

Korona zapski

STATISTIKA

M. Perman

Pembelajaran semester 2020

1.3. Centralni limitni izrek

Pogosto nas zanimajo porazdelitve
vseh slučajnih spremenljivk.

Analično se da porazdelitve
izračunati le v nekaj primerih.

Omejili se bomo na primere, ko
bodo X_1, X_2, \dots neodvisne, enakomerno
porazdeljene, vsota pa bo

$$S_n = X_1 + X_2 + \dots + X_n.$$

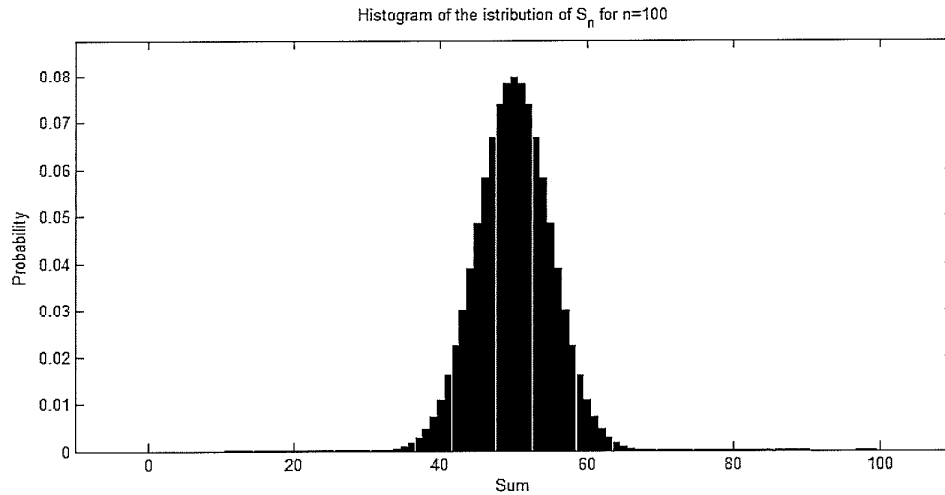
Oglejmo si

nekaj primerov histogramov vsot

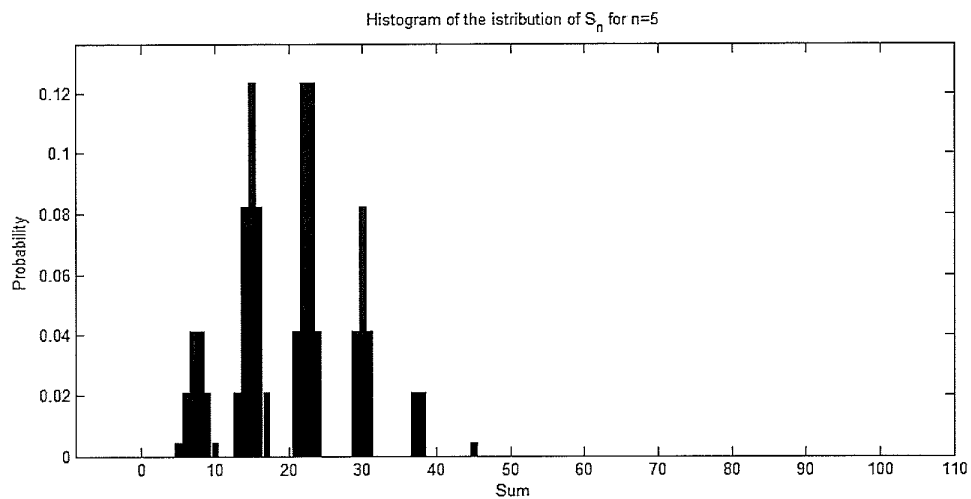
S_n za celoštevilске X_1, X_2, \dots

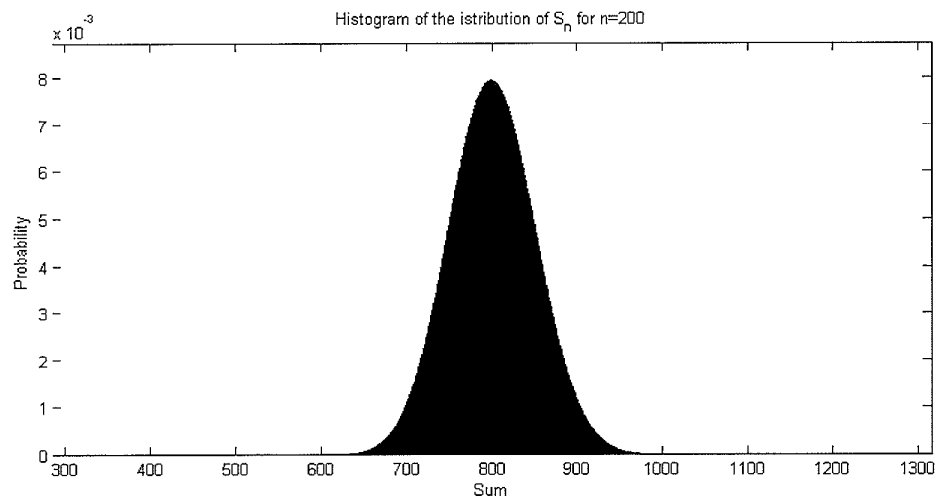
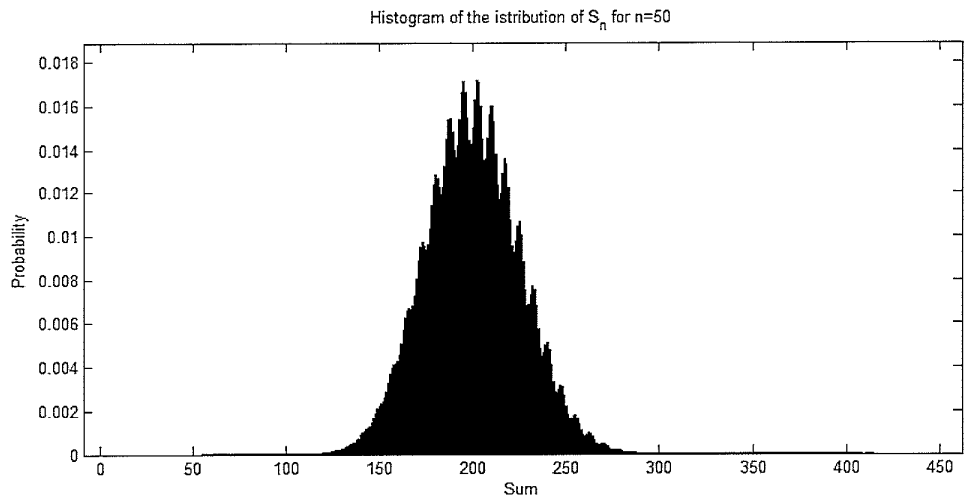
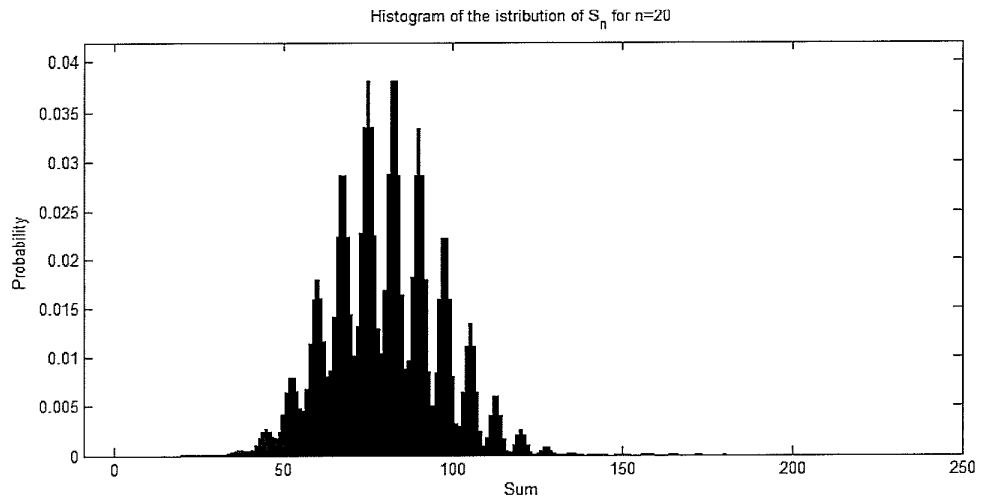
We will look at a few examples of distributions of S_n for different distributions of X_1 and different n .

1. Let $P(X_1 = 0) = P(X_1 = 1) = \frac{1}{2}$. Take $n = 100$. Let $S_n = X_1 + \dots + X_n$. The histogram of the distribution of S_n is:

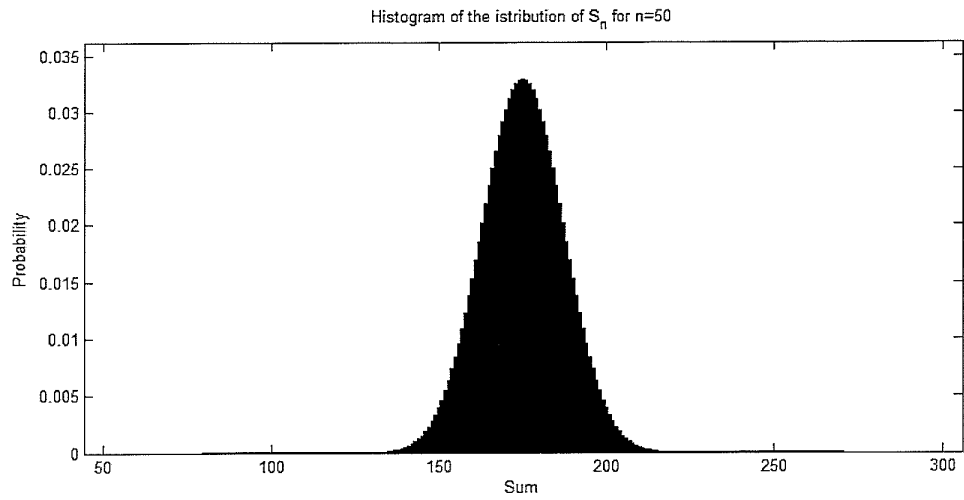


2. Let $P(X_1 = 1) = P(X_1 = 2) = P(X_1 = 9) = \frac{1}{3}$. Let $n = 5, 20, 50, 200$.

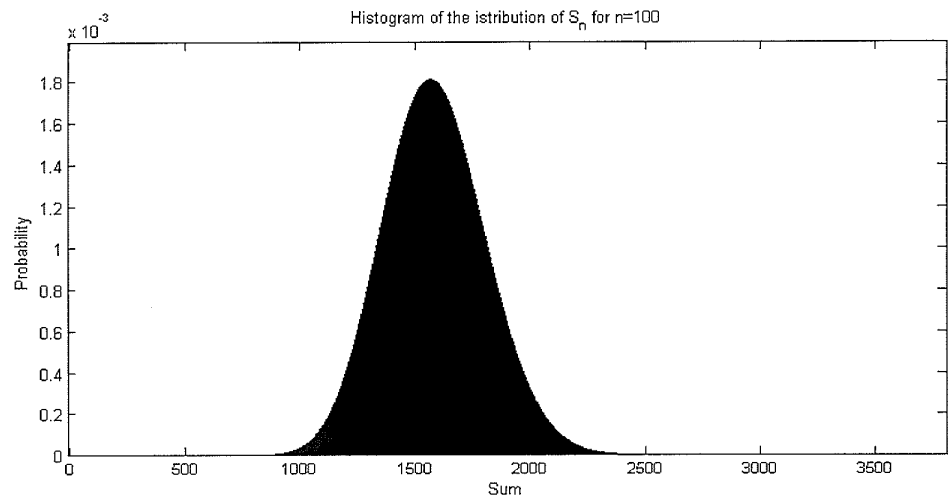




3. Take $P(X_1 = k) = 1/6$ for $k = 1, 2, \dots, 6$. Take $n = 50$.



4. Take $P(X_1 = 2^k) = 1/7$ for $k = 0, 1, 2, 3, 4, 5, 6$. Take $n = 100$.



~~1.2 Centralni limitni izrek~~

Centralni limitni izrek je eden od ključnih izrekov za statistiko. Kot smo večkrat opazili, se vsote neodvisnih, enako porazdeljenih slučajnih spremenljivk naravno pojavijo v številnih situacijah. Formulacija centralnega limitnega izreka je naslednja:

IZREK 1.10: Naj bodo X_1, X_2, \dots med sabo neodvisne, enako porazdeljene slučajne spremenljivke. Predpostavimo, da $E(|X_1|) < \infty$ in $\text{var}(X_1) < \infty$. Naj bo $S_n = X_1 + X_2 + \dots + X_n$. Za poljuben $x \in \mathbb{R}$ velja

$$\lim_{n \rightarrow \infty} P \left(\frac{S_n - E(S_n)}{\sqrt{\text{var}(S_n)}} \leq x \right) = \Phi(x),$$

kjer je Φ porazdelitvena funkcija standardizirano normalne porazdelitve.

Preden se lotimo formalnega dokaza, dodajmo še nekaj komentarjev.

1. Odštevanje $E(S_n)$ in deljenje z $\sqrt{\text{var}(S_n)}$ ima za posledico, da ima kvocient, ki nastopa v formulaciji izreka, pričakovano vrednost 0 in varianco 1. Temu se pogosto reče, da smo vsoto S_n standardizirali.
2. Iz Izreka 1.10 sledi, da za $\alpha < \beta$ velja

$$\lim_{n \rightarrow \infty} P \left(\alpha \leq \frac{S_n - E(S_n)}{\sqrt{\text{var}(S_n)}} \leq \beta \right) = \Phi(\beta) - \Phi(\alpha).$$

3. Izrek 1.10 v vsej splošnosti dokažemo s pomočjo karakterističnih funkcij, ki so za potrebe verjetnosti prilagojene verzije Fourierove transformacije. Elementarni dokaz, ki ga bomo navedli, zahteva dodatno predpostavko, da je $E(|X_1 - E(X_1)|^3) < \infty$. Ta predpostavka bo v vseh primerih praktične uporabe izpolnjena.

Dokaz centralnega limitnega izreka temelji na Lindeberg-Bergströmovi neenačbi.

IZREK 1.11: Naj bodo X_1, \dots, X_n med sabo neodvisne slučajne spremenljivke z $E(X_i) = 0$ in $\text{var}(X_1 + \dots + X_n) = 1$. Označimo $X := X_1 + \dots + X_n$ in

naj bo $Z \sim N(0, 1)$. Naj bo f trikrat zvezno odvedljiva funkcija na \mathbb{R} z $|f|, |f'|, |f''|, |f'''| \leq M < \infty$. Velja

$$|E(f(X)) - E(f(Z))| \leq \frac{1}{6} \left(1 + \sqrt{\frac{8}{\pi}} \right) M (E(|X_1|^3) + \dots + E(|X_n|^3)).$$

DOKAZ: Naj bodo Z_1, \dots, Z_n porazdeljene normalno z istimi pričakovanimi vrednostmi in variancami kot X_i ter neodvisne med seboj in od X_1, \dots, X_n . Taka izbira slučajnih spremenljivk je vedno možna. Privzamemo lahko, da je $Z = Z_1 + \dots + Z_n$. Zapišimo

$$E(f(X)) - E(f(Z)) = a_1 + \dots + a_n,$$

kjer je

$$a_i = E[f(X_1 + \dots + X_i + Z_{i+1} + \dots + Z_n)] - E[f(X_1 + \dots + X_{i-1} + Z_i + \dots + Z_n)].$$

V definiciji a_1 je v odštevanju vsota vseh Z_i , za $i = n$ pa je v prvem členu v a_n vsota vseh X_i . Funkcijo f lahko razvijemo po Taylorju. Za vsak $x \in \mathbb{R}$ velja

$$f(x+h) = f(x) + f'(x)h + \frac{1}{2}f''(x)h^2 + r$$

kjer je $|r| \leq \frac{1}{6}M|h|^3$. Če zdaj oba člena v a_i razvijemo okoli:

$$Y_i := X_1 + \dots + X_{i-1} + Z_{i+1} + \dots + Z_n$$

dobimo

$$a_i = E[f'(Y_i)(X_i - Z_i)] + \frac{1}{2}E[f''(Y_i)(X_i^2 - Z_i^2)] + r_i.$$

kjer je

$$|r_i| \leq \frac{1}{6}M(E(|X_i|^3) + E(|Z_i|^3))$$

Zaradi neodvisnosti X_i, Y_i, Z_i in predpostavk o enakosti pričakovanih vrednosti in varianc, sta pričakovani vrednosti v razvoju za a_i enaki 0, zato je $a_i = r_i$. Z uporabo Jensenove neenakosti za konveksno funkcijo $\phi(x) = x^{3/2}$ na $(0, \infty)$ in elementarnim dejstvom, da je

$$E(|Z_i|^3) = \sqrt{\frac{8}{\pi}} \text{var}(Z_i)^{3/2},$$

ki ga preverimo z integriranjem, lahko ocenimo

$$E(|Z_i|^3) = \sqrt{\frac{8}{\pi}} \text{var}(Z_i)^{3/2} = \sqrt{\frac{8}{\pi}} E(|X_i|^2)^{3/2} \leq \sqrt{\frac{8}{\pi}} E(|X_i|^3).$$

Končno lahko ocenimo

$$|E(f(X)) - E(f(Z))| \leq \frac{1}{6} \left(1 + \sqrt{\frac{8}{\pi}}\right) M(E(|X_1|^3) + \dots + E(|X_n|^3))$$

KOMENTAR: V Izreku 1.11 je dovolj privzeti samo $|f'''| \leq M < \infty$, vendar smo privzeli več, da ni treba utemeljevati obstoja pričakovanih vrednosti, ki nastopajo v dokazu izreka.

V posebnem primeru, ko so vse slučajne spremenljivke enako porazdeljene, nadomestimo X_i z

$$\tilde{X}_i = \frac{X_i - E(X_i)}{\sqrt{\text{var}(S_n)}}.$$

Te nove slučajne spremenljivke ustrezajo pogojem Izreka 1.11 in velja

$$\tilde{X}_1 + \dots + \tilde{X}_n = \frac{S_n - E(S_n)}{\sqrt{\text{var}(S_n)}}.$$

Vstavimo in zaradi $\text{var}(S_n) = n\text{var}(X_1)$ sledi

$$\left| E \left(f \left(\frac{S_n - E(S_n)}{\sqrt{\text{var}(S_n)}} \right) \right) - E(f(Z)) \right| \leq \frac{1}{6} \left(1 + \sqrt{\frac{8}{\pi}}\right) \frac{ME(|X_1 - E(X_1)|^3)}{\sqrt{n}\text{var}(X_1)^{3/2}}. \quad (1)$$

V zanimivem primeru, ko je $\text{var}(X_1) > 0$, gre desna stran zgornje neenačbe proti 0, ko $n \rightarrow \infty$.

Za dokaz centralnega limitnega izreka potrebujemo še nekaj Analize 1. Označimo z $\chi_{(-\infty, x]}$ karakteristično funkcijo intervala $(-\infty, x]$. Za poljubna $x \in \mathbb{R}$ in $\delta > 0$ obstajata trikrat zvezno odvedljivi funkciji $f_{-\delta}$ in f_{δ} z vrednostmi na $[0, 1]$, taki da veljajo neenakosti

$$\chi_{(-\infty, x-\delta]} \leq f_{-\delta} \leq \chi_{(-\infty, x]} \leq f_{\delta} \leq \chi_{(-\infty, x+\delta]}.$$

Ker sta $f_{\pm\delta}$ nekonstantni samo na $[x - \delta, x + \delta]$, imata zaradi zvezne odvedljivosti omejene prve tri odvode.

DOKAZ IZREKA 1.10: Izrek bomo, kot rečeno, dokazali pod dodatno predpostavko, da je $E(|X_i - E(X_i)|^3) < \infty$. Označimo

$$\tilde{S}_n = \frac{S_n - E(S_n)}{\sqrt{\text{var}(S_n)}}.$$

Fiksirajmo $x \in \mathbb{R}$. Naj bo $\epsilon > 0$. Ker je Φ zvezna, obstaja $\delta > 0$, da iz $|x - y|$ sledi $|\Phi(x) - \Phi(y)| < \epsilon$. Ker je

$$E \left[\chi_{(-\infty, x]}(\tilde{S}_n) \right] = P(\tilde{S}_n \leq x),$$

dobimo neenačbi

$$E \left[f_{-\delta}(\tilde{S}_n) \right] \leq P(\tilde{S}_n \leq x) \leq E \left[f_{\delta}(\tilde{S}_n) \right].$$

Iz posledice (1) sledi, da je

$$\lim_{n \rightarrow \infty} E \left[f_{\pm\delta}(\tilde{S}_n) \right] = E \left[f_{\pm\delta}(Z) \right],$$

kjer je $Z \sim N(0, 1)$. Sledi

$$E \left[f_{-\delta}(Z) \right] \leq \liminf_{n \rightarrow \infty} P \left(\tilde{S}_n \leq x \right) \leq \limsup_{n \rightarrow \infty} P \left(\tilde{S}_n \leq x \right) \leq E \left[f_{\delta}(Z) \right].$$

Iz neenačb za funkciji $f_{\pm\delta}$ sledi še

$$E \left[f_{-\delta}(Z) \right] \geq P(Z \leq x - \delta) \quad \text{in} \quad E \left[f_{\delta}(Z) \right] \leq P(Z \leq x + \delta).$$

Neenačbe združimo v

$$\Phi(x - \delta) \leq \liminf_{n \rightarrow \infty} P \left(\tilde{S}_n \leq x \right) \leq \limsup_{n \rightarrow \infty} P \left(\tilde{S}_n \leq x \right) \leq \Phi(x + \delta).$$

Limsup in liminf se razlikujeta kvečjemu za 2ϵ . Ker je bil ϵ poljuben, limita obstaja in je enaka $\Phi(x)$.

Tipičen primer uporabe centralnega limitnega izreka je, da aproksimiramo verjetnosti $P(a \leq S_n \leq b)$, kjer je S_n vsota med sabo neodvisnih, enako porazdeljenih slučajnih spremenljivk. Aproksimacijo dobimo kot

$$\begin{aligned} P(a \leq S_n \leq b) &= P \left(\frac{a - E(S_n)}{\sqrt{\text{var}(S_n)}} \leq \frac{S_n - E(S_n)}{\sqrt{\text{var}(S_n)}} \leq \frac{b - E(S_n)}{\sqrt{\text{var}(S_n)}} \right) \\ &\approx \Phi \left(\frac{b - E(S_n)}{\sqrt{\text{var}(S_n)}} \right) - \Phi \left(\frac{a - E(S_n)}{\sqrt{\text{var}(S_n)}} \right). \end{aligned}$$

Za velike n torej aproksimiramo člen zaporedja z limito zaporedja. Če so slučajne spremenljivke X_1, X_2, \dots celoštevilske, lahko aproksimacijo še nekoliko izboljšamo s *korekcijo za zveznost*, kar pomeni, da a nadomestimo z $a - \frac{1}{2}$ in b z $b + \frac{1}{2}$.

Kot pri vsaki aproksimaciji se postavi vprašanje napake. Na to vprašanje odgovarja izrek, ki ga tukaj navajamo brez dokaza.¹

IZREK 1.12 (Berry-Esséen) Naj veljajo enake oznake in predpostavke kot v izreku 1.10. Označimo $\rho = E(|X_1 - E(X_1)|^3)$. Velja neenačba

$$\sup_{x \in \mathbb{R}} \left| P \left(\frac{S_n - E(S_n)}{\sqrt{\text{var}(S_n)}} \leq x \right) - \Phi(x) \right| \leq \frac{C\rho}{\sqrt{n}\text{var}(X_1)^{3/2}},$$

kjer je C univerzalna konstanta.

