \eta_m$  független minták.
- Nullhipotézis:  $H_0 : \mu_\xi = \mu_\eta$ .
- Próbastatisztika, kritikus érték: nem tanuljuk.

A tanult statisztikai módszerek (becslések, tesztek) mögött matematikai tételek állnak. Ezek a tételek garantálják, hogy megfelelő feltételek mellett ezek a módszerek jól működnek. (Például a becslések erősen konzisztensek, a tesztek betartják az előírt elsőfajú hibát.) De mennyire hatékonyak ezek a módszerek akkor, ha a szükséges feltételek nem teljesülnek?

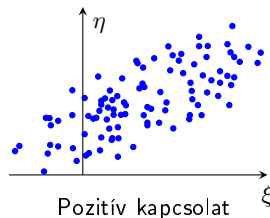
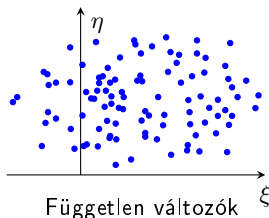
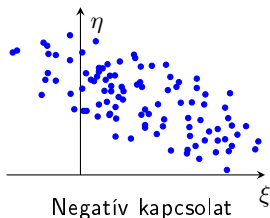
Egy statisztikai módszer **robosztus (robust)** egy feltételre nézve, ha a módszer a feltétel elhagyásával is jól alkalmazható. Például:

- A tanult  $t$ -próbák és a Welch-próba robusztus a normalitási feltételre nézve, ezek a próbák nem normális eloszlásra is alkalmazhatóak, ha a mintaméret elég nagy. (50 mintaelem elég szokott lenni.)
- A kétmintás  $t$ -próba nem robusztus a szórásfeltételre nézve. Ha a szórások nem azonosak, akkor a Welch-próbát kell alkalmazni.

# Kovariancia és korreláció

Eddig azt vizsgáltuk, hogy egy  $\xi$  mennyiségnek milyen az eloszlása egy populáción belül. A továbbiakban két mennyiség ( $\xi$  és  $\eta$ ) együttes viselkedésével foglalkozunk. Főleg az a kérdés, hogy milyen irányú és milyen erősségű kapcsolat van a két változó között. A fontosabb esetek:

- Pozitív irányú kapcsolat: a  $\xi$  és az  $\eta$  (jellemzően) azonos irányba mozdul el. (Például: testmagasság és testsúly.)
- Negatív irányú kapcsolat:  $\xi$  és  $\eta$  (jellemzően) egymással ellentétes irányba mozog.
- Független változók: nincs kapcsolat az értékek között.



Legyen  $\xi$  és  $\eta$  valószínűségi változó. Ekkor a két változó **kovarianciája (covariance)** illetve **korrelációs együtthatója (correlation coefficient)**:

$$\text{Cov}(\xi, \eta) = E\left([\xi - E(\xi)][\eta - E(\eta)]\right), \quad \text{corr}(\xi, \eta) = \frac{\text{Cov}(\xi, \eta)}{D(\xi)D(\eta)}.$$

Ha a kovariancia (és ezáltal a korrelációs együttható) értéke nulla, akkor azt mondjuk, hogy a két változó **korrelálatlan (uncorrelated)**.

A kovariancia és a korrelációs együttható fontosabb tulajdonságai:

- A kovariancia és a korrelációs együttható a teljes populációt jellemzi valamilyen (de vajon milyen?) szempontból.
- Lehetséges értékek:  $\text{Cov}(\xi, \eta) \in \mathbb{R}$ ,  $\text{corr}(\xi, \eta) \in [-1, +1]$ .
- Szimmetria:  $\text{Cov}(\xi, \eta) = \text{Cov}(\eta, \xi)$ ,  $\text{corr}(\xi, \eta) = \text{corr}(\eta, \xi)$ .
- Ha  $\xi$  és  $\eta$  függetlenek, akkor korrelálatlanok is. Viszont a korrelálatlanságból nem következik a függetlenség.
- Ha a két változó együttesen normális eloszlású, akkor a függetlenség és a korrelálatlanság ekvivalens fogalmak.

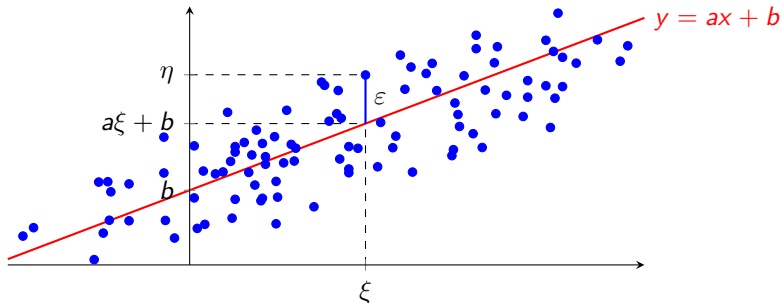


Lineáris regresszió a teljes populáción: szeretnénk megérteni, hogy az  $\eta$  változó értéke milyen módon alakul ki. Vegyük a következő reprezentációt:

$$\eta = (a\xi + b) + \varepsilon = \text{predikciós tag} + \text{hibatag}.$$

Elnevezések és modellfeltevések:

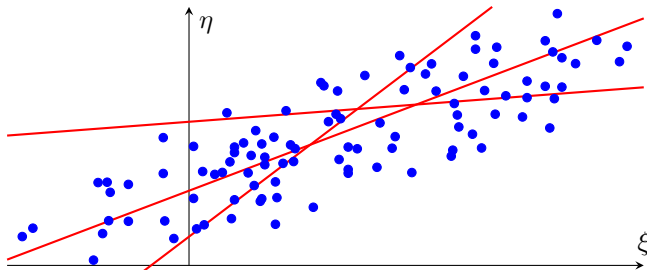
- $\xi$  a **magyarázó változó**,  $\eta$  a **függő változó**,
- $\varepsilon$  az egyedre jellemző **hibatag**, ami független a  $\xi$  változótól,
- $E(\varepsilon) = \text{átlagos hiba} = 0$ .



A regressziós modell lehetséges alkalmazásai:

- Becslés az  $\eta$  változóra: ha a hibatag kicsi, akkor  $\eta \approx a\xi + b$ .
- Megérteni, hogy milyen hatások szerint alakul ki az  $\eta$  változó.

Probléma: több olyan egyenes is létezhet, ami teljesíti a modellfeltevéseket.



Mi a  $D(\varepsilon)$  szórárs jelentése a populációra nézve?

$$D(\varepsilon) = \varepsilon \text{ átlagos eltérése a } 0 \text{ várható értéktől}$$

$$= \text{az abszolút hiba átlagos értéke}$$

Cél: azt az egyenest keressük, melynél  $D(\varepsilon)$  minimális, ugyanis ez az egyenes illeszkedik a legjobban a ponthalmazra.

