

II Uvod u statistiku

Statistika se bavi

(i) prikupljanjem podataka ("dizajn pokusa")

(ii) prikupljenim podacima

("deskriptivna statistika")

(iii) zaključivanjem na osnovu podataka ("statističko zaključivanje")

(1) Populacija i uzorak

Populacija - skupina jedinici (ili stvari)
o kojima želimo nešto znati

↳ za danu populaciju, tipično nas
zanimaju jedna ili više varijabli, npr.

- skup svih punoljetnih ljudi u RH,
a varijabla je stranaka koju osoba
podržava

- sve batenje za automobile proizvedene u nekoj tvornici, a vanjake je trajnije batenje.
- svi studenti na PAF-KO, a vanjake su visine i težina.

Općenito, vanjake mogu biti

(i) Kategorijalne (ili kvalitativne)

np. spol (M/Ž), stranka,

pušač (DA/NE)

(ne možemo ih prizvano predstaviti na nekoj numeričkoj skali)

iii) Numeričke (ili kvantitativne)

(unjednat ili različito broj)

neprekidne

- trajanje batanje
- mjesečno plaće
- visina, težina

diskretne

- broj studenata
- broj položenih kolegija u semestru

Uzorek - grupa jedinici iz populacije
na kojima "mjerimo" promatrani
varijablu (ukupno n EM)

Uzorek (mjerjenja) - oprežena mjerenja,
ostataka će biti

"mali" x -ovi



$x_1, \dots, x_n \in \mathbb{R}$ (ili \mathbb{R}^d)

(x_i - mjerenje za i -tu osobu)

Intuicijom, cilj stat. zaključivanja je
 donijeti zaključke o karakteristikama
 promatrane varijable u celej
 populaciji samo na temelju
 uzerke x_1, \dots, x_n .

↳ Šećejnost: Ako je

X_i = vrijednost varijable za šećejnu
 određenu osobu u populaciji,

te $F_i = F_X$, = pretpostavljamo

da su podaci x_1, \dots, x_n konkretne

realizacije mjd šećejnih varijabli

x_1, \dots, x_n t.d. $F_{X_i} = F$, $\forall i$.

Každemo do je X_1, \dots, X_n slučajno
uzorak za X (ili iz F).

Formalno, caji \rightarrow test, zaključivanje je
iz x_1, \dots, x_n zaključiti nešto o F .

Načel Definicija slučajnog uzorka odgovara
ideji da želimo x_1, \dots, x_n koji su

"reprezentativni" za cele populaciju.

↳ Loši primjeri

- želimo nas % ljudi a em koji najviše
za Dinamo, a uzorak uzememo u Splitu.
- želimo nas prosječan broj otkaza iznajmljivača
u Beogradu, a uzorak uzememo teksto da
masovno ispituje mlade na oglasnom
tržištu. \rightarrow Problem?

2) Deskriptivna statistika

→ primjeri u R-u

Histogramom

Podatke x_1, \dots, x_n grupiramo u k razreda

$$I_1 = [a_0, a_1), I_2 = [a_1, a_2), \dots, I_k = [a_{k-1}, a_k)$$

$$\begin{array}{|l} \downarrow \\ \hline f_1 = \# x_i\text{-ova koji} \\ \text{se nalaze u } I_1 \\ \hline v_1 = \frac{\%}{10} \text{ --- } \\ \hline = f_1/n \end{array}$$