Najboljša do zdaj znana ocena za univerzalno konstanto je $C < 0,4748$.²

PRIMER: Vzemimo neodvisne X_1, X_2, \dots s $P(X_i = 0) = P(X_i = 1) = \frac{1}{2}$. Vemo, da je $S_n \sim \text{Bin}(n, \frac{1}{2})$. Velja, da je $E(S_n) = n/2$ in $\text{var}(S_n) = n/4$. Aproksimacija z normalno porazdelitvijo nam da za $n = 10.000$, da je

$$P(4950 \leq S_n \leq 5050) = P(-1 \leq \frac{S_n - E(S_n)}{\sqrt{\text{var}(S_n)}} \leq 1) \approx \Phi(1) - \Phi(-1).$$

Za funkcijo Φ ni analitičnega izraza, zato uporabimo ali tabele ali pa statistične programe, ki imajo to funkcijo vgrajeno. Program R nam da $\Phi(1) - \Phi(-1) = 0,6827$. Program R ima vgrajeno tudi binomsko porazdelitev, tako da je točna verjetnost $0,6875$. Če uporabimo še korekcijo za zveznost, kar pomeni da spodnji meji odštejemo $\frac{1}{2}$ in zgornji prištejemo $\frac{1}{2}$, je rezultat aproksimacije $0,6875$, ker je točno na štiri decimalke!

PRIMER: Naj bodo spet X_1, X_2, \dots neodvisne in tokrat naj bo $P(X_1 = 1) = P(X_1 = 2) = P(X_1 = 9) = \frac{1}{3}$. Vzemimo $n = 300$. Izračunamo, da je

$$E(S_{300}) = 300 \cdot 4 \quad \text{in} \quad \text{var}(S_{300}) = 300 \cdot \frac{38}{3}.$$

¹Dokaz je, recimo, v Y. S. Chow, H. Teicher, Probability Theory: Independence, Interchangeability, Martingales, 3rd Edition, Springer 2003.

²Shevtsova, I., On the accuracy of the normal approximation for sums of independent symmetric random variables. (Russian) Dokl. Akad. Nauk 443 (2012), no. 6, 671-676.

Aproksimiramo

$$\begin{aligned} & P(1100 \leq S_{300} \leq 1300) \\ &= P\left(\frac{1100 - 1200}{\sqrt{3800}} \leq \frac{S_n - E(S_n)}{\sqrt{\text{var}(S_n)}} \leq \frac{1300 - 1200}{\sqrt{3800}}\right) \\ &\approx \Phi(1,6222) - \Phi(-1,6222) \\ &= 0,8952. \end{aligned}$$

Izračun s hitro Fourierovo transformacijo da rezultat 0,8970. Če uporabimo še korekcijo za zveznost, je rezultat aproksimacije z normalno porazdelitvijo enak 0,8970!

Centralni limitni izrek je najbolj znan primer izreka, ki govori o konvergenci v porazdelitvi. Ideja je, da zaporedje slučajnih spremenljivk konvergira v porazdelitvi, če so njihove porazdelitve vedno bolj podobne neki končni porazdelitvi. Kot iztočnico za matematično definicijo vzemimo Izrek 1.11.

DEFINICIJA: Zaporedje slučajnih spremenljivk Y_1, Y_2, \dots konvergira v porazdelitvi proti slučajni spremenljivki Y , če za vsako omejeno zvezno funkcijo $f: \mathbb{R} \rightarrow \mathbb{R}$ velja

$$E[f(Y_n)] \rightarrow E[f(Y)],$$

ko $n \rightarrow \infty$. Matematična oznaka je

$$Y_n \xrightarrow{d} Y,$$

ko $n \rightarrow \infty$.

V definiciji vzamemo omejene funkcije, da ni težav z obstojem pričakovanih vrednosti. Zvezne funkcije pa izberemo zato, da v porazdelitvi konvergirajo zaporedja konvergentnih konstantnih slučajnih spremenljivk. Da si nekoliko bolje predstavimo, kaj pomeni zgornja definicija, dokažimo naslednji izrek.

IZREK 1.13: Naj bo $(Y_n)_{n \geq 1}$ zaporedje slučajnih spremenljivk. Zaporedje konvergira v porazdelitvi proti slučajni spremenljivki Y , če in samo če za vsak $x \in \mathbb{R}$, v katerem je F_Y zvezna, velja

$$F_{Y_n}(x) \rightarrow F_Y(x),$$

ko $n \rightarrow \infty$.

DOKAZ: Naj bo F_Y zvezna v točki x . Naj bo $\epsilon > 0$. Izberimo $\delta > 0$, tako da bo za $|x - y| < \delta$ veljalo $|F_Y(x) - F_Y(y)| < \epsilon$. Vzemimo enaki funkciji $f_{\pm\delta}$ kot v Izreku 1.11. Ker sta funkciji zvezni, iz predpostavke sledi

$$E(f_{-\delta}(Y)) \leq \liminf_{n \rightarrow \infty} P(Y_n \leq x) \leq \limsup_{n \rightarrow \infty} P(Y_n \leq x) \leq E(f_{\delta}(Y)).$$

Nadalje iz izbire funkcij sledi

$$P(Y \leq x - \delta) \leq \liminf_{n \rightarrow \infty} P(Y_n \leq x) \leq \limsup_{n \rightarrow \infty} P(Y_n \leq x) \leq P(Y \leq x + \delta).$$

Ker je bil ϵ poljuben, iz zveznosti porazdelitvene funkcije v x sledi $F_{Y_n}(x) \rightarrow F_Y(x)$, ko $n \rightarrow \infty$.

Obratno predpostavimo, da $F_{Y_n}(x) \rightarrow F_Y(x)$, ko $n \rightarrow \infty$ za vsako točko zveznosti F_Y . Naj bo $\epsilon > 0$. Ker je F_Y nepadajoča, so nezveznosti lahko samo skoki, ki jih je največ števno mnogo. Naj bo f omejena zvezna funkcija. Označimo supremum te funkcije z M . Obstajata točki zveznosti $a < b$, da je $F_Y(a) < \epsilon$ in $F_Y(b) > 1 - \epsilon$. Zaradi zveznosti na $[a, b]$ lahko najdemo particijo $a = x_0 < x_1 < \dots < x_n = b$, tako da so vse točke v particiji točke zveznosti F_Y in funkcija f na $[x_k, x_{k+1}]$ variira za manj kot za ϵ . Funkcijo f aproksimirajmo s funkcijo

$$f^\epsilon(x) = \begin{cases} 0 & \text{za } x \leq a \\ f(x_k) & \text{za } x \in (x_{k-1}, x_k] \text{ za } k = 1, 2, \dots, n. \\ 0 & \text{za } x > b. \end{cases}$$

Velja

$$|E(f^\epsilon(Y)) - E(f(Y))| \leq \epsilon + 2M\epsilon,$$

kot se hitro prepričamo. Po drugi strani je

$$E[f^\epsilon(Y_n)] = \sum_{k=1}^n f(x_k) (F_{Y_n}(x_k) - F_{Y_n}(x_{k-1})).$$

Podobna vsota velja tudi za $E[f^\epsilon(Y)]$. Vsota je končna, zato po predpostavki obstaja dovolj velik N , da bo za $n \geq N$

$$|E[f^\epsilon(Y_n)] - E[f^\epsilon(Y)]| \leq \epsilon.$$

Po potrebi povečamo N , da bo za $n \geq N$ veljalo $F_{Y_n}(a) < \epsilon$ in $F_{Y_n}(b) > 1 - \epsilon$. Podobno kot za Y za $n \geq N$ velja

$$|E(f^\epsilon(Y_n)) - E(f(Y_n))| \leq \epsilon + 2M\epsilon,$$

Ko trikotniško združimo nenakosti, bo za $n \geq N$ veljalo

$$|E[f(Y_n)] - E[f(Y)]| \leq 3\epsilon + 4M\epsilon.$$

Ker je bil $\epsilon > 0$ poljuben, je izrek dokazan.

S tem izrekom lahko centralni limitni izrek, glede na to, da je Φ povesod zvezna, napišemo kot

$$\frac{S_n - E(S_n)}{\sqrt{\text{var}(S_n)}} \xrightarrow{d} Z,$$

kjer je $Z \sim N(0, 1)$. Definicija konvergence v porazdelitvi pa ima tudi druge prednosti, na samo da lahko elegantno povemo centralni limitni izrek.

Splošnejšo definicijo konvergence v porazdelitvi lahko uporabimo *mutatis mutandis* za vektorje.

DEFINICIJA: Zaporedje slučajnih vektorjev $\mathbf{Y}_1, \mathbf{Y}_2, \dots$ dimenzije r konvergi-
ra v porazdelitvi proti slučajnemu vektorju \mathbf{Y} , če za vsako omejeno zvezno
funkcijo $f: \mathbb{R}^r \rightarrow \mathbb{R}$ velja

$$E[f(\mathbf{Y}_n)] \rightarrow E[f(\mathbf{Y})],$$

ko $n \rightarrow \infty$. Matematična oznaka je

$$\mathbf{Y}_n \xrightarrow{d} \mathbf{Y},$$

ko $n \rightarrow \infty$.

Centralni limitni izrek ima tudi vektorsko varianto, ki jo tukaj samo nava-
jamo. Dokaz z milimi dodatnimi predpostavkami je podoben dokazu v eni
dimenziji, le pisanja in ocenjevanja je nekaj več.

IZREK 1.14: Naj bodo $\mathbf{Y}_1, \mathbf{Y}_2, \dots$ med sabo neodvisni, enako porazdeljeni
slučajni vektorji. Predpostavimo, da obstajata $E(\mathbf{Y}_1)$ in $\Sigma = \text{var}(\mathbf{Y}_1)$.
Označimo $\mathbf{S}_n = \mathbf{Y}_1 + \dots + \mathbf{Y}_n$. Velja

$$\frac{\mathbf{S}_n - E(\mathbf{S}_n)}{\sqrt{n}} \xrightarrow{d} N(0, \Sigma),$$

ko $n \rightarrow \infty$.

Za namene statistike je pomembna naslednja posledica definicij.

IZREK 1.15: Naj $\mathbf{Y}_n \xrightarrow{d} \mathbf{Y}$, ko $n \rightarrow \infty$, kjer so vsi vektorji r -dimenzionalni. Naj bo $\phi: \mathbb{R}^r \rightarrow \mathbb{R}$ zvezna funkcija. Potem

$$\phi(\mathbf{Y}_n) \xrightarrow{d} \phi(\mathbf{Y}),$$

ko $n \rightarrow \infty$.

DOKAZ: Dokaz sledi iz preproste opazke, da je za vsako omejeno zvezno funkcijo f kompozitum $f \circ \phi$ omejena zvezna funkcija in trditev sledi iz definicije konvergence v porazdelitvi za vektorje.

PRIMER: Multinomska porazdelitev nastane pri metanju kroglic v škatle. Recimo, da je škatel r , meti so med sabo neodvisni, i -to škatlo pa zadenemo z verjetnostjo p_i . Označimo z $N_i(n)$ število kroglic v škatli i po n metih. Definirajmo slučajni vektor $\mathbf{X}^{(k)}$ po komponentah s predpisom

$$\mathbf{X}_i^{(k)} = \begin{cases} 1 & \text{če } k\text{-ta kroglica pade v } i\text{-to škatlo} \\ 0 & \text{sicer.} \end{cases}$$

Zaradi neodvisnosti metov so tudi vektorji $\mathbf{X}^{(1)}, \mathbf{X}^{(2)}, \dots$ med seboj neodvisni in enako porazdeljeni. S temi definicijami je

$$\begin{pmatrix} N_1(n) \\ N_2(n) \\ \vdots \\ N_r(n) \end{pmatrix} = \mathbf{X}^{(1)} + \mathbf{X}^{(2)} + \dots + \mathbf{X}^{(n)}.$$

Iz definicij tudi sledi, da je

$$E(\mathbf{X}^{(1)}) = \mathbf{p},$$

kjer označimo $\mathbf{p} = (p_1, \dots, p_r)^T$. Ker je $\mathbf{X}_i^{(1)}$ indikatorska slučajna spremenljivka, je

$$\text{var}(\mathbf{X}_i^{(1)}) = p_i(1 - p_i).$$

Opazimo, da za $i \neq j$ velja $\mathbf{X}_i^{(1)}\mathbf{X}_j^{(1)} = 0$, zato je

$$\text{cov}(\mathbf{X}_i^{(1)}, \mathbf{X}_j^{(1)}) = -E(\mathbf{X}_i^{(1)})E(\mathbf{X}_j^{(1)}) = -p_i p_j.$$

Sledi

$$\Sigma = \text{var}(\mathbf{X}^{(1)}) = \text{diag}(p_1, p_2, \dots, p_r) - \mathbf{p}\mathbf{p}^T.$$

Zapišimo

$$\begin{pmatrix} \frac{N_1(n) - np_1}{\sqrt{np_1}} \\ \frac{N_2(n) - np_2}{\sqrt{np_2}} \\ \vdots \\ \frac{N_r(n) - np_r}{\sqrt{np_r}} \end{pmatrix} = \frac{\mathbf{A}(\mathbf{X}^{(1)} + \mathbf{X}^{(2)} + \dots + \mathbf{X}^{(n)} - n\mathbf{p})}{\sqrt{n}},$$

kjer je

$$\mathbf{A} = \text{diag}\left(\frac{1}{\sqrt{p_1}}, \dots, \frac{1}{\sqrt{p_r}}\right).$$

Po izrekih 1.14 in 1.15 velja, da

$$\begin{pmatrix} \frac{N_1(n) - np_1}{\sqrt{np_1}} \\ \frac{N_2(n) - np_2}{\sqrt{np_2}} \\ \vdots \\ \frac{N_r(n) - np_r}{\sqrt{np_r}} \end{pmatrix} \xrightarrow{d} \mathbf{A}\mathbf{X},$$

ko $n \rightarrow \infty$, kjer je $\mathbf{X} \sim N(\mathbf{0}, \Sigma)$. Porazdelitev vektorja $\mathbf{A}\mathbf{X}$ je tudi večrazsežna normalna in sicer $N(\mathbf{0}, \mathbf{A}\Sigma\mathbf{A}^T)$. Z množenjem matrik sledi

$$\mathbf{A}\Sigma\mathbf{A}^T = \mathbf{I} - \mathbf{q}\mathbf{q}^T,$$

kjer je

$$\mathbf{q} = (\sqrt{p_1}, \dots, \sqrt{p_r})^T.$$

Norma vektorja \mathbf{q} je enaka 1, zato je $\mathbf{q}\mathbf{q}^T$ idempotentna matrika in s tem tudi $\mathbf{I} - \mathbf{q}\mathbf{q}^T$ idempotentna matrika z rangom $r - 1$.

Če je $\mathbf{Z} \sim N(\mathbf{0}, \mathbf{I})$ in je \mathbf{H} idempotentna matrika, je $\mathbf{W} = \mathbf{H}\mathbf{Z} \sim N(\mathbf{0}, \mathbf{H})$. Sledi, da je

$$W_1^2 + \dots + W_r^2 = \mathbf{W}^T\mathbf{W} = (\mathbf{H}\mathbf{Z})^T\mathbf{H}\mathbf{Z} = \mathbf{Z}^T\mathbf{H}\mathbf{Z} \sim \chi^2(\text{rang}(\mathbf{H})).$$

Sklepamo lahko, da velja

$$\chi^2(n) = \sum_{i=1}^r \frac{(N_i(n) - np_i)^2}{np_i} \xrightarrow{d} U \sim \chi^2(r - 1),$$

ko $n \rightarrow \infty$. Ta limitni rezultat ima široko uporabo v statistiki.

2. VZORČENJE

2.1 Uvodni primeri

~~4.1~~ ~~Uvodni primeri~~

~~4.1.1~~ PLEBISCIT 1990

Plebiscit o neodvisnosti Slovenije decembra 1990 je bil prelomni dogodek pri osamosvajanju. Ko je takrat jeseni naraščala napetost, so mnogi nestrpno pričakovali rezultate predplebiscitnih anket. Ena od odmevnejših je bila SJM 90 (Slovensko javno mnenje 90), ki so jo izpeljali na Fakulteti za družbene vede na Univerzi v Ljubljani.

Za kaj pravzaprav gre pri taki anketi? Pred plebiscitom je nemogoče ugotoviti mnenje vsakega volivca, zato se je treba zateči k izbiranju manjšega števila enot iz populacije volivcev. Izbrani skupini v statističnem žargonu pravimo *vzorec*. Anketa SJM 90 je bila zasnovana na izbiri vzorca velikosti 2074 volivcev, od katerih se jih je 1306 nedvoumno izreklo tako za samostojnost kot za odcepitev Slovenije. Natančni rezultati ankete so v spodnji tabeli ¹.

		SAMOSTOJNOST		
		DA	NE	DRUGO
ODCEPITEV	DA	1306	11	34
	NE	183	125	63
	DRUGO	110	12	230

Tabela 4.1: Tabela rezultatov SJM90

Podatki iz vzorca kažejo, da se je velika večina volivcev iz vzorca izrekla za neodvisnost Slovenije. Iz tega bi sklepali, da se bo tudi večina vseh volilcev v Sloveniji odločila za neodvisnost. Ta sklep je sedaj, ko je rezultat plebiscita že zdavnaj znan, očiten. Poskusimo pa se postaviti v čas jeseni 1990. Vzorec je zajel samo neznamenit delež celotnega volilnega telesa, komajda nekaj več kot 0,1% vseh upravičencev. Sceptiki bi se gotovo vprašali, ali lahko na podlagi tako neznatnega vzorca zanesljivo sklepamo o

¹Vir: SJM90

volji celotne populacije volivcev. Lahko si zamislimo, da bi se v vzorcu znašlo ali preveč privržencev neodvisnosti ali preveč nasprotnikov. V takem primeru bi bila napoved izida seveda napačna. Ali nam statistika lahko pomaga, da presodimo zanesljivost napovedi na podlagi vzorčnih podatkov? Odgovor je pritrdilen.

Prvi korak pri razmišljanju o zanesljivosti ocen na podlagi vzorca mora biti opis načina izbire vzorca ali, kot se temu pravi v statistiki, vzorčni načrt. Pri ustrezno izpeljanih anketah je postopek izbire vzorca natanko predpisan. Poglejmo si način vzorčenja pri anketah SJM, ki je bil uporabljen tudi za anketo o plebiscitu.

Okvir vzorčenja pri anketah SJM je centralni register prebivalstva, ki ga vzdržuje Statistični urad RS. Gre preprosto za primerno urejen seznam vseh prebivalcev Slovenije, iz katerega zlahka dobimo tudi seznam vseh volivcev. Ta seznam je za potrebe vzorčenja po geografskem ključu razdeljen na manjše dele po 4200 volivcev. Tem skupinam pravimo *primarne vzorčne enote*. Vsaka od teh manjših skupin po 4200 volivcev je nadalje razdeljena na skupine po 100, spet po takem ključu, da volivci v podskupini živijo v skupnosti, kot je naselje ali vas. Tem manjšim skupinam pravimo *sekundarne vzorčne enote*. Izbira vzorca SJM poteka v treh korakih:

1. Na prvem koraku anketarji naključno izberejo 140 primarnih vzorčnih enot, torej 140 skupin po 4200 volivcev.
2. Na drugem koraku so v vsaki na prvem koraku izbrani primarni vzorčni enoti izbrane tri sekundarne vzorčne enote, torej tri skupine po 100 volivcev.
3. Na tretjem koraku nato anketarji v vsaki izbrani sekundarni vzorčni enoti naključno izberejo 5 volivcev.

Če izračunamo, je na koncu v vzorec izbranih $140 \cdot 3 \cdot 5 = 2100$ volivcev. Nato je seveda treba z izbranimi volivci stopiti v stik in jih povprašati po njihovem mnenju. Izvajalci ankete na terenu imajo napotek, da od vsake izbrane osebe poskusijo dobiti odgovor. Če prvi poskus ni uspešen, anketarji poskušajo znova do največ petkrat, če je to potrebno. Če so vsi poskusi neuspešni, anketarji neznan odgovor obravnavajo kot manjkajoči podatek. Torej, način vzorčenja je vnaprej določen in anketarji se ga morajo držati. Drug pomemben dejavnik pri vzorčenju za SJM je, da je izbira enot na vsakem koraku naključna. Ni vnaprej določeno, katere skupine ali podskupine ali

nazadnje volivci bodo izbrani, temveč so vse cnote izbrane kot na "loteriji". Gotovo je na tem mestu utemeljeno vprašanje, zakaj je treba uvesti tovrstno naključno izbiro. Razlog je v odpravljanju vsakršne pristranosti pri izbiri vzorca. Nekako tako, kot da bi pred jemanjem vzorca kapljice tekočine iz steklenice le-to dobro pretresli. Vzorec iz dobro pretresene steklenice mnogo bolje pokaže njeno vsebino, kot pa če steklenice ne bi pretresli.

Kot bomo videli v nadaljevanju, način vzorčenja bistveno vpliva na zanesljivost ocen. Na podlagi opisa vzorčenja lahko domnevamo, da so bili rezultati ankete SJM 90 precej zanesljivi, saj je bila "steklenica zelo dobro pretresena". V naslednjem razdelku si bomo ogledali, kako natančneje opišemo, do kolikšne mere lahko verjamemo, da vzorčni rezultati zanesljivo odražajo stanje v celotni populaciji.

~~4.1.1~~ INDEKS CEN ŽIVLJENJSKIH POTREBŠČIN

Statistični urad Republike Slovenije vsak mesec objavi indeks cen življenjskih potrebščin, ki je eden od osnovnih indikatorjev inflacije. Osnovi za izračun tega indeksa sta:

1. Seznam najvažnejših predmetov in storitev gospodinske porabe in podatki o višini stroškov za posamezne postavke. Osnova za določitev obsega te porabe je anketa o porabi v gospodinjstvih, ki jo izvaja Statistični urad RS.
2. Povprečna porast cen na drobno za posamezne postavke s seznama.

Formula, ki jo Statistični urad RS uporablja za izračun tega indeksa, je

$$I = \frac{\sum_{i=1}^m (p_{1i}/p_{0i}) \cdot w_{0i}}{\sum_{i=1}^m w_{0i}} \cdot 100\%.$$

Pri tem je kvocient p_{1i}/p_{0i} porast cene na drobno za predmet ali storitev i s seznama, utež w_{0i} pa je povprečni delež stroškov, ki je namenjen za dani predmet ali storitev glede na celotne izdatke v gospodinjstvu. Za primer navedimo, da gospodinjstva v Sloveniji za hrano v povprečju namenijo 23,1%² celotne mesečne porabe.

²Vir: Statistični letopis 1996, str. 234

Zgornji formuli ne bomo posvetili posebne pozornosti. Omenimo le, da so uteži potrebne zato, ker morajo imeti predmeti ali storitve s seznama, za katere gospodinjstva namenijo večji delež izdatkov, pri določanju rasti življenjskih stroškov večjo težo. Podražitev hrane ima na življenjsko raven večji vpliv kot, recimo, zvišanje cen frizerskih storitev.

Če želimo zgornjo formulo uporabiti, moramo priti na dan z dejanskimi številkami. Med drugim moramo za posamezne predmete ali storitve s seznama ugotoviti povprečne deleže porabe v gospodinjstvih. Ker bi bilo težko spremljati porabo v vseh gospodinjstvih, si Statistični urad RS pomaga z vzorčenjem. Gospodinjstva v Sloveniji so popisana in urejena v popisne okoliše. Za potrebe izbire vzorca so popisni okoliši razdeljeni v 6 podskupin glede na lokacijo in tip. Te podskupine v statistiki imenujejo *stratumi*.

Vzorčenje poteka v dveh korakih:

1. Najprej v vsakem stratumu naključno izberemo popisne okoliše.
2. Na drugem koraku v vsakem izbranem popisnem okolišu naključno izberemo 5 gospodinjstev.

V vzorec je zajetih 3270 gospodinjstev. Število popisnih okolišev, izbranih na prvem koraku, je tako, da v vsakem stratumu na koncu postopka izbire zajamemo 0,5% gospodinjstev. Statistični urad RS izbranim gospodinjstvom razdeli vprašalnike za vodenje evidence o porabi. Anketarji urada zberejo izpolnjene evidenčne vprašalnike in na podlagi tako zbranih podatkov določijo deleže porabe za posamezne predmete ali storitve. Prej omenjenih 23,1% stroškov za hrano je ena od tako dobljenih ocen.