Minimalizáljuk a hibatag szórását az  $a$  és  $b$  változóknban!

$$D^2(\varepsilon) = D^2(\eta - a\xi - b) = D^2(\eta) + a^2D^2(\xi) - 2aD(\xi)D(\eta) \text{corr}(\xi, \eta)$$

Deriválással:

$$0 = \frac{\partial D^2(\varepsilon)}{\partial a} = 2aD^2(\xi) - 2D(\xi)D(\eta) \text{corr}(\xi, \eta),$$

amiből következik, hogy

$$a = \frac{\text{corr}(\xi, \eta) D(\eta)}{D(\xi)}.$$

Azt is tudjuk, hogy  $E(\varepsilon) = 0$ , tehát

$$0 = E(\varepsilon) = E(\eta - a\xi - b) = E(\eta) - aE(\xi) - b,$$

amiből

$$b = E(\eta) - aE(\xi) = E(\eta) - \frac{\text{corr}(\xi, \eta) D(\eta)}{D(\xi)} E(\xi).$$

**Feladat:** A Tisza és a Maros vízhozama számunkra ismeretlen eloszlást követ. Azt tudjuk, hogy a Maros torkolata felett a Tisza vízhozamának a várható értéke  $660 \text{ m}^3/\text{s}$ , szórása  $160 \text{ m}^3/\text{s}$ , míg a Maros vízhozamának a várható értéke  $200 \text{ m}^3/\text{s}$ , szórása  $50 \text{ m}^3/\text{s}$ . A korrelációs együttható  $0,8$ . Egy adott napok a Tisza vízhozama  $800 \text{ m}^3/\text{s}$ . Milyen becslést adhatunk a Maros vízhozamára?

Legyen  $\xi$  és  $\eta$  a Tisza illetve a Maros vízhozama torkolat felett. Ekkor

$$E(\xi) = 660, \quad D(\xi) = 160, \quad E(\eta) = 200, \quad D(\eta) = 50, \quad \text{corr}(\xi, \eta) = 0,8.$$

Az előző oldalakon kapott formulák alkalmazásával kapjuk, hogy

$$a = \frac{\text{corr}(\xi, \eta)D(\eta)}{D(\xi)} = 0,25, \quad b = E(\eta) - aE(\xi) = 35.$$

Tehát a két vízhozamra az alábbi regressziós modell írható fel:

$$\eta = a\xi + b + \text{hibatag} = 0,25\xi + 35 + \text{hibatag} \approx 0,25\xi + 35.$$

A mai napon  $\xi = 800$ , tehát  $\eta \approx 0,25 \cdot 800 + 35 = 235$ .

Viszsgáljuk meg a varianciákat! A tagok függetlensége miatt

$$D^2(\eta) = D^2(a\xi + b + \varepsilon) = D^2(a\xi + b) + D^2(\varepsilon).$$

Egy kis számolás után (amit nem részletezünk) a predikciós tag varianciája

$$D^2(a\xi + b) = \text{corr}^2(\xi, \eta)D^2(\eta).$$

Ebből következik, hogy a hibatag varianciája

$$D^2(\varepsilon) = D^2(\eta) - D^2(a\xi + b) = [1 - \text{corr}^2(\xi, \eta)]D^2(\eta).$$

Foglaljuk össze, hogy mit kaptunk:

- A predikciós tag varianciája a teljes varianciának  $\text{corr}^2(\xi, \eta)$  hányada. „Ekkora mértékben magyarázza a  $\xi$  változó az  $\eta$  értékét.”
- A hibatag varianciája a teljes varianciának  $1 - \text{corr}^2(\xi, \eta)$  hányada. „Ekkora mértékben magyarázza az  $\varepsilon$  változó az  $\eta$  értékét.”

**Feladat:** Mekkora a predikciós tag illetve a hibtag várható értéke illetve szórása a vízhozamos feladatban? Milyen arányban magyarázza a Tisza vízhozama a Maros vízhozamát?

A modell:  $\eta = (0,25\xi + 35) + \varepsilon = \text{predikciós tag} + \text{hibtag}$ .

A predikciós tag várható értéke:

$$E(\text{predikciós tag}) = 0,25 E(\xi) + 35 = 0,25 \cdot 600 + 35 = 200.$$

A varianciák:

$$D^2(\text{predikciós tag}) = D^2(a\xi + b) = \text{corr}^2(\xi, \eta) D^2(\eta) = 0,8^2 \cdot 50^2 = 1600,$$

$$D^2(\text{hibtag}) = D^2(\varepsilon) = [1 - \text{corr}^2(\xi, \eta)] D^2(\eta) = [1 - 0,8^2] \cdot 50^2 = 900.$$

Foglaljuk össze ezeket egy táblázatban:

Hatás	E	Var	D	Var%
Tisza ( $a\xi + b$ )	200	1600	40	64%
hibtag ( $\varepsilon$ )	0	900	30	36%
Maros ( $\eta$ )	200	2500	50	100%

A Tisza 64%, a hibtag 36% arányban magyarázza a Maros vízhozamát.

Milyen módon jellemzi a korrelációs együttható a  $\xi$  és az  $\eta$  változó kapcsolatát? Induljunk ki a következő két összefüggésből:

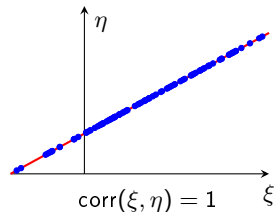
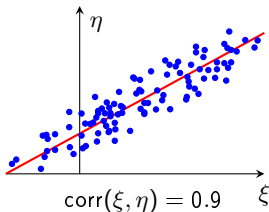
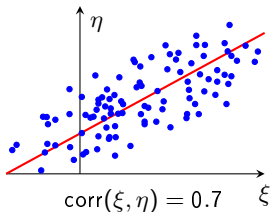
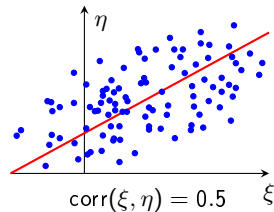
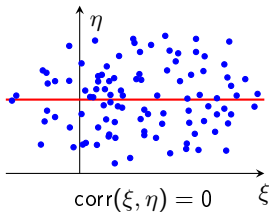
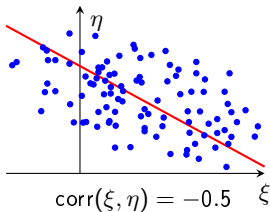
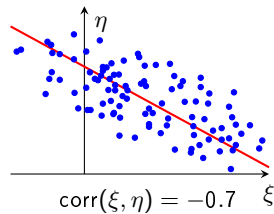
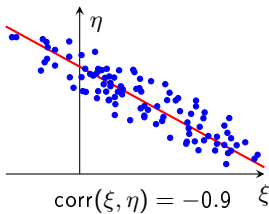
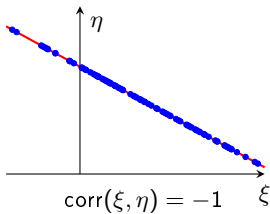
$$a = \frac{\text{corr}(\xi, \eta)D(\eta)}{D(\xi)}, \quad D^2(\varepsilon) = [1 - \text{corr}^2(\xi, \eta)]D^2(\eta).$$

A regressziós egyenes meredeksége ( $a$ ) alapján:

- Ha  $\text{corr}(\xi, \eta) > 0$ , akkor  $a > 0$ , tehát a változók között pozitív irányú kapcsolat van.
- Ha  $\text{corr}(\xi, \eta) < 0$ , akkor  $a < 0$ , tehát a változók között negatív irányú kapcsolat van.

