



Valószínűségszámítás és Statisztika

8. előadás
2017. november 12.

Példa

- Milyen valószínűséggel születik fiúgyermek?
- 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.
- Igaz-e, hogy a valószínűség 0,5? És 0,1?

$$X_i = \begin{cases} 1, & i.\text{fiú} \\ 0, & i.\text{lány} \end{cases} \Rightarrow$$

$$P(X_i = 1) = p, n = 2.644.757, \xi = \frac{\sum_{i=1}^n X_i}{n} \Rightarrow$$

$$EX_i = p, D^2 X_i = p(1-p), P\left(\frac{\sum_{i=1}^n X_i - nEX_1}{DX_1 \sqrt{n}} < x\right) \sim \Phi(x) \Rightarrow$$

$$P\left(-u < \sqrt{\frac{n}{p(1-p)}} (\xi - p) < u\right) \sim 2\Phi(u) - 1$$

$$p = 0.5 \Rightarrow \sqrt{\frac{n}{p(1-p)}} (\xi - p) = 37$$

$$u = 4 \Rightarrow 2\Phi(u) - 1 = 0,999936$$

$$p(1-p) \leq \frac{1}{4} \Rightarrow$$

$$2\Phi(u) - 1 \sim P\left(-u < \sqrt{\frac{n}{p(1-p)}} (\xi - p) < u\right) \leq$$

$$\leq P\left(-u < 2\sqrt{n} (\xi - p) < u\right) =$$

$$= P\left(\frac{-u}{2\sqrt{n}} < (\xi - p) < \frac{u}{2\sqrt{n}}\right) = P\left(\xi - \frac{u}{2\sqrt{n}} < p < \xi + \frac{u}{2\sqrt{n}}\right)$$

Esetünkben 0,9973 valószínűséggel $0,5132 \leq p \leq 0,5150$

Statisztikai mező

$$(\Omega, \mathcal{A}, P_{\vartheta}), \vartheta \in \Theta$$

statisztikai mező, ha Θ paraméterhalmaz
és $(\Omega, \mathcal{A}, P_{\vartheta})$ minden paraméter
esetén valószínűségi mező.

Egy érmédobás modellje

- Nem ismerjük a fejdobás valószínűségét:

$$\Omega = \{F, I\}, A = \{\emptyset; \{F\}; \{I\}; \{F, I\}\},$$

$$P_p(\{F\}) = p, P_p(\{I\}) = 1 - p, p \in [0, 1].$$

Minta

Def.: A $\xi = \begin{pmatrix} \xi_1 \\ \vdots \\ \xi_n \end{pmatrix} : \Omega \rightarrow X \subseteq \mathbb{R}^n$ valószínűségi vektorváltozót

mintának nevezzük.

n : mintanagyság

ξ_i : i . mintaelem

Def.: minta realizációja: $\mathbf{x} = \begin{pmatrix} x_1 \\ \vdots \\ x_n \end{pmatrix}$ a konkrét megfigyelt számsorozat.

Mintatér

- Def: \mathfrak{X} mintatér: a minta lehetséges értékeinek halmaza. Elemei a mintaértékek.
- n -elemű valós minta esetén: $\mathfrak{X} = \mathbb{R}^n$
- n -elemű pozitív egész értékű minta esetén: $\mathfrak{X} = \mathbb{N}^n$
- Példa: egy biztosítónál 10 napon keresztül figyelték a bejelentett károk számát, ekkor $\mathfrak{X} = \mathbb{C}_0^n$

Egy benzinkútnál tankoló autók száma 5 napon keresztül

- Megfigyelések: 78, 89, 167, 90, 85
- Minta realizációja:
 $(78, 89, 167, 90, 85)^T$
- Mintanagyság: 5

A minták típusai

- Független minta: a mintaelemek függetlenek.
- Független azonos eloszlású minta: a mintaelemek függetlenek és azonos eloszlásúak.
- Diszkrét minta: a mintaelemek diszkrétek.
- Abszolút folytonos eloszlású minta: a mintaelemek abszolút folytonosak.