Seveda se moramo vprašati, do kolikšne mere lahko ocenam na podlagi vzorca zaupamo. Odgovor je odvisen od tega, kako dobro je bil vzorčni načrt premišljen in ali je samo vzorčenje bilo izvedeno z nadzorom. Element naključnosti je bistvena sestavina vzorčenja, saj nam izkušnje iz preteklosti in teoretična razmišljanja kažejo, da tako najbolje dosežemo nepristranost vzorčnih ocen, poleg tega pa edino strog vzorčni načrt omogoča presojo o tem, kako zanesljiva je ocena. Brez naključne izbire vzorca taka presoja ni možna. Velikost vzorca je izbrana tako, da so dobljeni vzorčni odstotki dovolj zanesljivi za uporabo v formuli za izračun indeksa cen življenjskih potrebščin.

Rezultat zgornje formule za december 1996 je bil 1,3%.

Naslednji primer, kjer zbiranje podatkov poteka z vzorčenjem, so raziskave uspešnosti šolskih sistemov. Učenci slovenskih osnovnih in srednjih šol so vključeni v mednarodne primerjalne raziskave znanja matematike, naravoslovja, računalništva, bralne pismenosti in drugih področij znanja. V Sloveniji izvaja te primerjalne raziskave Pedagoški inštitut v Ljubljani. Za opis postopka vzorčenja v teh raziskavah si izberimo zadnjo izmed raziskav TIMSS, ki je potekala od leta 1991 in je bila končana leta 1998. To raziskavo smo že srečali v prejšnjih poglavjih.

V jeziku prvega poglavja spadajo v populacijo, ki nas v tem primeru zanima, vsi učenci sedmih in osmih razredov slovenskih šol. V letu 1995, ko je potekalo dejansko vzorčenje, je bila ta populacija velika $N = 54.965$ učencev³. Iz praktičnih razlogov je bilo izključenih 310 učencev (večinoma učenci iz šol za slepe in slabovidne ter podobnih), tako da je bila na koncu populacija, iz katere je bil izbran vzorec, velika $N = 54.655$.

Zanima nas znanje matematike in naravoslovja pri učencih iz opisane populacije. Znanje je bilo merjeno z nalogami, ki so bile sestavljene posebej za TIMSS. Če se spet vrnemo k jeziku prvega poglavja, sta spremenljivki raven znanja matematike in raven znanja naravoslovja posameznega učenca. Ta ravni izrazimo kot število na posebnih lestvicah, ki so premišljene tako, da omogočajo smiselno primerjavo med državami, ki so bile vključene v raziskavo. Za nas je v tem trenutku pomembno to, da vsaki enoti v populaciji pripadeta neki vrednosti, ki pomenita znanje matematike oziroma naravoslovja.

Preden se lotimo opisa vzorčenja, se moramo seveda vprašati, zakaj je vzorčenje sploh potrebno. Populacija je dokaj velika, zato si ni težko predstavljati, koliko dela bi bilo z izvedbo preizkušanja znanja, ki bi zajela vse učence v populaciji, da o vrstoglavih stroških tako velikega projekta sploh ne govorimo. Take raziskave za celotno populacijo učencev ni mogoče izpeljati iz povsem praktičnih razlogov. Rešitev je v vzorčenju. Iz celotne populacije izberemo manjši vzorec učencev in povprečno raven znanja slovenskih sedmošolcev in osmošolcev ocenimo na podlagi znanja učencev iz tega vzorca. V raziskavi TIMSS je bilo v vzorec zajetih 5927 učencev.

³Vir: Pedagoški inštitut Ljubljana

Samo izbiranje vzorca je bilo mednarodno usklajeno in vnaprej predpisano. Vzorčenje ni potekalo z neposrednim izbiranjem učencev, temveč posredno v dveh korakih. Na prvem koraku so sodelavci Pedagoškega inštituta med 455 osnovnimi šolami v Sloveniji naključno izbrali 150 šol, in sicer tako, da so imele večje šole nekaj več verjetnosti, da bodo izbrane v vzorec. Tukaj besedo "naključno" uporabljamo še v nekoliko ohlapnem pomenu, kasneje pa bomo ta pojem tudi natančneje opredelili. Ko so bile šole izbrane, je bil naslednji korak izbiranje razreda. Izbira je bila spet naključna, in sicer taka, da je bil izbran po en 7. in en 8. razred na vsaki šoli, ki je bila izbrana na prvem koraku. V končnem vzorcu so bili zajeti vsi učenci iz izbranih razredov, torej zgoraj omenjenih 5927 učencev. Takemu načinu izbiranja vzorca statistiki pravijo "vzorčenje v skupinicah", ker izbiramo celotne skupinice enot, v tem primeru razrede. Vsi izbrani učenci so reševali izbrane naloge in na podlagi njihovih odgovorov je bila ocenjena povprečna raven znanja matematike oziroma naravoslovja vseh slovenskih sedmošolcev in osmošolcev.

Tudi tukaj si moramo zastaviti vprašanje o zanesljivosti dobljenih ocen. Raziskava je zajela samo okrog 10% vseh učencev iz opisane populacije in na podlagi njihovih rezultatov je bila potem ocenjena raven znanja za celotno populacijo. Do kolikšne mere je to utemeljeno? Kako opisati zanesljivost ocen? Na to vprašanje bomo odgovorili v razdelkih, ki sledijo.

~~4.1.1~~ PREDSEDNIŠKE VOLITVE V ZDA LETA 1936

Kot zadnji primer vzorčenja si oglejmo znamenito anketo, kjer so šle stvari pri napovedovanju izida volitev zelo hudo narobe. Pred predsedniškimi volitvami v ZDA leta 1936 je prestižna revija *Literary Digest* na podlagi obsežne ankete napovedala zmagovalca. Danes vemo, da je bil na omenjenih volitvah s precejšnjo večino izvoljen Franklin D. Roosevelt, revija pa je takrat napovedala, da bo zmagal njegov republikanski tekmeč Alfred Landon. V tabeli 4.2 so vsebovane napovedi revije *Literary Digest* in dejanski rezultati. Jasno je, da anketa ni mogla zajeti vseh volivcev v ZDA leta 1936, zato je bilo treba izbrati vzorec. Vzorec, ki so ga izbrali anketarji revije *Literary Digest*, je bil največji, kar so jih sploh kdaj izbrali, in je vseboval 2,4 milijona oseb. Kljub tako velikemu vzorcu pa se je napoved razlikovala od dejanskega rezultata za skoraj 20%!

	Napoved revije <i>Literary Digest</i>	Dejanski rezultat
F. D. Roosevelt	43%	62%
A. Landon	57%	38%

Tabela 4.2: Napovedi in rezultati predsedniških volitev v ZDA leta 1936

Zanimivo je seveda vprašanje, kaj je povzročilo tako zelo zgrešeno napoved. Če želimo dobiti odgovor na to vprašanje, si moramo ogledati, kako je bil vzorec izbran. Izvor za izbiro volivcev v vzorec so bili sezname naročnikov revije *Literary Digest*, telefonski imeniki, sezname članov elitnih klubov in podobno. Vprašalnike so poslali izbranim osebam po pošti, in sicer kar 10 milijonom, odgovorilo pa je samo 2,4 milijona naslovnikov, kar je že razlog za previdnost. Če pomislimo, da je bilo leto 1936 najbolj črno leto velike depresije in je bilo v ZDA 11 milijonov brezposelnih, nam takoj pade v oči, da od le-teh velika večina ni imela telefona in torej njihova imena niso bila v telefonskem imeniku ali na seznamu članov kakšnega elitnega kluba. Če še pomislimo, da republikansko stranko v ZDA po pravilu podpirajo bogatejši sloji, je eden od razlogov za slabe napovedi na dlani. Anketa je bila že vnaprej načrtovana tako, da so bili v vzorec zajeti tisti, ki so tudi med depresijo imeli kaj pod palcem. Velikost vzorca seveda ni pomagala, ker se je samo ponavljala ena in ista napaka. V vzorcu se je vedno znova znašlo več republikanskih volivcev kot pa podpornikov predsednika Roosevelta. V statističnem žargonu bi lahko rekli, da je bil vzorčni načrt pri tej anketi slab, z drugimi besedami, način izbire vzorca je bil slabo premišljen. Še bolj primerna izjava bi bila, da vzorčnega načrta sploh ni bilo. Revija *Literary Digest* se je kmalu po volitvah 1936 znašla v stečaju.

Kot zanimivost lahko povemo še to, da je v istem času mladi statistik George Gallup na podlagi vzorca velikosti 5000 oseb napovedal zmago F. D. Roosevelta s 56% glasov. Še bolj zanimivo je to, da je George Gallup napovedal tudi napoved revije *Literary Digest*, še preden jo je ta objavila. Izbral je vzorec velikosti 3000 izmed tistih, ki so prejeli vprašalnike, in na podlagi izbranega vzorca napovedal, da bo napoved 44% za Roosevelta. Res ne bi bilo treba izbirati vzorca velikosti 2,4 milijona!



Pogosto je nemogoče zbrati podatke o vrednostih spremenljivke za vse enote v populaciji. Zato iz populacije izberemo vzorec, ki zajema le njen manjši del. Iz vrednosti spremenljivke za enote v vzorcu potem ocenimo želene količine, na primer povprečno vrednost spremenljivke ali odstotek enot z določeno lastnostjo, za celotno populacijo. Če nas zanima povprečna vrednost spremenljivke za celotno populacijo, kot oceno za to količino vzamemo povprečno vrednost spremenljivke za enote v vzorcu. Če nas zanima odstotek enot v populaciji z neko lastnostjo, za oceno tega odstotka vzamemo odstotek enot v vzorcu, ki imajo izbrano lastnost.



Vzorčni načrt je vnaprej predpisan postopek izbiranja vzorca iz vnaprej določene in natančno opredeljene populacije. Če je izbira enot ali skupin enot naključna, potem takemu vzorčenju pravimo *verjetnostno vzorčenje*. Z vpeljavo naključnosti se najbolje izognemo pristranosti.

2.1.422 ENOSTAVNO SLUČAJNO VZORČENJE

4221 POJEM ENOSTAVNEGA SLUČAJNEGA VZORCA

Iz primerov v prejšnjem razdelku, posebej iz zadnjega, je razvidno, da je način izbire vzorca zelo pomemben. V tem razdelku si bomo ogledali najpreprostejši vzorčni načrt: *enostavno slučajno vzorčenje*. Čeprav se ta tip vzorčenja v dejanskih anketah redko uporablja, je sestavni del mnogih, tudi bolj zapletenih vzorčnih načrtov, poleg tega pa je pri njem najbolj razviden pojem standardne napake. Za okvir razmišljanja vzemimo populacijo velikosti N , iz katere želimo izbrati vzorec velikosti n . Zamislimo

mi, če imamo za vsako od N
enot v populaciji listek, mi ga
olamo v škatlo.

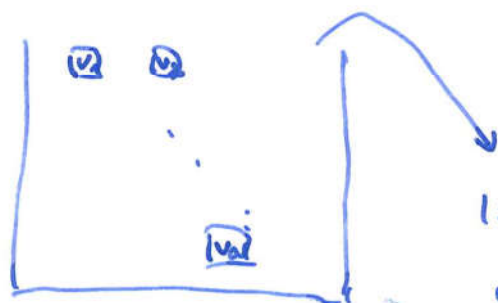
Preda izberemo enostavno slučajni vzorec, ki lahko predstavljamo, da imamo v katali listke, na katerih je zapisana vrednost. Če je listkov N , označimo vrednosti na njih z v_1, v_2, \dots, v_N . Definiramo populacijsko povprečje z

$$\mu = \sum_{k=1}^N v_k$$

in populacijsko varianco z

$$\sigma^2 = \frac{1}{N} \sum_{k=1}^N (v_k - \mu)^2.$$

Slika:



izberemo naključno
 v enot. Populacijsko

povprečje ocenimo z
 vzorčnim
 povprečjem.

Glavna ideja: ko izberemo vzorec in izračunamo vzorčno povprečje, smo kreivali stohastično spremenljivko \bar{X} .

Definirajmo

$$I_k = \begin{cases} 1, & \text{če je } k\text{-ta enota izbrana} \\ & \text{vzorec.} \\ 0, & \text{ničev.} \end{cases}$$

Velja $I_1 + \dots + I_N = n$. Zapišemo lahko

$$\bar{X} = \frac{1}{n} \sum_{k=1}^N v_k \cdot I_k$$

Vemo

$$\begin{aligned} E(\bar{X}) &= \frac{1}{n} \sum_{k=1}^N v_k \cdot E(I_k) \\ &= \frac{1}{n} \sum_{k=1}^N v_k P(I_k=1) \\ &= \frac{1}{n} \sum_{k=1}^N v_k \cdot \frac{n}{N} \\ &= \frac{1}{N} \sum_{k=1}^N v_k \\ &= \mu \end{aligned}$$

Zavada simetrije je $\text{cov}(I_k, I_l)$
enaka za vsa $k \neq l$. Računamo

$$\begin{aligned}\text{cov}(I_k, I_l) &= P(I_k=1, I_l=1) \\ &\quad - P(I_k=1)P(I_l=1) \\ &= \frac{\binom{N-2}{D-2}}{\binom{N}{n}} - \frac{n}{N} \cdot \frac{n}{N} \\ &= \frac{n(n-1)}{N(N-1)} - \frac{n^2}{N^2} \\ &= \frac{n}{N} \left[\frac{n-1}{N-1} - \frac{n}{N} \right] \\ &= \frac{n}{N} \left[\frac{(n-1)N - n(N-1)}{N(N-1)} \right] \\ &= \frac{n}{N} \left[\frac{-N+n}{N(N-1)} \right] \\ &= -\frac{n}{N} \cdot \frac{(N-n)}{N(N-1)}\end{aligned}$$

Velja

$$\text{var}(I_k) = \frac{n}{N} \cdot \left(1 - \frac{n}{N}\right) = \frac{n}{N} \cdot \frac{(N-n)}{N}$$

Racunamo

$$\text{var}(\bar{x})$$

$$= \frac{1}{n^2} \left[\sum_{k=1}^N v_k^2 \text{var}(I_k) + \sum_{\substack{k,l=1 \\ k \neq l}}^n v_k v_l \text{cov}(I_k, I_l) \right]$$

$$= \frac{1}{n^2} \left[\frac{n}{N} \cdot \frac{N-n}{2} \sum_{k=1}^N v_k^2 - \frac{n}{2} \cdot \frac{(N-n)}{N(N-1)} \sum_{\substack{k,l=1 \\ k \neq l}}^N v_k v_l \right]$$

$$= \frac{1}{n} \cdot \frac{N-n}{N^2} \cdot \left[\sum_{k=1}^N v_k^2 - \frac{1}{N-1} \sum_{\substack{k,l=1 \\ k \neq l}}^N v_k v_l \right]$$

$$= \frac{1}{n} \cdot \frac{N-n}{N^2} \left[\sum_{k=1}^N v_k^2 \left(1 + \frac{1}{N-1} - \frac{1}{N-1} \right) - \frac{1}{N-1} \sum_{\substack{k,l=1 \\ k \neq l}}^N v_k v_l \right]$$

$$= \frac{1}{n} \cdot \frac{N-n}{N^2} \left[\frac{N}{N-1} \sum_{k=1}^n v_k^2 - \frac{1}{N-1} \sum_{k=1}^n v_k^2 - \frac{1}{N-1} \sum_{\substack{k,l=1 \\ k \neq l}}^N v_k v_l \right]$$

Zavlačni nimeetrije je $E(X_k) = \mu$ in
 $\text{var}(X_k) = \sigma^2$ za vse $k = 1, \dots, n$.

Potrebujeemo še $\text{cov}(X_k, X_l)$.

Misljemo si belimo, da enote
izbiramo do zadnje. Pokem je

$X_1 + X_2 + \dots + X_n = \text{const.}$ in zato

$$\text{cov}(X_1, X_1 + X_2 + \dots + X_n) = 0$$

"

$$\text{cov}(X_1, X_1) + (n-1) \text{cov}(X_1, X_2) \Rightarrow$$

$$\text{cov}(X_1, X_2) = - \frac{\sigma^2}{n-1}$$

Sledi

$$\text{var}(\bar{X}) = \frac{1}{n^2} [n \text{var}(X_1) + n(n-1) \text{cov}(X_1, X_2)]$$

$$= \frac{1}{n} \left[\sigma^2 + (n-1) \left(- \frac{\sigma^2}{n-1} \right) \right]$$

$$= \frac{\sigma^2}{n} \left[1 + \left(- \frac{n-1}{n-1} \right) \right]$$

$$= \frac{\sigma^2}{n} \cdot \frac{n-1}{n-1}$$

$$= \frac{1}{n} \cdot \frac{N-n}{N^2} \left[\frac{N}{N-1} \sum_{k=1}^N v_k^2 - \frac{1}{N-1} \left(\sum_{k=1}^N v_k \right)^2 \right]$$

$$= \frac{1}{n} \cdot \frac{N-n}{N-1} \left[\frac{1}{N} \sum_{k=1}^N v_k^2 - \mu^2 \right]$$

$$= \frac{1}{n} \cdot \frac{N-n}{N-1} \cdot \sigma^2,$$

ker je $\sigma^2 = \frac{1}{N} \sum_{k=1}^N v_k^2 - \mu^2$.

Zračunajmo se nekoliko drugače.

Letno n predstavljamo, da izbiramo enote naključno eno po eno.

Vrednosti na izbranih listkih so slučajne spremenljivke X_1, X_2, \dots, X_n in

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

se. spr. X_1 ima vrednosti v_1, \dots, v_N + verjetnosti $\frac{1}{N}$, kar pomeni

$$E(X_1) = \mu \quad \text{in}$$

$$\text{var}(X_1) = \sigma^2.$$

Opomba: \bar{X} smo lahko napisali
na več načinov kot linearno
kombinacijo enostavnejših
stohastičnih spremenljivk.

Najpomembnejša opazka pri tem
drugem zapisu je, da je \bar{X}
vsota sh. spr. deljen z n . To
centralnem limitnem izreku bo
ta vsota približno normalno
porazdeljena. Torej porazdelitev \bar{X}
lahko aproksimiramo z normalno
porazdelitvijo!

Definicija: Porazdelitev \bar{X} večino
vsotina porazdelitev.

Definicija: Stohastični spremenljivki
 \bar{X} večino cenitka.

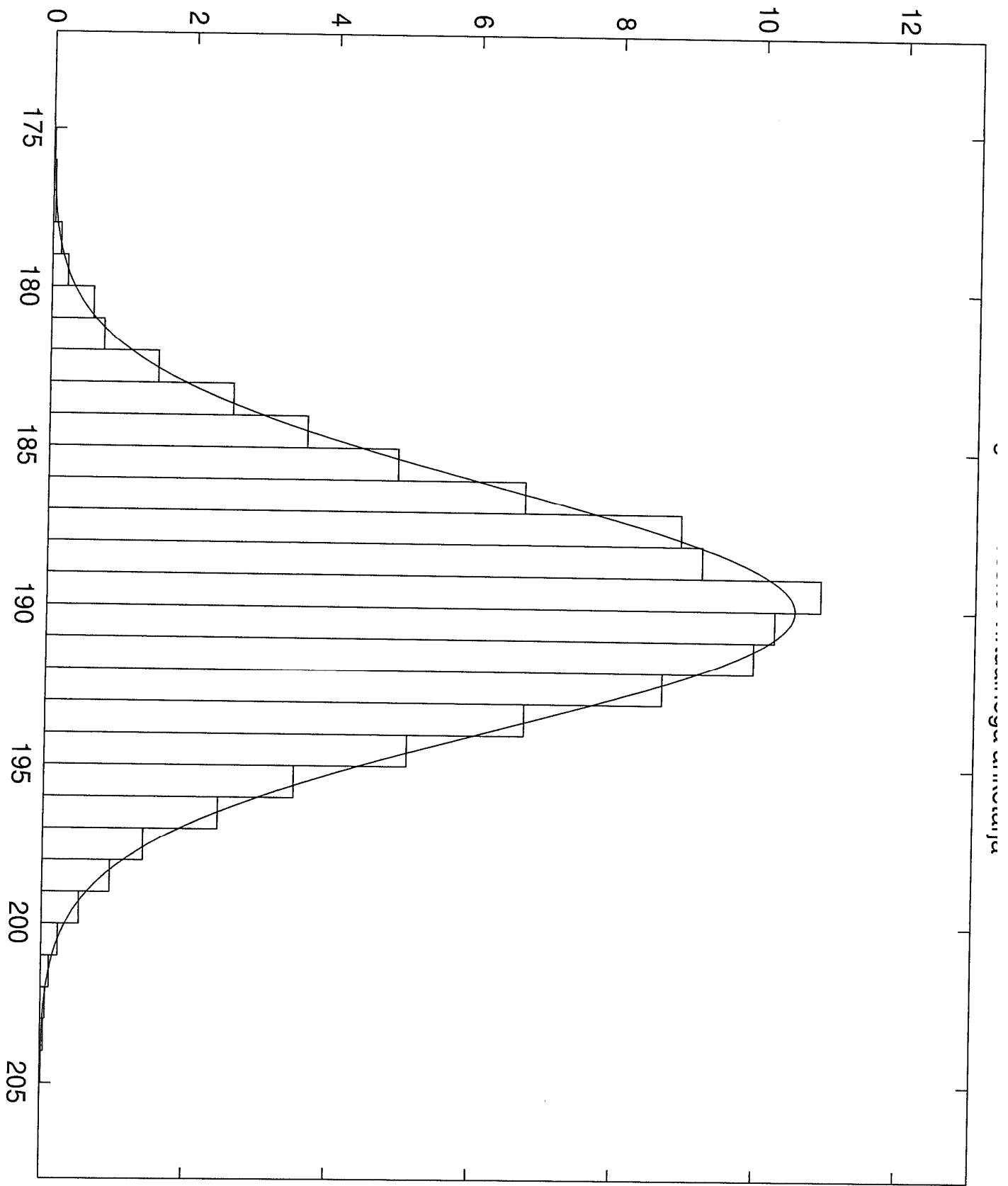
Zaduirek: Centralni limitni izrek
velja za neodvisne X_1, X_2, \dots

V vsoti $\frac{1}{n} \sum_{k=1}^n X_k$ so slučajne
spremenljivke kolektivne! Možna
sta dva odgovora:

(1) Odvisnost je šibka, dokler je
vzorec majhen delež populacije.
CLI še vedno dobro aproksimira
porazdelitev vsote.

(2) CLI velja tudi v primerih, ko
so slučajne spremenljivke odvisne,
vendar so formulacije in dokazi
bolj komplicirani.

SKLEP: Vzorecna porazdelitev \bar{X} je
približno normalna s parametroma
 μ in $\frac{\sigma^2}{n} \cdot \frac{N-n}{N-1}$. Na podlagi tega
lahko dajemo izjave o zaupljivosti
ocen.



Definicija: Cenilka je nepristrana
če je njena pričakovana vrednost
enaka parametru, ki ga ocenjujemo.

Pri vzorčju ne poznamo σ^2 .

Kako bi ta parameter ocenili?

V volkah imamo samo vzorec in
večemo

$$\hat{\sigma}^2 = \frac{1}{n} \sum_{k=1}^n (x_k - \bar{x})^2.$$

V statistiki s n. otuščimo cenilke

Rečunamo

$$\begin{aligned} E(\hat{\sigma}^2) &= \frac{1}{n} \sum_{k=1}^n E[(x_k - \bar{x})^2] \\ &= E[(x_1 - \bar{x})^2] \quad (\text{simetrija}) \\ &= \text{var}(x_1 - \bar{x}) \\ &= \text{var}(x_1) + \text{var}(\bar{x}) \\ &\quad - 2 \text{cov}(x_1, \frac{1}{n} \sum_{k=1}^n x_k) \end{aligned}$$

$$\begin{aligned}
&= \hat{\sigma}^2 + \frac{\hat{\sigma}^2}{n} \cdot \frac{N-n}{N-1} \\
&\quad - 2 \cdot \frac{1}{n} \left(\hat{\sigma}^2 - \frac{\hat{\sigma}^2 (n-1)}{N-1} \right) \\
&= \hat{\sigma}^2 \left[1 + \frac{N-n}{n(N-1)} - \frac{2}{n} \cdot \frac{N-n}{N-1} \right] \\
&= \hat{\sigma}^2 \left[1 - \frac{N-n}{n(N-1)} \right] \\
&= \hat{\sigma}^2 \frac{n(N-1) - N + n}{n(N-1)} \\
&= \hat{\sigma}^2 \cdot \frac{nN - N}{n(N-1)} \\
&= \hat{\sigma}^2 \cdot \frac{N(n-1)}{n(N-1)}
\end{aligned}$$

Cenilka $\hat{\sigma}^2$ ni nepristranska.
Vendar jo lahko popravimo v

$$\hat{\sigma}^2 = \frac{N-1}{N(n-1)} \sum_{k=1}^n (x_k - \bar{x})^2, \text{ ki}$$

je nepristranska cenilka $\hat{\sigma}^2$.

