

Sadržaj poglavlja

- Točkaste procjene parametara
- Kriterij najveće izglednosti

Novčić je bačen stotinu puta, pri čemu se pismo pojavilo 40 puta. Je li način bacanja bio korektan?

Ako se pri 100 bacanja kocke šestica pojavi 10 puta, je li ta kocka ispravna?

Pri ponavljanju identičnog pokusa slučajna varijabla poprimila je vrijednosti 6.1, 6.3, 6.4, 6.8, 7.2. Ako je njezina razdioba normalna, a parametri te razdiobe nepoznati, koji ćemo broj uzeti za očekivanje a koji za disperziju te varijable?

Slučajna varijabla poprima vrijednosti unutar intervala [0, 1]. Bilježenje rezultata dalo je podatke: 0.11, 0.14, 0.28, 0.44, 0.48, 0.68, 0.76. O razdiobi ove varijable nemamo nikakvu informaciju. S kojom sigurnošću možemo tvrditi da je riječ o jednolikoj razdiobi?

Na ova, i slična pitanja, odgovor daje *matematička statistika*.

10.1. Točkaste procjene parametara

10.1.1. Uvod

Predmet svakog statističkog proučavanja jest neki (masovni) skup, kojeg nazivamo **populacija** ili **generalni skup**. Populaciju mogu činiti na primjer

- stanovnici države, općine, mjesta;
- privredni potencijali države, regije, grada;
- proizvodnja neke tvornice u jednom danu, mjesecu ili godini

i slično. Podatak koji proučavamo u danoj populaciji nazivamo **obilježje**. Kod iste populacije možemo promatrati više obilježja. Npr. ako je u pitanju stanovništvo, možemo se zanimati za, recimo

- promjenu brojčanog stanja stanovništva tijekom godina;
- broj stanovnika prema starosnoj dobi;
- zaposlenost po vrstama zanimanja;
- školsku spremu

i za stotinjak drugih podataka. Promatramo li proizvodnju, obilježja mogu biti

- broj (količina) proizvedenih dobara u nekom vremenu;
- proizvodnja po vrstama proizvoda;
- profit;
- broj (postotak) škartnih proizvoda u ukupnoj proizvodnji.

Statistički se mogu pratiti i mnoge druge pojave. Tako na primjer, analiziraju se

- metereološke prilike,
- učestalost i vrsta bolesti,
- ispitivanje kupovne moći, tržišta i slično.

U modelu matematičke statistike, populacija čini skup Ω . Obilježje je opisano vrijednošću slučajne varijable X . Osnovni problem matematičke statistike je u određivanju razdiobe varijable X , ili pak nekih njezinih numeričkih karakteraistika.

Statistika se može baviti proučavanjem podataka koji **točno** opisuju stanje u svakoj populaciji. Ti se podaci dobivaju uglavnom popisom, redovitim evidencijama i praćenjima. Tako na primjer, svake desete godine se organiziraju popisi cjelokupnog stanovništva države. Analiziranjem i prikazivanjem takvih podataka bavi se tzv. **deskriptivna statistika**.

Vrlo često je nemoguće (i nepotrebno!) statistički obraditi čitavu populaciju. Djelom zbog toga što je ona prevelika da bi se taj postupak mogao sprovesti ili da bi bio isplativ. Drugi mogući razlog jest što se u nekim postupcima ispitivanja (recimo u kontroli kvalitete proizvodnje) u samom postupku ispitivanja uništava taj element populacije. Zamislimo na primjer ispitivanje duljine života žarulje!

U tom slučaju se proučava samo jedan mali dio populacije koji nazivamo **uzorak**. Na osnovu tog uzorka, donosimo potom sud o čitavoj populaciji. Predmet **matematičke statistike** jest statistička obrada uzorka: način odabira uzorka (da bi on dobro predstavljao čitavu populaciju) analiza obilježja u uzorku i procjena u kojoj su mjeri ti rezultati vjerodostojni za čitavu populaciju. Kako zaključci u ovom slučaju ne mogu nikad biti apsolutno sigurni (oni se donose uvijek s nekim stupnjem vjerojatnosti), matematička statistika se izražava i koristi metodama teorije vjerojatnosti.

Primjer 10.1.

Rezultati općih izbora postaju poznati (i službeni) kad se zna za glas svakog birača, tj. tek nakon što se obradi čitava populacija. Međutim, mnogo prije toga se rezultati mogu predvidjeti (s velikom dozom sigurnosti) na osnovu glasanja nekog dobro izabranog uzorka, koji može biti po veličini i 10 000 puta manji od čitave populacije!

Primjer 10.2.

Jedna proizvodna traka proizvodi otpornike. Dozvoljena granica škarta je 2%. Kako ćemo kontrolirati je li proizvodnja ispravna, t.j. je li postotak škarta unutar tih granica?

Bilo bi nerazumno, i skoro nemoguće za ovakav tip proizvoda, kontrolirati čitavu proizvodnju. Umjesto toga, uzimamo relativno maleni uzorak, odabran na pogodan način: recimo, svaki stoti proizvod. Ako je broj škartnih proizvoda u tom uzorku veći od određene granice, uz veliku dozu sigurnosti možemo zaključiti da je broj škartova u čitavoj populaciji veći od 2%, tj. da je došlo do grešaka u proizvodnji koje treba ispraviti. Kolika je ta dozvoljena granica škartnih proizvoda u uzorku i kolika je sigurnost u našem zaključku, to je predmet izučavanja matematičke statistike.

10.1.2. Populacija. Uzorak

Upoznajmo se s oznakama i temeljnim pojmovima matematičke statistike. Sa X ćemo označiti slučajnu varijablu koja će biti predmet proučavanja. Nju ćemo zvati **populacija**. Njezinu funkciju distribucije označavat ćemo sa F , funkciju gustoće (ako postoji) sa f , očekivanje s a i disperziju sa σ^2 .

U ovisnosti o problemu koji promatramo, neki parametri $\vartheta_1, \vartheta_2, \dots$ u ovoj razdiobi mogu biti nepoznati. Najčešći zadatak matematičke statistike jest dati odgovarajuću procjenu za te parametre. Ta se procjena postiže na temelju poznatih realizacija x_1, x_2, \dots, x_n slučajne varijable X .