A hibatag varianciája alapján:

- Ha  $\text{corr}(\xi, \eta) \approx \pm 1$ , akkor  $D(\varepsilon) \approx 0$ . Ebben az esetben jó az illeszkedés a regressziós egyeneshez, a változók közötti kapcsolat erős.
- Ha  $\text{corr}(\xi, \eta) \approx 0$ , akkor a  $D(\varepsilon)$  szórás nagy. Ebben az esetben nem jó az illeszkedés a regressziós egyeneshez, a változók között gyenge a kapcsolat (vagy akár függetlenek is).





# Statisztikai lineáris regresszió

Az előző részben két mennyiség ( $\xi$  és  $\eta$ ) kapcsolatát vizsgáltuk a teljes populáción belül. Probléma: nem ismerjük a változók elméleti várható értékét, szórását és korrelációs együtthatóját. Megoldás: egy statisztikai minta alapján mindent becsülni fogunk.

Tekintsünk összetartozó mintákat a  $\xi$  és  $\eta$  változókra:

- $\xi_1, \dots, \xi_n$  független megfigyelések  $\xi$ -re,
- $\eta_1, \dots, \eta_n$  az  $\eta$  mennyiség értékei ugyanezen egyedeknél.

**Empirikus kovariancia (sample covariance) és empirikus korrelációs együttható (sample correlation coefficient):**

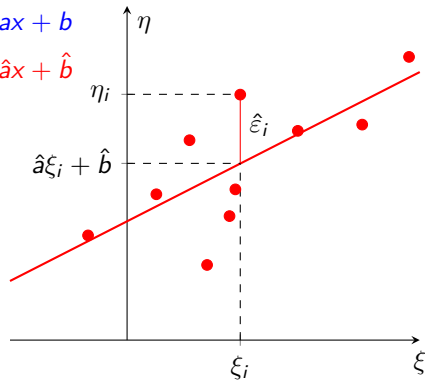
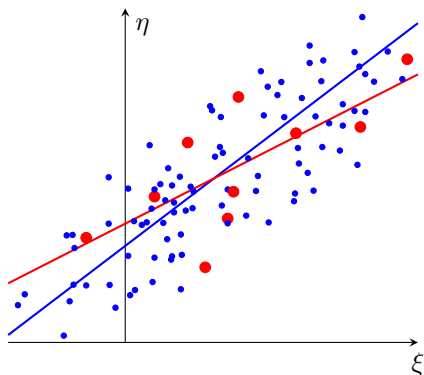
$$\text{Cov}_n(\xi, \eta) = \frac{1}{n-1} \sum_{i=1}^n (\xi_i - \bar{\xi})(\eta_i - \bar{\eta}), \quad \text{corr}_n(\xi, \eta) = \frac{\text{Cov}_n(\xi, \eta)}{D_n^*(\xi) D_n^*(\eta)}.$$

Lineáris regresszió a teljes populáción:

$$\eta = (a\xi + b) + \varepsilon = \text{predikciós tag} + \text{hibatag.}$$

Lineáris regresszió a mintaelemeken:

$$\eta_i = (\hat{a}\xi_i + \hat{b}) + \hat{\varepsilon}_i = \text{predikciós tag} + \text{reziduális}, \quad i = 1, \dots, n.$$



Az egyenest a **legkisebb négyzetes becslés (least squares estimation)** alkalmazásával kapjuk meg: keressük azon  $\hat{a}$  és  $\hat{b}$  számokat, melyekre

$$S(\hat{a}, \hat{b}) = \sum_{i=1}^n \hat{\varepsilon}_i^2 = \sum_{i=1}^n (\eta_i - \hat{a}\xi_i - \hat{b})^2 \longrightarrow \min.$$

Az összeget  $\hat{a}$  és  $\hat{b}$  szerint deriválva:

$$0 = \frac{\partial S}{\partial \hat{a}} = -2 \sum_{i=1}^n (\eta_i - \hat{a}\xi_i - \hat{b})\xi_i,$$

$$0 = \frac{\partial S}{\partial \hat{b}} = -2 \sum_{i=1}^n (\eta_i - \hat{a}\xi_i - \hat{b}).$$

Az egyenletrendszer megoldása:

$$\hat{a} = \frac{\text{corr}_n(\xi, \eta) D_n^*(\eta)}{D_n^*(\xi)}, \quad \hat{b} = \bar{\eta} - \frac{\text{corr}_n(\xi, \eta) D_n^*(\eta)}{D_n^*(\xi)} \bar{\xi}.$$

Vegyük észre: lényegében elemenként becsültünk mindent  $a$  és  $b$  formulájában.

A kapott becslések erősen konzisztensek, tehát  $n \rightarrow \infty$  esetén

$$\text{Cov}_n(\xi, \eta) \rightarrow \text{Cov}(\xi, \eta), \quad \text{corr}_n(\xi, \eta) \rightarrow \text{corr}(\xi, \eta), \quad \hat{a} \rightarrow a, \quad \hat{b} \rightarrow b.$$

Mennyire jó az illeszkedés a regressziós egyeneshez? A teljes populáción az illeszkedés „jóságát” a  $\text{corr}^2(\xi, \eta)$  mennyiség számszerűsíti. Becslése a minta alapján:  $R\text{-squared} = \text{corr}_n^2(\xi, \eta)$ . Tulajdonságai:

- $0 \leq R\text{-squared} \leq 1$ ,
- Annál jobb az illeszkedés, minél nagyobb az  $R\text{-squared}$  értéke,
- Ha  $R\text{-squared} \leq 0,5$ , akkor nagyon rossz az illeszkedés a regressziós egyeneshez, nagyok a hibatagok, ezért a modell alapján nem érdemes becsléseket végezni.

## Korrelációs teszt

Cél a  $\xi$  és az  $\eta$  valószínűségi változó függetlenségének tesztelése összetartozó minták alapján.

- Feltevés:  $\xi$  és  $\eta$  együttesen normális eloszlású.
- Nullhipotézis:  $H_0$  :  $\xi$  és  $\eta$  függetlenek.
- Próbastatisztika: egy ronda formula, amiben szerepel az empirikus korreláció.

Megjegyzések:

- Ez a teszt nem robusztus a normalitásfeltételre nézve: nagy minta esetén szimmetrikus eloszlásra még alkalmazható, de dőlt eloszlásra ne használjuk!  
Mit alkalmazzunk dőlt eloszlás esetén? Spearman-féle korreláció.
- Mi a jelentősége a függetlenségvizsgálatnak a regressziós modellben?  
Ha a próba nem veti el  $\xi$  és  $\eta$  függetlenségét, akkor nem érdemes regressziót végezni.

# Variansciaanalízis

Korábban kétmintás t-próbával teszteltük a várható értékek azonosságát részpopuláción belül. Most több részcsoporthoz fogunk vizsgálni egyszerre. Véletlenszerűen kiválasztva egy egyedet tekintünk két változót:

- $\xi$  = az egyed melyik részcsoporthoz esik (diszkrét változó),
- $\eta$  = egy vizsgált mennyiség (diszkrét vagy folytonos változó).