Še nekaj terminologije:

ko govorimo o cenilcih, je to učeloma pred izbiranjem vzorca.

Takrat še ne vemo njihovih vrednosti

ko je vzorec izbran, označimo

konkretne vrednosti z x_1, x_2, \dots, x_n ,

poprečje z \bar{x} in govorimo o

ocenah.

Za presojo kvalitete ocen pa

je nedodajna vzročna porazdelitev

Definicija: količini $\sqrt{\text{var}(\bar{x})}$ rečemo

standardna napaka cenilke \bar{x} in

označimo z $se(\bar{x})$.

Glede ugotovitosti cenilke potem

lahko rečemo: z verjetnostjo ~ 0.68

se zmotimo za manj kot $se(\bar{x})$,

z verjetnostjo 0.95 za manj kot

$1.96 \times se(\bar{x})$, z verjetnostjo ~ 0.99

za manj kot $2.56 \times se(\bar{x})$

ker je $se(\bar{x}) = \frac{s}{\sqrt{n}} \cdot \sqrt{\frac{N-n}{N-1}}$, moramo
tudi s posteriorni udeležiti $\frac{1}{s}$.

Definicija: Faktorju $\sqrt{\frac{N-n}{N-1}}$ rečemo
popravni faktor za kovčnost.

Opomba: Če bi vzorčili brez
z vračanjem, popravnega faktorja
ne bi bilo.

2.4. Intervali zaupanja

Intervali zaupanja so grafična metoda
za prikaz točnosti vzorčnih ocen.

Ideja je preprosta: če ocenjujemo neko
količino, za oceno dobimo število.
To število lahko "razpihujemo" v
interval

Če je z_α tako izbrano, da je za

$$z \sim N(0,1) \quad P(-z_\alpha \leq z \leq z_\alpha) = \alpha,$$

potem interval $\bar{x} \pm z_\alpha \times se(\bar{x})$

pokriva pravo vrednost parametra μ ,

če in samo če se ocena \bar{x} od μ

razlikuje za manj kot $z_\alpha \times se(\bar{x})$,

iz razprave o vročini povzdeletni pa

vedno, da je ta verjetnost

aproximativno α . Izjava se vedno

velja, če namesto z v izrazu

$z \times se(\bar{x})$ z oceno $\hat{\sigma}$.

Opomba: O verjetnostih lahko spet

govorimo pred izbravo vzorca. Če

rečemo, da bomo izbrali vzorec,

ocenili μ z \bar{x} , ocenil

standardno napako z $\frac{\hat{\sigma}}{\sqrt{n}} \cdot \sqrt{\frac{N-n}{N-1}}$,

potem bo interval zaupanja

$$\bar{x} \pm z_\alpha \times se(\bar{x})$$

pokril μ z verjetnostjo približno α .

Po izbiri vzorca v verjetnosti k ne govorimo več in so intervali zaupanja le grafičen način prikaza točnosti vzorčnih ocen.

2.5 Primeri drugih vzorčnih načrtov

Kot smo videli, v praksi enostavno slučajno vzorčenje ni v uporabi.

Večinoma se uporablja večfazno vzorčenje kot v slovenskem javnem mnenju. Ogledali si bomo dva od možnih vzorčnih načrtov in izpeljati ustrezne cenilke. V realnih primerih so lahko cenilke dokaj zapletene.

Stratificirano vzorčenje

Ločja je, da populacijo razdelimo na podskupine, ki jim večemo stratumi. Recimo, da je populacija velikosti N in jo

razdelimo v k podskupin velikosti:

N_1, N_2, \dots, N_k , tako da je $N = N_1 + \dots + N_k$.

Iz vsake podskupine = stratum
izberemo enostavni slučajni vzorec
velikosti n_1, n_2, \dots, n_k . Predpostavljamo
da so postopki izbire vzorcev v
stratumih neodvisni.

Oznake :

μ = populacijsko povprečje

σ^2 = populacijska varianca

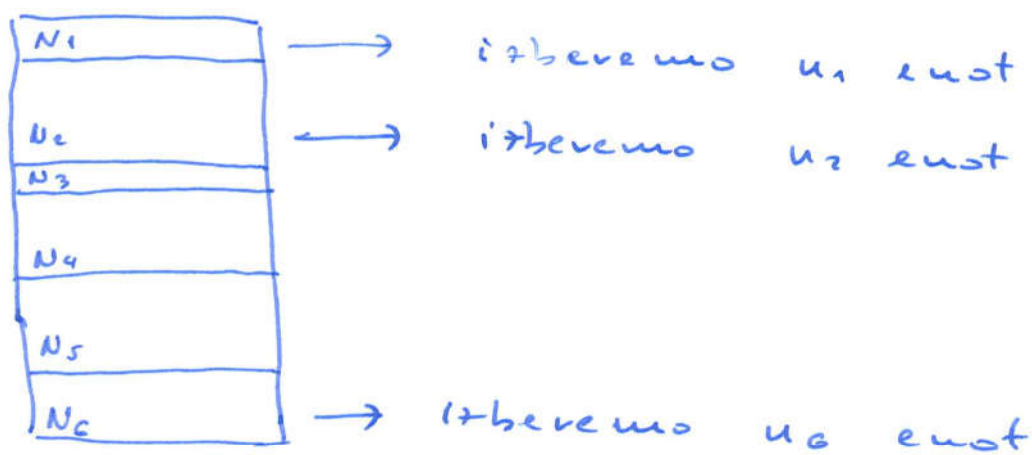
μ_k = populacijsko povprečje
v k -tem stratumu

σ_k^2 = populacijska varianca
v k -tem stratumu

$w_k = \frac{N_k}{N}$ = delež k -tega
stratuma v
celotni
populaciji.

Zakaj bi stratificirali? Če je znotraj stratumov varianca občutno manjša kot je σ^2 , bodo lažje končne ocene bolj točne.

Slika :



\bar{S}_e nekaj oznak:

\bar{x}_k = vzorčno povprečje za enostavni slučajni vzorec v k -tem stratumu.

$\hat{\sigma}_k^2$ = nepristranska cenilka σ_k^2 na podlagi enostavnega slučajnega vzorca

Kako bi ocenili μ ? Opazimo,
da je

$$\mu = \sum_{k=1}^K \frac{N_k}{N} \cdot \mu_k = \sum_{k=1}^K w_k \cdot \mu_k$$

Razumno je reči, da bo cenilka

$$\bar{X} = \sum_{k=1}^K w_k \cdot \bar{X}_k$$

cenilka μ za
stratificirano
vzorec

Zaradi linearnosti je

$$\begin{aligned} E(\bar{X}) &= \sum_{k=1}^K w_k \cdot E(X_k) \\ &= \sum_{k=1}^K w_k \cdot \mu_k \\ &= \mu \end{aligned}$$

Zgornja cenilka za μ je tvoj

nepristranska. Iz predpostavke sledi,

da so $\bar{X}_1, \bar{X}_2, \dots, \bar{X}_K$ neodvisne.

Zaradi neodvisnosti je

$$\begin{aligned} \text{var}(\bar{X}) &= \sum_{k=1}^K w_k^2 \cdot \text{var}(\bar{X}_k) \\ &= \sum_{k=1}^K w_k^2 \cdot \frac{\hat{\sigma}_k^2}{n_k} \cdot \frac{N_k - n_k}{N_k - 1} \end{aligned}$$

Ker lahko $\hat{\sigma}_k^2$ ocenimo nepistransko z $\hat{\sigma}_k^2$, je

$$\sum_{k=1}^K w_k^2 \cdot \frac{\hat{\sigma}_k^2}{n_k} \cdot \frac{N_k - n_k}{N_k - 1}$$

nepistranska celinka $\text{var}(\bar{X})$. Standardni napako ocenimo z

$$\hat{se}(\bar{X}) = \left(\sum_{k=1}^K w_k^2 \frac{\hat{\sigma}_k^2}{n_k} \cdot \frac{N_k - n_k}{N_k - 1} \right)^{1/2}$$

Kako pa je z normalnostjo vzorčne porazdelitve? $\bar{X}_1, \bar{X}_2, \dots, \bar{X}_K$ so približno normalno porazdeljena.

Zato lahko pričakujemo, da bodo tudi njihove lineerne kombinacije aproksimativno normalne

Opomba: Sledeje dejstvo se da dokazati,
vendar je dokaz dokaj zahteven.

Pri stratifikiranem vzorčevju se
pojavi še naslednje vprašanje. Recimo,
da lahko izberemo vzorec velikosti n .
Če imamo stratumne velikosti N_1, \dots, N_k ,
se moramo odločiti, kako velike
vzorce bomo izbrali iz posameznih
stratumov. Kako izbrati, da bo
standardna napaka cenilke \bar{x} čim
manjša? Najh. moramo n_1, \dots, n_k , da
bo $n = n_1 + n_2 + \dots + n_k$ in bo

$$\text{var}(\bar{x}) = \sum_{k=1}^k \frac{\sigma_k^2}{n_k} \cdot \frac{N_k - n_k}{N_k - 1} = W_k^2$$

čim manjša?

V večini praktičnih situacij so korekturni faktorji $\frac{n_k - u_k}{n_k - 1}$ praktično 1 in jih zamearimo. Rešujemo problem vezanega ekstrema:

$$f(u_1, u_2, \dots, u_k) = \sum_{k=1}^k \frac{z_k^2}{n_k} \cdot w_k^2$$

z vezjo $u_1 + u_2 + \dots + u_k = n$. Sestavimo Lagrangeovo funkcijo

$$F(u_1, u_2, \dots, u_k) = f(u_1, \dots, u_k) - \lambda(u_1 + u_2 + \dots + u_k - n)$$

Parcialne odvode izenačimo z 0:

$$\frac{\partial F}{\partial u_k} = -\frac{z_k^2}{n_k^2} w_k^2 - \lambda = 0$$

Sledi

$$-\lambda n_k^2 = z_k^2 w_k^2$$

Iz pogoja $u_1 + u_2 + \dots + u_k = n$ potem sledi

$$n_k = n \cdot \frac{z_k \cdot w_k}{\sum_{j=1}^k z_j \cdot w_j}$$

Našli smo optimalno razmestitev.

Opombe:

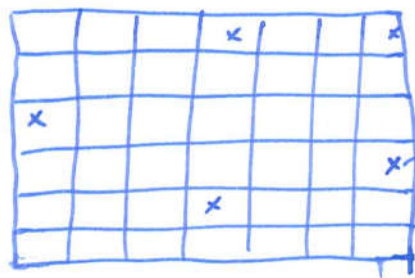
- (i) Na videz je stvar nesmisel, ker za optimalno razmestitev potrebujemo z_1^2, \dots, z_k^2 . Vendar te količine pogosto lahko vsaj približno ocenimo na podlagi bivših vzorcev (kot to vidimo delo statistični urad) ali na podlagi najširih pilotskih vzorcev. Točnost lahko s tem hitreje izboljšamo.
- (ii) Velikosti ne v optimalni rešitvi niso v splošnem cela števila, zato jih zaokrožimo. Če ne bi zaneмали korekturnega faktorja bi dobili rešitev, v kateri w_k zamejamo z $w_k \cdot \sqrt{\frac{N_k}{N_k - 1}}$, kar je v primerjavi z zaokrožanjem zane mavljivo.

Vzorčenje v skupinicah

Recimo, da je populacija velikosti N razdeljena na k enako velikih podskupin velikosti M , tako da je $N = k \cdot M$. Vzorec izberemo na naslednji način:

- (i) izberemo enostavni slučajni vzorec skupin velikosti k .
- (ii) v vzemi izbrani skupini izberemo enostavni slučajni vzorec velikosti m .

Slika:



S kvitcem so označene izbrane podskupine vzorce velikosti m

Oceniti želimo populacijsko povprečje μ in izračunati ter oceniti standardno napako.

Predpostavljamo, da so postopki
izbiranja vzorca v izbranih podskupinah
 neodvisni.

Oznake: $\mu_k =$ populacijsko povprečje
v k -ti skupini.

$\sigma_k^2 =$ populacijska varianca
v k -ti skupini

$\mu =$ populacijsko povprečje

$$\sigma_b^2 = \frac{1}{K} \sum_{k=1}^K (\mu_k - \mu)^2$$

$\bar{x}_k =$ vzorčno povprečje v
 k -ti skupini.

$$I_k = \begin{cases} 1, & \text{če } k\text{-to skupino} \\ & \text{izberemo} \\ 0, & \text{sicer.} \end{cases}$$

It računov pri enostavnem slučajnem

vzorčenju vemo, da je $E(I_k) = \frac{k}{K}$

in $\text{cov}(I_k, I_l) = 0$, $\forall k \neq l$.

$$= -\frac{k}{K} \left(1 - \frac{k}{K}\right) \cdot \frac{1}{K-1}$$

Kandidatka za cenilno μ je

$$\bar{X} = \frac{1}{k} \sum_{j=1}^k \bar{X}_j \cdot I_j$$

= povprečje ocen μ_k za
izbrane podskupine.

Opomba: Pišemo $\bar{X}_1, \dots, \bar{X}_k$,
čeprav ne bomo jemali vzorcev iz
vseh podskupin.

Vendar kot sh. spr. te le
ne obstajajo kot matematični objekti.

Iz predpostavke v bolj matematični
obliki sledi, da sta vektorja
 $(\bar{X}_1, \dots, \bar{X}_k)$ in (I_1, I_2, \dots, I_k) neodvisna.

Sledi, kov so $\bar{X}_1, \dots, \bar{X}_k$ neodvisni

$$E(\bar{X} | I_1, \dots, I_k) = \frac{1}{k} \sum_{j=1}^k \mu_j \cdot I_j$$

$$\text{kov}(\bar{X} | I_1, \dots, I_k) = \frac{1}{k^2} \sum_{j=1}^k I_j \text{var}(\bar{X}_j)$$

Рачунамо

$$\begin{aligned} E(\bar{X}) &= \frac{1}{k} \cdot \sum_{j=1}^k a_j \cdot E(I_j) \\ &= \frac{1}{k} \sum_{j=1}^k a_j \\ &= \mu. \end{aligned}$$

Свилене \bar{X} је непуистранска.

Рачунамо

$$\begin{aligned} \text{var}(\bar{X}) &= E(\text{var}(\bar{X} | I_1, \dots, I_k)) \\ &\quad + \text{var}(E(\bar{X} | I_1, \dots, I_k)) \\ &= \frac{1}{k^2} \cdot \sum_{j=1}^k E(I_j) \text{var}(\bar{X}_j) \\ &\quad + \frac{1}{k^2} \text{var}\left(\sum_{j=1}^k a_j I_j\right) \\ &= \frac{1}{k \cdot k} \sum_{j=1}^k \text{var}(\bar{X}_j) \\ &\quad + \frac{1}{k^2} \sum_{j=1}^k a_j^2 \cdot \frac{k}{k} \left(1 - \frac{k}{k}\right) \end{aligned}$$

$$+ \frac{1}{k^2} \sum_{i \neq j} \sigma_i \mu_j \frac{k}{k} \left(1 - \frac{k}{k}\right) \left(-\frac{1}{k-1}\right)$$

$$= \frac{1}{k \cdot k} \sum_{j=1}^k \frac{\sigma_j^2}{m} \cdot \frac{M-m}{M-1}$$

$$+ \frac{1}{k \cdot k} \left(1 - \frac{k}{k}\right)$$

$$\times \left[\sum_{j=1}^k \mu_j^2 - \sum_{i \neq j} \sigma_i \mu_j \left(\frac{1}{k-1}\right) \right]$$

$$= \frac{1}{k \cdot k} \cdot \sum_{j=1}^k \frac{\sigma_j^2}{m} \cdot \frac{M-m}{M-1}$$

$$+ \frac{k-k}{k \cdot k^2} \left[\sum_{j=1}^k \left(1 + \frac{1}{k-1} - \frac{1}{k-1}\right) \mu_j^2 - \sum_{i \neq j} \sigma_i \mu_j \frac{1}{k-1} \right]$$

$$= \quad -u- \quad$$

$$+ \frac{k-k}{k \cdot k^2} \cdot \left[\frac{k}{k-1} \sum_{j=1}^k \mu_j^2 - \frac{1}{k-1} \left(\sum_{j=1}^k \sigma_j\right)^2 \right]$$

$$= \frac{1}{k \cdot k} \sum_{j=1}^k \frac{\sigma_j^2}{m} \cdot \frac{M-m}{M-1}$$

$$+ \frac{1}{k} \cdot \frac{k-k}{k-1} \cdot \frac{1}{k} \sum_{j=1}^k (\sigma_j - \mu)^2$$

$$= \frac{1}{k \cdot k} \sum_{j=1}^k \frac{z_j^2}{m} \cdot \frac{M-m}{M-1}$$

$$+ \frac{1}{k} \cdot \frac{k-k}{k-1} \cdot \sigma_b^2$$

Če želimo presojati točnost cenilke \bar{x} moramo tudi oceniti količini

$$s^2 = \frac{1}{k} \sum_{j=1}^k z_j^2 \quad \text{in} \quad \sigma_b^2.$$

~~III~~ ⇒

Za s^2 hitro najdemo nepristransko cenilko:

$$\hat{s}^2 = \frac{1}{k} \sum_{j=1}^k \hat{\sigma}_j^2 \cdot I_j, \quad \text{kjer je}$$

$\hat{\sigma}_j^2$ nepristranska cenilka σ_j^2 . Za

σ_b^2 pa je nekoliko bolj komplicirano.

2c priručnik

$$\sigma_b^2 = \frac{1}{K} \sum_{j=1}^k a_j^2 - \mu^2$$

1 de ja :

$$\hat{\sigma}_b^2 = \frac{1}{K} \sum_{j=1}^k \bar{x}_j^2 \cdot I_j - \bar{x}^2$$

Vemo: $E(\bar{x}_j^2) = \text{var}(\bar{x}_j) + \mu_j^2$

$$E(\bar{x}^2) = \text{var}(\bar{x}) + \mu^2$$

Računamo

$$E(\hat{\sigma}_b^2) = \sigma_b^2 + \frac{1}{K} \sum \text{var}(\bar{x}_j) - \text{var}(\bar{x})$$

$$= \sigma_b^2 + f \cdot \frac{1}{m} \cdot \frac{M-m}{M-1}$$

$$= \frac{1}{K} f \frac{M-m}{m(M-1)}$$

$$= \frac{1}{K} \frac{k-k}{k-1} \sigma_b^2$$

Preoblikujemo v

$$E(\hat{\sigma}_b^2) = \sigma_b^2 \cdot \frac{K(K-1)}{K(K-1)} + \mu \cdot \frac{K-1}{K} \cdot \frac{M-m}{m(M-1)}$$

Torej lahko $\hat{\sigma}_b^2$ popravimo v nepristransko cenilko:

$$\begin{aligned} \tilde{\sigma}_b^2 &= \left[\hat{\sigma}_b^2 - \hat{\mu} \cdot \frac{K-1}{K} \cdot \frac{M-m}{m(M-1)} \right] \cdot \frac{K(K-1)}{K(K-1)} \\ &= \frac{K(K-1)}{K(K-1)} \hat{\sigma}_b^2 - \hat{\mu} \cdot \frac{(M-m)(K-1)}{m(M-1) \cdot K} \end{aligned}$$

Ostane še vprašanje aproksimativne normalnosti. Lahko si mislimo, da skupine izbiramo eno po eno obklev jih ne izberemo K . V vsaki tudi izberemo enostavni slučajni vzorec in dobimo slučajne sp. $\bar{X}_1, \bar{X}_2, \dots, \bar{X}_K$.

Cenilka je $\frac{1}{K}(\bar{X}_1 + \dots + \bar{X}_K)$ in vse $\bar{X}_1, \dots, \bar{X}_K$ so enako porazdeljene. Kot pri enostavnem sl. vzorčju je potem vzorčna povprečitev aproksimativno normalna

POUZETEK

1. Osnova za vzorčenje je vzoreni učrt. Vzorenje je verjetnostno, če lahko za vsak možen vzorec vnajprej povernemo verjetnost, da bo izbran.
2. Ocena. katerekoli količine $(\mu, \sigma^2, \tau, \sigma^2_\tau)$ je neka funkcija vrednosti statistične spremenljivke na izbranih enotah.
3. Na postopek ocenjevanja lahko matematično gledamo kot na "ustvarjanje" slučajne spremenljivke cenilke. Vsa informacija o točnosti je vsebovana v tem, kar znamo povedati o porazdelitvi cenilke.

3. Ocenjevanje parametrov

3.1. Pojem statističnega modela

Statistika se ukvarja z analizo podatkov.

Eno od bistvenih orodij za ta namen so statistični modeli. Statistični model je opis mehanizma, za katerega merimo, da je generiral podatke.

Tak mehanizem praviloma vsebuje naključnost, zato so statistični modeli formulirani v jeziku slučajnih spremenljivk, porazdelitev, neodvisnosti ali v jeziku pogojnih porazdelitev.

Oglejmo si nekaj primerov.

Primer: Linearna regresija

Podatki so oblike $y_i, x_{i1}, x_{i2}, \dots, x_{im}$

z $i = 1, 2, \dots, n$;

Tipično veckamo, da je y_i odziv i -te osebe, vrednosti $x_{i1}, x_{i2}, \dots, x_{in}$ pa kovariante. Tipično bi vadili "pojasnili" spremenljivko y_i kot funkcijo kovariat $x_{i1}, x_{i2}, \dots, x_{in}$, vendar funkcijteka odvisnost ni "točna", tako da novo vlogo igra še naključje.

Prvi primer verjetneje je zgodovinsko ta, da je F. Galton želel pojasniti telesno višino sina (y) s telesno višino očeta (x). Imel je podatke:

y_1	x_1
y_2	x_2
\vdots	\vdots
y_n	x_n

Statistični model pravi: pravi $(x_i, y_i), \dots, (x_n, y_n)$ "vstanejo" kot med sabo neodvisni slučajni vektorji iz neke porazdelitve $(x_i, y_i), (x_1, y_2), \dots, (x_n, y_n)$.

Predpostavljamo, da velja

$$Y_i = \alpha + \beta X_i + \varepsilon_i,$$

kjer so $\varepsilon_1, \varepsilon_2, \dots, \varepsilon_n$ neodvisne in enako porazdeljene slučajne spremenljivke z $E(\varepsilon_i) = 0$ in $\text{var}(\varepsilon_i) = \sigma^2$ za $i = 1, 2, \dots, n$.

V grobem torej rečemo, da je telesna višina sina linearna funkcija telesne višine očeta + slučajne člen. Na drugačen način bi lahko zapisali:

$$E(Y_i | X_i) = \alpha + \beta X_i$$

$$\text{var}(Y_i | X_i) = \sigma^2$$

Torej, statistični model je predpostavka da so (X_i, Y_i) generirani iz neke porazdelitve, o kateri nekaj predpostavljamo.

Primer : LOGISTRČNA REGRESIJA

Banke tipično morajo razvrščati kvaliteto kreditotjemalcev v razrede. To različna bančna zakonodaja.

Podatki iz preteklosti so oblike

$$\begin{array}{l} y_1 \quad x_{11} \quad x_{12} \quad \dots \quad x_{1m} \\ y_2 \quad x_{21} \quad x_{22} \quad \dots \quad x_{2m} \\ \vdots \\ y_n \quad x_{n1} \quad x_{n2} \quad \dots \quad x_{nm} \end{array}$$

Tu tem je $y_i \in \{0, 1\}$ indikator, ali je kreditotjemalec nehal odplačevati oblj. $x_{i1}, x_{i2}, \dots, x_{im}$ pa "kovariante", ki so starost, delovno doba, dohodki, znanstveni stan, tip zaposlitve, Vprašanje: ali lahko iz kovariet napovemo verjetnost, da kreditotjemalec ne bo vrnil kredita?

Logistična regresija je ena od možnosti.