Informacije o nepoznatoj razdiobi populacije X dobivamo samo na temelju realizacija te slučajne varijable.

Uzorak
<p>Neka je X slučajna varijabla s razdiobom F. Za slučajne varijable X_1, \dots, X_n kažemo da su nezavisne kopije slučajne varijable X, ako one imaju svojstva:</p> <ol style="list-style-type: none">međusobno su nezavisne, imaju razdiobu identičnu razdiobi slučajne varijable X. <p>Tako dobivenu n-torku slučajnih varijabli (X_1, \dots, X_n) nazivamo uzorak. Ako je x_1 je realizacija varijable X_1, x_2 realizacija varijable X_2 i t.d., tada se (x_1, \dots, x_n) naziva vrijednost ili realizacija uzorka (X_1, \dots, X_n).</p> <p>Broj n označava veličinu (dimenziju ili volumen) uzorka.</p>

Možemo zamisliti da varijable X_1, \dots, X_n opisuju ponašanje slučajne varijable X pri ponavljanju stohastičkog eksperimenta u nepromijenjenim uvjetima.

Radi jednostavnosti, pretpostavimo za sada da je u razdiobi varijable X poznat jedan parametar ϑ . Funkciju gustoće varijable X označavat ćemo s $f_\vartheta(x)$ ili s $f(\vartheta, x)$, jer ona ovisi o tom nepoznatom parametru ϑ .

Vrijednost parametra ϑ trebamo procijeniti na temelju realizacija x_1, x_2, \dots, x_n varijable X . Bit će definirana funkcija

$$\hat{\vartheta} = g(x_1, x_2, \dots, x_n)$$

koja će dati procjenu $\hat{\vartheta}$ parametra ϑ . Ta procjena ovisi, dakle, o realizacijama x_1, x_2, \dots, x_n . Realizacije su slučajne, pa je prirodno da će se pri ponavljanju pokusa pojaviti neka druga n -torka, a onda i druga vrijednost za procjenu $\hat{\vartheta}$. Zato je normalna situacija da procjena $\hat{\vartheta}$ nije jednaka nepoznatom parametru ϑ . (Jedan od zadataka matematičke statistike jest da pruži *mjeru sigurnosti* za točnost ove procjene.)

Budući da su x_1, x_2, \dots, x_n realizacije slučajnih varijabli X_1, X_2, \dots, X_n , onda će i $\hat{\vartheta}$ biti realizacija slučajne varijable

$$\Theta := g(X_1, X_2, \dots, X_n).$$

Statistika, procjenitelj i procjena
<p>Slučajna varijabla</p> $\Theta := g(X_1, X_2, \dots, X_n).$ <p>naziva se statistika. Statistikom nazivamo <i>svaku</i> funkciju koja ovisi o uzorku X_1, X_2, \dots, X_n, a ne ovisi (eksplicitno) o nepoznatom parametru.</p> <p>Neka je ϑ nepoznati parametar u populaciji X. Za statistiku (1) kažemo da je procjenitelj parametra ϑ. Vrijednost te statistike</p> $\hat{\vartheta} = g(x_1, x_2, \dots, x_n)$ <p>nazivamo procjenom parametra ϑ.</p>

Prema tome, procjenitelj je slučajna varijabla. Nakon realizacije uzorka, vrijednost procjenitelja daje nam procjenu nepoznatog parametra.

10.1.3. Statistika za procjenu očekivanja

Želimo procijeniti nepoznato očekivanje a populacije X . Prirodno je onda odabrati statistiku

$$\overline{X} := \frac{X_1 + X_2 + \dots + X_n}{n}.$$

Ona se naziva **sredina** uzorka.

Označimo nepoznato očekivanje i disperziju populacije X :

$$\begin{cases} E(X) = a, \\ D(X) = \sigma^2. \end{cases}$$

Varijabla \overline{X} je slučajna. Izračunajmo njezino očekivanje i disperziju! Prema svojstvima očekivanja, vrijedi

$$\begin{aligned} E(\overline{X}) &= E\left(\frac{X_1 + X_2 + \dots + X_n}{n}\right) \\ &= \frac{1}{n} \left[E(X_1) + E(X_2) + \dots + E(X_n) \right] = a. \end{aligned}$$

Varijable X_1, \dots, X_n su nezavisne, pa je

$$\begin{aligned} D(\overline{X}) &= D\left(\frac{X_1 + X_2 + \dots + X_n}{n}\right) \\ &= \frac{1}{n^2} \left[D(X_1) + D(X_2) + \dots + D(X_n) \right] = \frac{\sigma^2}{n}. \end{aligned}$$

Procjena očekivanja
<p>Nepoznato očekivanje a populacije X procjenjujemo pomoću sredine uzorka:</p> $\overline{X} = \frac{1}{n} \sum_{i=1}^n X_i.$ <p>Za tu slučajnu varijablu vrijedi</p> $E(\overline{X}) = a, \quad D(\overline{X}) = \frac{\sigma^2}{n},$ <p>gdje je σ^2 varianca (disperzija) populacije.</p>

Primjećujemo da je disperzija statistike \overline{X} obrnuto proporcionalna veličini uzorka. Ako je uzorak dovoljno velik, vrijednosti varijable \overline{X} bit će koncentrirane oko srednje vrijednosti $E(\overline{X}) = a$. Zato je jasno da će \overline{X} biti dobra procjena za a . O kvaliteti te procjene bit će više riječi u nastavku. ◀

10.1.4. Nepristrani procjenitelji

Među svim statistikama želimo odabrati one koje su, po nekim kriterijima, bolje od drugih. Zato ćemo izdvojiti neka poželjna svojstva statistika te dati kriterij za usporedbu različitih statistika.

U prethodnom primjeru, statistika \overline{X} za parametar a imala je svojstvo:

$$E(\overline{X}) = a.$$

Dakle, *očekivanje statistike podudara se s vrijednošću parametra*. Statistike koje posjeduju to poželjno svojstvo nazvat ćemo posebnim imenom.

