# 2. tétel

Nagy Dániel 2019. június 11.

#### **Kivonat**

Bootstrap módszerek. A maximum likelihood módszer. Hipotézis tesztelés. Extrém statisztikák. Post hoc analízis. Regresszió. Függetlenségvizsgálat. Egzakt tesztek.

#### 1. Bevezetés

### 1.1. Valószínűségszámítás alapfogalmak

- Eseménytér (ez egy abstrakt fogalom):  $\Omega = \{\omega_1, \omega_2, ..., \omega_n\}$  pl. kockadobás esetén  $\Omega = \{\omega_1 = \text{"lest dobok"}, \omega_2 = \text{"2est dobok"}, \omega_3 = \text{"párosat dobok"}...\}$
- Valószínűségi változó:  $X: \Omega \to \mathbb{R}$  pl. kockadobás esetén  $X(\omega_1) = 1, X(\omega_2) = 2, ...$
- Valószínűség: P egy mérték, amely  $\Omega$  részhalmazaihoz számot rendel:
  - $-P:\mathcal{P}(\Omega)\to\mathbb{R}$
  - $-P(\Omega) = 1$  és  $P(\emptyset) = 0$
  - $-0 \le P(A) \le 1 \ \forall A \in \Omega$
  - Ha  $A_1, A_2, \dots$  diszjunkt részhalmazai  $\Omega$ -nak, akkor

$$P\left(\bigcup_{i=1}^{\infty} A_i\right) = \sum_{i=1}^{\infty} P(A_i)$$

- Hasznos összefüggések:
  - $-P(A \cup B) = P(A) + P(B) P(A \cap B)$
  - Két esemény független  $\Longleftrightarrow P(A\cap B) = P(A)P(B)$
  - $-P(A|B) = \frac{P(A \cap B)}{P(B)}$
  - Teljes valószínűség: Ha $A_1,A_2,\dots$  az  $\Omega$ egy felosztása, akkor

$$P(B) = \sum_{k} P(B|A_k)P(A_k)$$

– Bayes-tétel: Ha  $A_1, A_2, \dots$  az  $\Omega$  egy felosztása, akkor

$$P(A_k|B) = \frac{P(B|A_k)P(A_k)}{P(B)} = \frac{P(B|A_k)P(A_k)}{\sum_{j} P(B|A_j)P(A_j)}$$

• Eloszlásfüggvény (CDF - cumulative distribution function):

$$F_X(x) = P(X < x) = P(\{\omega \in \Omega | X(\omega) < x\})$$

diszkrét esetben

$$F_X(x) = P(X = x) = P(\{\omega \in \Omega | X(\omega) = x\})$$

2

Ha az X változó F eloszlást követ, akkor így jelöljük:  $X \sim F$ .

• Sűrűségfüggvény (PDF - Probability density function): Ha az X változó eloszlásfüggvénye  $F_X(x)$ , akkor a sűrűségfüggvény definíciója

$$F_X(x) = \int_{-\infty}^x \rho_X(\xi) d\xi \iff P(a \le X(\omega) \le b) = \int_a^b \rho_X(x) dx$$

Megjegyzés: sűrűségfüggvénye csak folytonos eloszlású valószínűségi változónak van.

#### • Várható érték

folytonos eset 
$$E(X) = \langle X \rangle = \int_{-\infty}^{\infty} x \rho(x) dx$$
  
diszkrét eset  $E(X) = \langle X \rangle = \sum_{k} x_{k} p_{k} = \sum_{k} x_{k} P(X = x_{k})$ 

• Várható értékre vonatkozó azonosságok:

– Ha 
$$Y=g(X)\Rightarrow E(Y)=E(g(X))=\int\limits_{-\infty}^{\infty}g(x)\rho(x)\mathrm{d}x$$
 –  $E\left(\sum_{k}a_{k}X_{k}\right)=\sum_{k}a_{k}E(X_{k})$  – Ha  $X_{1},X_{2},...$  független változók, akkor  $E\left(\prod_{k}X_{k}\right)=\prod_{k}E(X_{k})$ 

• Variancia (szórásnégyzet)