Predpostavka je: ~~vektori~~ vektorji

$(y_i, x_{i1}, \dots, x_{im})$ „nastanejo“ kot

med sebo neodvisni vektorji

$(y_i, x_{i1}, \dots, x_{im})$ + enako porazdelitvijo.

Predpostavimo, da je

$$P(y_i = 1 \mid x_{i1}, x_{i2}, \dots, x_{im})$$

$$= \frac{e^{\beta_0 + \sum_{j=1}^m \beta_j x_{ij}}}{1 + e^{\beta_0 + \sum_{j=1}^m \beta_j x_{ij}}}$$

za $i = 1, 2, \dots, n$. Konstantam

$\beta_0, \beta_1, \dots, \beta_m$ rečemo konstante

v modelu ali parametri.

Parametre ocenimo na podlagi

znanih podatkov. ko imamo

konkretnega kreditnega

s kovariatami x_1, x_2, \dots, x_m in ko
imamo ocenjene parametre
 $\hat{\beta}_0, \hat{\beta}_1, \dots, \hat{\beta}_m$ lahko ocenimo
verjetnost, da kreditno ualec ne
bo vrnil kredita kot

$$\hat{p}_i = \frac{e^{\hat{\beta}_0 + \sum_{j=1}^m \hat{\beta}_j x_j}}{1 + e^{\hat{\beta}_0 + \sum_{j=1}^m \hat{\beta}_j x_j}}.$$

Torej, predpostavljamo, da so
podatki nastali na določen
način naključno in predpostavljau
določene lastnosti. Teh pozdelite
Ali lahko trdimo, da je model
točen? Tipično ne, lahko pa
izvedemo določena preverjanja.

V tem poglavju bomo obravnavali
kako oceniti parametre in kaj
lahko večemo o točnosti ocen.

3.2. Cenilne, vzorčne porazdelitve, intervali zaupanja

Primer: Recimo, da imamo
podatke x_1, x_2, \dots, x_n . Podatkom
pogosto večemo opazovane vrednosti.
Predpostavimo, da so podatki
nastali kot slučajne spremenljivke
 x_1, x_2, \dots, x_n , za katere predpostavljamo
 neodvisnost in enako $N(\mu, \sigma^2)$
porazdeljenost. Parametro μ in σ^2
ne poznamo, predpostavljamo le
 $(\mu, \sigma^2) \in \mathbb{R} \times (0, \infty)$.

Kako bi ocenili μ in σ^2 !

Idelja:

$$\hat{\mu} = \frac{1}{n} \sum_{i=1}^n x_i$$

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

Oceni $\hat{\mu}$ in $\hat{\sigma}^2$ smo dobili
tako, da smo "sprožili"
slučajni spremenljivki

$$\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i \quad \text{in}$$

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

Postava je podobna, kot pri
vzorecju. Preden podatki
"nastanejo" so slučajne
spremenljivke, le da pri
vzorecju vemo, da jih sprožimo
z uvedbo vzorečnega načrta,
tukaj pa je "sprožanje" nevolno

bolj abstraktno. Kaj znamo povedati
o parametrovih \bar{X} in S^2 ?

$$E(\bar{X}) = \frac{1}{n} \sum_{i=1}^n E(X_i) = \mu$$

$$\begin{aligned} \text{var}(\bar{X}) &= \frac{1}{n^2} \text{var}(X_1 + \dots + X_n) \\ &= \frac{\sigma^2}{n} \end{aligned}$$

Kot pri vzorčju bomo slučajni
spremenljivici \bar{X} rekli cenilka.

Številu $\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$ večemo ocene.

Oznako $\hat{\mu}$ bomo uporabljali za
eno ali drugo, vaba pa bo
jasna iz konteksta.

Cenilka $\hat{\mu} = \bar{X}$ je nepristranska.

Poleg tega je $\text{se}(\bar{X}) = \frac{\sigma}{\sqrt{n}}$ in
 \bar{X} je normalna.

Рачунамо \bar{x}

$$\begin{aligned} E(S^2) &= \frac{1}{n} \sum_{i=1}^n E[(x_i - \bar{x})^2] \\ &= E[(x_1 - \bar{x})^2] \quad (\text{симетрија}) \\ &= \text{var}(x_1 - \bar{x}) \\ &= \text{var}(x_1) + \text{var}(\bar{x}) \\ &\quad - 2 \text{cov}(x_1, \frac{1}{n} \sum_{i=1}^n x_i) \\ &= \sigma^2 + \frac{\sigma^2}{n} - \frac{2}{n} \cdot \sigma^2 \\ &= \frac{n-1}{n} \cdot \sigma^2 \end{aligned}$$

Секундарна S^2 ни непуистративна,
Затога па ја користишемо $\hat{\sigma}^2$
непуистративна ~~непуистративна~~ секундарна

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

Kako pa je s porazdelitvijo $\hat{\sigma}^2$?

17 poglavja o večvarstveni
normalni porazdelitvi vemo,
da je

$$\sum_{i=1}^n (x_i - \bar{x})^2 \sim \sigma^2 \chi^2(n-1)$$

$\chi^2(n)$ porazdelitev je vsota med
sebo neodvisnih slučajnih
spremenljivk \Rightarrow po centralnem
limitnem izreku aproksimativno
normalna. Rečemo lahko torej,
da je $\hat{\sigma}^2$ približno normalna
porazdeljena s parametroma
 σ^2 in $2\sigma^4(n-1)$, kar je
varianca $\chi^2(n)$ porazdelitve $2n$.
Z drugimi besedami

$$\sqrt{n} (\hat{\sigma}^2 - \sigma^2) \stackrel{\text{aprox.}}{\sim} N\left(0, 2\sigma^4 \frac{(n-1)}{n}\right)$$

Potrebujemo nekaj terminologije.

Prijeteli bomo, da so podatki ali

opazovane vrednosti x_1, x_2, \dots, x_n

"nastale" kot slučajni vektor

$\underline{X} = (X_1, X_2, \dots, X_n)$ + gostoto

$f_{\underline{X}}(\underline{x}, \underline{\theta})$, kjer je $\underline{\theta} \in \Theta \subseteq \mathbb{R}^m$ vektor
možnih parametrov.

Opomba: Zaenkrat ne privzamemo, da
so X_1, X_2, \dots, X_n neodvisne.

Opomba: Za diskretne vektorje
 \underline{X} bodo od parametra $\underline{\theta}$ odvisne
verjetnosti $p(\underline{x}, \underline{\theta}) = P_{\underline{\theta}}(\underline{X} = \underline{x})$.

Definicije:

(i) Funkciji $\hat{\underline{\theta}} = \hat{\underline{\theta}}(x_1, x_2, \dots, x_n)$, ki
"ocenjuje" $\underline{\theta}$ večemo ocena.

(ii) Situacijumu vektorju

$$\hat{\underline{\theta}} = \hat{\underline{\theta}}(x_1, x_2, \dots, x_n) = \hat{\underline{\theta}}(\underline{x})$$

rečemo ceniška parametra $\underline{\theta}$.

Opomba: k komponente pišemo $\hat{\theta}_j$.

(iii) Ceniška $\hat{\underline{\theta}}$ je nepristranska, če je

$$E_{\underline{\theta}}(\hat{\underline{\theta}}) = \underline{\theta}.$$

(iv) Srednja kvadratična napaka ceniške $\hat{\theta}$ je količina

$$MSE(\hat{\theta}) = E[(\hat{\theta} - \theta)^2],$$

če je parameter euodimensionalu.

(v) k euodimensionalno $\hat{\theta}$ je standardna napaka količina

$$se(\hat{\theta}) = \sqrt{\text{var}_{\underline{\theta}}(\hat{\theta})}.$$

Primer: Recimo, da so x_1, x_2, \dots, x_n
vzorec iz Poissonove porazdelitve
 $P_0(\lambda)$. Povprečje

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

je neprirodna skalar celilka λ , ker je
 $E(X_i) = \lambda$. Velja tudi

$$\text{var}(\bar{X}) = \frac{1}{n} \cdot \text{var}(X_i) = \frac{\lambda}{n},$$

tovej je $se(\bar{X}) = \frac{\sqrt{\lambda}}{\sqrt{n}}$. Po
centralnem limitnem izreku je
 \bar{X} aproksimativno normalno,
varianco pa lahko ocenimo z
 $\frac{\hat{\lambda}}{n}$. Spet lahko rečemo

$$\sqrt{n}(\hat{\lambda} - \lambda) \stackrel{\text{aprox}}{\sim} N(0, \lambda)$$

Primer: Opazovane vrednosti $\underline{x}_1, \underline{x}_2, \dots, \underline{x}_n$ so lahko tudi vektorji. Privzemimo da so nastali kot neodvisni, enako porazdeljeni normalni vektorji, torej $\underline{x}_1, \underline{x}_2, \dots, \underline{x}_n \rightarrow \underline{x}_k \sim N(\underline{\mu}, \underline{\Sigma})$.

Recimo

$$\hat{\underline{\mu}} = \frac{1}{n} \sum_{k=1}^n \underline{x}_k$$

$$\hat{\underline{\Sigma}} = \frac{1}{n-1} \sum_{k=1}^n (\underline{x}_k - \bar{\underline{x}})(\underline{x}_k - \bar{\underline{x}})^T$$

Računamo

$$E(\hat{\underline{\mu}}) = \underline{\mu} \quad (\text{nepristopna})$$

$$\begin{aligned} E(\hat{\underline{\Sigma}}) &= \frac{1}{n-1} \sum_{k=1}^n E[(\underline{x}_k - \bar{\underline{x}})(\underline{x}_k - \bar{\underline{x}})^T] \\ &= \frac{n}{n-1} E[(\underline{x}_1 - \bar{\underline{x}})(\underline{x}_1 - \bar{\underline{x}})^T] \\ &= \frac{n}{n-1} E[\underline{x}_1 \underline{x}_1^T - \bar{\underline{x}} \underline{x}_1^T \\ &\quad - \underline{x}_1 \bar{\underline{x}}^T + \bar{\underline{x}} \bar{\underline{x}}^T] \end{aligned}$$

$$\begin{aligned}
&= \frac{n}{n-1} \left[\underline{\Sigma} + \underline{\mu} \cdot \underline{\mu}^T \right. \\
&\quad - \text{cov}(\underline{\bar{x}}, \underline{x}_1) + \underline{\mu} \underline{\mu}^T \\
&\quad - \text{cov}(\underline{x}_1, \underline{\bar{x}}) - \underline{\mu} \underline{\mu}^T \\
&\quad \left. + \text{cov}(\underline{\bar{x}}, \underline{\bar{x}}) + \underline{\mu} \cdot \underline{\mu}^T \right]
\end{aligned}$$

$$= \frac{n}{n-1} \left[\underline{\Sigma} - \frac{1}{n} \underline{\Sigma} - \frac{1}{n} \underline{\Sigma} + \frac{1}{n} \underline{\Sigma} \right]$$

$$= \frac{n}{n-1} \cdot \frac{n-1}{n} \underline{\Sigma}$$

$$= \underline{\underline{\Sigma}}$$

Predlagana $\underline{\underline{\hat{\Sigma}}}$ je nepristvarna.

Po vektorski verziji centralnega limitnega izreka je

$$\sqrt{n} (\underline{\hat{\mu}} - \underline{\mu}) \overset{\text{aprox.}}{\sim} N(0, \underline{\underline{\Sigma}})$$

3.3 . Metoda največjega verjetja

Primer: Predpostavimo, da so opazovane vrednosti x_1, x_2, \dots, x_n vzorec iz Weibullove gostote. To pomeni, da so "nastale" kot med sabo neodvisne, enako porazdeljene slučajne spremenljivke z gostoto

$$f(x, \alpha, \beta) = \begin{cases} \frac{\alpha}{\beta} \left(\frac{x}{\beta}\right)^{\alpha-1} e^{-\left(\frac{x}{\beta}\right)^\alpha}, & x \geq 0 \\ 0, & \text{sicer.} \end{cases}$$

Pri tem je $\alpha > 0$ in $\beta > 0$.

Opomba: Abstraktno je $\underline{\theta} = (\alpha, \beta)$
in $\Theta = (0, \infty)^2$.

Kako bi ocenili α in β !

V tem primeru mi očitno, kako
oceniti parametre. Ugotovimo,
da je

$$F_X(x) = 1 - e^{-\left(\frac{x}{b}\right)^\alpha}$$

Če je $Y = \left(\frac{X}{b}\right)^\alpha$, je

$$\begin{aligned} P(Y \leq y) &= P\left(\left(\frac{X}{b}\right)^\alpha \leq y\right) \\ &= P(X \leq b \cdot y^{1/\alpha}) \\ &= 1 - e^{-\left(\frac{b \cdot y^{1/\alpha}}{b}\right)^\alpha} \\ &= 1 - e^{-y}, \end{aligned}$$

tovej je $Y \sim \exp(1)$. Sledi

$$\begin{aligned} E(X) &= E\left[b \cdot Y^{1/\alpha}\right] \\ &= b \cdot \int_0^\infty y^{1/\alpha} e^{-y} dy \\ &= b \cdot \Gamma\left(1 + \frac{1}{\alpha}\right) \end{aligned}$$

$$E(x^2) = b^2 \Gamma\left(1 + \frac{2}{\alpha}\right)$$

V tem primeru \bar{x} in izvedenke ne pomagajo. Morda bi lahko rekli

$$\bar{x} \approx b \Gamma\left(1 + \frac{1}{\alpha}\right) \text{ in}$$

$$\frac{1}{n} \sum_{k=1}^n x_k^2 \approx E(x^2) = b^2 \Gamma\left(1 + \frac{2}{\alpha}\right)$$

in izrazili α in b , vendar se to t. i. metoda momentov izkaže za slabotno.

Ideja: Vektor podatkov

$\underline{x} = (x_1, x_2, \dots, x_n)^T$ je nastal kot naključna izbira točka

$\underline{x} = (x_1, x_2, \dots, x_n)^T$ + gostota

$$f_{\underline{x}}(\underline{x}, \alpha, b) = \prod_{k=1}^n f_{x_k}(x_k, \alpha, b).$$

kje smo bolj verjetno izbrali \underline{x} ,
tam kjer je gostota velika ali
tam, kjer je gostota majhna?

Če se moramo izbirati, tam
kjer je gostota velika. Obrnimo
razmislek: za fiksne opazovane
vrednosti x_1, x_2, \dots, x_n si
oglejmo funkcijo

$$h(\alpha, \beta | \underline{x}) = \prod_{k=1}^n f_{X_k}(x_k, \alpha, \beta),$$

tovej funkcijo

$$(\alpha, \beta) \xrightarrow{L} \prod_{k=1}^n f_{X_k}(x_k, \alpha, \beta).$$

Če ima ta funkcija maksimum,
je argument maximuma dober
kandidat za $(\hat{\alpha}, \hat{\beta})$.

Ideji, da prikivamo \underline{x} in iščemo maksimum po $\underline{\theta}$ se imenuje metoda največjega verjetja.

(angl. maximum likelihood method)

Autor ideje je angleški statistik Sir Ronald A. Fischer (1890-1962).

V primeru Weibulove porazdelitve dobimo

$$L(\alpha, \beta | \underline{x}) = \frac{\alpha^n}{\beta^n} \prod_{k=1}^n \left(\frac{x_k}{\beta}\right)^{\alpha-1} e^{-\left(\frac{x_k}{\beta}\right)^\alpha}$$

ker je $x_k > 0$ lahko logaritmiramo in dobimo

$$\begin{aligned} \log L(\alpha, \beta | \underline{x}) &= n \log \alpha - n \log \beta \\ &+ (\alpha-1) \sum_{k=1}^n \log x_k - \sum_{k=1}^n \left(\frac{x_k}{\beta}\right)^\alpha \\ &- (\alpha-1) n \log \alpha \end{aligned}$$

Poenostavimo

$$\log L(\alpha, \beta | \underline{x}) = n \log \alpha - n\alpha \log \beta + (\alpha - 1) \sum_{k=1}^n \log x_k - \sum_{k=1}^n \left(\frac{x_k}{\beta}\right)^\alpha.$$

Parcialno odvajamo po α in β .

$$\begin{aligned} \frac{\partial \log L}{\partial \alpha} &= \frac{n}{\alpha} - n \log \beta + \sum_{k=1}^n \log x_k - \sum_{k=1}^n \left(\frac{x_k}{\beta}\right)^\alpha \log\left(\frac{x_k}{\beta}\right) \\ &= 0 \end{aligned}$$

$$\frac{\partial \log L}{\partial \beta} = -\frac{n\alpha}{\beta} + \sum_{k=1}^n \frac{x_k^\alpha}{\beta^{\alpha+1}} (-\alpha) = 0$$

Iz druge enačbe sledi:

$$\frac{1}{n} \sum_{k=1}^n \left(\frac{x_k}{\beta}\right)^\alpha = 1.$$

Prvo enačbo delimo + n in prevedimo:

$$\frac{1}{\alpha} = \frac{1}{n} \sum_{k=1}^n \log x_k - \frac{1}{n} \sum_{k=1}^n \left(\frac{x_k}{\beta}\right)^\alpha \log x_k$$

Evolutivnost rešitev enačb po metodi največjega verjetja.

Za logaritemsko funkcijo verjetja smo dobili

$$l(\alpha, b | \underline{x})$$

$$= n \log \alpha - n \log b + (\alpha - 1) \sum_{k=1}^n \log x_k - \sum_{k=1}^n \left(\frac{x_k}{b}\right)^\alpha$$

Parcialno odvajanje do enačbi

$$\frac{\partial l}{\partial \alpha} = \frac{n}{\alpha} - n \log b + \sum_{k=1}^n \log x_k - \sum_{k=1}^n \left(\frac{x_k}{b}\right)^\alpha \log\left(\frac{x_k}{b}\right) = 0$$

$$\frac{\partial l}{\partial b} = -\frac{n\alpha}{b} + \frac{\alpha}{b} \sum_{k=1}^n \left(\frac{x_k}{b}\right)^\alpha = 0$$

12) drugje enačbo sledi, da je

$$b^\alpha = \left(\sum_{k=1}^n x_k^\alpha\right) / n \quad \text{in} \quad \sum_{k=1}^n \left(\frac{x_k}{b}\right)^\alpha = 1.$$

Ustavimo v prvo enačbo in dobimo po deljenju $\frac{1}{\alpha}$

$$\frac{1}{\alpha} = -\frac{1}{n} \sum_{k=1}^n \log x_k + \frac{1}{n} \frac{\sum_{k=1}^n x_k^\alpha \log x_k}{\sum_{k=1}^n x_k^\alpha} = 0$$

Prepišemo

$$\frac{1}{\alpha} = \frac{\sum_{k=1}^n x_k^\alpha \log x_k}{\sum_{k=1}^n x_k^\alpha} = \frac{1}{n} \sum_{k=1}^n \log x_k = g(\alpha)$$

Računamo

$$g'(\alpha) = \frac{\sum_{k=1}^n x_k^\alpha \log^2 x_k \cdot \sum_{k=1}^n x_k^\alpha - \left(\sum_{k=1}^n x_k^\alpha \log x_k \right)^2}{\left(\sum_{k=1}^n x_k^\alpha \right)^2}$$

По Cauchy-Schwarzovi neenačbi je

$$\left(\sum_{k=1}^n x_k^\alpha \log x_k \right)^2 \leq \left(\sum_{k=1}^n x_k^\alpha \right) \cdot \left(\sum_{k=1}^n x_k^\alpha \log^2 x_k \right)$$

Opomba: Pišemo $x_k^\alpha = x_k^{d/2} \cdot x_k^{d/2}$.

Enaost dobimo le, če sta vektorja

$$(x_1^{d/2}, \dots, x_n^{d/2}) \text{ in } (x_1^{d/2} \log x_1, \dots, x_n^{d/2} \log x_n)$$

kolinearna. To je možno le, če je $x_1 = x_2 = \dots = x_n$, kar lahko izključimo.

Funkcija $g(\alpha)$ je na $(0, \infty)$ strogo naraščajoča.

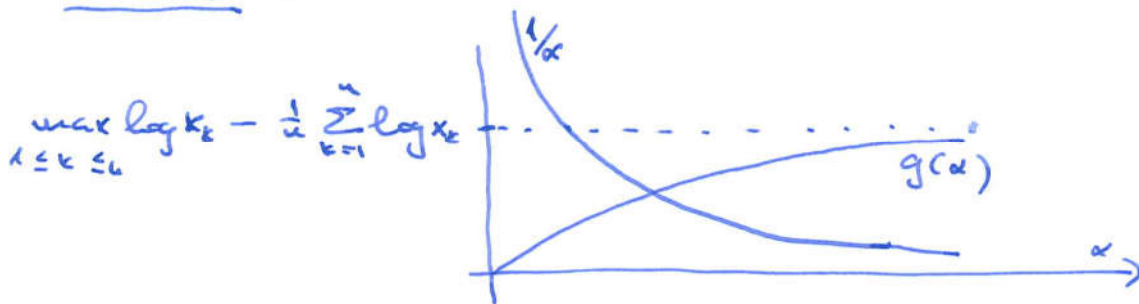
Volja $\lim_{x \rightarrow 0} g(x) = 0$ in

$$\lim_{x \rightarrow \infty} \frac{\sum_{k=1}^n x_k^x \log x_k}{\sum_{k=1}^n x_k^x} = \max_{1 \leq k \leq n} \log x_k$$

ker miso vsi x_1, x_2, \dots, x_n su tučni, je

$$\max_{1 \leq k \leq n} \log x_k > \frac{1}{n} \sum_{k=1}^n \log x_k.$$

Slika :



Ker učba $1/x = g(x)$ ima natanko eno rešitev > 0 . Iz tega sledi tudi enoličnost ζ .

Enočbe nimajo analitične vešitve.
Z nekaj truda lahko dokazemo, da
je vešitev evolične in obstaja, vendar
eksplicitne oblike \hat{x} in \hat{z} ne
poznamo. Kaj pa zdaj?

Oglejmo si simulacije.

3.4. Asimptotske lastnosti cenilek po metodi največjega verjetja

Iz simulacij izhaja, da utegnete
biti \hat{x} in \hat{z} aproksimativno
normalno porazdeljeni. Z izpeljavo
tega dejstva potrebujemo nekaj
pripravnih delov izrekov. Najprej
bo trditev sledila iz centralnega
limitnega izreka.

Izrek 3.1 : Naj ima slučajna
spremenljivka gostoto $f(x, \theta)$, kjer
je $\Theta \subseteq \mathbb{R}$ odprta množica. Definirjmo

$$\log_+ f(x, \theta) = \begin{cases} \log f(x, \theta) & \text{za} \\ & f(x, \theta) > 0; \\ 0 & \text{sicer.} \end{cases}$$

Predpostavimo, da povsod kjer je
treba lahko integrale s parametrom
odvajamo pod integralnim znakom.

Predpostavimo, da je $\theta \mapsto f(x, \theta)$
zvezo odvedljiva po θ za vsak
fiksni x . Definirjmo

$$Z = \frac{\partial}{\partial \theta} \log_+ f(x, \theta) \Big|_{\theta = \theta_0}$$

$$W = \frac{\partial^2}{\partial \theta^2} \log_+ f(x, \theta) \Big|_{\theta = \theta_0}$$

Velja

$$E(Z) = 0 \quad \text{in} \quad E(W) = - \text{var}(Z)$$

Доказ : Рассмотрим

$$E(z) = \int_{-\infty}^{\infty} \frac{\partial}{\partial \theta} \log_y f(x, \theta) \Big|_{\theta = \theta_0} \cdot f(x, \theta_0) dx$$

$$= \int_{-\infty}^{\infty} \frac{\partial}{\partial \theta} \log f(x, \theta) \Big|_{\theta = \theta_0} f(x, \theta_0) dx$$

{ $f(x, \theta_0) > 0$ }

$$= \int_{\{f(x, \theta_0) > 0\}} \frac{\frac{\partial}{\partial \theta} f(x, \theta) \Big|_{\theta = \theta_0}}{f(x, \theta_0)} f(x, \theta_0) dx$$

$$= \int_{\{f(x, \theta_0) > 0\}} \frac{\partial}{\partial \theta} f(x, \theta) \Big|_{\theta = \theta_0} dx$$

$$= \left(\frac{\partial}{\partial \theta} \int_{-\infty}^{\infty} f(x, \theta) dx \right) \Big|_{\theta = \theta_0}$$

$$= \left(\frac{\partial}{\partial \theta} 1 \right) \Big|_{\theta = \theta_0}$$

$$= 0 .$$

Podobno va enunciamo za w .

$$E(w) = \int_{\{f(x, \theta_0) > 0\}} \frac{\partial^2}{\partial \theta^2} \log f(x, \theta) \Big|_{\theta = \theta_0} f(x, \theta_0)$$

$$= \int_{\{f(x, \theta) > 0\}} \frac{\frac{\partial^2 f}{\partial \theta^2}(x, \theta_0) f(x, \theta_0) - \left(\frac{\partial}{\partial \theta} f(x, \theta_0)\right)^2}{f^2(x, \theta_0)} \times f(x, \theta_0) dx$$

$$= \int_{-\infty}^{\infty} \frac{\partial^2}{\partial \theta^2} f(x, \theta_0) dx - \int_{\{f(x, \theta) > 0\}} \left(\frac{\partial}{\partial \theta} \log f(x, \theta)\right)^2 \Big|_{\theta = \theta_0} \cdot f(x, \theta_0) dx$$

$$= \left(\frac{\partial^2}{\partial \theta^2} 1\right) \Big|_{\theta = \theta_0} - E(z^2)$$

$$= -E(z^2)$$

Ke- je $E(z) = 0$, je drugi del tudi trivno dokazati.

