

POROČILO SEMINARSKE NALOGE

STATISTIKA

SARA BIZJAK

UNIVERZA V LJUBLJANI
FAKULTETA ZA MATEMATIKO IN FIZIKO
ODDELEK ZA MATEMATIKO

JULIJ 2020

1. NALOGA

Podatki so vzeti iz datoteke `Kibergard`, kjer se nahajajo informacije o 43.886 družinah, ki stanujejo v mestu Kibergard. Za vsako družino so zabeleženi naslednji podatki:

- Tip družine (od 1 do 3)
- Število članov družine
- Število otrok v družini
- Skupni dohodek družine
- Mestna četrt, v kateri stanuje družina (od 1 do 4)
- Stopnja izobrazbe in vodje gospodinjstva (od 31 do 46: opisi v datoteki `z navodili`)

Nalogo sem reševala s pomočjo programa R. Koda, uporabljena za generiranje enostavnih slučajnih vzorcev in izračune, je dostopna v priloženi datoteki *naloga1.R*.

PRIMER A

Vzamemo enostavni slučajni vzorec 200 družin in na njegovi podlagi ocenimo delež družin v Kibergardu, v katerih vodja gospodinjstva nima srednješolske izobrazbe (niti poklicne niti splošne mature). Opisan delež znaša $p = 0.195$.

PRIMER B

Ocenimo standardno napako in postavimo 95% interval zaupanja. Standardno napako za delež izračunamo po formuli

$$\hat{se}(p) = \sqrt{\frac{p \cdot (1 - p)}{n - 1} \cdot \left(1 - \frac{n}{N}\right)},$$

kjer so $p = 0.195$, $n = 200$, $N = 43.886$.

Dobimo rezultat $\hat{se}(p) = 0.02802185$.

Interval zaupanja je enak: $[0.1400782, 0.2499218]$.

PRIMER C

Vzorčni delež in ocenjeno standardno napako primerjamo s populacijskim deležem in pravo standardno napako.

- Vzorčni delež: 0.195
- Populacijski delež: 0.2115025
- Razlika obeh deležev: 0.01650253
- Ocenjena standardna napaka (iz vzorca): 0,02802185
- Prava standardna napaka (iz celotne populacije): 0.02888282
- Razlika med ocenjeno in pravo standardno napako: 0.0008609634

Ker velja $0.2115025 \in [0.1400782, 0.2499218]$, interval zaupanja pokrije populacijski delež.

PRIMER D

GRAF

INTERVALI ZAUPANJA + KOLIKO JIH POKRIJE POPULACIJSKI DELEŽ

PRIMER E

Standardni odklon vzorčnih deležev za 100 prej dobljenih vzorcev je enak 0.02881085. Prava standardna napaka za vzorec velikosti 200 pa je 0.02888282. Razlikujeta se za $7.196778 \cdot 10^{-5}$.

PRIMER F

GRAF

INTERVALI ZAUPANJA + KOLIKO JIH POKRIJE POPULACIJSKI DELEŽ

2. NALOGA

3. NALOGA

Opazimo n neodvisnih realizacij zvezne porazdelitve z gostoto:

$$f(x \mid \sigma) = \begin{cases} \frac{\Gamma(2\alpha)}{(\Gamma(\alpha))^2} [x(1-x)]^{\alpha-1} & ; 0 < x < 1 \\ 0 & ; \text{sicer,} \end{cases}$$

kjer je $\alpha > 0$ neznan parameter. Če je X slučajna spremenljivka s to gostoto, se da izračunati:

$$E(X) = \frac{1}{2}, \quad \text{var}(X) = \frac{1}{4(2\alpha + 1)}.$$

PRIMER A

Določimo obliko porazdelitve v odvisnosti od α .

PRIMER B

Ocenimo α po metodi momentov.

PRIMER C

Poiščemo enačbo, ki določa cenilko po metodi največjega verjetja. Pogledamo, kdaj ta cenilka obstaja.

PRIMER D

Poiščemo asimptotično varianco cenilke po metodi največjega verjetja.

4. NALOGA

5. NALOGA

X in Y sta slučajni spremenljivki, za kateri velja:

- $E(X) = \mu_x$,
- $E(Y) = \mu_y$,
- $\text{var}(X) = \sigma_x^2$,
- $\text{var}(Y) = \sigma_y^2$,
- $\text{cov}(X, Y) = \sigma_{x,y}$.