Modellezzük az  $\eta$  mennyiséget a következő módon: ha az egyed a  $j$ . csoportba esik (tehát  $\xi = j$ ), akkor legyen

$$\eta = a_j + \varepsilon = \text{csoporthatás} + \text{egyedi hatás},$$

ahol

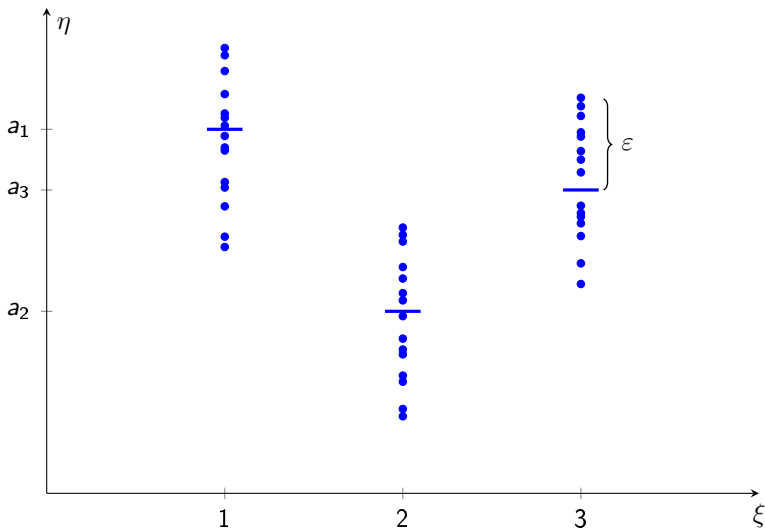
- $a_j$  rögzített valós szám, a  $j$ . csoport hatása,
- $E(\varepsilon) = 0$ , és  $\varepsilon$  független a  $\xi$  változótól.

Az  $\eta$  változó átlagos értéke a  $j$ . csoportban (feltételes várható érték):

$$E(\eta \mid \xi = j) = a_j + E(\varepsilon) = a_j.$$

Tehát a modell:

$$\eta = a_j + \varepsilon = \text{csoporthatás} + \text{egyedi hatás}$$

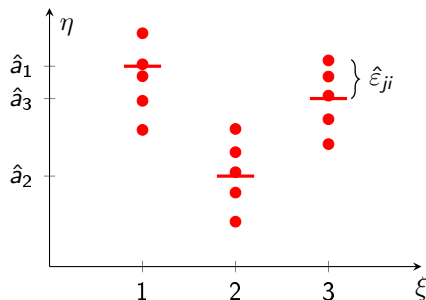
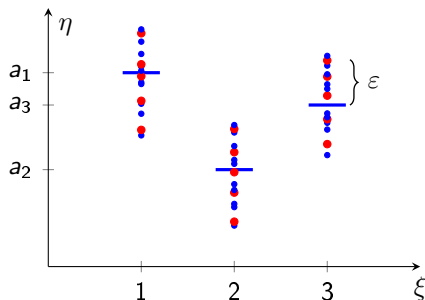


Adjunk becslést a populációátlagra és a csoporthatásokra. A teljes minta:

- minta az 1. részcsoporthra:  $\eta_{11}, \eta_{12}, \dots, \eta_{1n_1}$ ,
- minta a 2. részcsoporthra:  $\eta_{21}, \eta_{22}, \dots, \eta_{2n_2}$ ,
- ...
- minta az utolsó ( $r$ .) részcsoporthra:  $\eta_{r1}, \eta_{r2}, \dots, \eta_{rn_r}$ .

A teljes minta elemszáma:  $n = n_1 + n_2 + \dots + n_k$ . Mintaátlag:  $\bar{\eta}$ .

Mintaátlag a  $j$ . csoportban:  $\bar{\eta}_j = (\eta_{j1} + \eta_{j2} + \dots + \eta_{jn_j})/n_j$ .





Legyenek  $\hat{a}_1, \hat{a}_2, \dots, \hat{a}_k$  a becslések. A  $j$ . csoport egy tetszőleges egyedére vonatkozó **reziduális**:

$$\hat{\varepsilon}_{ji} = \eta_{ji} - \hat{a}_j \approx \eta_{ji} - a_j = \varepsilon_{ji}.$$

A paramétereket a legkisebb négyzetek módszerével fogjuk becsülni:

SSW = a mintaelemek reziduálisainak négyzetösszege

$$= \sum_{j=1}^r \sum_{i=1}^{n_j} \hat{\varepsilon}_{ji}^2 = \sum_{j=1}^r \sum_{i=1}^{n_j} [\eta_{ji} - \hat{a}_j]^2 \longrightarrow \min$$

Egy kis számolás után a következő becsléseket kapjuk:

Mennyiség	Elméleti érték	Becslés
$j$ . csoport átlaga	$a_j$	$\hat{a}_j = \bar{\eta}_j$
Egyedi hatás (hibatag)	$\varepsilon_{ji} = \eta_{ji} - a_j$	$\hat{\varepsilon}_{ji} = \eta_{ji} - \bar{\eta}_j$

Az előző oldalról:

$$\text{SST} = \text{sum of squares (total)} = \sum_{j=1}^r \sum_{i=1}^{n_j} (\eta_{ji} - \bar{\eta})^2$$

= milyen mértékben szóródnak az adatok a mintaátlag körül

Ez a szóródás két forrásból származik:

$$\text{SSW} = \text{sum of squares (within groups)} = \sum_{j=1}^r \sum_{i=1}^{n_j} (\eta_{ji} - \bar{\eta}_j)^2$$

= milyen mértékben szóródnak az adatok a csoportátlagok körül

$$\text{SSB} = \text{sum of squares (between groups)} = \sum_{j=1}^r \sum_{i=1}^{n_j} (\bar{\eta}_j - \bar{\eta})^2$$

= milyen mértékben szóródnak a csoportátlagok a mintaátlag körül

Megmutatható, hogy  $\text{SST} = \text{SSW} + \text{SSB}$ .

## Varianciaanalízis (Analysis of Variances, ANOVA)

A cél azt tesztelni, hogy minden csoportnak azonos a hatása, tehát a teljes populáción belül azonosak a csoportátlagok.

- Feltevések: az  $\eta$  változó minden csoporton belül normális eloszlást követ, és minden csoporton belül azonos a szórása.
- Nullhipotézis:  $H_0 : a_1 = a_2 = \dots = a_r$ .
- Próbat statisztika:  $F = \frac{SSB}{SSW} \frac{n-r}{r-1}$ .

A kapott értékeket az **ANOVA táblázatban** szoktuk összefoglalni:

Effect	Df	Sum Sq	Mean Sq	F value	Pr(>F)
$\xi$	$r - 1$	SSB	$\frac{SSB}{r-1}$	$\frac{SSB}{SSW} \frac{n-r}{r-1}$	p-value
Residuals	$n - r$	SSW	$\frac{SSW}{n-r}$		
Total	$n - 1$	SST			

Ha a csoportonkénti szórások azonosak, akkor  $D(\varepsilon) \approx SSW/(n - r)$ .

## Megjegyzések:

- Alapötlet: ha  $H_0$  igaz, akkor  $SSB \approx 0$ , tehát  $F \approx 0$ . Akkor fogadjuk el a nullhipotézist, ha  $F \leq c$ , ahol  $c$  a kritikus érték.
- A teszt robusztus a normalitásfeltételre nézve, de nem robusztus a szórásfeltételre nézve. Ha a szórások nem azonosak, akkor használjuk inkább a Welch-féle ANOVA tesztet.
- Két csoport ( $r = 2$ ) esetén: ANOVA = kétmintás t-próba.
- Azt szeretjük, ha közel ugyanannyi megfigyelés esik minden csoportba.