Opomba: Podoben račun velja za
parcialne odvode \neq več pisarje.

Če je $\underline{\theta} = (\theta_1, \theta_2, \dots, \theta_m)$ lahko
definiramo

$$z_i = \frac{\partial}{\partial \theta_i} \log_+ f(\underline{x}, \underline{\theta}) \Big|_{\underline{\theta} = \underline{\theta}_0}$$

$$w_{i,j} = \frac{\partial^2}{\partial \theta_i \partial \theta_j} \log_+ f(\underline{x}, \underline{\theta}) \Big|_{\underline{\theta} = \underline{\theta}_0}$$

Velja (pod podobnimi predpostavkami)

$$E(z_i) = 0 \quad \text{in}$$

$$E(w_{i,j}) = -\text{cov}(z_i, z_j)$$

Primer: Če je X Weibullova je

$$\begin{aligned}\log f(x, \alpha, \beta) &= \log \alpha - \log \beta \\ &\quad + (\alpha - 1) \log \frac{x}{\beta} \\ &\quad - \left(\frac{x}{\beta}\right)^\alpha\end{aligned}$$

Odvajamo

$$\frac{\partial}{\partial \alpha} \log f(x, \alpha, \beta)$$

$$= \frac{1}{\alpha} + \log \frac{x}{\beta} - \left(\frac{x}{\beta}\right)^\alpha \cdot \log \frac{x}{\beta}$$

Nadomestimo $x = Y$ in izračunamo

$$Z_1 = \frac{1}{\alpha} + \log \frac{Y}{\beta} - \left(\frac{Y}{\beta}\right)^\alpha \cdot \log \frac{Y}{\beta}$$

$$= \frac{1}{\alpha} + \frac{1}{\alpha} \log \left(\frac{Y}{\beta}\right)^\alpha - \frac{1}{\alpha} \left(\frac{Y}{\beta}\right)^\alpha \log \left(\frac{Y}{\beta}\right)$$

$$= \frac{1}{\alpha} + \frac{1}{\alpha} \log Y - \frac{1}{\alpha} Y \log Y$$

$$E(\log Y) = \int_0^{\infty} \log y \cdot e^{-y} dy$$

Vemo :

$$\Gamma(x) = \int_0^{\infty} u^{x-1} e^{-u} du \quad \Rightarrow$$

$$\Gamma'(x) = \int_0^{\infty} u^{x-1} \log u \cdot e^{-u} du$$

$$\Gamma''(x) = \int_0^{\infty} u^{x-1} (\log u)^2 e^{-u} du$$

Sledi

$$E(\log Y) = \Gamma'(1) \quad E[\log Y \cdot Y^2]$$

$$E(Y \log Y) = \Gamma'(2) \quad = \Gamma'(3)$$

$$\Gamma(x+1) = x \Gamma(x) \quad \Rightarrow$$

$$\Gamma'(x+1) = \Gamma(x) + x \Gamma'(x) \quad \Rightarrow$$

$$\Gamma'(2) = \Gamma(1) + \Gamma'(1) = 1 + \Gamma'(1)$$

Dobivamo

$$E(Z_1) = \frac{1}{\alpha} + \frac{1}{\alpha} \Gamma'(1) - \frac{1}{\alpha} (1 + \Gamma'(1))$$

$$= 0 !$$

Podobno je

$$\frac{\partial}{\partial b} \log(x, \alpha, b)$$

$$= -\frac{\alpha}{b} + \frac{\alpha}{2} \left(\frac{x}{b}\right)^{\alpha}, \text{ torej}$$

$$z_2 = \frac{\alpha}{2} \left(-1 + \left(\frac{x}{b}\right)^{\alpha}\right)$$

$$= \frac{\alpha}{b} (-1 + Y),$$

torej $E(z_2) = 0$.

Računamo še

$$\frac{\partial^2}{\partial x \partial b} \log f(x, \alpha, b)$$

$$= -\frac{1}{b} + \frac{1}{2} \left(\frac{x}{b}\right)^{\alpha} \left(1 + \log\left(\frac{x}{b}\right)^{\alpha}\right)$$

Zamejamo $x \neq X$ in

izračunamo pričakovano vrednost:

$$E(W_{1,2}) = E\left[-\frac{1}{b} + \frac{1}{2} Y (1 + \log Y)\right]$$

$$= \frac{1}{6} E[Y \log Y]$$

$$= \frac{1}{6} \Gamma'(2)$$

Previous :

$$\text{cov}(z_1, z_2)$$

$$= \text{cov}\left(\frac{1}{2} \log Y - \frac{1}{2} Y \log Y, \frac{\alpha}{6} Y\right)$$

$$= \frac{1}{6} \text{cov}(\log Y - Y \log Y, Y)$$

$$= \frac{1}{6} [E(Y \log Y - Y^2 \log Y)$$

$$- E[\log Y - Y \log Y] E(Y)]$$

$$= \frac{1}{6} [\Gamma'(2) - \Gamma'(3)$$

$$- \Gamma'(1) + \Gamma'(2)]$$

$$\Gamma'(3) = 1 + 2 \Gamma'(2)$$

$$= \frac{1}{6} [-1 - \Gamma'(1)] = -\frac{1}{6} \Gamma'(2) \checkmark$$

Za izpeljavo bomo uporabili
Taylorjevo vrsto. Aproximacija
bo potekala v več korakih.

Najprej nekaj terminologije in
predpostavke:

Definicija:

- (i) Funkcija $\underline{\theta} \mapsto f(\underline{x}, \underline{\theta})$
imenujemo funkcija verjetja
in označimo z $L(\underline{\theta} | \underline{x})$
(angl. likelihood function)
- (ii) Funkcija $\underline{\theta} \mapsto \log_+(L(\underline{\theta} | \underline{x}))$
imenujemo logaritmsko
funkcijo verjetja. Označi: $l(\underline{\theta} | \underline{x})$.

Pri izpeljavi se bomo omejili
na eno dimenzijo z $\theta \in (a, b)$.

Korak 1 : Predpostavimo, da so x_1, x_2, \dots, x_n uzevec iz gostote $f(x, \theta)$ in da vedno obstaja enolično doloceno maksimum $\hat{\theta}$.

Dvomba : "Vedno" pomeni z verjetnostjo 1. Fiksirujmo $\theta_0 \in (a, b)$ in predpostavimo, da so x_1, x_2, \dots, x_n generirani iz $f(x, \theta_0)$. Velja

$$l'(\hat{\theta} | \underline{x}) = 0$$

Po Taylorjevi formuli je

$$l'(\theta_0 | \underline{x}) = \underbrace{l'(\hat{\theta} | \underline{x})}_{=0} + (\theta_0 - \hat{\theta}) l''(\hat{\theta} | \underline{x}) + R(\underline{x}, \theta_0)$$

Korak 2 : Ignoriramo R i
zaprimamo

$$\hat{\theta} - \theta_0 \approx - \frac{l'(\theta_0 | \underline{x})}{l''(\hat{\theta} | \underline{x})}$$

Korako x_1, x_2, \dots, x_n nezavisne, je

$$l(\theta_0 | \underline{x}) = \sum_{k=1}^n \log_e f(x_k, \theta)$$

$$l'(\theta_0 | \underline{x}) = \sum_{k=1}^n \left. \frac{d}{d\theta} \log_e f(x_k, \theta) \right|_{\theta=\theta_0}$$

Nadomestimo uale x + vel. uimix:

$$l'(\theta_0 | \underline{x}) = \sum_{k=1}^n \underbrace{\left. \frac{d}{d\theta} \log_e f(x_k, \theta) \right|_{\theta=\theta_0}}_{z_k}$$

$$= \sum_{k=1}^n z_k$$

Izrek 3.1 nam pove, da je
 $E(z_k) = 0$. Poleg tega so
 z_1, z_2, \dots, z_n enako porazdeljene
in neodvisne.

Korak 3: Ker je $\hat{\theta}$ "blizu" θ_0
v imenovalcu $\hat{\theta}$ nadomestimo z
 θ_0 . Zapišemo

$$l''(\theta_0 | \underline{x}) = \sum_{k=1}^n \underbrace{\frac{d^2}{d\theta^2} \log f(x_k, \theta)}_{w_k} \Big|_{\theta=\theta_0}$$

Nadomestimo x z X .

$$l''(\theta_0 | \underline{x}) = \sum_{k=1}^n w_k.$$

Izrek 3.1 pove, da je $E(w_k) =$
 $= -\text{var}(z_k)$

Korak 4 : 1+ veljati smo

$$\sqrt{n} (\hat{\theta} - \theta_0) \stackrel{\sim}{=} \frac{\frac{1}{\sqrt{n}} \sum_{k=1}^n z_k}{\frac{1}{n} \sum_{k=1}^n w_k}$$

Ker je $E(z_k) = 0$, je v_0

centralnem limitnem izreku
sterca

približno normalno

porazdeljen. Imenovalec pa

je "približno" konstanta

$E(w_k)$. Če normalno porazdeljeno

shičjavo spremenljivko delimo

s konstanto, je porazdelitev še

vedno normalna. Tukaj

"približno" normalno delimo

s "skraj" konstanto.

Variance itevca je $\text{var}(z_1)$,

konstanta v imenovalec pa

$E(W_1) = -\text{var}(z_1)$. Torej je

$$\sqrt{n}(\hat{\theta} - \theta_0) \approx N\left(0, -\frac{1}{E(W_1)}\right)$$

Sklep: Karoli je θ_0 , je

$\sqrt{n}(\hat{\theta} - \theta_0)$ približno normalno porazdeljena (kot cevilka!).

Komentar: $\hat{\theta}$ ni nujno

nepistranska. Izpoljavo nam v

resnici da srednjo kvadratično

napako.

Definicija: Naj bo $f(x, \underline{\theta})$

gostota. Definiciramo

$$(\underline{\theta} = (\theta_1, \theta_2, \dots, \theta_m))$$

$$w_{i,j} = \frac{\partial^2}{\partial \theta_i \partial \theta_j} \log_+ f(x, \underline{\theta})$$

Matrisko

$$\underline{I}(\underline{\theta}) = \left(-E(w_{i,j}) \right)_{i,j=1}^m$$

imeenujemo Fisherjeva matrika informacije.

Komentarja:

(i) V eni dimenziji je $I(\underline{\theta})$ število odvisno od $\underline{\theta}$.

(ii) Če v primeru

učetvažsežnega $\underline{\theta}$ definiravamo

$$z_i = \frac{z}{\sigma_i} \log_+ f(\underline{x}_i, \underline{\theta})$$

dobivamo, da je

$$E(W_{i,j}) = -\text{cov}(z_i, z_j)$$

V splošnem imamo več parametrov. Enaki ukazi kot za eno dimenzijo z več pisavja nam dajejo

$$\sqrt{n} (\hat{\underline{\theta}} - \underline{\theta}_0) \xrightarrow{d} N(0, \underline{I}(\underline{\theta}_0)^{-1})$$

komentar: Strogo formulacija bom objavil na spletni strani.

Primer : za Weibullovo
gostoto je

$$l(\alpha, b | x) = \log \alpha - \alpha \log b + \\ + \alpha \log x - \left(\frac{x}{b}\right)^\alpha$$

Odvajamo

$$\frac{\partial^2 l}{\partial \alpha^2} = -\frac{1}{\alpha^2} - \left(\frac{x}{b}\right)^\alpha \log^2 \frac{x}{b}$$

$$= -\frac{1}{\alpha^2} - \frac{1}{\alpha^2} \left(\frac{x}{b}\right)^\alpha \log^2 \left(\frac{x}{b}\right)^\alpha$$

$$\frac{\partial^2 l}{\partial \alpha \partial b} = -\frac{1}{b} + \frac{1}{b} \left(\frac{x}{b}\right)^\alpha \left(1 + \log \left(\frac{x}{b}\right)^\alpha\right)$$

$$\frac{\partial^2 l}{\partial b^2} = -\frac{\alpha(1+\alpha)}{b^2} \left(\frac{x}{b}\right)^\alpha + \frac{\alpha}{b^2}$$

Умноживаем $Y = \left(\frac{x}{b}\right)^\alpha \sim \exp(\lambda)$

Добиваем

$$I_{11}(\alpha, b) = \frac{1}{\alpha^2} (1 + E[Y \log^2 Y])$$

$$I_{12}(\alpha, b) = -\frac{1}{b} (-1 + E[Y(1 + \log Y)])$$

$$I_{22}(\alpha, b) = \frac{\alpha(\alpha+1)}{b^2} E(Y) - \frac{\alpha}{b^2}$$

Следя

$$\underline{I}(\alpha, b) = \begin{pmatrix} \frac{1}{\alpha^2} (1 + \Gamma''(2)) & * \\ -\frac{1}{b} \Gamma'(2) & \frac{\alpha^2}{b^2} \end{pmatrix}$$

Шаг 3 :

$$\sqrt{n} ((\hat{\alpha}, \hat{b}) - (\alpha_0, b_0))$$

$$\sim N(0, \underline{I}^{-1}(\alpha_0, b_0))$$

V končni fazi zamenjamo netična α_0, β_0 z ocenama $\hat{\alpha}_0$ in $\hat{\beta}_0$.

S tem lahko vidimo

$$\widehat{\text{MSE}}(\hat{\alpha}) = \frac{I_{11}^{-1}(\hat{\alpha}, \hat{\beta})}{n} \quad \text{in}$$

$$\widehat{\text{MSE}}(\hat{\beta}) = \frac{I_{22}^{-1}(\hat{\alpha}, \hat{\beta})}{n}$$

Ti ocenici in dejstva, da sta $\hat{\alpha}$ in $\hat{\beta}$ aproksimativno normalni nam omogočata izjavo o točnosti vzorčnih ocenilk.

Zaključne opombe

- (i) Metoda največjega verjetja je najpogostejša, ne pa edina. Je pa najbolj uporabna in prožna metoda.
- (ii) Cenilka po metodi največjega verjetja niso nikoli nepristranske. Če je $\hat{\theta}_n$ cenilka na osnovi n opazovanih vrednosti, potem

$$P(|\hat{\theta}_n - \theta_0| > \varepsilon) \rightarrow 0, \text{ ko } n \rightarrow \infty$$

za vsak $\varepsilon > 0$. V otulkih

$$\hat{\theta}_n \xrightarrow{P} \theta_0. \text{ Rečemo, da je}$$

cenilka dosledna.

4. Preizkušanje domnev

4.1. Primeri in terminologija

Primer: Recimo, da želimo ugotoviti ali je ruketni cilindrični "poščen". To bi pomenilo, da so vsi izidi ... enako vezjetni, zaporedni izidi pa med sabo neodvisni.

Že preitkus potrebujemo podatke. Recimo, da imamo n izidov x_1, x_2, \dots, x_n . Matematično je predpostavka, da so x_1, x_2, \dots, x_n nastale kot slučajne spremenljivke. Privzeli bomo, da so te slučajne spremenljivke x_1, x_2, \dots, x_n med sabo neodvisne, ničev pa imajo

vse slučajne porazdelitve z možnimi vrednostmi v $\{0, 1, 2, \dots, 36\}$.

Predpostavljamo $P(X_i = l) = p_l$ za $l = 0, 1, \dots, 36$;

Kako ugotoviti ali je rulet ni cilindev pošten?

Upoštevanje: Ali bomo lahko trdili, da je, z gotovostjo na podlagi nekajinčih izidov? Ne!

Označimo z $N_i(u)$ število pojavitev izida i v u igrah. Vemo, da je

$$N_i(u) \sim \text{Bin}(u, p_i)$$

$$\text{za } i = 0, 1, \dots, 36;$$

Po tem, kar vemo o binomski
porazdelitvi: je

$$E[N_i(n)] = np_i$$

in

$$\text{var}(N_i(n)) = np_i(1-p_i)$$

Centralni limitni izrek pove,
da lahko pričakujemo, da

$$\text{bo } N_i(n) \approx np_i \pm 2.56 \times \sqrt{np_i(1-p_i)}$$

z verjetnostjo 0.99. To velja

za fiksni i , možnih izidov

pa je 37. V grobem lahko

rečemo, da za vsak izid i

pričakujemo, da bo $N_i(n)$

„nekje okrog“ pričakovane vrednosti,

vedno le za nekaj standardnih

odklonov stran.

loleja: Na nek način moramo "izmeriti", koliko se dejanski podatki razhajajo od tega, kar bi pričakovali od postregega cilindra.

Tipično v statistiki velikost odklona merimo s standardnimi odkloni in odklone merimo s kvadrati. To je vodilo angleškega statistika Karla Pearsona (1857-1936), da je definiral slučajno spremenljivko

$$\chi^2(n) = \sum_{i=0}^{36} \frac{(N_i(n) - n \cdot p_i)^2}{n \cdot p_i}$$

Ta količina nekako "mevi",
kako daleč so izidi od
"pričakovanih", če privzamemo,
da so verjetnosti p_0, p_1, \dots, p_{36} .

V konkretnem primeru je
 $p_0 = p_1 = \dots = p_{36}$.

Komentar: V imenovalcu niso
variance sl. spr. $N_i(u)$. Spustimo
jih iz matematičnih razlogov, da
ima $\chi^2(u)$ porazdelitev, ki jo
znamo opisati in aproksimirati.

Manjka nam še naslednje:

kde je ta meva vzbujanja
"prevelika", da bi se verjeli,

da je $p_0 = p_1 = \dots = p_{36}$.

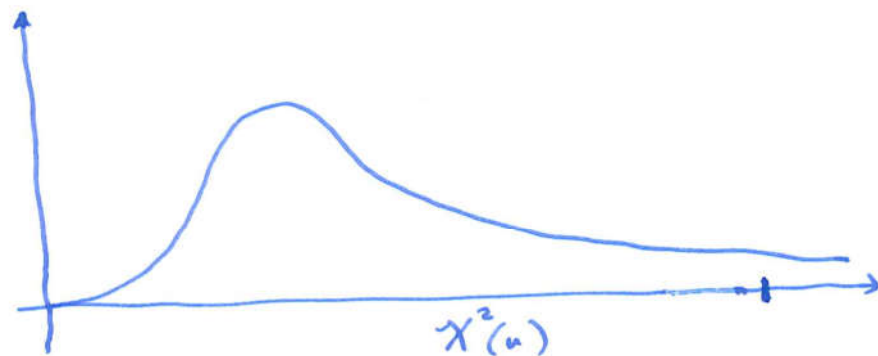
Oglejmo si simulacije. Zavedati se moramo, da je $\chi^2(n)$

stohastična spremenljivka. Iz poglavja o centralnem limitnem izreku vemo, da

$$\chi^2(n) \xrightarrow{d} \chi^2(36) = \Gamma(18, 1/2)$$

To pomeni, da imo za "končne" n $\chi^2(n)$ aproksimativno $\chi^2(36)$ porazdelitev.

Slika :



Statistični software da

$$P(X^2(n) \leq 50.99) \doteq 0.95$$

$$P(X^2(n) \leq 58.62) \doteq 0.99$$

$$P(X^2(n) \leq 67.99) \doteq 0.999$$

$$P(X^2(n) \leq 76.36) \doteq 0.9999$$

Sklep: Aproximativna porazdelitev nam pomaga, da precejimo ali in koliko je "mera razhajanja" prevelika. Kolaj bomo smatrali, da je verjetnost take velike $X^2(n)$ prevelika, pa je stvar presoje. Tipično se bomo odločili za nek prag c_α in rekli, da je razhajanje preveliko, če $X^2(n) \gg c_\alpha$

Nekaj terminologije in definicij

(i) Vedno bomo privzeli statistični model. Rekl. bomo, da so podatki \underline{x} nastali kot slučajni vektor \underline{x} s gostoto $f_{\underline{x}}(\underline{x}, \underline{\theta})$ ali v diskretnem primeru porazdelitvi $P(\underline{x} = \underline{x})$. Parameter $\underline{\theta}$ bo iz neke množice $\underline{\theta} \in \Theta$.

Primer: V primeru ruletnega cilindra je $\underline{x} = (x_1, x_2, \dots, x_n)$, kjer so komponente neodvisne in je parameter $\underline{\theta} = p \in \Delta$
z $\Delta = \{p : p_i \geq 0, \sum_i p_i = 1\}$.

(ii) Radi bi preizkusili tujitev
o parametra θ . Tipično ho-
to tujitev $\theta \in \Theta_0 \subseteq \Theta$. Tej
tujitvi bomo rekli ničelna
domneva in označili

$$H_0: \theta \in \Theta_0.$$

Alternativa je alternativna domneva

$$H_1: \theta \in \Theta \setminus \Theta_0.$$

Primer: V primeru valitnega
cilindra je Θ_0 singleton
 $\{(1/3, \dots, 1/3)\}$.

(iii) Med različnimi, torej
slučajni spremenljivimi, ki
med različnimi med
tem, kaj pričakujemo in

kar smo videli v podetkih,
rečemo testna statistika.

Tipično je to enovrstna
stohastična spremenljivka, ki
jo bomo označili s T .

Primer: Za valjeni cilindri
je testna statistika enaka $\chi^2(n)$.

(iv) Izbrati si moramo, kaj
bo za nas "majhna verjetnost".
Tej itkivi rečemo stopnja
tvegavanja $\alpha \in (0, 1)$.

Komentar: Zgodovinsko so
tipične itkive $\alpha = 0.05, 0.01$
in 0.001 .

(v) Meva razhajaja T bo
"prevelika", če bo padla
v kritično območje C_α .

Pri tem potrebujemo, da bo za
 $\underline{\theta} \in \Theta_0$

$$P_{\underline{\theta}}(T \in C_\alpha) \leq \alpha.$$

Komentar: Če je Θ_0 singleton,
lahko zahtevamo tudi enakost

$$P_{\underline{\theta}_0}(T \in C_\alpha) = \alpha.$$

Primer: Za valjni cilindri
in $X^2(u)$ je $C_\alpha = [C_\alpha, \infty)$ z
 $P(X^2(36) \geq C_\alpha) = \alpha$.

(vi) Če je $T \in C_\alpha$ večemo,
da ničelno domnevo zavrujemo
pri stopnji tveganja α .

Komentar: z zavrnitvijo nisimo
dokončno potrdili H_0 , le go
najboljših močeh smo se
odločili.

(vii) Funkciji

$$F(\underline{\theta}) = P_{\underline{\theta}}(T \in C_\alpha)$$

večemo moč testa. Za
 $\theta \neq \theta_0$ nam pove, s koliko
verjetnostjo se bomo
odločili prav v tem primeru

(viii) Če zavrnuemo H_0 pa ta
drži tagrašimo napako
I vrste. Če ne zavrnuemo
 H_0 pa ta ne drži, smo
zagrešili napako II vrste.

Prvo verjetnost poskušamo
nadzorovati, da je α ,
druga verjetnost, torej
verjetnost napake II vrste pa
je moč testa.