Nepristrani procjenitelji
<p>Za statistiku Θ kažemo da je nepristrani procjenitelj ili nepristrana statistika parametra ϑ, ukoliko vrijedi</p> $E(\Theta) = \vartheta.$

Kriterij nepristranosti svakako je poželjan, ali nije jedini odlučujući za odabir statistike. Upoznat ćemo primjere kod kojih pristrani procjenitelji mogu imati bolja svojstva od nepristranih. (Na primjer, njihova disperzija može biti manja.)

10.1.5. Usporedba statistika

Usporedba statistika
<p>Neka je (X_1, \dots, X_n) uzorak, ϑ nepoznati parametar te $\Theta_1(X_1, \dots, X_n)$, $\Theta_2(X_1, \dots, X_n)$ dvije nepristrane statistike za ϑ. Kažemo da je Θ_1 bolja (efikasnija) od Θ_2 ako je $D(\Theta_1) < D(\Theta_2)$.</p>

Još je jedno poželjno svojstvo koje bi dobra statistika trebala imati: povećanjem uzorka statistika mora davati sve bolju aproksimaciju nepoznatog parametra.

Valjane statistike
<p>Statistiku $\Theta_n = \Theta(X_1, X_2, \dots, X_n)$ nazivamo valjanom procjenom parametra ϑ ako za svaki $\varepsilon > 0$ slučajna varijabla Θ_n konvergira prema ϑ <i>po vjerojatnosti</i>:</p> $\lim_{n \rightarrow \infty} P(\Theta_n - \vartheta < \varepsilon) = 1.$

Teorem 10.1.

Da bi nepristrana statistika bila valjana, dovoljno je da joj disperzija teži u nulu (kad n teži u beskonačnost).

DOKAZ. Ta tvrdnja slijedi iz Čebiševljeve nejednakosti:

$$P(|\Theta_n - \vartheta| < \varepsilon) \geq 1 - \frac{E[(\Theta_n - \vartheta)^2]}{\varepsilon^2} = 1 - \frac{D(\Theta_n)}{\varepsilon^2} \rightarrow 1.$$

10.1.6. Procjena disperzije, uz poznato očekivanje

Pretpostavimo sad da nam je očekivanje populacije poznato, a disperzija σ^2 nije. Za procjenu disperzije bismo statistiku

$$D^2 = \frac{1}{n} \sum_{i=1}^n (X_i - a)^2.$$

Očekivanje ove statistike je:

$$E(D^2) = \frac{1}{n} \sum_{i=1}^n E(X_i - a)^2 = \frac{1}{n} \sum_{i=1}^n D(X_i) = D(X) = \sigma^2.$$

Dakle, ova je statistika nepristrana.

O kvaliteti procjene odlučivat će disperzija statistike. Zbog nezavisnosti i jednake distribuiranosti slučajnih varijabli X_i bit će:

$$\begin{aligned} D(D^2) &= \frac{1}{n^2} \sum_{i=1}^n D\left[(X_i - a)^2\right] = \frac{1}{n} D\left[(X - a)^2\right] \\ &= \frac{1}{n} \left(E\left[(X - a)^4\right] - \left[E(X - a)^2\right]^2 \right) = \frac{1}{n} \left(\mu_4 - \sigma^4 \right) \end{aligned} \quad (5)$$

Ovdje je

$$\mu_4 = E\left[(X - a)^4\right]$$

četvrti centralni moment populacije X .

Vidimo da disperzija statistike D^2 opada obrnuto proporcionalno veličini uzorka. Prema Teoremu 10.1, ova je statistika valjana.

10.1.7. Procjena disperzije, uz nepoznato očekivanje

Ako je očekivanje poznato, tada je statistika

$$D^2 = \frac{1}{n} \sum_{i=1}^n (X_i - a)^2$$

nepristrani procjenitelj za disperziju. Koju ćemo statistiku koristiti ako je i očekivanje a nepoznato? Prirodno je zamijeniti ga u ovoj formuli s \bar{X} . Tako dobivamo statistiku

$$\Theta = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2.$$

Provjerimo je li ona nepristrana. Njezino očekivanje je

$$E(\Theta) = E\left(\frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2\right) = \frac{1}{n} \sum_{i=1}^n E(X_i - \bar{X})^2.$$

Vrijedi $E(X_i - \bar{X}) = a - a = 0$, pa je $E(X_i - \bar{X})^2 = D(X_i - \bar{X})$. Sada je, zbog nezavisnosti varijabli X_1, X_2, \dots, X_n ,

$$\begin{aligned} E(\Theta) &= \frac{1}{n} \sum_{i=1}^n D(X_i - \bar{X}) \\ &= \frac{1}{n} \sum_{i=1}^n D\left[X_i - \frac{1}{n} \sum_{j=1}^n X_j\right] \\ &= \frac{1}{n} \sum_{i=1}^n D\left[\frac{n-1}{n} X_i - \frac{1}{n} \sum_{j \neq i} X_j\right] \\ &= \frac{1}{n} \sum_{i=1}^n \left[\left(\frac{n-1}{n}\right)^2 D(X_i) + \frac{1}{n^2} \sum_{j \neq i} D(X_j) \right] \\ &= \frac{1}{n} \cdot n \left[\left(\frac{n-1}{n}\right)^2 \sigma^2 + \frac{1}{n^2} \cdot (n-1) \sigma^2 \right] \\ &= \frac{n-1}{n} \sigma^2. \end{aligned}$$

Prema tome, očekivanje statistike Θ ne podudara se s parametrom σ^2 . Ovaj procjenitelj nije nepristran. Primjetimo ipak da se razlika očekivanja procjenitelja i parametra smanjuje povećavanjem veličine uzorka n .