Milyen módszerrel ellenőrizhető a csoportonkénti szórások egyenlősége?

- Levene-teszt: formálisan kell hozzá a csoportonkénti normalitás, de robusztus erre a feltételre nézve. (Opcióknál: median!)
- Bartlett-teszt: formálisan ehhez is kell a csoportonkénti normalitás, és nem robusztus erre a feltételre nézve. (Ezért kevésbé ajánlott.)
- F-próba: csak  $r = 2$  csoport esetén működik, és ez a legérzékenyebb a normalitásra.

Várható érték és szórás tesztelése független minták alapján, illetve a robusztusság a normalitásfeltételre nézve:

Csoportok száma	Várható érték teszt (azonos szórás)	Várható érték teszt (tetszőleges szórás)	Szórás teszt
$r = 2$	Kétmintás t-próba (robusztus)	Welch-próba (robusztus)	F-próba (nem robusztus!!!)
$r \geq 2$	ANOVA (robusztus)	Welch-féle ANOVA (robusztus)	Levene-teszt (robusztus)

# Lineáris modellek

A lineáris regresszió és az ANOVA azonos alapötletre épül: egy magyarázó változó megfelelő függvényével modellezzük egy függő változó értékét. Ezt több magyarázó változóval is meg lehet tenni:

- Két (vagy több) magyarázó változó:  $\xi$ ,  $\zeta$ .
- Függő változó:  $\eta$ .
- Az egyedre jellemző hibatag:  $\varepsilon$ , független a magyarázó változóktól.

Lássunk néhány ilyen modellt!

## **Többszörös lineáris regresszió (multivariate linear regression):**

Tipikusan akkor alkalmazzuk, ha a magyarázó változók folytonosak, de diszkrét változókra is lehet. A modell:

$$\eta = (a\xi + b\zeta + c) + \varepsilon = \text{predikciós tag} + \text{hibatag.}$$

Például: felnőtt embereket vizsgálunk,

$$\text{testsúly} = a \cdot \text{testmagasság} + b \cdot \text{derékbőség} + c + \varepsilon.$$

## Többszemponos varianciaanalízis (multi-factor ANOVA):

A magyarázó változók **faktorok (factors)**, tehát olyan diszkrét változók, melyek részpopulációkat definiálnak. Ha az egyedre  $\xi = i$  és  $\zeta = j$ , akkor

$$\eta = (m + a_i + b_j) + \varepsilon = \text{predikciós tag} + \text{hibatag.}$$

Jelölések:

- $m$  = populációátlag,
- $a_i$  = az első faktor szerint az  $i$ -edik részpopuláció hatása,
- $b_j$  = a második faktor szerint a  $j$ -edik részpopuláció hatása.

Például: felnőtt embereket vizsgálunk,

- $\xi$ : nem (1=férfi, 2=nő), csoporthatások:  $a_1, a_2$ ,
- $\zeta$ : életmód, szokott-e sportolni (1=soha, 2=időnként, 3=naponta), csoporthatások:  $b_1, b_2, b_3$ .

$$\text{testsúly} = \text{populációátlag} + \text{nem hatása} + \text{életmód hatása} + \varepsilon.$$

**Kevert modell (mixed model):** Akkor alkalmazzuk, ha a magyarázó változók között van faktor és folytonos változó is. Legyen  $\xi$  a folytonos változó és  $\zeta$  a faktor. A modellben ha  $\zeta = j$ , akkor

$$\eta = (a\xi + b_j + c) + \varepsilon = \text{predikciós tag} + \text{hibatag.}$$

Például: felnőtt embereket vizsgálunk,

$$\text{testsúly} = a \cdot \text{testmagasság} + \text{életmód hatása} + c + \varepsilon.$$

**Lineáris modell (linear model):** A magyarázó változók tetszőlegesen (folytonosak vagy diszkrét) lehetnek. A modellben rögzítettek a  $g$  és  $h$  függvények, melyek segítségével:

$$\eta = a \cdot g(\xi) + b \cdot h(\zeta) + c + \varepsilon.$$

A bemutatott modellek (többszörös lineáris regresszió, többszemponos ANOVA, kevert modell) mind felírhatóak ilyen alakban, tehát ezek speciális esetei a lineáris modellnek.



# Valószínűségek becslése és tesztelése

Feladat: Becsüljük meg egy statisztikai minta alapján azon egyedek arányát egy populációban, melyek rendelkeznek egy adott tulajdonsággal.

Válasszunk ki véletlenszerűen egy egyedet a populációból, és legyen:

$A$  = a kiválasztott egyed rendelkezik a vizsgált tulajdonsággal

Ekkor:  $P(A)$  = a vizsgált tulajdonság aránya a populációban (ismeretlen)

Vezessük be a következő valószínűségi változót:

$$\xi = \begin{cases} 1, & \text{ha a kiválasztott egyed rendelkezik a vizsgált tulajdonsággal,} \\ 0, & \text{ha nem rendelkezik.} \end{cases}$$

Ekkor:  $P(\xi = 1) = P(A)$  és  $P(\xi = 0) = 1 - P(A)$ , tehát

$$\begin{aligned} E(\xi) &= 0 \cdot P(\xi = 0) + 1 \cdot P(\xi = 1) = 0 + 1 \cdot P(A) = P(A) \\ &= \text{a vizsgált tulajdonság aránya a populációban.} \end{aligned}$$

Tehát igazából egy várható értéket kell megbecsülnünk!

Statisztikai minta: kiválasztunk  $n$  egyedet a populációból, és jelölje  $\xi_1, \dots, \xi_n$  a  $\xi$  változó értékét a kiválasztott egyedek esetében. Ekkor:

$$\text{a tulajdonság aránya a populációban} = P(A) = E(\xi) \approx \frac{\xi_1 + \dots + \xi_n}{n}.$$

Sőt, ez egy erősen konzisztens becslés, tehát  $\bar{\xi} \rightarrow P(A)$ , amint  $n \rightarrow \infty$ .

**Gyakoriság, tapasztalati gyakoriság (frequency):**  $k_A = \xi_1 + \dots + \xi_n$ .

A mintaelemek közül ennyi rendelkezik a vizsgált tulajdonsággal.

**Relatív gyakoriság (relative frequency):**  $k_A/n$ .

A mintaelemek ekkora hányada rendelkezik a vizsgált tulajdonsággal.

A relatív gyakoriság erősen konzisztens becslés: ha  $n \rightarrow \infty$ , akkor

$$k_A/n \rightarrow P(A) = \text{a tulajdonság aránya a teljes populációban}$$

**Feladat:** Megvizsgáltunk 200 japán nemzetiségű embert, közülük 84 esett az A vércsoportba. Adjunk becslést az A vércsoport arányára Japánban!

Az A vércsoport tapasztalati gyakorisága illetve relatív gyakorisága:

$$k_A = 84, \quad k_A/n = 84/200 = 42\% \approx \text{arány a populációban.}$$

Tudjuk: a tulajdonság aránya a populációban =  $P(A) = E(\xi)$ .