Primer: Recimo, da so
podatki x_1, x_2, \dots, x_n nastali
kot neodvisne, enako
porazdeljene slučajne
spremenljivke x_1, x_2, \dots, x_n

2 $X_i \sim N(\mu, \sigma^2)$. Vprašamo se

$$H_0: \mu = \mu_0$$

$$H_1: \mu \neq \mu_0$$

V tem primeru je

$$\Theta = \mathcal{L}(\mu, \sigma^2) : \mu \in \mathbb{R}, \sigma^2 > 0 \}$$

$$\Theta_0 = \mathcal{L}(\mu_0, \sigma^2) : \sigma^2 > 0 \}.$$

Vemo, da je \bar{X} nepristranska
cenilka μ in

$$S^2 = \frac{1}{n-1} \sum_{k=1}^n (X_k - \bar{X})^2$$

nepristranska cenilka σ^2 .

Po izloji, da razdalje v
statistični merini s
standardnimi odkloni

definiramo T tako, da
razdaljo med \bar{X} in μ_0
izmerimo s standardnim
odklonom \bar{X} , ki je $\frac{\hat{\sigma}}{\sqrt{n}}$.

Vendar $\hat{\sigma}$ ne poznamo, zato
razmerimo cenilko $\hat{\sigma}^2$.

Torej bi definirali

$$T = \frac{\bar{X} - \mu_0}{\sqrt{\hat{\sigma}^2/n}}$$
$$= \frac{\sqrt{n} (\bar{X} - \mu_0)}{\sqrt{\frac{1}{n-1} \sum_{k=1}^n (x_k - \bar{X})^2}}$$

Vemo, da sta \bar{X} in S^2
nezavisni in je T
povzeto deljena po t_{n-1} -
povzeto delitvi, če je $\mu = \mu_0$.

Še en kotček terminologije
(ix) Če je θ_0 singleton $\underline{\theta}_0$,
potem večemo porazdelitvi

T testna porazdelitev.

Enako velja, če je porazdelitev

T druga za vse $\theta \in \underline{\theta}_0$.

V zadnjem primeru porazdelitev

T ni odvisna od σ^2 in je

za vsak (μ_0, σ^2) druga.

Statistični software nam

za vsak $\alpha \in (0, 1)$ da c_α tako

da je $P_{\mu_0, \sigma^2} (|T| \geq c_\alpha) = \alpha$.

Kritično območje je torej

oblike $C_\alpha = (-\infty, -c_\alpha] \cup [c_\alpha, \infty)$.

Komentar: Trditev o porazdelitvi velja le, če je $\mu = \mu_0$, ker je v tem primeru

$$\sqrt{n}(\bar{X} - \mu_0) \sim N(0, \sigma^2)$$

Moč testa je odvisna od μ in σ kot funkcija

$$F(\mu, \sigma) = P\left(\left|\frac{\sqrt{n}(\bar{X} - \mu_0)}{\left(\frac{1}{n-1} \sum_{k=1}^n (X_k - \bar{X})^2\right)^{1/2}}\right| \geq c_\alpha\right)$$

za $X_1, \dots, X_n \sim N(\mu, \sigma^2)$. Moč je težko izračunljiva.

Komentar: Zgorajji preizkus

domneve je izumil

William Sealy Gossett (1876-1937)

znan pod vzdevkom Student.

(xi) ko imamo podatke x_1, x_2, \dots, x_n , statistični model in testno statistiko T , lahko izračunamo dejansko vrednost t . Če H_0 zavrnemo, če je t prevelika ali $|t|$ prevelika, lahko izračunamo

$$P_{H_0}(T \geq t)$$

ali

$$P_{H_0}(|T| \geq t).$$

Tej verjetnosti rečemo p -vrednost.

H_0 zavrnemo, če je p -vrednost manjša od α .

Komentar: p -vrednost je eden od najbolj zlorabljenih konceptov v statistiki.

Kako pa v splošnem prideemo do testnih statistik T .

Kakšna splošna metoda?

4.2 Test s kvocientom verjetij in Nilssonov izrek

Kot smo videli, je tipična situacija

$$H_0: \underline{\theta} \in \Theta_0 \quad \text{proti} \quad H_1: \underline{\theta} \in \Theta \setminus \Theta_0.$$

Recimo, da imajo podatki gostoto $f(x, \underline{\theta})$. Funkcija

$$\underline{\theta} \mapsto f(x, \underline{\theta})$$

za fiksen x interpretiramo kot
"verjetje": loka

Lahko izračunamo pri fiksnem \underline{x}

$$\max_{\underline{\theta} \in \Theta_0} f(\underline{x}, \underline{\theta}) \quad \text{ali} \quad \max_{\underline{\theta} \in \Theta} f(\underline{x}, \underline{\theta}).$$

Zapišimo drugače:

$$\max_{\underline{\theta} \in \Theta_0} L(\underline{\theta} | \underline{x}) \quad \text{in} \quad \max_{\underline{\theta} \in \Theta} L(\underline{\theta} | \underline{x}).$$

Če je drugi maksimum "dosti" večjiⁿ od prvega bi to pomenilo bolj $\underline{\theta} \in \Theta \setminus \Theta_0$, torej dokazno guardsko proti ničelni domnevi.

Definicija: kvocient

$$\Lambda = \frac{\max_{\underline{\theta} \in \Theta} L(\underline{\theta} | \underline{x})}{\max_{\underline{\theta} \in \Theta_0} L(\underline{\theta} | \underline{x})}$$

imenujemo Wilksova lambda statistika.

Komentarja:

(i) pisali smo max, čeprav vnaprej ni očitno, da bo tudi dosežen.

Preverjali bomo v konkretnih situacijah.

(ii) Statistika Λ je odvisna samo od podatkov (v okviru predpostavljenege statističnega modela). To dejstvo je pomembno, saj lahko odločamo o H_0 samo na podlagi podatkov.

Primer: Naj bodo X_1, X_2, \dots, X_n normalne + neodvisne in neodvisne in X_1, X_2, \dots, X_n utorec. v tem primeru je

$$L(\mu, \sigma^2 | \underline{x}) = \prod_{k=1}^n \frac{1}{\sqrt{2\pi} \sigma} e^{-\frac{1}{2\sigma^2} (x_k - \mu)^2}$$

Testiravamo $H_0: \mu = 0$ protiv

$H_1: \mu \neq 0$. Točnj je

$$\Theta = \{ (\mu, \sigma^2) : \mu \in \mathbb{R}, \sigma^2 > 0 \}$$

in

$$\Theta_0 = \{ (\mu, \sigma^2) : \sigma^2 > 0 \}$$

Za izvođenje maksimuma logaritmiramo

$$\begin{aligned} \ell(\mu, \sigma^2 | \underline{x}) &= \frac{n}{2} \log 2\pi - n \log \sigma \\ &\quad - \frac{1}{2\sigma^2} \sum_{k=1}^n (x_k - \mu)^2 \end{aligned}$$

Izvođenjem parcijalno odvođa

$$\begin{aligned} \frac{\partial \ell}{\partial \mu} &= \frac{1}{2\sigma^2} \cdot 2 \cdot \sum_{k=1}^n (x_k - \mu) = 0 \\ &\Rightarrow \mu = \bar{x} \end{aligned}$$

in

$$\frac{\partial \ell}{\partial \sigma} = -\frac{n}{2} + \frac{1}{\sigma^3} \sum_{k=1}^n (x_k - \mu)^2 = 0 \Rightarrow$$

$$\hat{\sigma}^2 = \frac{1}{n} \sum_{k=1}^n (x_k - \bar{x})^2$$

Če je $\mu = 0$, potem maksimiramo

$$l(0, \sigma^2 | \underline{x}) = \frac{n}{2} \log 2\pi - n \log \sigma - \frac{1}{2\sigma^2} \sum_{k=1}^n x_k^2$$

po σ . Odvajamo in sledi

$$-\frac{n}{\sigma} + \frac{1}{\sigma^3} \sum_{k=1}^n x_k^2 = 0 \quad \Rightarrow$$

$$\hat{\sigma}^2 = \frac{1}{n} \sum_{k=1}^n x_k^2$$

Po izločaji bomo zavrnili H_0 , če bo Λ "prevelika" oziroma, če bo $\log \Lambda$ prevelika..

Definicija: količini

$$\lambda = 2 \log \Lambda$$

rečemo Wilksova logaritemska lambda statistika.

17va ču u u u u

$$l(\hat{\mu}, \hat{\sigma}^2 | \underline{x})$$

$$= \frac{n}{2} \log 2\pi - n \log \hat{\sigma} - \frac{1}{2\hat{\sigma}^2} \sum_{k=1}^n (\underline{x}_k - \bar{x})^2$$

$$= \frac{n}{2} \log 2\pi - n \log \hat{\sigma} - \frac{n}{2}$$

iu

$$l(0, \tilde{\sigma}^2 | \underline{x}) = \frac{n}{2} \log 2\pi - n \log \tilde{\sigma} - \frac{n}{2}$$

Stedi

$$\lambda - 2 \log \Lambda$$

$$= 2 \left\{ \frac{n}{2} \log 2\pi - n \log \hat{\sigma} - \frac{n}{2} - \frac{n}{2} \log 2\pi + n \log \tilde{\sigma} + \frac{n}{2} \right\}$$
$$= n \log \tilde{\sigma} - n \log \hat{\sigma}$$

Нително доменува завршено, де
 во $\lambda \geq \lambda_0$ за нек λ_0 . То
 помени

$$2n \log \frac{\hat{\sigma}^2}{\sigma^2} \geq \lambda_0 \quad \text{али}$$

$$\log \frac{\hat{\sigma}^2}{\sigma^2} \leq -\frac{\lambda_0}{2n}$$

$$\begin{aligned} \frac{\hat{\sigma}^2}{\sigma^2} &= \frac{\frac{1}{n} \sum_{k=1}^n (x_k - \bar{x})^2}{\frac{1}{n} \sum_{k=1}^n x_k^2} \\ &= \frac{\frac{1}{n} \sum_{k=1}^n x_k^2 - \bar{x}^2}{\frac{1}{n} \sum_{k=1}^n x_k^2} \\ &= 1 - \frac{n \bar{x}^2}{\sum_{k=1}^n x_k^2} \end{aligned}$$

Нително доменува завршено, де
 је

$$\frac{n \bar{x}^2}{\sum_{k=1}^n x_k^2} \geq c_\alpha$$

По Cauchy-Schwarz-у је

$$\left(\sum_{k=1}^n x_k \cdot 1 \right)^2 \leq \left(\sum_{k=1}^n x_k^2 \right) \cdot n \quad \text{али}$$

$$n \bar{x}^2 \leq \sum_{k=1}^n x_k^2,$$

зато то $C_\alpha < 1$. Се означава
приметно

$$n \bar{x}^2 \geq C_\alpha \sum_{k=1}^n x_k^2 \quad \text{али}$$

$$\begin{aligned} (1 - C_\alpha) \bar{x}^2 \cdot n &\geq C_\alpha \left(\sum_{k=1}^n x_k - n \bar{x} \right)^2 \\ &\geq C_\alpha \sum_{k=1}^n (x_k - \bar{x})^2 \end{aligned}$$

али

$$\frac{\bar{x}^2 \cdot n}{\sum_{k=1}^n (x_k - \bar{x})^2} \geq \frac{C_\alpha}{1 - C_\alpha}.$$

али

$$\frac{|\bar{x}| \cdot \sqrt{n}}{\sqrt{\sum_{k=1}^n (x_k - \bar{x})^2}} \geq \sqrt{\frac{C_\alpha}{1 - C_\alpha}}.$$

Torej je $\lambda \geq \lambda_0$ enako kot
 $|T| \geq t_\alpha$ za T statistiko iz
 prejšnjega vzorca. Poleg tega
 je porazdelitev λ neodvisna
 od z . Ponovno smo predelali
 t-test!

Komentar: t_n -porazdelitev je
 definirana kot

$$T = \frac{z}{\sqrt{\frac{1}{n} u}}, \quad \text{kjer sta}$$

z, u neodvisni in $u \sim \chi^2(n), z \sim N(0,1)$

Za velike n je $\frac{1}{n} \cdot u \sim 1$, torej

bo $T^2 \sim z^2 \sim \chi^2(1)$.

Primer: Vrnilo se k primeru
 vulete. Opatovane vrednosti x_1, x_2, \dots
 x_n so neodvisne in enako
 porazdeljene na $\{0, 1, \dots, 365\}$.

Vemo, da je

$$\Theta = \left\{ (p_0, \dots, p_{36}) : p_i \geq 0, \sum_{i=0}^{36} p_i = 1 \right\}.$$

Recimo, da je $\Theta_0 = \{p^0\}$. Vemo, da je

$$l(p|x) = \sum_{k=0}^{36} n_k \log p_k,$$

kjer je n_k število pojavitev izida k .

Maksimum po Θ najdemo z

Lagrangeovo metodo: maksimiziramo

$$\sum_{k=0}^{36} n_k \log p_k \quad p_k \text{ pogju}$$

$$\sum_{k=0}^{36} p_k = 1.$$

Definiramo

$$F(p) = \sum_{k=0}^{36} n_k \log p_k - \lambda \left(\sum_{k=0}^{36} p_k - 1 \right),$$

$$\frac{\partial F}{\partial p_k} = \frac{n_k}{p_k} - \lambda = 0$$

Sledi, da je $p_k = c \cdot n_k$ za
neko konstanto. Ker se p_k
seštejejo v 1, je

$$\hat{p}_k = \frac{n_k}{n}$$

Torej je

$$L(\hat{P} | \underline{x}) = \sum_{k=0}^{36} n_k \cdot \log \frac{n_k}{n}.$$

Ker je Θ_0 singleton, nam ni
treba računati maksimuma.

$$L(P^0 | \underline{x}) = \sum_{k=0}^{36} n_k \log p_k^0.$$

Sledi, da je

$$\lambda = 2 \sum_{k=0}^{36} n_k \cdot \log \frac{n_k}{n} - 2 \sum_{k=0}^{36} n_k \cdot \log p_k^0$$

Ali z mano kaj veči o porazdelitvi

λ , če H_0 drži?

Funkciju $f(x) = x \log(x/x_0)$
razvijemo okolo točke x_0 .

$$f(x) \approx f(x_0) + f'(x_0)(x-x_0) + \frac{1}{2} f''(x_0)(x-x_0)^2$$

Računamo

$$f'(x_0) = 1$$

$$f''(x_0) = \frac{2}{x_0} - \frac{1}{x_0} = \frac{1}{x_0}$$

Prepisemo

$$\lambda = 2n \sum_{k=0}^{36} \frac{n_k}{n} \cdot \log \frac{n_k/n}{p_k^0}$$

$$= 2n \cdot \sum_{k=0}^{36} \hat{p}_k \log \frac{\hat{p}_k}{p_k^0}$$

Če n_0 dosti, bo \hat{p}_k blizu p_k^0
in bo Taylorjeva aproksimacija
dobra.

Torej bo $\chi_0 = p_k^0$

$$\lambda \approx 2n \left\{ \sum_{k=0}^{36} (\hat{p}_k - p_k^0)^2 + \frac{1}{2} \sum_{k=0}^{36} (\hat{p}_k - p_k^0)^2 / p_k^0 \right\}$$

$$= \sum_{k=0}^{36} \frac{(n\hat{p}_k - np_k^0)^2}{n \cdot p_k^0}$$

$$= \sum_{k=0}^{36} \frac{(n_k - np_k^0)^2}{n \cdot p_k^0}$$

$$= \chi^2 !$$

Torej nam da kvocient verjetij

χ^2 statistiko! Vemo tudi, da

ima v primeru, ko H_0 drži,

χ^2 aproksimativno $\chi^2(36)$

porazdelitev.

V obeh primerih smo ugotovili, da ima λ aproksimativno χ^2 porazdelitev. V obeh primerih nam je ideja vrnila statistiko, ki jih potujemo, torej deluje dobro.

Kaj pa v splošnem? Oglejmo si preprost primer, ko je

$$\Theta = (a, b) \text{ in } \Theta_0 = \{\theta_0\} \text{ z } \theta_0 \in (a, b).$$

Predpostavimo, da so opazovane vrednosti x_1, x_2, \dots, x_n vzorce iz porazdelitve z gostoto $f(x, \theta)$.

Naj bo $\hat{\theta}$ ocena θ po metodi največjega verjetja. Po

Taylorju bo

$$2l(\theta_0 | \underline{x}) - 2l(\hat{\theta} | \underline{x})$$

$$\approx \underbrace{2l'(\hat{\theta} | \underline{x})}_{=0 \text{ po def.}} + 2 \cdot \frac{1}{2} l''(\hat{\theta} | \underline{x}) (\theta_0 - \hat{\theta})^2$$

$$= \frac{1}{n} l''(\hat{\theta} | \underline{x}) \cdot [\sqrt{n}(\theta_0 - \hat{\theta})]^2$$

Vemo, da je v primeru, ko H_0 drži

$$\frac{1}{n} l''(\hat{\theta} | \underline{x})$$

$$\approx \frac{1}{n} l''(\theta_0 | \underline{x})$$

$$\approx -\text{var}(z_1) = E(w_1)$$

7 otuokami + 3. poglavja.

Ampak $\sqrt{n}(\hat{\theta} - \theta_0)$ je približno normalno porazdeljena s parametroma 0 in $\frac{1}{\text{var}(z_1)}$.

To pomeni

$$\lambda \approx \left[\sqrt{n} (\hat{\theta}_1 - \theta_0) \cdot \sqrt{\text{var}(\hat{\theta}_1)} \right]^2$$

Ampak izraz v ogledih oklepjih
ima približno standardizirano

normalno porazdelitev. To

pomeni, da ima λ aproksimativno

$\chi^2_{(1)}$ porazdelitev.

Ta grob račun nakazuje, da bi

utegnila imeti λ v splošnem

χ^2 porazdelitev.

Izrek 4.1 (Wilks) Naj bodo x_1, x_2, \dots, x_n

vzorec iz porazdelitve z gostoto $f(\underline{x}, \theta)$

z $\theta \in \Theta \subseteq \mathbb{R}^p$ odprta množica. Naj

bo $g: \Theta \rightarrow \mathbb{R}^k$ zvezna parcialno

odvedljiva z $\text{rang } Dg = k$ na \mathbb{R} .

Naj bo

$$H_0: g(\underline{\theta}) = \underline{\gamma}_0 \quad \text{in}$$

$$H_1: g(\underline{\theta}) \neq \underline{\gamma}_0.$$

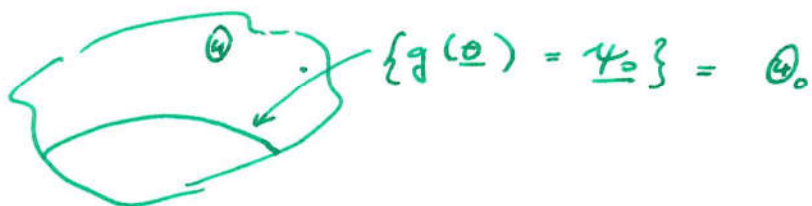
Naj bodo izpoljeni vsi pogoji izreka
o asimptotski normalnosti cenik
po metodi največjega verjetja. Potem
velja

$$\lambda \xrightarrow{d} \chi^2(p-k), \quad \text{ko } n \rightarrow \infty.$$

Komentarji:

(i) $g(\underline{\theta}) = \underline{\gamma}_0$ pomeni, da $\underline{\theta}$ leži
na gladni ploskvi v Θ .

Slika :



(ii) Tipično lakuo H_0 in H_1
preformuliramo v

$$H_0: \underline{\theta} \in \underline{\Theta}_0 \quad \text{in} \quad H_1: \underline{\theta} \in \underline{\Theta} \setminus \underline{\Theta}_0.$$

Tipično je $\underline{\Theta}_0$ neka podmnožitevost.

Rečemo lakuo, da je λ

povazdelitev približno $\chi^2(r)$,

kjer je $r = \dim \underline{\Theta} - \dim \underline{\Theta}_0$.

(iii) Implikacija izreka je, da aproks.

povazdelitev λ ni odvisna od

$\underline{\theta}_0 \in \underline{\Theta}_0$ oziroma od $\underline{\theta}_0$, ta

katerega je $\underline{g}(\underline{\theta}_0) = \underline{\gamma}_0$.

Zato je λ ustrezna statistika.

REGRESIJA

4.1. Uvodni primeri

Uvodni primeri so v elementarnih kvadratih objavljeni na spletni strani predmeta.

4.2 Ocenjevanje parametrov, izrek

Gauss-Markova

Iz primerov razberemo, da bomo prišli, da so podatki y_1, y_2, \dots, y_n nastali kot

$$y_k = \alpha + \beta x_k + \varepsilon_k,$$

kjer so x_1, x_2, \dots, x_n konstante, za katere privzamemo, da so znane. Parametra α in β sta neznanata, za slučajne spremenljivke $\varepsilon_1, \varepsilon_2, \dots, \varepsilon_n$ pa privzamemo, da je

$$E(\varepsilon_k) = 0, \quad \text{cov}(\varepsilon_k) = \sigma^2 \quad \text{za}$$

vse $k = 1, 2, \dots, n$ in

$$\text{cov}(\varepsilon_k, \varepsilon_l) = 0 \quad \text{za } k \neq l.$$

Komentar: Zgorajne predpostavke večino predpostavke standardne linearne regresije.

V matričnem obliki lahko zapišemo:

$$\underline{X} = \begin{pmatrix} 1 & x_1 \\ 1 & x_2 \\ \vdots & \vdots \\ 1 & x_n \end{pmatrix} \quad \underline{Y} = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix} \quad \underline{\varepsilon} = \begin{pmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \vdots \\ \varepsilon_n \end{pmatrix}$$

$$\underline{\varepsilon} = \begin{pmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \vdots \\ \varepsilon_n \end{pmatrix} \quad \text{S tem označujemo so}$$

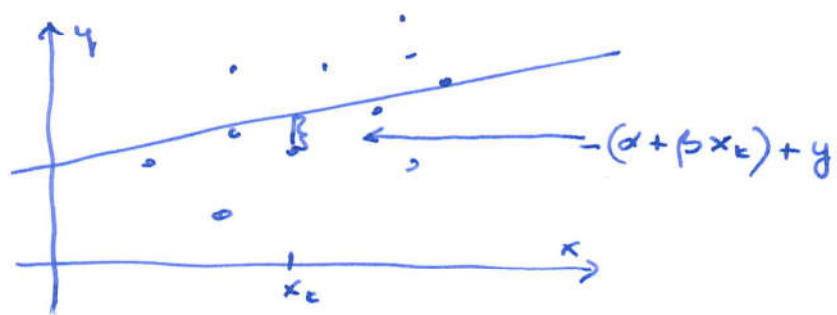
predpostavke $E(\underline{\varepsilon}) = 0$, $\text{var}(\underline{\varepsilon}) = \sigma^2 \underline{I}$.

Predpostavljamo

$$\underline{Y} = \underline{X} \cdot \begin{pmatrix} \alpha \\ \beta \end{pmatrix} + \underline{\varepsilon}$$

1. delja ocenjevanja:

Slika:



Vsaka točka ima navpično razdaljo do premice $\alpha + \beta x$.

Iščemo premico, ki se bo čim bolj
"prilegala" točkam. Sledimo

Gausso in v nadalje kvadriramo in
seštejemo. Iščemo premico, za
katere je vsota kvadratov vzdaljš
najmanjša možna. Definiavimo

$$s(\alpha, \beta) = \sum_{k=1}^n (y_k - \alpha - \beta x_k)^2$$

Parcialno odvajamo po α in β in
izenačimo z 0.

$$\frac{\partial s}{\partial \alpha} = -2 \sum_{k=1}^n (y_k - \alpha - \beta x_k) = 0$$

$$\frac{\partial s}{\partial \beta} = -2 \sum_{k=1}^n (y_k - \alpha - \beta x_k) x_k = 0$$

Pohiimo linearni enačbi za α
in β . Zapišimo v matrični
obliki.

$$n\alpha + \beta \sum_{k=1}^n x_k = \sum_{k=1}^n y_k$$

$$\alpha \cdot \sum_{k=1}^n x_k + \beta \sum_{k=1}^n x_k^2 = \sum_{k=1}^n x_k y_k$$

ali

$$\underbrace{\begin{pmatrix} n & \sum_{k=1}^n x_k \\ \sum_{k=1}^n x_k & \sum_{k=1}^n x_k^2 \end{pmatrix}}_{\underline{X}^T \cdot \underline{X}} \begin{pmatrix} \alpha \\ \beta \end{pmatrix} = \underbrace{\begin{pmatrix} \sum_{k=1}^n y_k \\ \sum_{k=1}^n x_k y_k \end{pmatrix}}_{\underline{X}^T \cdot \underline{y}}$$

Predpostavimo, da ima \underline{X} poln rang.