Ha  $E(X) = \mu$ , akkor a szórásnégyzet a változó és a várható értéke közötti különbség négyzetének várható értéke:

$$\sigma^2(X) = V(X) = E((X - \mu)^2) = \langle (X - \mu)^2 \rangle = \langle X^2 \rangle - \mu^2$$

• Ha  $X_1, X_2, ...$  függetlenek, akkor

$$\sigma^2 \left( \sum_k (a_k X_k + b_k) \right) = \sum_k a_k^2 \sigma^2(X_k)$$

• Szórás (standard deviation) definíciója:

$$\sigma(X) = \sqrt{\sigma^2(X)} = \sqrt{\langle X^2 \rangle - \langle X \rangle^2}$$

Minta

Matematikailag egy statisztikai minta megfelel N darab azonos eloszlású, független (iid) változónak egy adott F eloszlásból.

• Minta átlaga:  $\overline{X} = \frac{1}{N} \sum_{i=1}^{N} X_i \ (p_k$ -t a relatív gyakorisággal közelítjük)

- Minta varianciája:  $s^2 = \frac{1}{N-1} \sum_{i=1}^{N} (X_i \overline{X})^2$ , standard hibája  $SE = \sqrt{s^2}$ . A nevezőben az N-1 faktor az ún. Bessel-korrekció [1].
- Megjegyzés: Ha egy teljes populáció esetén  $E(X) = \mu$  és  $V(X) = \sigma^2$ , attól még általában  $\overline{X} \neq \mu$  illetve  $s^2 \neq \sigma^2$ .
- Egy minta esetében  $\overline{X}$ ,  $s^2$ , SE maguk is valószínűségi változók, hiszen minden mintavételezés esetén más-más értéket vehetnek fel. Ezért van értelme arról beszélni, hogy pl.  $s^2$  értéke milyen eloszlást követ. Ha a minta (mérési pontok) iid változók, és  $E(X_i) = \mu$ ,  $V(X_i) = \sigma^2$ , akkor

$$E(\overline{X}) = \mu$$

$$V(\overline{X}) = \sigma^2/N$$

$$E(s^2) = \sigma^2$$

### 1.2. Statisztikai következtetés (inference)

- Az alapprobléma: van egy adathalmaz, ami tartalmazza a méréseket. Ezek  $X_1, X_2, ..., X_N \sim F$  független, azonos F eloszlást követő valószínűségi változók.
- A statisztikai következtetés feladata, hogy a minta alapján meghatározzuk az F eloszlásfüggvényt. Ezzel ekvivalens, ha F helyett a  $\rho$  sűrűségfüggvényt határozzuk meg.
- $\bullet$  Ehhez használhatunk parametrikus és nem-parametrikus modelleket. A parametrikus modell egy olyan  $\mathcal{F}$  halmaz, ami a lehetséges PDF-eket tartalmazza:

$$\mathcal{F} = \{ \rho(x|\theta) : \theta \in \Theta \},\$$

ahol $\Theta$ a lehetséges paraméterek halmaza. Pl. ha normális eloszlást feltételezünk, akkor a parametrikus modell

$$\mathcal{F} = \left\{ \rho(x|\mu, \sigma) = \frac{1}{\sigma\sqrt{2\pi}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right) : \mu \in \mathbb{R}, \sigma > 0 \right\},\,$$

a feladat pedig  $\mu$  és  $\sigma$  meghatározása. Nem-parametrikus modellek azok, amelyeket nem lehet véges számú valós paraméterrel definiálni, pl.  $\mathcal{F} = \{az \text{ összes létező PDF}\}.$ 

# 2. Bootstrap módszerek

#### 2.1. Jackknife módszer

## 3. Maximum likelihood

A maximum likelihood módszer egy olyan becslési eljárás, amelynek segítségével egy parametrikus modell paramétereinek értékét próbáljuk a minta alapján meghatározni. Ehhez felírjuk

az ún. likelihood-függvényt, ami azt fejezi ki, hogy a mért adatok esetén mekkora a valószínűsége a  $\theta$  paramétereknek. Ha a változó elposzlása ismert, akkor ezzel megadható a likelihood függvény:

diszkrét változóra: 
$$\mathcal{L}(\theta) = P(X = x | \theta)$$
 folytonos változóra:  $\mathcal{L}(\theta) = \rho(x | \theta)$ 