Eloszláscsaládok

$$F_g(\mathbf{s}) = P_g(\xi_1 < s_1, \dots, \xi_n < s_n)$$

Független minta esetén:

$$F_g(\mathbf{s}) = \prod_{i=1}^n P_g(\xi_i < s_i)$$

Független azonos eloszlású minta esetén:

$$F_g(\mathbf{s}) = \prod_{i=1}^n P_g(\xi_i < s_i) = \prod_{i=1}^n F_g(s_i)$$

Jelölések:

E_g : várható érték P_g esetén,

D_g : szórás P_g esetén,

f_g : sűrűségfüggvény P_g esetén (absz. folyt. minta)

$p_g(s) = P_g(\xi_i = s)$ (diszkrét minta)

Példák

- Egy érmédobás. Fej esetén 1-et írunk, írás esetén 0-át.

$$p_p(k) = P_p(\xi_1 = k) = \begin{cases} p & k = 1 \\ 1 - p & k = 0 \end{cases} = p^k (1 - p)^{1-k}$$

Benzinkutas példa. Azt feltételezzük, hogy megfigyeléseink független, azonos eloszlású Poissonok.

$$p_\lambda(k) = P_\lambda(\xi_i = k) = \lambda^k e^{-\lambda} / k!, \quad k = 0, 1, 2, \dots$$

Statisztikák

Def.: Statisztika: a minta függvénye.

$$T : X \rightarrow R^k$$

Def'.: Statisztika:

$T(\xi)$, ha $T : X \rightarrow R^k$ függvény.

Példák

Tapasztalati momentumok:

$$X = R^n,$$

$$\text{mintaközép: } T(\mathbf{x}) = \bar{x} = \frac{\sum_{i=1}^n x_i}{n}, \quad T(\xi) = \bar{\xi} = \frac{\sum_{i=1}^n \xi_i}{n},$$

$$\text{tapasztalati } k. \text{ momentum: } T(\mathbf{x}) = \frac{\sum_{i=1}^n x_i^k}{n}, \quad T(\xi) = \frac{\sum_{i=1}^n \xi_i^k}{n}.$$

Tapasztalati szórásnégyzet

$$X = R^n,$$

$$T(\mathbf{x}) = \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n}, \quad T(\xi) = s^2 = \frac{\sum_{i=1}^n (\xi_i - \bar{\xi})^2}{n}$$

Rendezett minta

- A ξ_1, \dots, ξ_n minta elemeit nagyság szerint sorbarendezeve kapjuk az $\xi_1^{(n)} \leq \xi_2^{(n)} \leq \dots \leq \xi_n^{(n)}$ rendezett mintát.
- Ez n -dimenziós statisztika
- Mostantól: a ξ_1, \dots, ξ_n minta elemei független, azonos eloszlásúak.
- Ha feltesszük, hogy a közös eloszlásuk abszolút folytonos, akkor felírható a rendezett minta k -adik elemének, $\xi_k^{(n)}$ -nek a sűrűségfüggvénye. (gyakorlat)
- Spec.: minimum, maximum.
- Def.: minta terjedelme: $\xi_n^{(n)} - \xi_1^{(n)}$

Tapasztalati eloszlásfüggvény

- Előző előadáson már szerepelt.
- Tapasztalati eloszlás eloszlásfüggvénye: tapasztalati eloszlásfüggvény:

