

2. tétel

Nagy Dániel

2019. június 8.

Kivonat

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

1. Bevezetés

1.1. Valószínűesszámitás alapfogalmak

- **Eseménytér** (ez egy absztrakt fogalom): $\Omega = \{\omega_1, \omega_2, \dots, \omega_n\}$ pl. kockadobás esetén $\Omega = \{\omega_1 = \text{"1est dobok"}, \omega_2 = \text{"2est dobok"}, \omega_3 = \text{"párosat dobok"} \dots\}$
- **Valószínűségi változó**: $X : \Omega \rightarrow \mathbb{R}$ pl. kockadobás esetén $X(\omega_1) = 1, X(\omega_2) = 2, \dots$
- **Valószínűség**: P egy mérték, amely Ω részhalmazaihoz számot rendel:

- $P : \mathcal{P}(\Omega) \rightarrow \mathbb{R}$
- $P(\Omega) = 1$ és $P(\emptyset) = 0$
- $0 \leq P(A) \leq 1 \ \forall A \in \Omega$
- Ha A_1, A_2, \dots diszjunkt részhalmazai Ω -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 $\iff P(A \cap B) = P(A)P(B)$
- $P(A|B) = \frac{P(A \cap B)}{P(B)}$
- Teljes valószínűség: Ha A_1, A_2, \dots az Ω egy felosztása, akkor

$$P(B) = \sum_k P(B|A_k)P(A_k)$$

- Bayes-tétel: Ha A_1, A_2, \dots az Ω 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

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

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

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

Megjegyzés: sűrűségfüggvénye csak folytonos eloszlású valószínűségi változónak van.

- **Várható érték**

$$\text{folytonos eset } E(X) = \langle X \rangle = \int_{-\infty}^{\infty} x \rho(x) dx$$

$$\text{diszkrét eset } E(X) = \langle X \rangle = \sum_k x_k p_k = \sum_k x_k P(X = x_k)$$

- **Várható értékre vonatkozó azonosságok:**

$$- \text{ Ha } Y = g(X) \Rightarrow E(Y) = E(g(X)) = \int_{-\infty}^{\infty} g(x) \rho(x) dx$$

$$- E\left(\sum_k a_k X_k\right) = \sum_k a_k E(X_k)$$

$$- \text{ Ha } X_1, X_2, \dots \text{ független változók, akkor } E\left(\prod_k X_k\right) = \prod_k E(X_k)$$

- **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, \dots 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:** $\bar{X} = \frac{1}{N} \sum_{i=1}^N X_i$ (p_k -t a relatív gyakorisággal közelítjük)

- **Minta varianciája:** $s^2 = \frac{1}{N-1} \sum_{i=1}^N (X_i - \bar{X})^2$, standard hibája $SE = \sqrt{s^2}$. A nevezőben az $N-1$ faktor az ún. Bessel-korrekciónak [1].
- **Megjegyzés:** Ha egy teljes populáció esetén $E(X) = \mu$ és $V(X) = \sigma^2$, attól még általában $\bar{X} \neq \mu$ illetve $s^2 \neq \sigma^2$.
- Egy minta esetében \bar{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

$$\begin{aligned} E(\bar{X}) &= \mu \\ V(\bar{X}) &= \sigma^2/N \\ E(s^2) &= \sigma^2 \end{aligned}$$

1.2. Statisztikai következtetés (inference)

- Az alapprobléma: van egy adathalmaz, ami tartalmazza a méréseket. Ezek $X_1, X_2, \dots, X_N \sim F$ 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 ρ sűrűségfüggvényt határozzuk meg.
- 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 Θ a lehetséges paraméterek halmaza. Pl. ha normális eloszlást feltételezünk, akkor a parametrikus modell

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

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

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 θ paramétereknek:

$$\mathcal{L}(\theta|x) = P(\theta|X = x)$$

A gyakorlatban sokszor a log-likelihood függvényt használjuk:

$$\ell(\theta|x) = \ln \mathcal{L}(\theta|x)$$

A maximum likelihood módszer lényege, hogy megkeressük azt a θ paramétert, ami a likelihood függvényt maximalizálja:

$$\hat{\theta}_{MLE} = \operatorname{argmax}_{\theta \in \Theta} \mathcal{L}(\theta|x)$$

4. Extrém statisztikák

5. Post-hoc analízis

6. Regresszió

7. Hipotézis tesztelés

8. z-teszt, t-test

8.1. Konfidenciaintervallumok

8.2. Függetlenségvizsgálat, χ^2 -próba

Hivatkozások

[1] Bessel's correction https://en.wikipedia.org/wiki/Bessel%27s_correction.