

# Pričakovana vrednost, varianca

Nataša Kejžar

## Povzetek

Spoznali boste kako:

- generiramo vzorce iz končnih in neskončnih populacij
- računamo kvantile v populaciji in na vzorcu
- narišemo gostoto, verjetnostno funkcijo, porazdelitveno funkcijo
- s simulacijami preverjamo pristranskost cenilke: ugotovili smo, da je povprečna ocena natančnejša, če naredimo večje število simulacij

Nekaj osnovnih verjetnostnih izrazov, ki se v statistiki pogosto uporabljajo

## Pričakovana vrednost

- izračun po definiciji

$$E(X) = \sum_{i=1}^{\infty} x_i p(x_i)$$
$$E(X) = \int_{-\infty}^{\infty} x f(x) dx$$

- lastnosti

$$E(aX) = a \cdot E(X)$$
$$E(X + Y) = E(X) + E(Y)$$

- Jensenova neenakost: za konkavno funkcijo  $f$  velja (za konveksno pa ravno obratno)
- <https://www.statlect.com/fundamentals-of-probability/Jensen-inequality>

$$f(E(X)) \geq E(f(X))$$

## Varianca

- po definiciji

$$\text{var}(X) = E[(X - E(X))^2]$$

- lastnosti

$$\begin{aligned}
var(X) &= E(X^2) - E(X)^2 \\
var(aX) &= a^2 \cdot var(X) \\
var(X + Y) &= var(X) + var(Y) + 2 \cdot cov(X, Y) \\
var(X) &= E_Y(var(X|Y)) + var_Y(E(X|Y))
\end{aligned}$$

- izračun variance za diskretne spremenljivke (za izračun  $E(X)$  glej zgoraj)

$$var(X) = \sum_{i=1}^{\infty} x_i^2 p(x_i) - E(X)^2$$

## Kovarianca

- po definiciji

$$cov(X, Y) = E[(X - E(X))(Y - E(Y))]$$

- lastnosti

$$cov(X, Y + Z) = cov(Y + Z, X) = cov(X, Y) + cov(X, Z)$$

## Naloge

- Pomoč za porazdelitve v R: `help(Distributions)`. Za vzorec velikosti 1000 iz porazdelitve  $N(120, 30^2)$  narišite:
  - verjetnostno porazdelitev (`hist`, parameter `probability`, `breaks`, ali s knjižnico `ggplot2`)
    - Kaj je gostota porazdelitve? Kaj dobimo, če narišemo rezultate funkcije `density`?
    - Na histogram dodajte še pričakovano vrednost za to spremenljivko (`abline`, parameter `v`; ali `geom_vline`)
    - Na histogram dodajte še gostoto normalne porazdelitve (uporabite funkciji `dnorm` in `curve`, parameter `add=TRUE`; ali `stat_function`)
  - porazdelitveno funkcijo (`ecdf`) Ponovite isto za vzorec velikosti 10. Komentirajte opažanja.
- Povprečna porodna teža novorojenčka v Sloveniji je 3.300 gramov, večina donošenih novorojenčkov ob rojstvu tehta med 2.500 in 4.100 grami.
  - Če velja, da je teža novorojenčkov v populaciji normalno porazdeljena in da sta spodnja in zgornja meja 5. in 95. percentil, z R izračunajte, kolikšen je standardni odklon teže novorojenčkov.
  - Izračunajte interkvartilni razmik za težo.
  - V katerem percentilu se nahaja novorojenček, ki je težak 2900 g?
  - Iz teoretične porazdelitve izberite vzorec 5, 50, 500 enot in na vzorcu izračunajte percentil za novorojenčka s 2900 grami. Primerjajte in komentirajte rezultata.
  - Za izbran vzorec s 500 enotami izračunajte 90% **interval zaupanja za povprečje**. Primerjajte ga z intervalom, ki so ga poročali na začetku naloge. Komentirajte.
- Teoretične nalogice:
  - pokažite s pomočjo definicije za varianco, da drži  $var(X) = E(X^2) - E(X)^2$
  - pokažite s pomočjo vsote varianc, da drži  $var(2X) = 4 \cdot var(X)$
  - pokažite, da varianco končne populacije (velikosti  $N$ ) lahko zapišemo tudi kot

$$var(X) = \frac{1}{2N^2} \sum_{i=1}^N \sum_{j=1}^N (x_i - x_j)^2.$$

- Verjetnostna porazdelitev končne populacije, ki vas zanima, je predstavljena v spodnji tabeli pod nalogo. V njej je 1500 enot.
  - Izračunajte pričakovano vrednost populacije v R.
    - po definiciji
    - tako da zapišete populacijo v vektor (funkcija `rep`)
  - Izračunajte varianco populacije v R. Zakaj izračun s funkcijo `var` da drugačen rezultat?
  - Iz populacije bi radi izbrali vzorec velikosti 15 s ponavljanjem. Zapišite kodo v R. (`sample`)
  - Iz populacije bi radi izbrali vzorec velikosti 15 brez ponavljanja. Zapišite kodo v R.