Opazimo X in želimo napovedati Y .

PRIMER A

Poiščemo napoved, ki je oblike $\hat{Y} = \alpha + \beta \cdot X$, kjer α in β izberemo tako, da je srednja kvadratična napaka $E[(Y - \hat{Y})^2]$ minimalna.

Uporabimo namig

$$E[(Y - \hat{Y})^2] = [E(Y) - E(\hat{Y})]^2 + \text{var}(Y - \hat{Y}). \quad (1)$$

Ker sta oba člena desne strani enačbe večja od 0, je dovolj, da poiščemo vrednosti α in β , ki minimizirata ta dva člena – potem bo najmanjša možna tudi njuna vsota.

Poglejmo si najprej prvi člen enačbe. Vrednost enačbe

$$\left[E(Y) - E(\hat{Y}) \right]^2 = [\mu_y - \alpha - \beta\mu_x]^2$$

bo najmanjša, ko bo $\mu_y - \alpha - \beta \cdot \mu_x = 0$. To pa bo res natanko takrat, ko bo $\alpha = \mu_y - \beta\mu_x$.

Poglejmo si še drugi člen enačbe. Vidimo, da je

$$\begin{aligned} \text{var}(Y - \hat{Y}) &= \text{var}(Y - \alpha - \beta X) = \text{var}(Y - \beta X) = \text{var}(Y) - 2\beta \text{cov}(X, Y) + \beta^2 \text{var}(X) \\ &= \sigma_y^2 - 2\beta \sigma_{x,y} + \beta^2 \sigma_x^2 \end{aligned}$$

funkcija spremenljivke β . Za izračun minimuma enačbo najprej odvajamo po β in enačimo z 0.

$$\frac{\partial}{\partial \beta} (\text{var}(Y - \hat{Y})) = -2\sigma_{x,y} + 2\beta\sigma_x^2 = 0$$

To bo res, ko bo $\beta = \frac{\sigma_{x,y}}{\sigma_x^2}$.

Vrednosti α in β , pri katerih bo izraz $E \left[(Y - \hat{Y})^2 \right]$ minimalen, sta

$$\alpha = \mu_y - \mu_x \frac{\sigma_{x,y}}{\sigma_x^2} \quad \text{in} \quad \beta = \frac{\sigma_{x,y}}{\sigma_x^2}$$

PRIMER B

Pokažemo, da se pri tako izbranih koeficientih *determinacijski koeficient* (kvadrat korelacijskega koeficienta) izraža v obliki

$$r_{x,y}^2 = 1 - \frac{\text{var}(Y - \hat{Y})}{\text{var}(Y)}.$$

Spomnimo se najprej formule za *korelacijski koeficient*:

$$\text{corr}(X, Y) = \frac{\text{cov}(X, Y)}{\sqrt{\text{var}(X)\text{var}(Y)}} = r_{x,y}.$$

Ker je determinacijski koeficient enak kvadratu korelacije, je enak

$$r_{x,y}^2 = \frac{\text{cov}(X, Y)^2}{\text{var}(X)\text{var}(Y)}.$$

Dobimo enakost

$$1 - \frac{\text{var}(Y - \hat{Y})}{\text{var}(Y)} = \frac{\text{cov}(X, Y)^2}{\text{var}(X)\text{var}(Y)}$$

in če upoštevamo še vrednosti α in β iz prvega dela naloge:

$$1 - \frac{\text{var}(Y - \hat{Y})}{\text{var}(Y)} = \frac{\text{var}(Y) - \text{var}(Y - \hat{Y})}{\text{var}(Y)} = \frac{\sigma_y^2 - \sigma_y^2 + \frac{\sigma_{x,y}^2}{\sigma_x^2}}{\sigma_y^2} = \frac{\sigma_{x,y}^2}{\sigma_x^2 \sigma_y^2} = \frac{\text{cov}(X, Y)^2}{\text{var}(X) \text{var}(Y)} = r_{x,y}^2.$$

Pokazali smo, da je determinacijski koeficient res enak $1 - \frac{\text{var}(Y - \hat{Y})}{\text{var}(Y)}$.