

Matematikai statisztika

Informatika alapszak, "B" spec.

Zempléni András

Valószínűségelméleti és Statisztika Tanszék
Matematikai Intézet
Természettudományi Kar
Eötvös Loránd Tudományegyetem

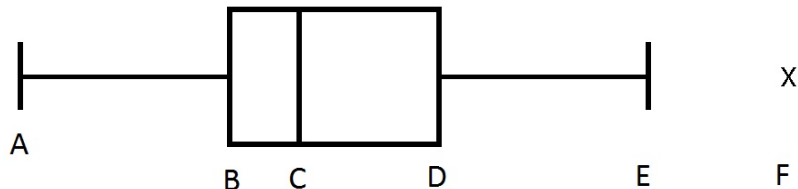
Honlap: zempleni.elte.hu

E-mail: andras.zempleni@ttk.elte.hu

Szoba: D 3-310

2. előadás

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



A betűk a következő értékeket jelentik:

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

ahol $IQR = Q_3 - Q_1$ az interkvartilis terjedelem

- **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.

- Statisztikai mező: $(\Omega, \mathcal{A}, P_{\vartheta}) : \vartheta \in \Theta$
- Paraméterter: Θ . Ez lehet egydimenziós, de akár végtelen dimenziós is
- Minta: $X = (X_1, \dots, X_n)$ független, azonos F eloszlású valószínűségi változók
- Mintater $\mathcal{X} : \mathbb{R}^n$ azon része, ahova a mintaelemek eshetnek
- A mintaelemek eloszlása ismeretlen, de paraméterezhető: $F \rightarrow F_{\vartheta}$
- Példák:
 - Poisson eloszlású minta, ekkor $\vartheta \rightarrow \lambda \in \Theta = (0; \infty)$
 - normális eloszlású minta, ekkor $\vartheta \rightarrow (\mu, \sigma) \in \Theta = (-\infty; \infty) \times (0; \infty) \subset \mathbb{R}^2$
 - F -ről nem tudunk semmit, ekkor Θ végtelen dimenziós. De ekkor is lehet egydimenziós paramétereket értelmezni, például várható érték, szórás

Motiváció – becslésmélet

Az Asus kicseréli táblagépeit, amennyiben a vevők 8-nál több pixelhibát jelentenek be vásárlástól számítva 3 napon belül. A Samsung már egyetlen, 3 napon belül bejelentett pixelhiba esetén is új készüléket biztosít. A Sony-nál legalább 2 pixelhiba esetén jár új táblagép.

Hogyan tudnánk megbecsülni, hogy a gyártónak éves szinten milyen mértékű vesztesége származik ezekből a cserékből?

- Kulcskérdés: mi az esélye, hogy egy, a gyártósorról véletlenszerűen leemelt készüléket pixelhiba miatt ki kell cserélni?
- Ha X a pixelhibák száma, akkor a kérdéses valószínűség például a Sony-nál: $P(X \geq 2)$
- Milyen eloszlású lehet X (Poisson?) → *illeszkedésvizsgálat*
- Ha tudom, hogy Poisson-eloszlású, akkor hogyan becsüljem meg a paramétert? → *pontbecslés*
- Milyen intervallumban lesz "nagy" valószínűséggel a becsült paraméter? → *intervallumbecslés*
- Ezután készíthető a kérdéses valószínűségre intervallumbecslés, abból pedig egy intervallumbecslés a várható veszteségre.

- Legyen $X = (X_1, \dots, X_n)$ i.i.d. minta egy ϑ valós paraméterű eloszláscsaládból. $T : \mathcal{X} \rightarrow \mathbb{R}$ becslés ϑ -ra.
- Tulajdonságai:
 - Torzítatlanság: $E_{\vartheta} T(X) = \vartheta$ minden $\vartheta \in \Theta$ paraméterre
 - Aszimptotikus torzítatlanság: $E_{\vartheta} T_n(X) \rightarrow \vartheta$ (ha $n \rightarrow \infty$) minden $\vartheta \in \Theta$ paraméterre
 - Konzisztencia: $T_n(X) \rightarrow \vartheta$ sztochasztikusan (ha $n \rightarrow \infty$) minden $\vartheta \in \Theta$ paraméterre (ez a gyenge, erős, ha 1 vszű a konv.)
- Megj.: A konzisztenciához elégséges, hogy T_n aszimptotikusan torzítatlan legyen és $D^2(T_n) \rightarrow 0$