Međutim, množenjem s konstantnim faktorom $\frac{n}{n-1}$ ovaj se procjenitelj može učiniti nepristranim:

Procjene disperzije
Ako je očekivanje a populacije X poznato, nepristrana procjena nepoznate disperzije σ^2 računa se formulom
$D^2 := \frac{1}{n} \sum_{i=1}^n (X_i - a)^2. \quad (6)$
Ako su očekivanje a i disperzija σ^2 populacije X nepoznati, onda se nepristrani procjenitelj za disperziju računa formulom
$S^2 := \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2. \quad (7)$

Statistika S^2 je valjana, jer joj disperzija teži k nuli. Vrijedi naime

$$\begin{aligned} D(S^2) &= E[(S^2 - \sigma^2)^2] \\ &= E(S^4) - 2\sigma^2 E(S^2) + \sigma^4 = E(S^4) - \sigma^4. \end{aligned}$$

Iz prikaza

$$S^2 = \frac{1}{n} \sum_{i=1}^n (X_i - a)^2 - \frac{1}{n(n-1)} \sum_{\substack{i,j=1 \\ i \neq j}}^n (X_i - a)(X_j - a)$$

nakon kvadriranja ovog izraza i računanja očekivanja svakog člana, dobivamo izraz sličan (5):

$$D(S^2) = \frac{1}{n} \left(\mu_4 - \frac{n-3}{n-1} \sigma^4 \right). \quad (8)$$

Primjer 10.3.

Da bi se utvrdila preciznost mjernog geodetskog uređaja koji nema sistematske pogreške, načinjeno je šest mjerenja. Dobiveni su rezultati (u metrima): 3540, 3582, 3555, 3578, 3564, 3548. Odredi nepristranu procjenu za varijancu, u slučajevima (a) ako je poznato da iznos mjerene veličine iznosi 3560 m, (b) ako nije poznat iznos mjerene veličine.

► (a) U ovom je slučaju poznato očekivanje slučajne varijable, jer ono mora biti jednako mjerenoj vrijednosti (zbog odsustva sistematske pogreške): $a = 3560$. Zato procjenu za varijancu računamo ovako:

$$\hat{d}^2 = \frac{1}{n} \sum_{i=1}^n (x_i - a)^2 = 232.17 \text{ m}^2.$$

(b) Očekivanje je nepoznato, pa ga računamo iz uzorka:

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i = 3561.17 \text{ m}.$$

Nepristranu procjenu varijance računamo ovako:

$$\hat{s}^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2 = 276.97 \text{ m}^2. \quad \blacktriangleleft$$

10.1.8. Uporaba džepnog računala

Formula

$$\hat{s}^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2$$

nije najprikladnija za račun džepnim računalom. Ona zahtjeva izvođenje točno $5n+1$ operacija (za računalu s inverznom notacijom, inače je broj neznatno veći). Pod operacijom se smatra svako unošenje podataka ili njihov poziv iz memorije, te svaka funkcijska ili aritmetička operacija.

Transformirajmo ovaj izraz na sljedeći način:

$$\hat{s}^2 = \frac{1}{n-1} \left(\sum_{i=1}^n x_i^2 - 2\bar{x} \sum_{i=1}^n x_i + n\bar{x}^2 \right) = \frac{1}{n-1} \left(\sum_{i=1}^n x_i^2 - n\bar{x}^2 \right)$$

Sad je nužno napraviti točno $3n+5$ operacija. Zato ćemo procjenu disperzije računati ovom formulom.

Praktički svi džepni kalkulatori imaju ugrađene elementarne statističke funkcije. Neki među njima specijalizirani su upravo za rješavanje statističkih zadataka. Na različitim računalima mogu postojati različiti načini korištenja tih funkcija, ali zajednički principi mogu se opisati ovako.

Niz podataka x_1, x_2, \dots, x_n unosi se posebnom tipkom, obično označenom s $\boxed{\Sigma}$. Na koncu unosa, u posebnim registrima spremjeni su sljedeći podatci:

- volumen uzorka n ;
- zbroj elemenata uzorka, $\sum x_i$;
- zbroj kvadrata elemenata uzorka, $\sum x_i^2$.

U posebnim su registrima također spremljene izračunate statističke funkcije. Na boljim računalima, pozivi tih registrara nalaze se na posebnim tipkama označenim s $\boxed{\bar{x}}$, \boxed{s} i $\boxed{\sigma}$.

10.1.9. Računanje s grupiranim podatcima

Podatci dani u uzorku vrlo su često grupirani u *razrede*. Uzorak tada ima ovakav oblik

x_1	n_1
x_2	n_2
\vdots	\vdots
x_r	n_r

Ovdje je $n = n_1 + \dots + n_r$ volumen uzorka. Sredina i disperzija uzorka računaju se tada formulama

$$\begin{aligned} \bar{x} &= \frac{1}{n} \sum_{i=1}^r n_i x_i, \\ \hat{s} &= \frac{1}{n-1} \sum_{i=1}^r n_i (x_i - \bar{x})^2 = \frac{1}{n-1} \left(\sum_{i=1}^r n_i x_i^2 - n\bar{x}^2 \right). \end{aligned}$$

Primjer 10.4.

Odredimo procjenu za očekivanje i disperziju na temelju uzorka normalne populacije:

x_i	2560	2600	2620	2650	2700
n_i	2	3	10	4	1

► Volumen uzorka je

$$n = n_1 + \dots + n_5 = 20.$$

Računanje očekivanja i disperzije olakšano je ako vrijednosti slučajne varijable translateramo za isti iznos C . Tu je C po volji odabrani broj. Pri tom vrijedi

$$E(X) = C + E(X - C), \quad D(X) = D(X - C)$$

(U računu koji slijedi koristit ćemo samo prvo svojstvo.) Za pogodnu konstantu u ovom primjeru možemo uzeti $C = 2620$:

$$\begin{aligned} \bar{x} &= 2620 + \frac{1}{n} \sum_{i=1}^5 n_i (x_i - 2620) \\ &= 2620 + \frac{1}{20} \left[-60 \cdot 2 + (-20) \cdot 3 + 30 \cdot 4 + 80 \right] = 2620 + 1 = 2621. \end{aligned}$$

$$\hat{s}^2 = \frac{1}{n-1} \left(\sum_{i=1}^5 n_i x_i^2 - n\bar{x}^2 \right) = 967.4 \quad \blacktriangleleft$$

Primjer 10.5. ■ Određivanje varijance na temelju zadanog uzorka varijanci

Pri kontroli kvalitete nekog proizvoda, ispituje se varijanca na kontrolnim uzorcima tijekom svakog dana. Dobivene su vrijednosti $s_1^2, s_2^2, \dots, s_k^2$, na temelju uzoraka veličina n_1, n_2, \dots, n_k . Kako ćemo odrediti procjenu varijance ove populacije?