$$\downarrow$$
$$f_2$$

$$v_2$$

$$\downarrow$$
$$f_k$$

$$v_k$$

→ iznad razreda I_j crtkom = produktiv

visine

$$\frac{v_j}{a_j - a_{j-1}}$$

vidi nap.
prije Pr. 6.2



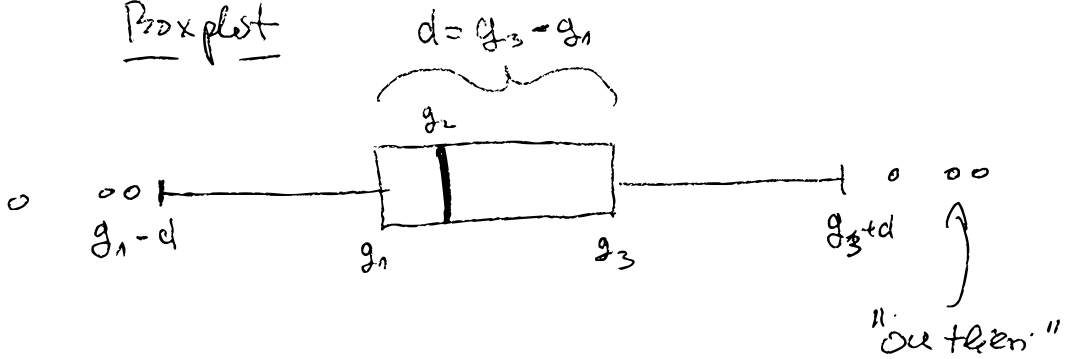
ukupna površina

jednaka 1

, te je histogram

zapravo procjenitelj za gustocu od X .

Boxplot



g_1, g_2, g_3 t.i.d.

- 50% podatku veče, a 50% manje od g_2
(tar. median uvek)
- 75% —||— , a 25% —||— g_1
- 25% —||— , a 75% —||— g_3

⇒ npr. 50% podatku se nalazi
u $[g_1, g_3]$

3. Testiranje statističkih hipoteza

[spada pod stat. zaključivanje]

Prisjetimo se, da slučajnu varijablu X
▷ neparametrom tipom distribucije F ,
iz uzorka X_1, \dots, X_n želimo nešto
zaključiti o F .

Tipično pretpostavljamo da F ovisi
o neparametrom parametru θ pa
se problem sveli na ujedruvu procenu.

$$\text{m.p.} \cdot X \sim \begin{pmatrix} 0 & 1 \\ \mu & \sigma^2 \end{pmatrix}, \theta = \mu$$

$$\cdot X \sim N(\mu, \sigma^2), \theta = (\mu, \sigma^2)$$

Elementi statističkog testa

(i) nul-hipoteza (H_0) i alternativna hipoteza (H_1)

Želimo utvrditi imamo li na temelju
podataka dovoljno dokaza da bismo

odbačili H_0 i prihvatili H_1 .

mpr. Standardni lijek pomaže u 60% slučajeva, a
želimo provjeriti je li novi lijek efikasniji.

→ model je $X \sim \binom{0}{1} \binom{p}{p}$ ← poboljšanje
od novog
lijeka

$H_0: p = 0.6$ (tj. nema poboljšanja)

$H_1: p > 0.6$ (ima poboljšanja)

↳ tipično, H_1 predočava ovisno o 50 želimu

= provjeriti.

(ii) Dvije vrste pogreške

odbacujemo H_0 ,
a ona je istinita
("pogreška prve vrste")

ne odbacujemo H_0 ,
a ona nije istinita
("pogreška druge vrste")

Za mali α (tipično $\alpha = 0.05, 0.01$), zahtijevamo
da vrijedi:

$P(\text{odbacujemo } H_0 \mid H_0) = \alpha$
vjerovatnost pogreške 1. vrste

→ tzv.
nesigurna
znakejnost
testa

[tipična je greška 1. vrste "skuplje", npr. odobrenjem manje cijene iako nije bogat od strane; osuđivanjem nevinu osobu.]

(iii) Testna statistika: $T = f(X_1, \dots, X_n)$

Kritična područja: skup $C_\alpha \subseteq \mathbb{R}$ t.d.

$$P(T \in C_\alpha \mid H_0) = \alpha$$

tada ćemo odbacivati H_0 [= vjerojatnost greške prve vrste]

F -je f i skup C_α biramo tako da znamo odrediti razdjelak od T pod uvjetom da vrijedi H_0 , te tako da je

$$P(T \in C_\alpha \mid H_0)$$

što veće.

u našem primjeru, tipična je utvrda

$$T = \frac{\bar{X}_n - 0.6}{\sqrt{0.04}}$$

gdje je

$$\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i$$

($X_i = 1$ i-ta osoba doživjela pobačaj)

Ako injekci H_0 , tj. X_1, \dots, X_n su njeki i.i.d.

$$X_i \sim \begin{pmatrix} 0 & 1 \\ 0.4 & 0.6 \end{pmatrix},$$

pa CGT-a znamo da je za velike n ,

$$(a) \quad T = \frac{S_n - n \cdot 0.6}{\sqrt{n \cdot 0.6 \cdot 0.4}} \stackrel{d}{\approx} N(0, 1) \quad (S_n = \sum_{i=1}^n X_i)$$

Uo čemu,

\bar{X}_n = % ljudi s pokazivačem nakon korištenja novog lijeka.

