

Matematika Statisztika

2. előadás

Arató Miklós

2022.09.19.

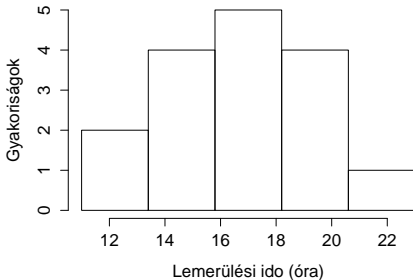
Hisztogram – Ha a mennyiségi ismerv folytonos vagy sok ismervérték van, akkor alkalmas módon osztályokat képezünk, majd minden egyes adatot pontosan egy osztályhoz rendeljük. A hisztogram az osztályok gyakoriságait ábrázolja.

- javaslat az osztályok számára:

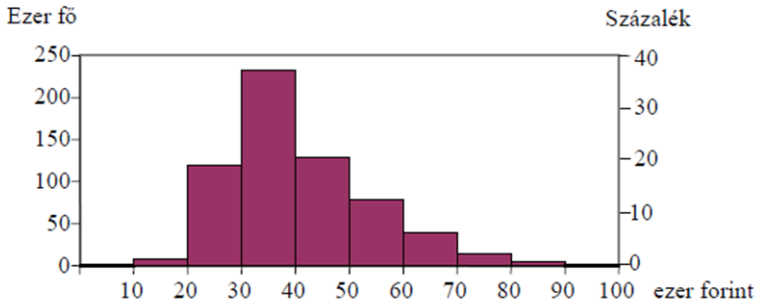
$$k = \lfloor \log_2 n \rfloor$$

- ha azonos hosszúságú (h) osztályközöket akarunk létrehozni, akkor $h = \frac{x_n^* - x_1^*}{k}$
- az f_i gyakoriságokat ábrázoljuk a függőleges tengelyen
- sűrűséghisztogramnál a $g_i = \frac{f_i}{n}$ relatív gyakoriságokat ábrázoljuk a függőleges tengelyen

- **ha az osztályközök különböző hosszúságúak, akkor a gyakoriságokat egy közös hosszra kell arányosítani**

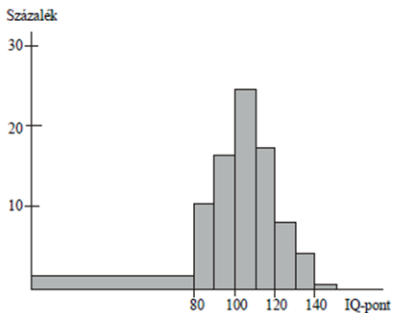
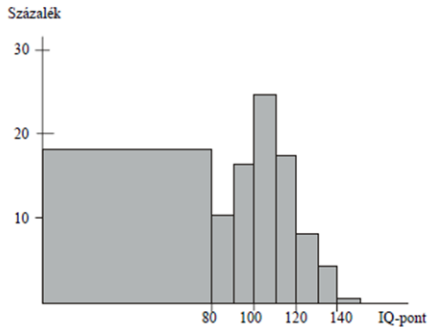


21. ábra. A férfi népesség megoszlása az öregségi nyugdíjak nagysága szerint
2000. január 1.



Adatforrás: Magyar statisztikai évkönyv, 1999 (2000). Központi Statisztikai Hivatal, Budapest.

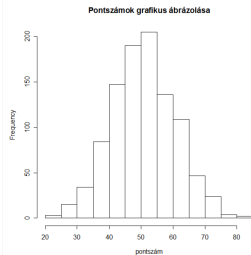
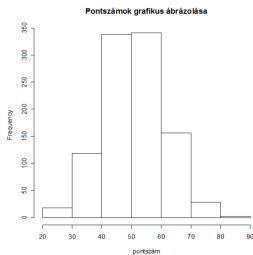
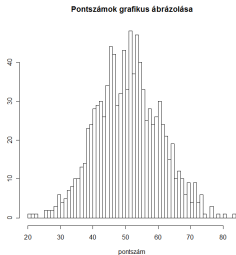
15. ábra. Egy népességcsoport megoszlása IQ-pontok szerint



15. táblázat. Egy népességcsoport megoszlása IQ-pontok szerint

IQ-pont	Százalék	IQ-pont	Százalék	IQ-pont	Százalék
-80	18,2	101-110	24,7	131-140	4,1
81-90	10,8	111-120	17,3	141-150	0,5
91-100	16,2	121-130	8,2	Összesen	100,0

Hisztogram (folyt.)