► Trebamo odrediti nepristrani procjenitelj za nepoznatu varijancu σ^2 . Izabrat ćemo statistiku

$$\hat{\Theta} = \frac{a_1 S_1^2 + a_2 S_2^2 + \dots + a_k S_k^2}{A}$$

gdje su a_1, a_2, \dots, a_k i A konstante koje treba odrediti.

Očekivanje ove statistike je

$$E(\hat{\Theta}) = \frac{a_1 E(S_1^2) + \dots + a_k E(S_k^2)}{A} = \frac{a_1 + \dots + a_k}{A} \cdot \sigma^2$$

Statistika će biti nepristrana ako je $A = a_1 + \dots + a_k$. Konstante a_1, \dots, a_k možemo birati po volji, ali je prirodno da one odgovaraju veličinama pojedinih uzoraka. Tako će dnevne procjene temeljene na većem uzorku imati veću težinu u konačnoj procjeni. Prema tome, tražena procjena je

$$\hat{\vartheta} = \frac{n_1 s_1^2 + n_2 s_2^2 + \dots + n_k s_k^2}{n_1 + n_2 + \dots + n_k}. \quad \blacktriangleleft$$

10.1.10. * Nepristrana procjena standardnog odstupanja

Pokazali smo da je

$$D^2 = \frac{1}{n} \sum_{i=1}^n (X_i - a)^2$$

nepristrana procjena disperzije σ^2 . **Standardno odstupanje (devijacija)** definira se kao korijen disperzije, $\sigma = \sqrt{\sigma^2}$. Logično je postaviti pitanje: je li veličina

$$D = \sqrt{D^2} = \left(\frac{1}{n} \sum_{i=1}^n (X_i - a)^2 \right)^{1/2}$$

nepristrana procjena za standardno odstupanje?

Može izgledati neobično, ali odgovor je negativan. Razlog tome je što funkcija drugog korijena "jače skuplja" velike brojeve od malih. Za bilo koju nedegeneriranu pozitivnu slučajnu varijablu Y općenito vrijedi

$$E(Y) < \sqrt{E(Y^2)}.$$

(Ova nejednakost slijedi iz Cauchy-Schwarz-Bunjakovskijeve nejednakosti.) Zato je

$$E(D) < \sqrt{E(D^2)} = \sqrt{\sigma^2} = \sigma.$$

Nepristranu procjenu za standardno odstupanje praktički je nemoguće utvrditi u općem slučaju, za bilo koju distribuciju populacije X . Ako populacija ima normalnu razdiobu, onda se može dokazati da nepristrana procjena glasi

$$\tilde{D} = k_{n+1} \sqrt{\frac{1}{n} \sum_{i=1}^n (X_i - a)^2}, \quad (9)$$

pri čemu se koeficijent k_n računa formulom

$$k_n = \sqrt{\frac{n-1}{2}} \cdot \frac{\Gamma\left(\frac{n-1}{2}\right)}{\Gamma\left(\frac{n}{2}\right)}. \quad (10)$$

Na isti se način dobiva nepristrana procjena standardnog odstupanja ukoliko očekivanje nije poznato:

$$\tilde{S} = k_n \sqrt{\frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2},$$

Za velike vrijednosti od n može se koristiti jednostavnija formula:

$$\tilde{S} = \sqrt{\frac{1}{n-1.45} \sum_{i=1}^n (X_i - \bar{X})^2},$$

Za male vrijednosti od n može biti korisna sljedeća tablica:

n	k_n	n	k_n	n	k_n
3	1.1284	10	1.0280	30	1.0087
4	1.0853	12	1.0230	35	1.0072
5	1.0640	15	1.0181	40	1.0064
6	1.0506	20	1.0134	45	1.0056
7	1.0423	25	1.0104	50	1.0051

Međuvrijednosti se mogu utvrditi interpolacijom.

U većini zadataka i teoriji koja slijedi, kao procjenu za odstupanje ćemo ipak, jednostavnosti radi, koristiti korijen varijance. Izuzetak će biti zadatci u kojima se eksplicitno traži nepristrana procjena standardnog odstupanja.

10.2. Kriterij najveće izglednosti

Pretpostavimo da nam je u razdiobi slučajne varijable X nepoznata vrijednost jednog parametra ϑ . Označimo sa $f(\vartheta, x)$ zakon razdiobe te slučajne varijable,

$$f(\vartheta, x) = P_{\vartheta}(\{X = x\}), \quad x \in S,$$

ako je X diskretnog tipa, odnosno, neka je

$$f(\vartheta, x) = f_{\vartheta}(x),$$

funkcija gustoće, ako je X neprekidna slučajna varijabla. Indeks ϑ označava da se vjerojatnost događaja $\{X = x\}$ i vrijednost funkcije gustoće $f(x)$ varijable X računaju uz pretpostavku da je nepoznata vrijednost parametra, o kojem ovise te vrijednosti, jednaka ϑ . Tu nepoznatu vrijednost ćemo pokušati procijeniti iz vrijednosti uzorka (x_1, \dots, x_n) .

Kriterij najveće izglednosti

Neka je x_1, x_2, \dots, x_n realizacija uzorka populacije X , čija funkcija gustoće $f(\vartheta, x)$ ovisi o nepoznatom parametru ϑ . **Funkcija izglednosti** definira se kao umnožak

$$L(\vartheta, x_1, \dots, x_n) := f(\vartheta, x_1)f(\vartheta, x_2) \cdots f(\vartheta, x_n). \quad (11)$$

Za procjenu parametra ϑ uzimamo onu vrijednost $\hat{\vartheta}$ za koju funkcija izglednosti poprima globalni maksimum.