A gyakorlatban sokszor a log-likelihood függvényt vagy az átlagolt log-likelihood függvényt használjuk:

$$\ell(\theta|x) = \ln \mathcal{L}(\theta|x)$$
$$\hat{\ell}(\theta|x) = \frac{1}{N} \ln \mathcal{L}(\theta|x)$$

A maximum likelihood módszer lényege, hogy megkeressük azt a  $\theta$  paramétert, ami a likelihood függvényt maximalizálja:

$$\hat{\theta}_{MLE} = \operatorname*{argmax}_{\theta \in \Theta} \mathcal{L}(\theta | x_1, x_2, ..., x_N)$$

Példa: normális eloszlás paraméterei

$$\rho(x|\mu,\sigma^2) = \frac{1}{\sigma\sqrt{2\pi}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right)$$

Ezért egy N elemű minta esetén azt feltételezve, hogy a mintát egy normális eloszlást követő populációból vesszük, az N minta sűrűségfüggvénye:

$$\rho(x_1, x_2, ..., x_N | \mu, \sigma^2) = \prod_{i=1}^N \rho(x_i | \mu, \sigma^2) = \prod_{i=1}^N \frac{1}{\sigma \sqrt{2\pi}} \exp\left(-\frac{(x_i - \mu)^2}{2\sigma^2}\right)$$
$$= \left(\frac{1}{2\pi\sigma^2}\right)^{n/2} \exp\left(-\frac{\sum_{i=1}^N (x_i - \mu)^2}{2\sigma^2}\right) = \mathcal{L}(\mu, \sigma^2)$$

A log-likelihood függvény pedig

$$\ell(\mu, \sigma^2) = \ln \mathcal{L}(\mu, \sigma^2) = -\frac{N}{2} \ln(2\pi\sigma^2) - \frac{1}{2\sigma^2} \sum_{i=1}^{N} (x_i - \mu)^2$$

Ahhoz, hogy megkapjuk a  $\hat{\mu}_{MLE}$  és  $\hat{\sigma}_{MLE}$  becsült paramétereket, az alábbi két egyenletet kell megoldani:

$$\frac{\partial \ell(\mu, \sigma^2)}{\partial \mu} = 0$$
$$\frac{\partial \ell(\mu, \sigma^2)}{\partial \sigma} = 0$$

Ha ezeket megoldjuk, az jön ki, hogy  $\hat{\mu}_{MLE} = \overline{x} = \frac{1}{N} \sum_{i=1}^{N} x_i$  és  $\hat{\sigma}_{MLE}^2 = \frac{1}{N} \sum_{i=1}^{N} (x_i - \mu)^2$ .

#### 4. Extrém statisztikák

#### 5. Post-hoc analízis

# 6. Regresszió

Tegyük fel, hogy a megfigyelt adathalmaz  $\{(x_1, y_1), (x_2, y_2), ..., (x_N, y_N)\}$ . Ekkor az X a független változó (feature variable, predictor, regressor), Y pedig a függő változó (outcome, response variable). Az r(X) regressziós függvény az Y várható értéke X függvényében:

$$r(x) = E(Y|X=x)$$

- Ha r(x) megadható véges számú valós paraméterrel, akkor **parametrikus regresszió**ról beszélünk.
- ullet Ha Y meghatárzása a cél, ismert X esetén, akkor **predikció**ról beszzélünk.
- $\bullet$  HaYdiszkrét (pl. kutya vagy macska látható a képen), akkor **klasszifikáció**ról beszélünk.
- Ha a cél az r(x) görbe meghatározása, akkor **regresszió**ról vagy **görbeillesztés**ről beszélünk.