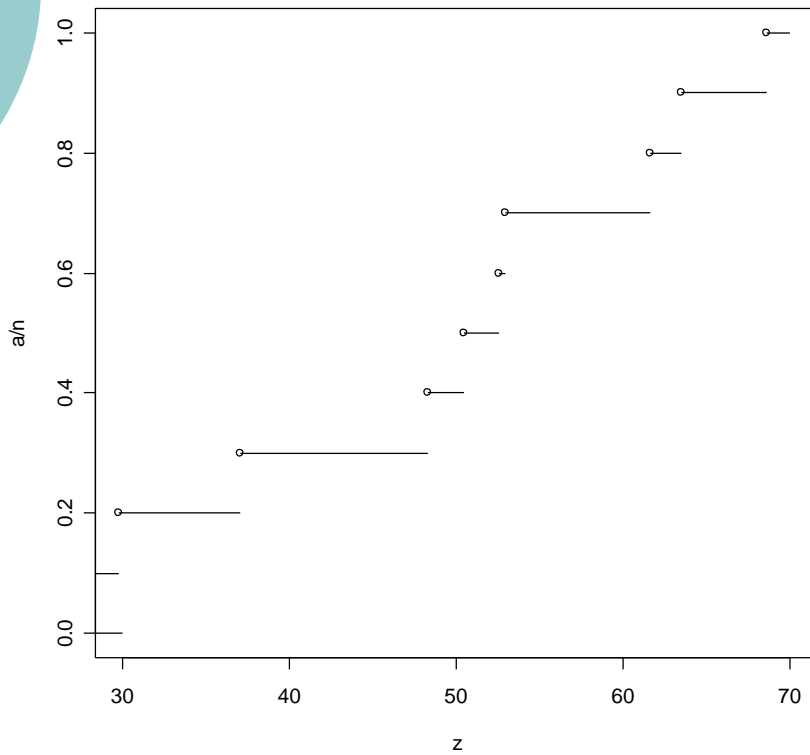
$$F_n(z) = \frac{1}{n} \sum_{i=1}^n \chi \{ \xi_i < z \}$$

$$F_n(z) = \frac{k}{n}, \text{ ha } \xi_k^{(n)} < z \leq \xi_{k+1}^{(n)}, \xi_0^{(n)} = -\infty, \xi_{n+1}^{(n)} = \infty$$

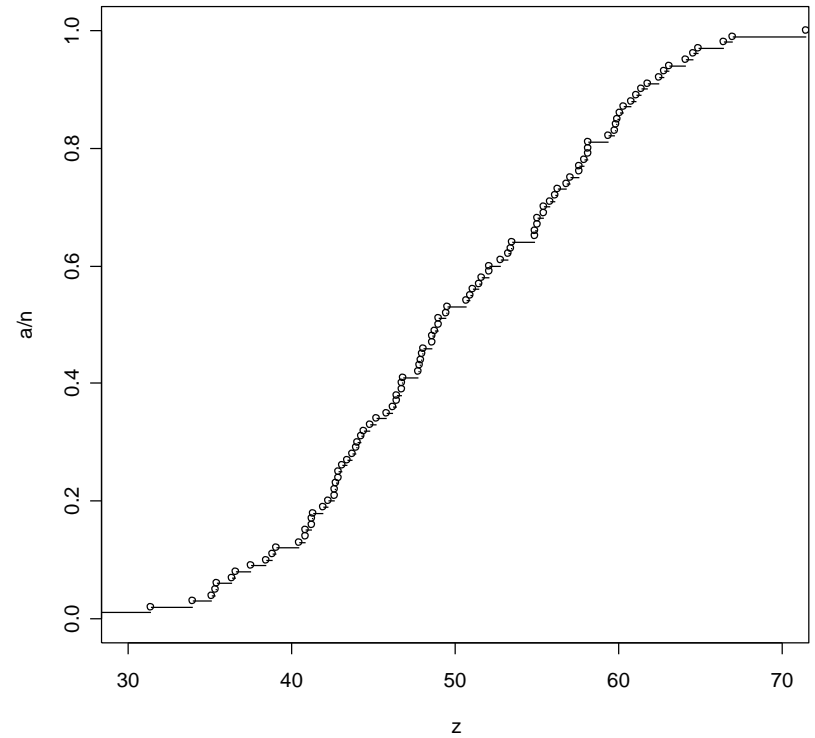
Mintaátlag éppen ennek az eloszlásnak a várható értéke.