Zašto se ova funkcija naziva *funkcija izglednosti*? Za zadani x , vrijednost $f(\vartheta, x)$ opisuje vjerojatnost da slučajna varijabla poprimi vjerojatnost u okolišu broja x . Zato umnožak (11) predstavlja vjerojatnost da uzorak (X_1, X_2, \dots, X_n) poprimi vrijednost u okolišu od (x_1, x_2, \dots, x_n) . Postavlja se pitanje: za koju će vrijednost parametra ϑ ta vjerojatnost biti najveća? Na taj način dobivamo *kriterij za odabir procjene parametra* ϑ . Za ϑ ćemo odabrati onaj parametar koji maksimizira funkciju izglednosti. Na taj način *maksimiziramo vjerojatnost pojavljivanja* uzorka koji se ostvario!

Opravljanje ovog uvjeta je intuitivna pretpostavka da onom uzorku koji se stvarno realizira trebamo dati prednost u odnosu na one koji se nisu ostvarili. Nakon što nam je poznata vrijednost uzorka (x_1, \dots, x_n) , uzimamo onu vrijednost parametra ϑ za koju ta realizacija ima najveću vjerojatnost pojavljivanja, veću nego bilo koja druga realizacija.

Primjer 10.6. ■ Procjena parametra eksponencijalne razdiobe

Vrijeme X ispravnog rada uređaja kojem se karakteristike ne mijenjaju vremenom, dobro je opisano eksponencijalnom razdiobom, s gustoćom

$$f(x) = \lambda e^{-\lambda x}, \quad x > 0.$$

Ovdje je λ nepoznati parametar razdiobe. Bilježeni su rezultati na probnom uzorku i dobiven niz x_1, x_2, \dots, x_n . Na temelju tih rezultata, koristeći se kriterijem najveće izglednosti, treba procijeniti očekivanje varijable X .

► U ovom je primjeru

$$\begin{aligned} L(\lambda, x_1, \dots, x_n) &= f(\lambda, x_1)f(\lambda, x_2) \cdots f(\lambda, x_n) \\ &= \lambda e^{-\lambda x_1} \cdot \lambda e^{-\lambda x_2} \cdots \lambda e^{-\lambda x_n} \\ &= \lambda^n e^{-\lambda z}. \end{aligned}$$

Tu smo označili $z = \sum_{i=1}^n x_i$. Izračunajmo maksimum ove funkcije:

$$\frac{\partial L}{\partial \lambda} = \lambda^{n-1} e^{-\lambda z} (-\lambda z + n).$$

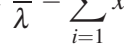
Ova se derivacija poništava kad je $\lambda = 0$ ili kad je $-\lambda z + n = 0$. Prva je mogućnost besmislena, a iz druge slijedi

$$\hat{\lambda} = \frac{n}{z} = \frac{n}{x_1 + x_2 + \dots + x_n} = \frac{1}{\bar{x}}.$$

Za eksponencijalnu funkciju je poznato da vrijedi $E(X) = 1/\lambda$. zato je procjena za očekivanje

$$\bar{x} = \frac{1}{\hat{\lambda}} = \frac{x_1 + x_2 + \dots + x_n}{n}.$$

Ovaj je rezultat u skladu sa standardnom procjenom za očekivanje slučajne varijable. ◀



Korisno je primjetiti da je funkcija izglednosti L uvijek pozitivna, pa je stoga definirana i funkcija $\ln L$. S obzirom da vrijedi

$$(\ln L)' = \frac{L'}{L},$$

ova funkcija poprima maksimum u istim točkama kao i L . Vrlo je često nju praktičnije derivirati nego funkciju L . U prethodnom primjeru je

$$\begin{aligned} \ln L &= n \ln \lambda - \lambda \sum_{i=1}^n x_i, \\ \frac{\partial \ln L}{\partial \lambda} &= \frac{n}{\lambda} - \sum_{i=1}^n x_i = 0, \end{aligned}$$

pa slijedi isti zaključak kao prije.

10.2.1. Procjena vjerojatnosti događaja

Događaj A ima nepoznatu vjerojatnost realiziranja p . Kako ćemo procijeniti tu vjerojatnost?

Više je mogućih odgovora. Navedimo dva najjednostavnija.

Primjer 10.7. ■ Procjena vjerojatnosti korištenjem relativne frekvencije

Pokus u kojem se događaj A može ostvariti ponavljamo n puta, pri nepromijenjenim uvjetima. U svakom ponavljanju bilježimo je li se događaj zbilo ili nije. Na taj je način dobiven uzorak x_1, x_2, \dots, x_n , pri čemu je $x_k = 0$ ako se A nije ostvarilo, a $x_k = 1$ ako se A ostvarilo. Na temelju tog uzorka treba procijeniti vjerojatnost p .

► Rezultat pokusa prate indikatorske slučajne varijable koje su nezavisne kopije slučajne varijable

$$X \sim \begin{pmatrix} 0 & 1 \\ q & p \end{pmatrix}.$$

Ova je razdioba zadana vjerojatnostima (ovisnim o nepoznatom parametru p):

$$f(p, x) = p^x (1-p)^{1-x}, \quad x = 0 \text{ ili } 1.$$

Funkcija izglednosti je

$$L(p, x_1, \dots, x_n) = \prod_{k=1}^n p^{x_k} (1-p)^{1-x_k}.$$

Zato je

$$\begin{aligned} \ln L &= \sum_{k=1}^n \left[x_k \ln p + (1-x_k) \ln(1-p) \right], \\ \frac{\partial \ln L}{\partial p} &= \sum_{k=1}^n \left(\frac{x_k}{p} - \frac{1-x_k}{1-p} \right) = \frac{m}{p} - \frac{n-m}{1-p} = 0. \end{aligned} \quad (12)$$

Tu smo s m označili

$$m = x_1 + \dots + x_n,$$

a taj je zbroj jednak broju pojavljivanja događaja A u n pokusa.