A várható értékre vonatkozó módszerek alkalmazásával lehetőség van:

- Konfidencia intervallumot adni a  $P(A)$  valószínűségre.
- Tesztelni a  $P(A)$  valószínűség értékét. Legyen  $p \in [0, 1]$  tetszőleges hipotetikus valószínűség, és tekintsük az alábbi nullhipotézist:

$$H_0 : P(A) = p \quad \text{tehát} \quad H_0 : E(\xi) = p.$$

Ez a nullhipotézis tesztelhető t-próbával.

**FONTOS:** Most a  $\xi$  háttérváltozó nem normális eloszlást követ, emiatt ezek a módszerek csak nagy mintaméretre működnek. Tipikusan legyen  $n \geq 20$ , de inkább  $n \geq 50$ .

Mit tegyünk, ha csak kevés mintaelemünk van?

- Ne használjunk t-próbát!
- Alkalmazzuk a **binomiális próbát**, ugyanis ez tetszőleges  $n$  esetén alkalmazható. (Ezt a próbát nem tanuljuk.)

**Feladat:** Megvizsgáltunk 200 japán nemzetiségű embert, közülük 84 esett az A vércsoportba. Teszteljük azt a nullhipotézist, hogy a japán emberek körében az A vércsoport aránya 40%! Adjunk 95% megbízhatóságú konfidencia intervallumot erre az arányra!

Vezessük be a következő valószínűségi változót:

$$\xi = \begin{cases} 1, & \text{ha a kiválasztott ember az A vércsoportba esik,} \\ 0, & \text{ha nem oda esik.} \end{cases}$$

Nullhipotézis:  $H_0 : E(\xi) = 0,4$ .

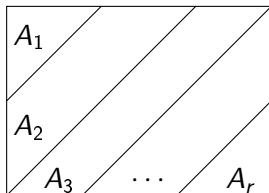
- Statisztikai minta  $(\xi_1, \dots, \xi_{200})$ : 84 db 1-es és 116 db 0-ás érték.
- Mintaátlag:  $\bar{\xi} = 0,42$ .
- Korrigált empirikus szórás:  $D_n^*(\xi) = 0,495$ .
- Standard hiba:  $SE = D_n^*(\xi)/\sqrt{n} = 0,035$ .
- Próbat statisztika:  $t = (\bar{\xi} - 0,4)/SE = 0,57$ .
- Kritikus érték:  $c = \Phi_{199}^{-1}(0,975) = 1,97$ .
- Döntés:  $|t| \leq c$ , ezért a nullhipotézist elfogadjuk.
- Konfidencia intervallum:  $[\bar{\xi} - c SE, \bar{\xi} + c SE] = [0,35, 0,49]$ .

## Mit tesztelhetünk még hasonló módszerrel?

- Két vagy több részpopuláción belüli arányok összehasonlítása független minták alapján. Például:
  - $\xi_1, \dots, \xi_{200}$ : 200 japán ember közül ki esik az A vércsoportba (0/1)
  - $\eta_1, \dots, \eta_{100}$ : 100 magyar ember közül ki esik az A vércsoportba (0/1) $H_0$  : a japánoknál és a magyaroknál azonos az A vércsoport aránya  
Tesz: Welch-próba vagy ANOVA
- Egy populáción belüli két arány összehasonlítása összetartozó minták alapján. Például:
  - $\xi_1, \dots, \xi_{200}$ : 200 japán ember közül ki esik az A vércsoportba (0/1)
  - $\eta_1, \dots, \eta_{200}$ : ugyanezen emberek közül ki esik a B vércsoportba (0/1) $H_0$  : a japánoknál azonos az A és a B vércsoport aránya  
Tesz: páros  $t$ -próba.

**FONTOS:** kell a nagy minta! Legalább 50 megfigyelés kell minden egyes változóra. Kis mintaelemszám esetén olyan tesztek kell keresni, melyek speciálisan arányokra vannak kitalálva.

Tegyük fel, hogy a populáció valamely szempont szerint több részcsoportha bontható fel. Ezek a részcsoporthok együttesen lefedik az összes egyedet, és minden csoportban van egyed. Feladat: teszteljük a részcsoporthok arányát.



Legyen  $r \geq 2$  a részcsoporthok száma. Véletlenszerűen kiválasztunk egy egyedet, és legyen

$A_i =$  a kiválasztott egyed az  $i$ -edik részcsoporthba esik,  $i = 1, \dots, r$ .

Válasszunk ki véletlenszerűen  $n$  egyedet a populációból. Ekkor:

- $k_i =$  az  $i$ -edik csoport gyakorisága a mintában,
- $k_i/n =$  az  $i$ -edik csoport relatív gyakorisága a mintában.

Becslés:  $k_i/n \approx P(A_i) =$  az  $i$ -edik részcsoporth aránya.

Tekintsünk hipotetikus valószínűségeket (részarányokat):  $p_1, \dots, p_m > 0$ , ahol  $p_1 + \dots + p_m = 1$ . A célunk a következő nullhipotézist tesztelni:

$$H_0 : P(A_i) = p_i \quad \text{minden } i \text{ részcsoport esetén.}$$

Ötlet: ha igaz a nullhipotézis, akkor az  $n$  elemű mintában körülbelül  $np_i$  olyan egyednek kell lennie, mely az  $i$ -edik csoportba tartozik.

### $\chi^2$ -próba (khinégzet-próba) valószínűségek tesztelésére

A fenti nullhipotézis tesztelhető a következő próbastatisztikával:

$$\chi^2 = \frac{(k_1 - np_1)^2}{np_1} + \dots + \frac{(k_r - np_r)^2}{np_r}$$

A kritikus érték a  $\chi^2$ -eloszlás táblázatából kereshető ki.

**Feltétel:** nagy minta,  $n \geq 5 / \min(p_1, \dots, p_r)$ .

Vegyük észre, a fenti összegben minden tag a következő módon áll elő:

$$\frac{(\text{tapasztalati gyakoriság} - \text{várt gyakoriság})^2}{\text{várt gyakoriság}}$$

**Feladat:** Megvizsgáltunk 200 japán nemzetiségű embert. Közülük rendre 62, 84, 38 illetve 24 esett a 0, az A, a B és az AB vércsoportba. Adjunk becslést a vércsoportok arányára a teljes populáción belül! Teszteljük 90%-os szignifikancia szinten azt a nullhipotézist, hogy a japán emberek körében a vércsoportok aránya rendre 30%, 40%, 20% illetve 10%!

	0	A	B	AB	össz.
Tapasztalati gyakoriság ( $k_i$ )	62	84	38	24	200
Relatív gyakoriság ( $k_i/n$ )	0,31	0,42	0,19	0,12	1
Hipotetikus valószínűség ( $p_i$ )	0,3	0,4	0,2	0,1	1
Várt gyakoriság ( $np_i$ )	60	80	40	20	200

Próbastatisztika:

$$\chi^2 = \frac{(62 - 60)^2}{60} + \frac{(84 - 80)^2}{80} + \frac{(38 - 40)^2}{40} + \frac{(24 - 20)^2}{20} = 1,17.$$

A kritikus érték a  $\chi^2$ -eloszlás táblázatából:  $c = 6,25$ .