- **Tapasztalati eloszlás:** minden megfigyeléshez azonos, $\frac{1}{n}$ súlyt rendelünk \Rightarrow ez egy diszkrét eloszlás
- A mintaátlag éppen ennek a várható értéke
- A tapasztalati eloszlás eloszlásfüggvényét hívjuk **tapasztalati eloszlásfüggvénynek**, ami egy tiszta ugrófüggvény, értéke minden mintaelem helyén $\frac{1}{n}$ nagyságot ugrik felfelé.
A tapasztalati eloszlásfüggvény az x helyen:

$$\frac{I(x_1 < x) + I(x_2 < x) + \dots + I(x_n < x)}{n} = \frac{\sum_{i=1}^n I(x_i < x)}{n}$$

Azt mutatja meg, hogy a mintaelemek hányad része kisebb x -nél.

Középértékek számítása

Adott az n elemű $\underline{x} = (x_1, x_2, \dots, x_n)$ tapasztalati minta; osztályközös gyakorisági sor esetén k jelöli az osztályok számát, x_i az osztályközöket, f_i pedig a gyakoriságokat.

Mintaátlag: az adatok átlagos értéke

- Számítása közvetlenül az adatokból: $\bar{x} = \frac{\sum_{i=1}^n x_i}{n}$
- Számítása osztályközös gyakorisági sorból: $\bar{x} = \frac{\sum_{i=1}^k f_i x_i}{n}$

Módusz: a legtöbbször előforduló ismérték

- Számítása osztályközös gyakorisági sorból:

$$Mo = x_{mo,a} + \frac{d_a}{d_a + d_f} \cdot h_{mo}, \text{ ahol}$$

- a móduszt tartalmazó osztályköz: amelyben egységnyi osztályköz hosszra a legnagyobb gyakoriság jut (\rightsquigarrow *korrigált gyakoriságok!*)
- $x_{mo,a}$: a móduszt tartalmazó osztályköz alsó értéke
- h_{mo} : a móduszt tartalmazó osztályköz hossza
- d_a : a móduszt tartalmazó osztályköz korrigált gyakorisága mínusz a móduszt közvetlenül megelőző osztályköz korrigált gyakorisága
- d_f : a móduszt tartalmazó osztályköz korrigált gyakorisága mínusz a móduszt közvetlenül követő osztályköz korrigált gyakorisága

Középértékek számítása

Adott az n elemű $\underline{x} = (x_1, x_2, \dots, x_n)$ tapasztalati minta; osztályközös gyakorisági sor esetén k jelöli az osztályok számát, x_i az osztályközöket, f_i pedig a gyakoriságokat.

Mintaátlag: az adatok átlagos értéke

- Számítása közvetlenül az adatokból: $\bar{x} = \frac{\sum_{i=1}^n x_i}{n}$
- Számítása osztályközös gyakorisági sorból: $\bar{x} = \frac{\sum_{i=1}^k f_i x_i}{n}$

Módusz: a legtöbbször előforduló ismérték

- Számítása osztályközös gyakorisági sorból:

$$Mo = x_{mo,a} + \frac{d_a}{d_a + d_f} \cdot h_{mo}, \text{ ahol}$$

- a móduszt tartalmazó osztályköz: amelyikben egységnyi osztályköz hosszra a legnagyobb gyakoriság jut (\rightsquigarrow *korrigált gyakoriságok!*)
- $x_{mo,a}$: a móduszt tartalmazó osztályköz alsó értéke
- h_{mo} : a móduszt tartalmazó osztályköz hossza
- d_a : a móduszt tartalmazó osztályköz korrigált gyakorisága mínusz a móduszt közvetlenül megelőző osztályköz korrigált gyakorisága
- d_f : a móduszt tartalmazó osztályköz korrigált gyakorisága mínusz a móduszt közvetlenül követő osztályköz korrigált gyakorisága

Középértékek számítása

Jelölje $x_1^* \leq x_2^* \leq \dots \leq x_n^*$ a rendezett tapasztalati mintát.