Sređivanjem jednakosti (12) dobivamo $p = \frac{m}{n}$. ◀

Primjer 10.8. ■ Procjena vjerojatnosti korištenjem geometrijske razdiobe

Ponavljamo pri nepromijenjenim uvjetima pokus u kojem se događaj A može ostvariti i bilježimo broj pokusa kad se to dogodilo. Čitav se postupak ponavlja n puta. Na taj je način dobiven uzorak x_1, x_2, \dots, x_n . Na temelju tog uzorka treba procijeniti vjerojatnost p .

► Slučajna varijabla koja opisuje pojavljivanje događaja A ima geometrijsku razdiobu s parametrom p . Ona je zadana vjerojatnostima

$$f(p, x) = p(1-p)^{x-1}, \quad x = 1, 2, 3, \dots$$

Funkcija izglednosti je (nepoznati parametar i dalje označavamo s p):

$$L(p, x_1, \dots, x_n) = p^n (1-p)^{x_1-1} \cdots (1-p)^{x_n-1}.$$

Sada imamo

$$\begin{aligned} \ln L &= n \ln p + \ln(1-p) \left(\sum_{i=1}^n x_i - n \right), \\ \frac{\partial \ln L}{\partial p} &= \frac{n}{p} - \frac{\sum_{i=1}^n x_i - n}{1-p} = 0. \end{aligned}$$

Iz posljednje jednakosti, nakon sređivanja, dobivamo

$$\hat{p} = \frac{n}{x_1 + x_2 + \dots + x_n} = \frac{1}{\bar{x}}. \quad \blacktriangleleft$$

Primjer 10.9. ■ Procjena vjerojatnosti

Postotak bijelih kuglica u kutiji je nepoznat. Zagrabili smo n kuglica i pobrojali m bijelih. Kolika je procjena za postotak bijelih kuglica?

► Taj postotak jednak je vjerojatnosti da izvučena kuglica iz kutije bude bijela. Neka je p ta vjerojatnost. Slučajna varijabla X koju promatramo je broj bijelih kuglica u uzorku veličine n . Njezina je razdioba $X \sim B(n, p)$. Zato vrijedi

$$f(p, x) = \binom{n}{x} p^x (1-p)^{n-x}, \quad x = 0, 1, \dots, n.$$

Odavde je

$$\begin{aligned} \frac{\partial}{\partial p} \ln f(p, x) &= \frac{\partial}{\partial p} [x \ln p + (n-x) \ln(1-p)] \\ &= \frac{x}{p} - \frac{n-x}{1-p} = 0 \implies \hat{p} = \frac{x}{n} \end{aligned}$$

Ako se u uzorku pojavilo m bijelih kuglica, onda je najbolja procjena $\hat{p} = \frac{m}{n}$. ◀

Primjer 10.10. ■ Procjena parametra Poissonove razdiobe

Neka X ima Poissonovu razdiobu, $X \sim \mathcal{P}(\lambda)$, λ nepoznat. Procijenimo vrijednost od λ .

► Sada je

$$f(\lambda, x) = P(X = x) = \frac{\lambda^x}{x!} e^{-\lambda}.$$

Funkcija izglednosti je

$$\begin{aligned} L(\lambda, x_1, \dots, x_n) &= \frac{\lambda^{x_1 + \dots + x_n}}{x_1! \cdots x_n!} e^{-n\lambda}, \\ \ln L(\lambda, x_1, \dots, x_n) &= -n\lambda + (x_1 + \dots + x_n) \ln \lambda - \sum \ln(x_i!), \\ \frac{\partial \ln L}{\partial \lambda} &= -n + \frac{x_1 + \dots + x_n}{\lambda} = 0, \\ \hat{\lambda} &= \frac{x_1 + \dots + x_n}{n} = \bar{x}. \quad \blacktriangleleft \end{aligned}$$



Kriterijom najveće izglednosti možemo odrediti i više od jednog nepoznatog parametra. Ako funkcija izglednosti ima oblik

$$L(\vartheta_1, \dots, \vartheta_s, x_1, \dots, x_n) = f(\vartheta_1, \dots, \vartheta_s, x_1) \cdots f(\vartheta_1, \dots, \vartheta_s, x_n)$$

onda nepoznate parametre $\vartheta_1, \dots, \vartheta_s$ dobivamo iz uvjeta

$$\frac{\partial L(\vartheta_1, \dots, \vartheta_s, x_1, \dots, x_n)}{\partial \vartheta_i} = 0, \quad i = 1, 2, \dots, s. \quad (13)$$

Primjer 10.11. ■ Procjena parametara normalne razdiobe

Slučajna varijabla X ima normalnu razdiobu $\mathcal{N}(a, \sigma^2)$ s nepoznatim i očekivanjem a i disperzijom σ^2 . Odredimo procjenu tih parametara koristeći kriterij najveće izglednosti.

► Pripadna funkcija izglednosti je

$$L(a, \sigma) = L(a, \sigma, x_1, \dots, x_n) = \frac{1}{(2\pi\sigma^2)^{n/2}} \exp\left\{-\frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - a)^2\right\}.$$

Logaritam ove funkcije je

$$\ln L(a, \sigma) = -\frac{n}{2} \ln(2\pi\sigma^2) - \frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - a)^2.$$

Ekstrem dobivamo ako vrijedi:

$$\frac{\partial}{\partial a} \ln L(a, \sigma) = \frac{1}{2\sigma^2} \cdot 2 \sum_{i=1}^n (x_i - a) = 0$$

i odavde

$$\begin{aligned} \hat{a} &= \frac{1}{n} \sum_{i=1}^n x_i, \\ \frac{\partial}{\partial \sigma} \ln L(a, \sigma) &= -\frac{n}{2} \cdot \frac{2}{\sigma} + \frac{1}{\sigma^3} \sum_{i=1}^n (x_i - a)^2 = 0 \end{aligned}$$

i odavde

$$\hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \hat{a})^2.$$

Primjer 10.12. ■ Odabir statistike

Zadana je jednolika razdioba na intervalu $[0, c]$, c nepoznat. Vrijednost od c trebali bismo procijeniti na osnovu uzorka: nekoliko na sreću odabranih brojeva iz tog intervala. Koju statistiku je prikladno upotrebiti?