Most  $|\chi^2| \leq c$ , tehát a nullhipotézist elfogadjuk.



# Függetlenségvizsgálat

Kérdés: milyen módon tesztelhető egy  $\xi$  és egy  $\eta$  mennyiség függetlensége?

Válasz: attól függ, hogy milyen változókról van szó...

- Korrelációs teszt a Pearson-féle korrelációs együtthatóval:
  - Csak normális eloszlású változókra.
  - A korrelációs együttható segítségével tényleg a függetlenséget teszteli.
- Korrelációs teszt a Spearman-féle korrelációs együtthatóval:
  - Folytonos eloszlású változókra, nem csak normálisra.
  - A korrelációs együttható segítségével tényleg a függetlenséget teszteli.
- ANOVA, kétmintás t-próba, Welch-próba:
  - Egyik változó diszkrét (csoportokat definiál), a másik folytonos.
  - Azt teszteli, hogy nincs csoporthatás, tehát az  $\eta$  változó várható értékére nem hat a  $\xi$  változó. Ez jóval kevesebb, mint a függetlenség!
- $\chi^2$ -próba (khinégyzet-próba) függetlenségre:
  - Nem azonos a valószínűségeknel tanult  $\chi^2$ -próbával!
  - Mindkét változó diszkrét, a függetlenséget teszteli.
  - Kell a nagy mintaméret!

Adott két mennyiség a populációban,  $\xi$  és  $\eta$ . Összetartozó minták:

- $\xi_1, \dots, \xi_n$ : a  $\xi$  mennyiség értéke  $n$  kiválasztott egyednél,
- $\eta_1, \dots, \eta_n$ : az  $\eta$  mennyiség értéke ugyanezen egyedeknél.

$H_0$ : a két mennyiség független egymástól.

Tesztelési módszer a változók típusa szerint:

	$\eta$ diszkrét	$\eta$ normális	$\eta$ folytonos
$\xi$ diszkrét	$\chi^2$ -próba (nagy minta!)	ANOVA	ANOVA (nagy minta!)
$\xi$ normális	ANOVA	Pearson-korreláció	Spearman-korreláció (nagy minta!)
$\xi$ folytonos	ANOVA (nagy minta!)	Spearman-korreláció (nagy minta!)	Spearman-korreláció (nagy minta!)

# Illeszkedésvizsgálat

Feladat: becsüljük meg és tesztljük le egy  $\xi$  mennyiség eloszlását a teljes populációban a  $\xi_1, \dots, \xi_n$  megfigyelések alapján!

Legyen  $\xi$  **diszkrét**, és legyen  $R_\xi = \{x_1, \dots, x_k\}$  az értékkészlete.  
Feladat: becsüljük meg és tesztljük a  $P(\xi = x_i)$  valószínűségeket.

Becslés:

$P(\xi = x_i) \approx$  relatív gyakoriság = az  $x_i$  érték aránya a mintában

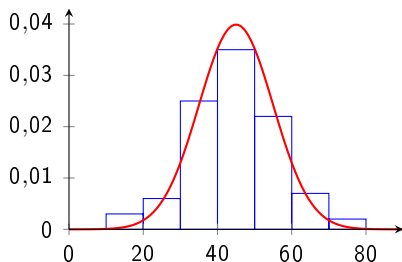
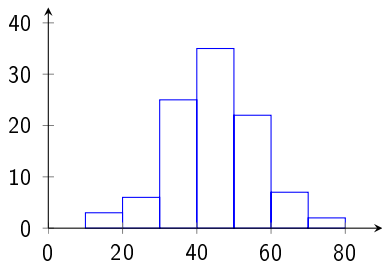
Tesztlés: adottak  $p_1, \dots, p_k > 0$  hipotetikus arányok,  $p_1 + \dots + p_k = 1$ .

$H_0 : P(\xi = x_i) = p_i$  minden  $i$ -re

Tesztlési módszer:  $\chi^2$ -próba a valószínűségekre. Probléma megoldva.

A továbbiakban csak azzal az esettel foglalkozunk, amikor  $\xi$  folytonos eloszlású egy ismeretlen  $f_\xi$  sűrűségfüggvénnyel.

Sűrűségfüggvény becslése **hisztogrammal**: bontsuk fel a számegyenest azonos (mondjuk  $h > 0$ ) hosszúságú intervallumokra. Minden intervallumra állítsunk egy olyan magas oszlopot, ahány elem esik az adott intervallumba. Az így kapott hisztogrammnak még nincs sok köze a sűrűségfüggvényhez.



Osszuk minden oszlop magasságát  $nh$ -val! Az így kapott új hisztogramm összterülete pontosan 1 lesz.

**Tétel.** Ha  $n \rightarrow \infty$  és  $h \rightarrow 0$ , akkor az átskálázott hisztogramm konvergál az ismeretlen sűrűségfüggvényhez. Ilyen módon grafikus becslést adhatunk a sűrűségfüggvényre.

A legtöbb statisztikai program nem olyan módon ábrázolja a boxplotot, ahogyan azt korábban tanultuk. Általában felmérnek a doboz aljára és tetejére  $1.5 \cdot \text{IQR}$  távolságot, és a bajúszt csak eddig ábrázolják. Az ezen kívül eső megfigyeléseket **outlier** értékeknek nevezzük, ezek egyesével vannak ábrázolva a boxploton.

Normalitásvizsgálat **boxplot** segítségével: amennyiben a minta normális eloszlásból jön, akkor a boxplotnak két speciális tulajdonsága van:

- A boxplot körülbelül szimmetrikus a mediánra.
- Az outlier értékek a teljes minta legfeljebb 1%-át teszik ki. Az ennél több outlier arra utal, hogy a minta nem normális eloszlásból jön. Ezt a tulajdonságot jól szemlélteti ez a [Wikipedia ábra](#).

Mi lehet még az oka a sok outlier értéknek? Mérési hibák, jegyzőkönyvezési hibák, stb. Az adatelemzés során az outlier értékeket külön-külön meg szokták vizsgálni, nem hibából származnak-e.

Rendezzük növekvő sorrendbe a megfigyelt értékeket:  $\xi_1^* \leq \dots \leq \xi_n^*$ . Ezt nevezzük **rendezett mintának**. Legyen  $q_\alpha$  a  $\xi$  változó  $\alpha$ -kvantilise.

**Tétel.** Nagy mintaelemszám esetén  $\xi_i^* \approx q_{i/(n+1)}$  minden  $i$ -re.

**Q-Q plot:** Koordináta-rendszerben ábrázoljuk a  $(q_{i/(n+1)}, \xi_i^*)$  pontokat.

Eloszlásvizsgálat Q-Q plot segítségével: kíváncsiak vagyunk arra, hogy a minta egy adott (például normális) eloszlásból származik-e.