Medián: azon ismérték, amelynél ugyanannyi kisebb vagy egyenlő, mint nagyobb vagy egyenlő ismérték fordul elő a mintában (a "közepső" elem)

- Számítása közvetlenül az adatokból:

$$\text{Me} = \begin{cases} x_{\frac{n+1}{2}}^*, & \text{ha } n \text{ páratlan} \\ \frac{x_{\frac{n}{2}}^* + x_{\frac{n}{2}+1}^*}{2}, & \text{ha } n \text{ páros} \end{cases}$$

- Számítása osztályközös gyakorisági sorból – két lépésben lineáris interpolációval:

1. Melyik osztályközben van a medián: azon i , amire $f'_{i-1} \leq \frac{n}{2}$ és $f'_i \geq \frac{n}{2}$

2. $\text{Me} = x_{i,a} + \frac{\frac{n}{2} - f'_{i-1}}{f'_i} \cdot h_i$, ahol

- $x_{i,a}$: a mediánt tartalmazó osztályköz alsó értéke
- h_i : a mediánt tartalmazó osztályköz hossza
- f'_{i-1} : a mediánt közvetlenül megelőző osztályköz kumulált gyakorisága
- f'_i : a mediánt tartalmazó osztályköz gyakorisága

Tapasztalati kvantilisok számítása

Tapasztalati y -kvantilis: azon ismértérték, amelynél a mintaelemek y -ad része kisebb vagy egyenlő, míg $(1 - y)$ -ad része nagyobb vagy egyenlő, $0 < y < 1$

Számítása nem egyértelmű, mi mindig az egyik interpolációs módszert alkalmazzuk két lépésben:

1. hányadik mintaelem a keresett kvantilis \rightsquigarrow sorszám: $s := (n + 1)y$
2. lineáris interpolációval a kvantilis kiszámítása

- Számítása közvetlenül az adatokból

1. Sorszám: $s = e + t$ (e : egészrész, t : törtrész)

2. $q_y = x_e^* + t(x_{e+1}^* - x_e^*)$

- Számítása osztályközös gyakorisági sorból – két lépésben lineáris interpolációval:

1. Melyik osztályközben van az s -edik elem: jelölje ezt i , azaz $f'_{i-1} \leq s$ és $f'_i \geq s$

2. $q_y = x_{i,a} + \frac{s-f'_{i-1}}{f'_i} h_i$, ahol

$x_{i,a}$, h_i , f'_{i-1} és f'_i ugyanazokat jelöli, mint az előző fólia alján, csak az adott y -kvantilisre vonatkozóan

Tapasztalati kvantilisok számítása

Tapasztalati y -kvantilis: azon ismértérték, amelynél a mintaelemek y -ad része kisebb vagy egyenlő, míg $(1 - y)$ -ad része nagyobb vagy egyenlő, $0 < y < 1$

Számítása nem egyértelmű, mi mindig az egyik interpolációs módszert alkalmazzuk két lépésben:

1. hányadik mintaelem a keresett kvantilis \rightsquigarrow sorszám: $s := (n + 1)y$
2. lineáris interpolációval a kvantilis kiszámítása

- Számítása közvetlenül az adatokból

1. Sorszám: $s = e + t$ (e : egészrész, t : törtrész)

2. $q_y = x_e^* + t(x_{e+1}^* - x_e^*)$

- Számítása osztályközös gyakorisági sorból – két lépésben lineáris interpolációval:

1. Melyik osztályközben van az s -edik elem: jelölje ezt i , azaz $f'_{i-1} \leq s$ és $f'_i \geq s$

2. $q_y = x_{i,a} + \frac{s-f'_{i-1}}{f'_i} h_i$, ahol

$x_{i,a}$, h_i , f'_{i-1} és f'_i ugyanazokat jelöli, mint az előző fólia alján, csak az adott y -kvantilisre vonatkozóan

A szakirodalomban a tapasztalati és az elméleti értékek között nem tesznek különbséget, mindegyiket nagy betűvel írják (ami néha meglehetősen zavaró...). Jelölje q_y a tapasztalati y -kvantilist.