## 7. Hipotézistesztelés

A hipotézistesztelés lényege, hogy a rendelkezésre álló adatok alapján egy feltevés (hipotézis) igazságtartalmára akarunk kijelentést tenni. Fontos, hogy a feltevés a teljes populációra vonatkozik, tehát egy hipotézis elfogadásának vagy elutasításának mindig van valamennyi bizonytalansága. Példák: 1. Dobunk egy dobókockával 100-szor, majd feltesszük azt a hipotézist, hogy a dobott számok egyenletes eloszlást követnek. 2. Megnézzük 1000 ember jövedelmét majd feltételezzük, hogy a jövedelem olyan lognormális eloszlást követ, amelyre  $\mu=150000 {\rm HUF}$  és  $\sigma=50000 {\rm HUF}$ .

Formális definíció:

- A  $\Theta$  paraméter-teret felosztjuk  $\Theta_0$  és  $\Theta_1$ -re.  $(\Theta_0 \cup \Theta_1 = \Theta, \Theta_0 \cap \Theta_1 = \emptyset)$
- A nullhipotézis  $H_0: \theta \in \Theta_0$
- Az alternatív hipotézis  $H_1: \theta \in \Theta_1$
- $\bullet$  Az összegyűjtött  $X \in \mathcal{X}$ adatok alapján akarjuk a hipotézist eldönteni.
- Az elvetési régió (rejection region) egy  $R \subset \mathcal{X}$  halmaz, amelyre

$$X \in R \Rightarrow H_0$$
-t elvetjük és elfogadjuk  $H_1$ -et  $X \notin R \Rightarrow H_0$ -t elfogadjuk

6

• Egy tesztfüggvény (test statistic) egy  $T: \mathcal{X} \to \mathbb{R}$  függvény, amelyre az elvetési régió így írható:

$$R = \{ x \in \mathcal{X} | T(x) > c \},$$

ahol c a teszt kritikus értéke.

- A hipotézis tesztelés lényege keresni egy olyan T-t és c-t, amivel a legjobb (legkevésbé káros) döntést hozhatjuk.
- Egyoldalú teszt: Ha X a rendelkezésre álló adathalmaz, és  $\theta \in \Theta$  egy paraméter (pl. átlag, variancia, stb), akkor egy egyoldalú teszt a következőt jelenti:

$$H_0: \theta = \theta_0 \text{ és } H_1: \theta < \theta_0 \text{ (vagy } > \theta_0)$$

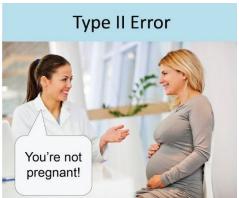
• Kettős hipotézisteszt: Ha X a rendelkezésre álló adathalmaz, és  $\theta \in \Theta$  egy paraméter (pl. átlag, variancia, stb), akkor egy kettős teszt a következőt jelenti:

$$H_0: \theta = \theta_0 \text{ és } H_1: \theta \neq \theta_0$$

#### • Hibák

	$H_0$ igaz	$H_0$ hamis
Elfogadjuk $H_0$ -t	OK	type 2 hiba (false positive)
Elvetjük $H_0$ -t	type 1 hiba (false negative)	OK





1. ábra. A két hibatípus.

• Statisztikai erő: Egy kettős hipotézisteszt statisztikai ereje nem más, mint az a valószínűség, hogy a teszt helyesen veti el a nullhipotézist, amikor az alternatív hipotézis igaz:

statisztikai erő = 
$$\beta = P(H_0 \text{ elutasítva} | H_1 \text{ igaz})$$

• szignifikancia-szint Egy statisztikai teszt  $\alpha$  szignifikancia szintű, ha a nullhipotézis elvetésének a valószínűsége  $\alpha$  feltéve, hogy a nullhipotézis igaz.

$$\alpha = P(H_0$$
-t elvetjük |  $H_0$  igaz)

p-érték: A p-érték azt mutatja meg, hogy ha igaz a nullhipotézis, akkor mekkora valószínűséggel kapunk olyan eredményt, amely legalább annyira extrém, mint a mért eredmény.
 Pl. Ha egy minta átlaga X̄, akkor a p = 5% azt jelenti, hogy ha a populáció átlaga μ, akkor 5% valószínűséggel mérhetek legalább | X̄ - μ| eltérést. A p-értéket így is definiálhatjuk:

$$p = P(\text{ezt az eredményt mérem} \mid H_0 \text{ igaz})$$

### 7.1. Egymintás *u*-próba [2]

(angolul z-test [3])