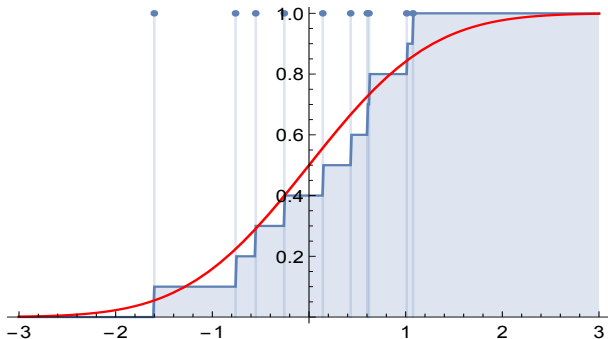
- Kiszámoljuk a kérdéses eloszlás kvantiliseit:  $q_{1/(n+1)}, \dots, q_{n/(n+1)}$ .
- Ábrázoljuk a Q-Q plotot.
- Ha a minta a kérdéses eloszlásból származik, akkor  $\xi_i^* \approx q_{i/(n+1)}$  minden  $i$ -re, tehát minden pont az  $x = y$  egyenes közelébe esik.
- Ha a minta nem a kérdéses eloszlásból jön, akkor  $\xi_i^* \not\approx q_{i/(n+1)}$  bizonyos megfigyelésekre, tehát egyes pontok nem illeszkednek az  $x = y$  egyeneshez.
- Hátrány: nem egzakt módszer, hanem szubjektív döntés.
- Előny: azt is le lehet olvasni az ábráról, hogy a minta mely érték-tartományban milyen mértékben illeszkedik a kérdéses eloszláshoz.

## Empirikus eloszlásfüggvény (sample distribution function):

$F_n(x)$  = az  $x$  értéknél kisebb elemek relatív gyakorisága a mintában

Tulajdonságok:

- Minden  $x$  valós számra  $F_n(x) \in [0, 1]$ .
- Az  $F_n$  függvény egy lépcsős függvény. A mintaelemeknél van ugrása, és minden mintaelemnél  $1/n$  az ugrás nagysága.
- Nagy mintaméret esetén  $F(x) = P(\xi < x) \approx F_n(x)$ .



**A matematikai statisztika alaptétele:** Jelölje  $F_\xi$  a minta valódi eloszlásfüggvényét. Ekkor

$$\max_{x \in \mathbb{R}} |F_n(x) - F_\xi(x)| \rightarrow 0, \quad n \rightarrow \infty.$$

A tétel következményei:

- A matematikai statisztikának van értelme: a mintában van elég információ ahhoz, hogy mindent meg tudjunk becsülni.
- Tesztelhetjük az eloszlásfüggvényt.

## Kolmogorov–Szmirnov-próba

Cél: teszteljük egy  $\xi$  változó eloszlásfüggvényét egy minta alapján. Legyen  $F_\xi$  az igazi eloszlásfüggvény (ismeretlen) és  $F_0$  egy tetszőleges hipotetikus eloszlásfüggvény. Nullhipotézis:  $H_0 : F_\xi = F_0$ . Próbastatisztika:

$$D_n = \max_{x \in \mathbb{R}} |F_n(x) - F_0(x)|.$$

Akkor fogadjuk el a nullhipotézist, ha  $D_n \approx 0$ .



## A Kolmogorov–Szmirnov-próba tulajdonsága:

- Csak konkrét eloszlást lehet vele tesztelni. Például:  
 $H_0$  : a  $\xi$  változó normális eloszlású 0 várható értékkel és 1 szórással.
- Tetszőleges eloszlás esetén alkalmazható, de főleg folytonos változókra szokták alkalmazni.
- Hátránya: kis mintaelemszám esetén alacsony az ereje, a hamis nullhipotéziseket is elfogadja.

## Cél a normalitás tesztelése:

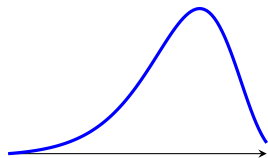
$H_0$  : a  $\xi$  változó normális eloszlást követ valamilyen paraméterekkel

## Fontosabb tesztelési módszerek:

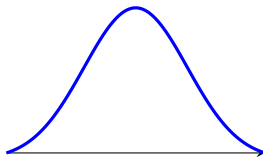
- **Lilliefors-teszt:** Becslést végez a várható értékre és a szórássra a minta alapján, majd azt teszteli a K–Sz-próbával, hogy a minta normális eloszlást követ  $\mu = E_n(\xi)$  és  $\sigma = D_n^*(\xi)$  paraméterekkel.
- **Shapiro–Wilk-teszt:** Sokak szerint ez a legjobb normalitásteszt, kis mintaelemszám esetén is magas az ereje.

**Skewness (ferdeség):** egy olyan statisztikai mutatószám, mely a minta alapján jellemzi a sűrűségfüggvény szimmetriáját. Tulajdonságai:

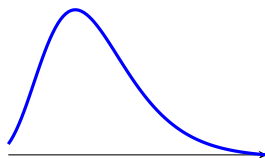
- skewness  $\approx 0$ : (közel) szimmetrikus sűrűségfüggvény
- skewness  $> 0$ : jobbra ferde sűrűségfüggvény
- skewness  $< 0$ : balra ferde sűrűségfüggvény



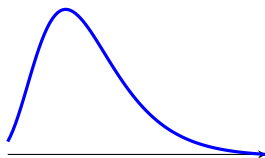
skewness = -1



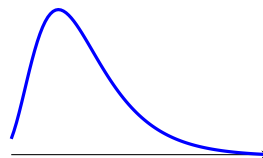
skewness = 0



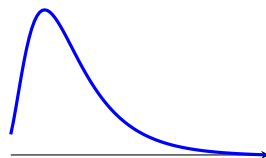
skewness = 1



skewness = 1.2



skewness = 1.5



skewness = 2

## Néhány további gondolat

Az adatelemzés akkor hatékony, ha a statisztikai minta jól reprezentálja a teljes populációt. Hogyan kaphatunk **reprezentatív mintát**?

Véletlenszerű mintavételezés: véletlenszerűen kiválasztunk egyedet a teljes populációból.

- Előny: egyszerűen és olcsón megvalósítható (biztos?)
- Hátrány: a reprezentativitáshoz nagy mintaméretre lehet szükség.

Irányítottan összeállított minta: a teljes populáción belüli arányokat figyelembe véve magunk állítunk össze egy mintát.

- Előny: kisebb mintaméret, mint véletlenszerű mintavételezésnél.
- Hátrány: előzetes ismeretekre van szükség a populációról; bonyolult és gyakran drága.
- Gyakran súlyozást alkalmaznak a populációarányok reprezentálásához.

Egyes statisztika programok (SPSS, R Commander) típusokba sorolják a változókat, és csak azokat az eljárásokat engedik futtatni, amik megfelelnek az adott típusnak. Milyen típusokról tanultunk eddig:

- Folytonos változó: mindig valós szám az értéke, az értékkészlete egy intervallumon.
- Diszkrét változó: valós szám vagy szöveg is lehet az értéke, de az értékkészlete véges.

Még egy kifejezés:

- Faktor: csoportokat definiáló változó. Mindig diszkrét.

Milyen típusokba sorolják egyes programok a változókat:

- Skálaváltozó: értelmezhetőek a matematikai műveletek (összeadás, átlagolás). Például: testmagasság, utódok száma, vizsgajegy(?).
- Ordinális változó: nem értelmezhetőek a matematikai műveletek, de van rendezés az értékek között. Például: rendfokozatok, ordinális skálák.
- Nominális változó: nem értelmezhetőek a matematikai műveletek és rendezés sincs az értékek között. Például: nem, nemzetiség.