Definíció. [Likelihood függvény] $L(\vartheta; x) = f_{\vartheta}(x) = \prod_{i=1}^n f_{\vartheta}(x_i)$, ha az eloszlás

folytonos és $L(\vartheta; x) = P_{\vartheta}(X = x) = \prod_{i=1}^n P_{\vartheta}(X_i = x_i)$, ha az eloszlás diszkrét

Definíció. [Log-likelihood függvény] $\ell(\vartheta; x) = \ln(L(\vartheta; x))$

Fontos becslések tulajdonságai

Tétel. Legyen X_1, \dots, X_n i.i.d. minta egy ϑ paraméterű eloszláscsaládból, $h: \mathbb{R} \rightarrow \mathbb{R}$ (mérhető) függvény. Tegyük fel, hogy a táblázatban szereplő összes várható érték/szórás létezik minden ϑ esetén.

Mit be- csülünk? $g(\vartheta)$	Mivel becsüljük? $T_n(X)$	Torzí- tatlan?	Aszimptotikusan torzítatlan?	Gyengén/ erősen konzisztens?
$E_{\vartheta}X_1$	$\bar{X} = \frac{\sum_{i=1}^n X_i}{n}$	igen	igen	igen
$D_{\vartheta}^2 X_1$	$S_n^2 = \frac{\sum_{i=1}^n (X_i - \bar{X})^2}{n}$	nem	igen	igen
$D_{\vartheta}^2 X_1$	$(S_n^*)^2 = \frac{\sum_{i=1}^n (X_i - \bar{X})^2}{n-1}$	igen	igen	igen
$F_{\vartheta}(x)$	$F_n(x) = \frac{\sum_{i=1}^n I(X_i < x)}{n}$	igen	igen	igen
$E_{\vartheta}h(X_1)$	$\frac{\sum_{i=1}^n h(X_i)}{n}$	igen	igen	igen

- **Maximum likelihood módszer** (ML-módszer): Azt a paraméterértéket keressük, ahol a likelihood függvény a legnagyobb értéket veszi fel: $\max_{\vartheta} L(\vartheta; \mathbf{x})$.
- Ez nyilván megegyezik azzal a paraméterértékkel, ahol a log-likelihood függvény veszi fel a legnagyobb értéket, azaz: $\max_{\vartheta} \ell(\vartheta; \mathbf{x})$.
- Amennyiben a függvény deriválható ϑ szerint, akkor a maximumot kereshetjük a deriváltak segítségével
- Célszerűbb a likelihood függvény helyett a log-likelihood függvény maximumhelyét keresni.
- Ha ϑ 1 dimenziós, akkor $\partial_{\vartheta} \ell(\vartheta, \mathbf{x}) = 0$, míg ha $\vartheta = (\vartheta_1, \dots, \vartheta_p)$ p dimenziós, akkor $\partial_{\vartheta_i} \ell(\vartheta, \mathbf{x}) = 0$ megoldásából kapjuk a becslést.

Tétel. [ML-becslés invariáns tulajdonsága] Ha ϑ ML-becslése $\hat{\vartheta}$, akkor tetszőleges g függvény esetén $g(\vartheta)$ ML-becslése $g(\hat{\vartheta})$.

Momentum módszer:

- A mintából számítható tapasztalati momentumokat ($m_i := \frac{1}{n} \sum_j x_j^i$) egyenlővé tesszük az elméleti momentumokkal ($M_i(\vartheta) := E_{\vartheta} X^i$),
- mégpedig annyit, amennyiből a paramétereket meg tudjuk határozni. p darab ismeretlen paraméter esetén p ismeretlenes egyenletrendszert oldunk meg ϑ -ra: $M_1(\vartheta) = m_1, \dots, M_p(\vartheta) = m_p$ (megjegyzés: $m_1 = \bar{x}$)
- Ha valamelyik egyenlet nem ad információt a keresett paraméterre, akkor magasabb hatványokat nézünk (míg megoldható nem lesz az egyenletrendszer)

Példák: Határozzuk meg a momentum és a maximum likelihood becslést az alábbi eloszlások paraméterére:

- Poisson
- exponenciális
- Legyen X_1 Bin(2; p) eloszlású (egyelemű) minta, ahol $p \in (0; 1)$ ismeretlen paraméter. Adjunk X_1 segítségével torzítatlan becslést $g(p) = \frac{1}{p}$ -re!