Példa

normális eloszlás közelítése, $n=10$



normális eloszlás közelítése, $n=100$



Glivenko-Cantelli tétel ("statisztika alaptétele")

Tétel: ξ_1, \dots, ξ_n független, azonos F eloszlásfüggvényűek. Ekkor $\sup_z |F_n(z) - F(z)| \xrightarrow{n \rightarrow \infty} 0$ majdnem mindenütt (1 vszgel).

Biz.: Csak folytonos F eloszlásfüggvényekre látjuk be. Ebből következik, hogy tetszőleges pozitív egész N -hez léteznek olyan valós z_1, \dots, z_N számok, hogy

$$F(z_0) = 0, F(z_1) = \frac{1}{N}, \dots, F(z_i) = \frac{i}{N}, \dots, F(z_{N-1}) = \frac{N-1}{N}, F(z_N) = 1,$$

$$z_0 = -\infty, z_N = \infty.$$

Becslésemélet

- A minta eloszlásának ismeretlen paraméterét közelítjük a minta függvényével
- Def.: becslőfüggvény: $\hat{\mathcal{G}}: X \rightarrow \Theta$
- Def.: becslés: $\hat{\mathcal{G}}(\xi)$

- A becslések maguk is statisztikák.
Szubjektíven: olyan statisztikák, amik jól közelítik az ismeretlen paramétert.

Példa (Milyen valószínűséggel születik fiúgyermek?)

- Svájcban 1871 és 1900 között a 2.644.757 megszületett gyermekből 1.359.671 fiú és 1.285.086 lány volt.
- Ekkor $n = 2.644.757$, $\mathfrak{X} = \{0;1\}^n$.
- Fiúk relatív gyakorisága így 0,5141.
- Mik ennek a becslésnek a tulajdonságai?

$$X_i = \begin{cases} 1, & i.\text{fiú} \\ 0, & i.\text{lány} \end{cases} \Rightarrow$$

$$P_p (X_i = 1) = p, n = 2.644.757,$$

$$\hat{p} = \hat{p}(\mathbf{X}) = \frac{\sum_{i=1}^n X_i}{n} \Rightarrow$$

$$E_p \hat{p} = p, \hat{p} \xrightarrow{n \rightarrow \infty} p \text{ mm.}$$

Becslések tulajdonságai

- Def.: *Torzítatlanság*: A paraméter $\hat{\vartheta}(\xi)$ becslése torzítatlan, ha

$$E_{\vartheta} \left(\hat{\vartheta}(\xi) \right) = \vartheta, \quad \forall \vartheta \in \Theta.$$

- *Konzisztencia* $\hat{\vartheta}(\xi) \xrightarrow{\mathcal{P}} \vartheta$ sztochasztikusan ($n \rightarrow \infty$) minden paraméterértékre.
- Példák:
 - Valószínűség becslése relatív gyakorisággal.
 - Glivenko tétele: a tapasztalati eloszlásfüggvény egyenletesen is konvergál az elméleti eloszlásfüggvényhez.
 - Várható érték becslése mintaátlaggal

Konzisztencia

- Elégséges feltétel $E_{\mathcal{G}} \left(\hat{\mathcal{G}}_n(\xi) \right) \rightarrow \mathcal{G}$
(aszimptotikus torzítatlanság)
és

$$D_{\mathcal{G}}^2 \left(\hat{\mathcal{G}}_n(\xi) \right) \rightarrow 0$$

Példák

- Poisson eloszlás paraméterére: mintaátlag
- Exponenciális eloszlás paraméterére:
 - mintaátlag reciproka: aszimptotikusan torzítatlan, konzisztens
 - $n \cdot \min(X_1, \dots, X_n)$ torzítatlan, de nem konzisztens
- Szórásnégyzetre

Becslések összehasonlítása

- Melyik a jobb becslés?