X	1	2	3	4	5
p(X)	1/15	1/5	4/15	2/5	1/15

- Naj bo spremenljivka  $X$  porazdeljena po  $Bernoulli(\pi = 0.85)$ .
  - Generirajte 100 opazovanj. (uporabite npr. `runif` ali `sample`)
  - Izračunajte po definiciji  $E(X)$  in  $var(X)$ .
  - Vemo, da je vsota neodvisnih, enako porazdeljenih Bernoullijevih spremenljivk porazdeljena po

- binomski porazdelitvi ( $Y = \sum_{i=1}^n X_i \sim \text{Bin}(n, \pi)$ ). Pokažite to s simulacijo velikega vzorca.
- d. Izpeljite, da velja  $\text{var}(Y) = n \cdot \pi(1 - \pi)$ .
6. Radi bi primerjali povprečni teži dveh velikih skupin morskih prašičkov. Prva skupina so tisti, ki so hišni ljubljenci v Evropi, druga tisti iz ZDA. V ta namen zberemo enako velika vzorca.
- Kateri statistični test bi uporabili? Kaj lahko poveste o kovarianci teh dveh skupin?
  - Simulirajte populacijsko kovarianco v R (funkcija `cov`) za naslednji spremenljivki (predpostavite lahko, da sta teži v populaciji pribl. normalno porazdeljeni):
    - $X_{EU}$  in  $X_{ZDA}$  (koliko je populacijska kovarianca?)
    - $X_{EU}$  in  $X_{EU} - X_{ZDA}$  (kaj pričakujete?)
  - Izračunajte kovarianco teoretično.
7. Imate normalno porazdeljeno spremenljivko ( $X \sim N(0, 1)$ ). Grafično preverite rezultat iz Rice-a (6. poglavje), da je  $X^2 \sim \chi_{df=1}^2$ .
- izberite dovolj velik vzorec iz  $N(0, 1)$ .
  - narišite *gostoto* kvadriranih vrednosti
  - dodajte gostoto porazdelitve  $\chi_1^2$
8. Imate dve neodvisni spremenljivki  $X_1$  in  $X_2$ , ki sta porazdeljeni po  $\chi^2$  porazdelitvi z  $df_1$  in  $df_2$  stopinjami prostosti.
- Grafično preverite, da je  $X_1 + X_2 = X_3$ , kjer je  $X_3 \sim \chi_{df_1+df_2}^2$ .
  - Naj bosta sedaj  $X_1 \sim \chi_{df_1}^2$  in  $X_3 \sim \chi_{df_1+df_2}^2$  neodvisni spremenljivki. Ali mislite, da velja tudi  $X_3 - X_1 = X_2$ , kjer je  $X_2 \sim \chi_{df_2}^2$ ? Komentirajte zakaj.
    - Narišite porazdelitev spremenljivke  $X_3 - X_1$ .
  - Kako mislite, da je porazdeljena razlika dveh neodvisnih standardno normalno porazdeljenih spremenljivk? Zakaj (oblika, variabilnost)? Preverite grafično.
9. Generirajte vrednosti za  $N = 10$  enot veliko populacijo s kodo pod nalogo. Vse naslednje primere v R zakodirajte tako, da boste lahko kadarkoli spremenili vrednosti enot in tudi število enot.
- Izračunajte povprečje po formuli  $\mu = 1/N \sum_{i=1}^N x_i$ .
  - Posebej seštejte pozitivne in negativne odmike od povprečja.
  - Izračunajte varianco po obeh formulah, ki ste jih omenili na predavanjih.
  - Izmislite si poljubno vrednost  $a$  in izračunajte vsoto kvadriranih odmkov (VKO) od te vrednosti. Za izračun te količine naredite funkcijo.
  - Za različne vrednosti  $a$ jev (npr. med  $[4, 10]$ ) izračunajte vsoto kvadriranih odmkov in dvojice narišite na graf. Z rdečo označite vrednost  $N\sigma^2$  (pri kateri vrednosti  $a$  se pojavi?).

```
set.seed(1)
x = sample(x = 1:10, size=10, replace=TRUE)
```

10. Odgovorite:
- S kakšnim namenom za izračun variance kvadriramo odmike od povprečja?
  - Kaj se dogaja z VKO pri različnih vrednostih  $a$ ? Pri katerem  $a$  je v splošnem VKO najmanjši? Zakaj?
  - Ponovite prejšnjo vajo z  $N = 20$  enot, ki lahko zavzamejo vrednosti 0 ali 1. Kakšne vrednosti  $a$  boste vzeli za risanje grafa in zakaj?
11. Izračunajte delež vrednosti nad 132 za vzorec  $n = 1000$ , ki ga dobite iz porazdelitve  $N(120, 30^2)$ . Primerjajte dobljeni delež z verjetnostjo, ki jo izračunamo iz gostote (`pnorm`).
- Primerjajte rezultate, ki jih dobite iz vzorca iz 10 enot in vzorca iz 1000 enot. Kje so bistvene razlike in zakaj pride do njih?

12. Zapišite intuitivno cenilko za populacijsko povprečje.
- Pokažite teoretično, da je cenilka nepristranska.
  - Pokažite to s simulacijami na velikosti vzorca 10 iz  $N(120, 30^2)$ :
    - generirajte 1000 vzorcev (uporabite `for` zanko ali **replicate**); za vsak vzorec izračunajte oceno za povprečje in te ocene shranite v nov vektor
    - če je cenilka nepristranska, morajo ocene variirati okrog populacijskega povprečja; preverite
  - Narišite graf
  - Izračunajte *pričakovano vrednost*