- Egymintás *u*-próbával a minta alapján a populáció átlagára vonatkozó hipotézist lehet tesztelni.
- $\bullet$  A populációt normális eloszlásúnak feltételezzük, melynek ismerjük a  $\sigma$  szórását.
- $\bullet$  Ha a  $\sigma$  szórás nem ismert, akkor a u-próba helyett t-próbát kell végezni.
- Azt akarjuk megvizsgálni, hogy a minta alapján a  $\mu$  érték tekinthető-e a populáció átlagának.
- $\bullet$  A teszt elvégzéséhez kiszámoljuk a minta  $\overline{X}$  átlagát, majd ebből a próbastatisztikát:

$$u = \frac{\overline{X} - \mu}{\sigma / \sqrt{n}}$$

- Kihasználjuk azt a feltételt, hogy az u próbastatisztika standard normális eloszlást követ, és adott előre meghatározott p-értékre kiszámoljuk a  $u_{p/2} = \Phi^{-1}(1 p/2)$ -t.
- Ha  $|u| \geq u_{p/2}$ , akkor ez azt jelenti, hogy a mért  $\overline{X}$  és a feltételezett  $\mu$  érték között p szignifikancia-szint mellett az eltérés szignifikáns. Tehát elvetjük a nullhipotézist, miszerint a populáció átlaga  $\mu$ .
- Példa: Begyűjtünk 3 tonna almát, amiről tudjuk, hogy az almák tömegének szórása  $\sigma=40g$ . Találomra kiválasztunk 100db almát, amelyeket megmérve azt kapjuk, hogy a tömegek átlaga  $\overline{X}=150g$ . Kijelenthető-e 95%-os biztonsággal, hogy a begyűjtött almák átlagos tömege  $\mu=140g$ ? És 99%-os biztonsággal?

$$-H_0: \mu = 140g$$

$$-H_1: \mu \neq 140g$$

$$-p_1 = 0.05, p_2 = 0.01$$

$$-u_{p_1/2} = \Phi^{-1}(1 - p_1/2) = \Phi^{-1}(1 - 0.05/2) \approx 1.96$$

$$-u_{p_2/2} = \Phi^{-1}(1 - p_2/2) = \Phi^{-1}(1 - 0.01/2) \approx 2.58$$

$$-u = \frac{\overline{X} - \mu}{\sigma/\sqrt{n}} = \frac{150 - 140}{40/10} = 2.5$$

$$- \text{WTF}???}$$

### 7.2. Egymintás t-próba [4]

- $\bullet$  A vizsgált valószínűségi változó normális eloszlást követ, de nem ismerjük a  $\sigma$ szórást.
- A kérdés ugyanaz, mint az előbb, hogy a populáció átlaga megegyezik-e statisztikai szempontból a feltételezett  $\mu$  értékkel.
- ullet t-próba esetén a próbastatisztika a következő lesz:

$$t = \frac{\overline{X} - \mu}{s/n} \,,$$

ahol  $s=\sqrt{\frac{1}{n-1}\sum_{i=1}^n(X_i-\overline{X})^2}$  a minta szórása. A próbastatisztika ebben az esetben Student t-eloszlást követ, f=n-1 szabadsági fokkal.

- $\bullet\,$  Az előbbihez hasonlóan az adott pszignifikancia szinthez megkeressük a  $t_p$  értéket.
- Ha  $|t| \ge t_p$ , akkor a nullhipotézist elvetjük, mert a minta átlaga szignifikánsan eltér a  $\mu$  értéktől.
- Ellenkező esetben, vagyis ha  $|t| < t_p$ , a nullhipotézist megtartjuk.

# 8. Függetlenségvizsgálat, $\chi^2$ -próba

# Hivatkozások

- [1] Bessel's correction https://en.wikipedia.org/wiki/Bessel%27s\_correction.
- [2] Egymintás u-próba https://hu.wikipedia.org/wiki/Egymintás\_u-próba.
- [3] Z-test https://en.wikipedia.org/wiki/Z-test.
- [4] Egymintás t-próba https://hu.wikipedia.org/wiki/Egymintás\_t-próba.