► Recimo, zbog jednostavnijeg razmatranja, da je uzorak dao sljedeće vrijednosti

$$3, 5, 8, 1, 3, 2, 6, 2.$$

Neka je X slučajna varijabla: vrijednost na sreću odabranog broja iz intervala $[0, c]$. Njezina je funkcija gustoće

$$f(c, x) = \begin{cases} 1/c, & 0 \leq x \leq c, \\ 0, & \text{inače.} \end{cases}$$

Pri tom vrijedi $E(X) = \frac{c}{2}$, $D(X) = \frac{c^2}{12}$.

Odaberimo statistiku pomoću koje možemo odrediti c . Budući je $E(X) = c/2$, a znamo statistiku pomoću koje određujemo $E(X)$, onda možemo izabrati statistiku

$$\Theta_1 = 2\bar{X} = \frac{2}{n} \sum_{k=1}^n X_k.$$

Vrijedi

$$E(\Theta_1) = \frac{2}{n} \sum_{k=1}^n E(X_k) = \frac{2}{n} \cdot \frac{c}{2} \cdot n = c$$

pa je ova statistika nepristrana. Međutim, ona nije baš najsretnije odabrana. Ukoliko uzorak poprimi opisanu vrijednost, onda dobivamo sljedeću procjenu za c

$$\hat{c} = \frac{2}{8} (3 + 5 + 8 + 1 + 3 + 2 + 6 + 2) = 7.5$$

što je apsurd, budući se pojavila realizacija 8.

Pokušajmo odrediti prikladniju statistiku. Pogledajmo što će dati kriterij najveće izglednosti. Imamo

$$L(c, x_1, \dots, x_n) = \begin{cases} \left(\frac{1}{c}\right)^n, & x_k \leq c, \forall k, \\ 0, & \text{inače.} \end{cases}$$

Maksimum se postiže ako je c najmanji moguć, a to je za $\hat{c} = \max_{1 \leq k \leq n} x_k$.

Time se nameće statistika

$$Y = \max\{X_1, \dots, X_n\}.$$

U zadanom uzorku, dobili bismo procjenu $\hat{c} = 8$. Očividno je da niti ona nije posve zadovoljavajuća: ako je *zaista* $c = 8$, nevjerovatno je da se baš ta maksimalna vrijednost i izabere.

Pogledajmo je li ova statistika nepristrana. Odredimo njezinu razdiobu:

$$\begin{aligned} F_Y(y) &= P(Y < x) = P(X_1 < x, \dots, X_n < x) \\ &= P(X < x)^n = \left(\frac{x}{c}\right)^n, \quad 0 < x \leq c. \end{aligned}$$

Funkcija gustoće je (u ovisnosti o parametru c)

$$f(c, x) = \frac{n}{c^n} x^{n-1}, \quad 0 < x < c.$$

Odavde dobivamo očekivanje

$$E(Y) = \frac{n}{c^n} \int_0^c x \cdot x^{n-1} dx = \frac{n}{n+1} c.$$

Statistika nije nepristrana. Stoga činimo korekciju i promatramo drugu statistiku

$$\Theta_2 = \frac{n+1}{n} Y = \frac{n+1}{n} \max\{X_1, \dots, X_n\}$$

koja daje nepristranu procjenu za c . ◀

Usporedimo statistike Θ_1 i Θ_2 iz ovog primjera. Vrijedi

$$D(\Theta_1) = D(2\bar{X}) = 4D\left(\frac{X_1 + \dots + X_n}{n}\right) = 4 \cdot \frac{D(X)}{n} = \frac{c^2}{3n}.$$

Izračunajmo disperziju statistike Y :

$$E(Y^2) = \frac{n}{c^n} \int_0^c x^2 \cdot x^{n-1} dx = \frac{n}{n+2} c^2$$

i odavde

$$\begin{aligned} D(Y) &= E(Y^2) - E(Y)^2 = \left(\frac{n}{n+2} - \frac{n^2}{(n+1)^2}\right) c^2 \\ &= \frac{n}{(n+2)(n+1)^2} c^2. \end{aligned}$$

Sada je

$$D(\Theta_2) = D\left(\frac{n+1}{n} Y\right) = \frac{(n+1)^2}{n^2} \cdot \frac{n}{(n+2)(n+1)^2} c^2 = \frac{c^2}{n(n+2)}.$$

Čim je $n \geq 2$, statistika Θ_2 je efikasnija od statistike Θ_1 .

Primjer 10.13. ■ Nepristranost i valjanost statistike

Kao procjenu za nepoznato očekivanje $a = E(X)$ možemo uzeti statistike

$$\begin{aligned} \Theta_1 &= \bar{X} = \frac{X_1 + X_2 + \dots + X_n}{n}, \\ \Theta_2 &= X_1. \end{aligned}$$

Provjerimo da su obje procjene nepristrane. Koja među njima je valjana?

► Za statistiku Θ_1 znamo da je nepristrana. Isto vrijedi i za Θ_2 :

$$E(\Theta_2) = E(X_1) = a.$$

Provjerimo da je Θ_1 valjana. Iz poznatog svojstva

$$D(\Theta_1) = \frac{1}{n^2} \sum_{i=1}^n D(X_i) = \frac{\sigma^2}{n}$$

s pomoću Čebiševljeve nejednakosti dobivamo

$$P(|\Theta_1 - a| > \varepsilon) \leq \frac{D(\Theta_1)}{\varepsilon^2} = \frac{\sigma^2}{\varepsilon^2 n} \rightarrow 0$$

kad $n \rightarrow \infty$. Dakle, $\Theta_1 = \bar{X}$ je valjana procjena za očekivanje.

Za statistiku Θ_2 vrijedi pak

$$P(|\Theta_2 - \vartheta| > \varepsilon) = P(|X_1 - \vartheta| > \varepsilon) > 0$$

čim je razdioba od X_1 netrivialna. Zato Θ_2 nije valjana statistika. ◀