- tercilisek: $T_1 = q_{1/3}$, $T_2 = q_{2/3}$
- **kvartilisek:**
 - $Q_1 = q_{1/4}$ (alsó kvartilis)
 - $Q_2 = \mathbf{Me} = q_{2/4}$ (középső kvartilis vagy medián)
 - $Q_3 = q_{3/4}$ (felső kvartilis)
- kvintilisek: $K_1 = q_{1/5}$, $K_2 = q_{2/5}$, $K_3 = q_{3/5}$, $K_4 = q_{4/5}$
- decilisek: $D_i = q_{i/10}$, $i = 1, 2, \dots, 9$
- percentilisek: $P_i = q_{i/100}$, $i = 1, 2, \dots, 99$

Szóródási mutatók számítása

Terjedelem: $R = x_n^* - x_1^*$ (R =range)

Interkvartilis terjedelem: $IQR = Q_3 - Q_1$

Tapasztalati szórás: az átlagtól való átlagos eltérés abszolút mértékegységben

- Számítása közvetlenül az adatokból: $s_n = \sqrt{\frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n}}$
- Számítása osztályközös gyakorisági sorból: $s_n = \sqrt{\frac{\sum_{i=1}^k f_i (x_i - \bar{x})^2}{n}}$

Korrigált tapasztalati szórás: az átlagtól való átlagos eltérés abszolút mértékegységben

- Számítása közvetlenül az adatokból: $s_n^* = \sqrt{\frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n-1}}$
- Számítása osztályközös gyakorisági sorból: $s_n^* = \sqrt{\frac{\sum_{i=1}^k f_i (x_i - \bar{x})^2}{n-1}}$
- ezt "szeretjük" a legjobban, minden szoftver, programcsomag szórás számításánál ezt veszi alapértelmezettnek

Relatív szórás vagy **szórási együttható**: az átlagtól való átlagos eltérés százalékban; lehet a korrigált és a korrigálatlan tapasztalati szórásnégyzetből is számítani:

$$V = \frac{s_n^*}{\bar{x}} \text{ vagy } V = \frac{s_n}{\bar{x}}$$

Kevésbé gyakran használt, szóródást mérő mutatók:

- átlagos abszolút eltérés: $\frac{\sum_{i=1}^n |x_i - \bar{x}|}{n}$
- Gini-együttható: $G = \frac{1}{n(n-1)} \sum_{i=1}^n \sum_{j=1}^n |x_i - x_j|$.

A szórást ezeknél is választhatjuk a tapasztalati vagy a korrigált tapasztalati szórásnak egyaránt.

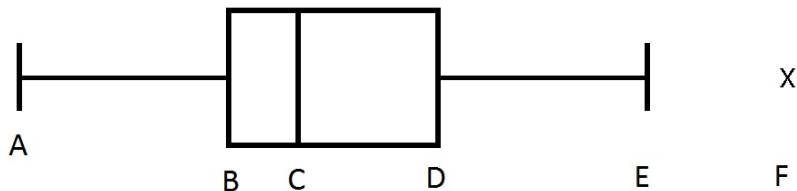
Tapasztalati ferdeség

- Számítása közvetlenül az adatokból: $\frac{\sum_{i=1}^n (x_i - \bar{x})^3}{(s_n)^3}$
- Számítása osztályközös gyakorisági sorból: $\frac{\sum_{i=1}^n f_i (x_i - \bar{x})^3}{(s_n)^3}$

Tapasztalati csúcsosság

- Számítása közvetlenül az adatokból: $\frac{\sum_{i=1}^n (x_i - \bar{x})^4}{(s_n)^4} - 3$
- Számítása osztályközös gyakorisági sorból: $\frac{\sum_{i=1}^n f_i (x_i - \bar{x})^4}{(s_n)^4} - 3$

Boxplot ábra (Box&Whiskers diagram) – ez fekvő, de lehet álló is



A betűk a következő értékeket jelentik:

- $A = \max\{x_1^*, Q_1 - 1,5 \cdot IQR\}$
- $B = Q_1$
- $C = Me$
- $D = Q_3$
- $E = \min\{x_n^*, Q_3 + 1,5 \cdot IQR\}$
- F : kieső érték (outlier) \rightsquigarrow azokat az adatpontokat tüntetjük fel, amik A -n vagy E -n kívülre esnek

ahol $IQR = Q_3 - Q_1$ az interkvartilis terjedelem

Milyen valószínűséggel születik fiúgyermek? (Példa)

- Svájcban 1871 és 1900 között a 2.644.757 megszületett gyermekből 1.359.671 fiú és 1.285.086 lány volt.
- Fiúk relatív gyakorisága így 0,5141.
- Hogyan becsüljük a születési valószínűségeket?
- Tudunk esetleg valamilyen intervallumot adni ezekre a valószínűségekre?
- Igaz-e, hogy a valószínűség 0,5? És 0,1?