$$X_i = \begin{cases} 1, & i.\text{fiú} \\ 0, & i.\text{lány} \end{cases}, P_p(X_i = 1) = p,$$

$$\hat{p}_1 = \frac{\sum_{i=1}^n X_i}{n}, \quad \hat{p}_2 = X_1, \text{ vagy}$$

$$\hat{p}_3 = \frac{\sum_{i=1}^{\lfloor n/2 \rfloor} X_i}{\lfloor n/2 \rfloor} ?$$

Becslések összehasonlítása (hatásos becslések)

- Torzítatlan becslésekre: T_1 hatásosabb becslése $h(\theta)$ -nak a T_2 -nél, ha

$$D_{\theta}^2(T_1(\underline{X})) \leq D_{\theta}^2(T_2(\underline{X}))$$

teljesül minden θ paraméterértékre.

Példa: a mintaátlag hatásosabb becslés a várható értékre minden alakú becslésnél.

$$\sum_{i=1}^n c_i X_i$$

Hatásos becslés

- Def.: A T torzítatlan becslés hatásos, ha minden más torzítatlan becslésnél hatásosabb.
- Miért a torzítatlanokra? Furcsa példa: azonosan 0-val becsüljük az ismeretlen paramétert.
- Ezért érdemes a hatásos becsléseket csak a torzítatlan becslések között keresni.
- Átlagos négyzetes eltérés:
$$E_{\theta} (T(\underline{X}) - \theta)^2$$

Hatásos becslés egyértelműsége

- **Áll.:** Amennyiben T_1 és T_2 hatásos becslései $h(\theta)$ -nak, akkor 1 valószínűséggel megegyeznek minden lehetséges paraméter esetén.

Becslési módszerek

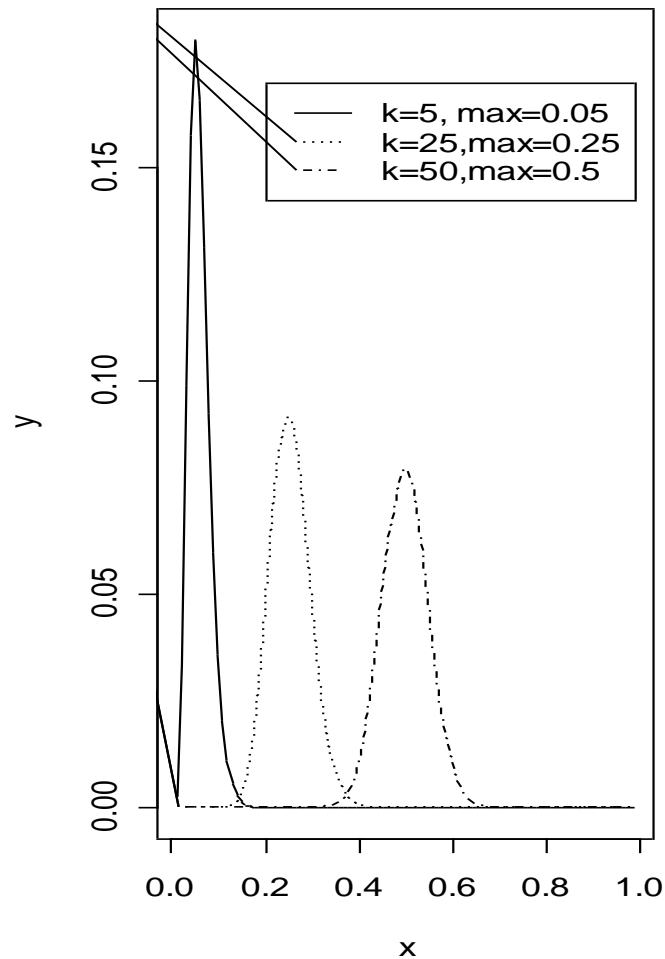
- Eddig: „ad hoc” módszerek
- Általános eljárás kellene
 - Példa: valószínűség becslése, n kísérletből. Jelölje k a sikeresek számát (X_i $i=1, \dots, n$ indikátorminta)