Ako je lijek zaista bolji, tj. ako $p > 0.6$, očekujemo da $\bar{X}_n > 0.6$, tj. $T > 0$.

$$C_d = [z_d, +\infty) \quad , \quad \text{gdje je}$$

$$z_d \in \mathbb{R} \text{ t.d. } \boxed{P(Z \geq z) = \alpha \text{ za } Z \sim N(0, 1).}$$

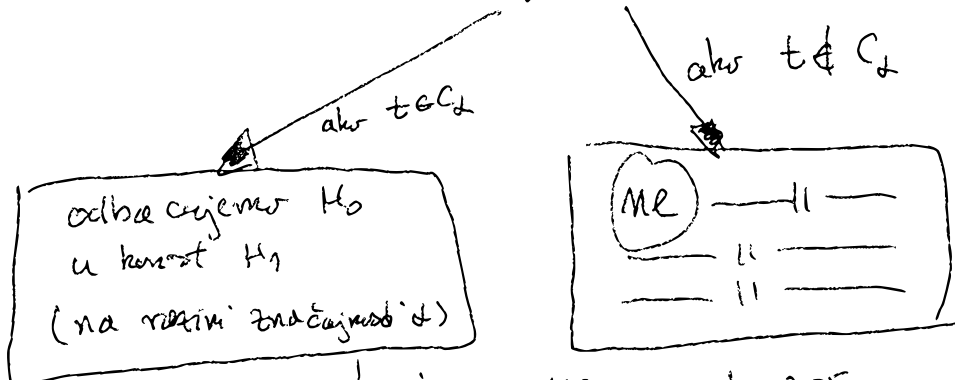
Uo čemu, zbog (a) zaista imamo

npr.
 $z_{0.05} = 1.65$

$$P(T \in C_d \mid H_0) \approx P(Z \geq z) = \alpha.$$

Nakon što smo odredili sve elemente, test dalje provodimo na sljedeći način:

Iz podataka x_1, \dots, x_n odredimo najprije testnu statistiku $t = f(x_1, \dots, x_n)$



u našem primjeru, ako je $n = 1000$ i $\alpha = 0.05$

($z_{0.05} = 1.65$),

• u slučaju $\bar{x}_n = 0.65$, imati brojnu

$$t = \frac{0.65 - 0.6}{\sqrt{0.6 \cdot 0.4}} \approx 3.23 > 1.65$$

\Rightarrow odbacili bi H_0 ("nema promjene")
na raz. znač. 0.05.

• u slučaju $\bar{x}_n = 0.61$,

$$t \approx 0.64 \leq 1.65$$

\Rightarrow na ovoj razin. znač. ne bi
odbacili H_0

(tj., čak i ako $\bar{x}_n > 0.6$, razlika nije velika te je moguće biti samo područje u odabranu uzorku)

Općenito, ako je $X \sim \begin{pmatrix} 0 \\ \sigma \\ \rho \end{pmatrix}$ za $\rho \in (0,1)$,

X_1, \dots, X_n uzorak za X , a želimo testirati

$H_0: \rho = \rho_0$ — nekada fiksnu vrijednost

u odnosu na

(a) $H_1: \rho > \rho_0$, ili

(b) $H_1: \rho < \rho_0$, ili

(c) $H_1: \rho \neq \rho_0$,

konstatiramo

$$T := \sqrt{n} \frac{\bar{X}_n - \rho_0}{\sqrt{\rho_0(1-\rho_0)}} \stackrel{d}{\approx} N(0,1),$$

uz H_0

te

(a) $C_\alpha = [z_{1-\alpha/2}, +\infty)$, (b) $C_\alpha = (-\infty, -z_{1-\alpha/2}]$,

(c) $C_\alpha = (-\infty, -z_{\alpha/2}] \cup [z_{\alpha/2}, +\infty)$

$z_{1-\alpha/2} = z_{\alpha/2}$

↳ u sve tri slučajeva imamo $P(T \in C_\alpha | H_0) \approx \alpha$.

Dop.] H_1 određivanog prije nego što smo uslijeli početke \rightarrow vidi zad. 8. os. i Dop. 9. os. ne uježbama

Dop.] Ako zaista ujedini to te preko puta parametara test, u prosjeku samo u jednako

100.0% slučajeva odbaciti H_0

(iako je točno) \rightarrow

[intuicija o razini značajnosti α]

[ostatak ne uježbama!]