$$P\left(\sum_{i=1}^n X_i = k\right) = \binom{n}{k} p^k (1-p)^{n-k}$$

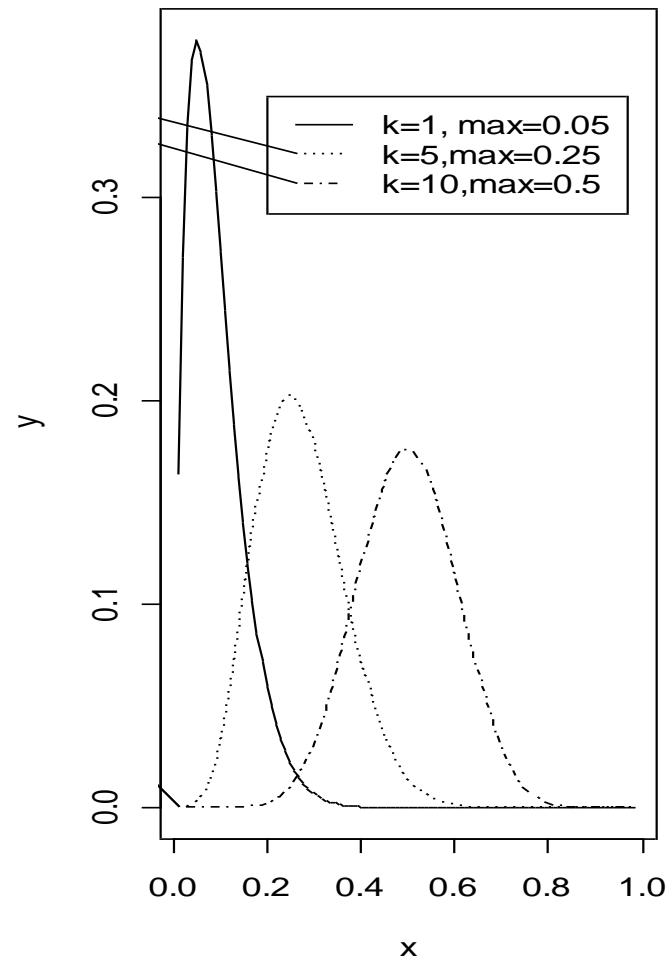
Most p függvényében nézzük, k rögzített (elnevezés: likelihood függvény).

A likelihood függvény maximumhelye logikus választás a valószínűség becslésének

likelihood függvény, $n=100$



likelihood függvény, $n=20$



A módszer általánosan

$$L(\theta; \underline{x}) = f_{\theta}(\underline{x}) = \prod_{i=1}^n f_{\theta}(x_i)$$

(a likelihood függvény) maximumhelye lesz a θ paraméter maximum likelihood becslése.

Ha a függvény deriválható, a loglikelihood függvény

$$l(\theta; \underline{x}) = \ln f_{\theta}(\underline{x}) = \sum_{i=1}^n \ln f_{\theta}(x_i)$$

maximumhelye deriválással

$$\frac{\partial}{\partial \theta} l(\theta; \underline{x}) = \frac{\partial}{\partial \theta} \ln f_{\theta}(\underline{x}) = \sum_{i=1}^n \frac{\partial}{\partial \theta} \ln f_{\theta}(x_i) = 0$$

megoldásaként megtalálható

Példák

- valószínűségre: relatív gyakoriság
- Poisson eloszlás paraméterére: x
- Exponenciális eloszlás paraméterére: $1/\bar{x}$
- Normális eloszlás várható értékére: \bar{x}
- A módszer többdimenzós paraméter becslésére is használható: $N(m, \sigma)$ esetén $(\bar{x}, \sum (x_i - \bar{x})^2 / n)$ a maximum likelihood becslés.