Potem bo rešitev zgorajje enačbe

$$\begin{pmatrix} \hat{\alpha} \\ \hat{\beta} \end{pmatrix} = (\underline{X}^T \underline{X})^{-1} \cdot \underline{X}^T \cdot \underline{y}$$

Rečemo, da sta $\hat{\alpha}$ in $\hat{\beta}$ oceni α in β po metodi najmanjših kvadratov. Opazimo se, da

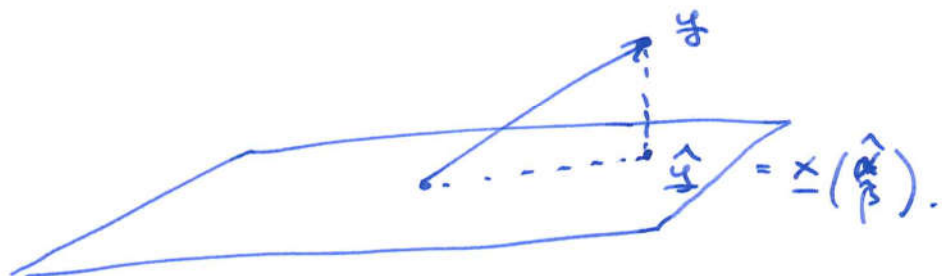
je

$$\underline{X} \begin{pmatrix} \hat{\alpha} \\ \hat{\beta} \end{pmatrix} = \underline{X} (\underline{X}^T \underline{X})^{-1} \underline{X}^T \underline{y}$$

Matrika $\underline{X}(\underline{X}^T \underline{X})^{-1} \underline{X}^T$ je ortogonalna
 projekcija na podprostor v \mathbb{R}^n ,
 ki ga ustvarjajo stolpca \underline{X} .
 Torej sta $\hat{\alpha}$ in $\hat{\beta}$ evklidovski
 rešitvi enačbe

$$\underline{X} \begin{pmatrix} \hat{\alpha} \\ \hat{\beta} \end{pmatrix} = \underline{X} (\underline{X}^T \underline{X})^{-1} \underline{X}^T \underline{y} = \hat{\underline{y}}$$

Slika :



Če zamenjamo \underline{y} z \underline{y} , dobimo
 slučajni vektor $\begin{pmatrix} \hat{\alpha} \\ \hat{\beta} \end{pmatrix}$, ki mu
 učimo cenilka. Iz matričnega
 zapisa hitro sledijo nekaj
 osnovnih trditev.

Če je
$$\begin{pmatrix} \hat{\alpha} \\ \hat{\beta} \end{pmatrix} = (\underline{X}^T \underline{X})^{-1} \underline{X}^T \underline{Y}$$
 je

$$E \left[\begin{pmatrix} \hat{\alpha} \\ \hat{\beta} \end{pmatrix} \right] = (\underline{X}^T \underline{X})^{-1} \underline{X}^T \underline{X} \cdot \begin{pmatrix} \alpha \\ \beta \end{pmatrix} \\ = \begin{pmatrix} \alpha \\ \beta \end{pmatrix}$$

Cenilka po metodi najmanjših kvadratov je nepristranska.

$$\text{var} \left(\begin{pmatrix} \hat{\alpha} \\ \hat{\beta} \end{pmatrix} \right) = (\underline{X}^T \underline{X})^{-1} \underline{X}^T \underbrace{\text{var}(\underline{Y})}_{\sigma^2 \mathbf{I}} \underline{X} (\underline{X}^T \underline{X})^{-1} \\ = \sigma^2 \cdot (\underline{X}^T \underline{X})^{-1}$$

Sled tovej

$$\text{var}(\hat{\alpha}) = \sigma^2 \cdot \frac{\sum_{k=1}^n x_k^2}{n \sum_{k=1}^n x_k^2 - \left(\sum_{k=1}^n x_k \right)^2}$$

$$\text{var}(\hat{\beta}) = \sigma^2 \cdot \frac{n}{n \sum_{k=1}^n x_k^2 - \left(\sum_{k=1}^n x_k \right)^2}$$

Ostane še uprati, kako oceniti σ^2 .

V splošnem obpustimo več stolpcev v \underline{x} . Podatki so oblike

$$\begin{array}{l} y_1 \quad x_{11} \dots x_{1m} \\ y_2 \quad x_{21} \quad \dots x_{2m} \\ \vdots \\ y_n \quad x_{n1} \quad \dots x_{nm} \end{array}$$

Pričakujemo, da so x_{kj} fiksne konstante, y_1, y_2, \dots, y_n pa so "nastali" kot

$$y_k = \beta_1 x_{k1} + \beta_2 x_{k2} + \dots + \beta_m x_{km} + \varepsilon_k.$$

Opomba: x_{k1} so lahko 1, ne zahtevamo pa tega nikoli.

Pišemo matrično:

$$\underline{x} = \begin{pmatrix} x_{11} & \dots & x_{1m} \\ x_{21} & \dots & x_{2m} \\ \vdots & & \vdots \\ x_{n1} & \dots & x_{nm} \end{pmatrix}, \quad \underline{y} = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix}$$

$$\underline{y} = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix}, \quad \underline{\beta} = \begin{pmatrix} \beta_1 \\ \beta_2 \\ \vdots \\ \beta_m \end{pmatrix}, \quad \underline{\varepsilon} = \begin{pmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \vdots \\ \varepsilon_n \end{pmatrix}.$$

Predpostavka standardnega regresijskega modela zapišemo kot

$$\underline{y} = \underline{X}\underline{\beta} + \underline{\varepsilon}$$

2

$$E(\underline{\varepsilon}) = 0 \quad \text{in} \quad \text{var}(\underline{\varepsilon}) = \sigma^2 \cdot \underline{I}$$

Opomba: ves čas predpostavljamo, da je \underline{X} matrika s polnim rangom.

Locija za ocenjevanje $\underline{\beta}$ je enaka kot prej. Iščeemo $\beta_1, \beta_2, \dots, \beta_m$, da bo

$$s(\underline{\beta}) = \sum_{k=1}^n (y_k - \beta_1 x_{k1} - \beta_2 x_{k2} - \dots - \beta_m x_{km})^2$$

najmanjša možna.

Parcialno odvajamo in dobimo

$$\frac{\partial S}{\partial \beta_i} = -2 \sum_{k=1}^n (y_k - \beta_1 x_{k1} - \dots - \beta_m x_{km}) x_{ki}$$

Izenačimo $\neq 0$ in zapišemo v matričnem obliki (pokujimo - 2):

$$(\underline{X}^T \cdot \underline{X}) \cdot \underline{\beta} = \underline{X}^T \cdot \underline{y}$$

Ocena $\underline{\beta}$ bo

$$\hat{\underline{\beta}} = (\underline{X}^T \underline{X})^{-1} \underline{X}^T \cdot \underline{y}$$

kot cenilka bo

$$\hat{\underline{\beta}} = (\underline{X}^T \underline{X})^{-1} \underline{X}^T \cdot \underline{y}$$

Izrek 5.1 : Naj bodo izpolnjene predpostavke standardnega linearnega regresijskega modela.

Ocenilka po metodi najmanjših kvadratov je nepristranska z

$$\text{var}(\hat{\underline{\beta}}) = \sigma^2 \cdot (\underline{X}^T \underline{X})^{-1}.$$

Dokaz: Enak kot prej.

Definicija: Vektorju $\hat{\underline{y}} = \underline{X} \cdot \hat{\underline{\beta}}$ večemo vektor prilagojenih vrednosti. Vektorju $\hat{\underline{\varepsilon}} = \underline{y} - \hat{\underline{y}}$ večemo vektor ostankov ali residualov.

Opomba: Tudi v več dimenzijah je $\hat{\underline{y}}$ ortogonalna projekcija \underline{y} na podprostor, ki ga najemajo stolpci \underline{X} .

Oceniti moramo še σ^2 .

Iz predpostavk bi na čelo ma

sledilo $\frac{1}{n} \sum_{k=1}^n \varepsilon_k^2 \approx \sigma^2$.

Vendar ε_k ne poznamo, poznamo
pa $\hat{\varepsilon}_k$. Izračunajmo

$$\begin{aligned}
 E\left(\sum_{k=1}^n \hat{\varepsilon}_k^2\right) &= E\left[\hat{\underline{\varepsilon}}^T \cdot \hat{\underline{\varepsilon}}\right] \\
 &= E\left[(\underline{y} - \underline{x} \cdot \hat{\underline{\beta}})^T (\underline{y} - \underline{x} \cdot \hat{\underline{\beta}})\right] \\
 &= E\left[(\underline{y} - \underline{x} \cdot (\underline{x}^T \underline{x})^{-1} \underline{x}^T \underline{y})^T (\underline{y} - \underline{x} \cdot (\underline{x}^T \underline{x})^{-1} \underline{x}^T \underline{y})\right] \\
 &= E\left[\underline{y}^T \cdot \underbrace{(\underline{I} - \underline{x} (\underline{x}^T \underline{x})^{-1} \underline{x}^T)}_{\text{Idempotentna matrika}} (\underline{I} - \underline{x} (\underline{x}^T \underline{x})^{-1} \underline{x}^T) \underline{y}\right] \\
 &= E\left[\underline{y}^T (\underline{I} - \underline{x} (\underline{x}^T \underline{x})^{-1} \underline{x}^T) \underline{y}\right] \\
 &= E\left[(\underline{y} - \underline{x} \underline{\beta})^T (\underline{I} - \underline{x} (\underline{x}^T \underline{x})^{-1} \underline{x}^T) (\underline{y} - \underline{x} \underline{\beta})\right] \\
 &\quad \text{ker je } (\underline{I} - \underline{x} (\underline{x}^T \underline{x})^{-1} \underline{x}^T) \underline{x} = 0 \\
 &= E\left[\underline{\varepsilon}^T (\underline{I} - \underline{x} (\underline{x}^T \underline{x})^{-1} \underline{x}^T) \underline{\varepsilon}\right] \\
 &= E\left[\text{sr}((\underline{I} - \underline{x} (\underline{x}^T \underline{x})^{-1} \underline{x}^T) \underline{\varepsilon} \cdot \underline{\varepsilon}^T)\right]
 \end{aligned}$$

$$= \text{SE} \left(\underbrace{(\underline{I} - \underline{X}(\underline{X}^T \underline{X})^{-1} \underline{X}^T)}_{\text{(linearnost)}} \underbrace{E(\underline{\varepsilon} \underline{\varepsilon}^T)}_{\text{var}(\underline{\varepsilon}) = \sigma^2 \cdot \underline{I}} \right)$$

$$= \text{SE} \left(\sigma^2 \cdot (\underline{I} - \underline{X}(\underline{X}^T \underline{X})^{-1} \underline{X}^T) \right)$$

$$= \sigma^2 \text{SE} \left((\underline{I} - \underline{X}(\underline{X}^T \underline{X})^{-1} \underline{X}^T) \right)$$

Za idempotentne matrice je sled
 eneke ranga, $\underline{I} - \underline{X}(\underline{X}^T \underline{X})^{-1} \underline{X}$ pa
 je ortogonalna projekcija na
 komplement podprostora, ki ga
 naprejo stolpci \underline{X} . Dimenzija
 tega podprostora je $n-m$.

Torej je

$$E\left(\sum_{k=1}^n \hat{\varepsilon}_k^2\right) = \sigma^2(n-m)$$

Izrek 4.2: Ceniška

$$\hat{\sigma}^2 = \frac{1}{n-m} \sum_{k=1}^n \hat{\varepsilon}_k^2$$

je nepristranska ceniška σ^2 .

Dokaz: Sumo \bar{z} .

Ali lahko kaj rečemo o tem, koliko je "dobra" cenička $\hat{\beta}$?

Ali morda obstaja boljša cenička?

Pri takih vprašanjih moramo definirati razred "tekmic".

Glede na to, da je cenička po metodi najmanjših kvadratov linearna funkcija \underline{y} se zdi smiselno, da najprej pogledamo linearne tekunice, potem pa še zahtevamo nepristranskost. Primerjali bomo

$\hat{\beta} = (\underline{X}^T \underline{X})^{-1} \cdot \underline{X}^T \cdot \underline{y}$ in tekunice oblike

$\tilde{\beta} = \underline{L} \cdot \underline{y}$, kjer je $\underline{L} (n \times m)$

matrica (ki je lahko odvisna od \underline{X}).

Ker zahtevamo nepristranskost, mora veljati

$$E(\tilde{\beta}) = E(\underline{L} \cdot \underline{y}) = \underline{L} \cdot \underline{X} \beta = \beta$$

To pomeni, da je $\underline{L} \cdot \underline{X} = \underline{I}$.

Izrek 4.3 (Gauss-Markov) Ceniška po metodi najmanjših kvadratov je najboljša med vsemi linearnimi nepristranskimi ceniškami $\hat{\beta}$.

Komentar: Kaj pomeni najboljša?

Pomen bomo razjasnili po dokazu.

Dokaz: Naj bo $\tilde{\beta} = \underline{L} \cdot \underline{Y}$ nepristranska.

Računamo

$$\begin{aligned} \text{var}(\tilde{\beta}) &= \text{var}(\tilde{\beta} - \hat{\beta} + \hat{\beta}) \\ &= \text{var}(\hat{\beta}) + \text{var}(\tilde{\beta} - \hat{\beta}) \\ &\quad + \text{cov}(\tilde{\beta} - \hat{\beta}, \hat{\beta}) \\ &\quad + \text{cov}(\hat{\beta}, \tilde{\beta} - \hat{\beta}) \end{aligned}$$

Posebej računamo

$$\text{cov}(\hat{\beta} - \tilde{\beta}, \hat{\beta})$$

$$= \text{cov}(\underline{L}\underline{y} - (\underline{X}^T\underline{X})^{-1}\underline{X}^T\underline{y}, (\underline{X}^T\underline{X})^{-1}\underline{X}^T\underline{y})$$

$$= (\underline{L} - (\underline{X}^T\underline{X})^{-1}\underline{X}^T) \underbrace{\text{var}(\underline{y})}_{\sigma^2 \cdot \underline{I}} \underline{X} \cdot (\underline{X}^T\underline{X})^{-1}$$

$$= \sigma^2 (\underline{L} - (\underline{X}^T\underline{X})^{-1}\underline{X}^T) \underline{X} (\underline{X}^T\underline{X})^{-1}$$

$$= \sigma^2 (\underline{L}\underline{X} - \underline{I}) (\underline{X}^T\underline{X})^{-1}$$

$$= \underline{0}$$

Podobno je $\text{cov}(\hat{\beta}, \hat{\beta} - \tilde{\beta}) = 0$.

Sledi

$$\text{var}(\tilde{\beta}) = \text{var}(\hat{\beta}) + \text{var}(\hat{\beta} - \tilde{\beta}).$$

Ure matrice so pozitivno
semidefinitne.

Definicija: Za pozitivno semidefinitu matrici \underline{A} in \underline{B} rečemo, da je $\underline{A} \geq \underline{B}$, če je $\underline{A} - \underline{B}$ pozitivno semidefinitna matrica.

V tem smislu lahko zapisemo

$$\text{var}(\hat{\beta}) \geq \text{var}(\hat{\beta}_i)$$

Med drugim to pomeni:

- (i) $\hat{\beta}_i$ je najboljša nepristranska linearna ocenilka β_i .
- (ii) če hočemo nepristransko oceniti $y = \underline{a}^T \cdot \beta$ za nek \underline{a} , je

$$\begin{aligned} \text{var}(\underline{a}^T \hat{\beta}) &= \underline{a}^T \cdot \text{var}(\hat{\beta}) \underline{a} \\ &\geq \underline{a}^T \text{var}(\hat{\beta}_i) \underline{a} \end{aligned}$$

Torej je $\hat{y} = \underline{a}^T \hat{\beta}_i$ najboljša ocenilka $y = \underline{a}^T \cdot \beta$.

(iii) Velja tudi, da $\hat{\beta}$ minimizira

$$E [(\hat{\beta} - \beta)^T (\tilde{\beta} - \beta)]$$

med vsemi nepristranskimi
linearnimi $\tilde{\beta}$. Računamo

$$E [(\hat{\beta} - \beta)^T (\tilde{\beta} - \beta)]$$

$$= E [(\tilde{\beta} - \beta + \hat{\beta} - \hat{\beta})^T (\tilde{\beta} - \beta + \hat{\beta} - \hat{\beta})]$$

$$= E [(\hat{\beta} - \hat{\beta})^T (\tilde{\beta} - \hat{\beta})]$$

$$+ E [(\hat{\beta} - \beta)^T (\hat{\beta} - \beta)]$$

$$+ E [(\tilde{\beta} - \hat{\beta})^T (\hat{\beta} - \beta)]$$

$$+ E [(\hat{\beta} - \beta)^T (\tilde{\beta} - \hat{\beta})]$$

Ampak

$$E [(\tilde{\beta} - \hat{\beta})^T (\hat{\beta} - \beta)]$$

$$= E [\text{se} ((\hat{\beta} - \beta) (\tilde{\beta} - \hat{\beta})^T)]$$

$$= \text{se} [\text{cov} (\hat{\beta}, \tilde{\beta} - \hat{\beta})]$$

$$= 0$$

Podobno dobimo za drugo
povzeto vrednost, zato je

$$E[(\tilde{\beta} - \beta)^T (\tilde{\beta} - \beta)] \\ \geq E[(\hat{\beta} - \beta)^T (\hat{\beta} - \beta)].$$

Definicija: Vektorju $\hat{y} = X(X^T X)^{-1} X^T \hat{\beta}$
rečemo vektor prilagojenih
vrednosti. (angl. fitted values).

Izrek Gauss-Markova lahko uveljavimo
fosplošno. Predpostavimo

$$\underline{y} = \underline{X}\underline{\beta} + \underline{\varepsilon}$$

$$\text{z } E(\underline{\varepsilon}) = 0, \quad \text{var}(\underline{\varepsilon}) = \sigma^2 \underline{\Sigma},$$

kjer je $\underline{\Sigma}$ znana pozitivno-definitna
matrica. Velja

$$\begin{aligned} \underline{\Sigma}^{-1/2} \underline{y} &= \underline{\Sigma}^{-1/2} \underline{X} \underline{\beta} + \underline{\Sigma}^{-1/2} \underline{\varepsilon} \\ \tilde{\underline{y}} &= \tilde{\underline{X}} \underline{\beta} + \tilde{\underline{\varepsilon}} \end{aligned}$$

$$\begin{aligned} \text{var}(\tilde{\underline{\varepsilon}}) &= \text{var}(\underline{\Sigma}^{-1/2} \underline{\varepsilon}) \\ &= \underline{\Sigma}^{-1/2} \sigma^2 \underline{\Sigma} \underline{\Sigma}^{-1/2} \\ &= \sigma^2 \underline{I}. \end{aligned}$$

Ozna $\underline{\beta}$ je $\hat{\underline{\beta}} = (\tilde{\underline{X}}^T \tilde{\underline{X}})^{-1} \tilde{\underline{X}}^T \tilde{\underline{y}}$

ali zapisano drugače

$$\hat{\underline{\beta}} = (\underline{X}^T \underline{\Sigma}^{-1} \underline{X})^{-1} \underline{X}^T \underline{\Sigma}^{-1} \underline{y}$$

Izrek 4.3a (Gauss-Markov sploščen)

Naj bo $\underline{Y} = \underline{X}\beta + \underline{\varepsilon}$ z $E(\underline{\varepsilon}) = 0$
in $\text{var}(\underline{\varepsilon}) = \sigma^2 \underline{\Sigma}$, kjer je $\underline{\Sigma}$ znana
matrika in σ^2 neznani parameter.

$$\text{Cenilka } \hat{\beta} = (\underline{X}^T \underline{\Sigma}^{-1} \underline{X})^{-1} \underline{X}^T \underline{\Sigma}^{-1} \underline{Y}$$

je najboljša nepristranska cenilka
 β , ki je linearna.

Dokaz: Cenilka $\hat{\beta}$ je linearna
in nepristranska. Recimo, da

bi obstajala boljša cenilka $\tilde{\beta} = k \cdot Y$.

Potem bi $\tilde{\beta}$ bila boljša cenilka

tudi za $\tilde{Y} = \tilde{X}\beta + \tilde{\varepsilon}$, kar pa velja
da ni.

Opomba: Lahko tudi na enak način
kot v dokazu izreka 4.3 dokažemo,
da je $\text{var}(\tilde{\beta}) \geq \text{var}(\hat{\beta})$.

Primer: Recimo, da je

$$y_{kl} = \alpha + \beta x_{kl} + \varepsilon_k + \varepsilon_{k,l}$$

za $k = 1, 2, \dots, K$ in za vsak k

je $l = 1, 2, \dots, L_k$. Privzamemo, da

so $\varepsilon_k, \varepsilon_{k,l}$ vse neodvisne z

$$E(\varepsilon_k) = 0, E(\varepsilon_{k,l}) = 0, \text{ var}(\varepsilon_k) = \sigma^2$$

in $\text{var}(\varepsilon_{k,l}) = \tau^2$. Privzamemo, da

je $\rho = \tau^2 / \sigma^2$ znano število.

Zapišemo

$$\begin{pmatrix} y_{11} \\ y_{12} \\ \vdots \\ y_{1L_1} \\ y_{21} \\ \vdots \\ y_{2L_2} \\ \vdots \\ y_{K1} \\ \vdots \\ y_{KLK} \end{pmatrix} = \begin{pmatrix} 1 & x_{11} \\ 1 & \vdots \\ \vdots & \vdots \\ 1 & x_{1L_1} \\ \vdots & \vdots \\ \vdots & \vdots \\ 1 & x_{K1} \\ \vdots & \vdots \\ 1 & x_{KLK} \end{pmatrix} \begin{pmatrix} \alpha \\ \beta \end{pmatrix} + \underbrace{\begin{pmatrix} \varepsilon_1 + \varepsilon_{11} \\ \vdots \\ \varepsilon_1 + \varepsilon_{1L_1} \\ \vdots \\ \varepsilon_K + \varepsilon_{K1} \\ \vdots \\ \varepsilon_K + \varepsilon_{KLK} \end{pmatrix}}_{\eta}$$

Ugotovimo

$$\text{cov}(\eta) = \left(\begin{array}{c|c} \begin{matrix} \sigma^2 + \tau^2 & & & \\ & \sigma^2 & & \\ & & \dots & \\ & & & \sigma^2 \end{matrix} & \\ \hline & \begin{matrix} & & & \\ & & & \\ & & & \\ & & & \end{matrix} \\ \hline & & & & & & & & & \end{array} \right) = \sigma^2 \underline{\Sigma}$$

$\underbrace{\hspace{10em}}_{L_k}$

Iz postavimo σ^2 in τ^2 v matrike postanejo oblike

$$\sigma^2 \begin{pmatrix} 1+\rho & 1 & \dots & 1 \\ 1 & 1+\rho & \dots & 1 \\ \vdots & \vdots & \ddots & \vdots \\ 1 & 1 & \dots & 1+\rho \end{pmatrix} = \sigma^2 (\rho \underline{I} + \underline{1}\underline{1}^T)$$

U smislu izreka 4.3 a) bo

$$\underline{\Sigma} = \begin{pmatrix} \square & & & \\ & \square & & \\ & & \square & \\ & & & \square \\ & & & & \square \end{pmatrix}$$

s matrikami oblike $\rho \underline{I} + \underline{1}\underline{1}^T$

1.2 postavimo δ^2 in $\bar{\lambda}$ katlice postanejo oblike

$$\delta^2 \begin{pmatrix} 1+\rho & 1 & 1 & 1 \\ 1 & 1+\rho & & 1 \\ 1 & & & 1+\rho \\ 1 & 1 & & 1+\rho \end{pmatrix} = \delta^2 (\rho I + \underline{1} \cdot \underline{1}^T)$$

V smislu izreka 4.3 a bo

$$\underline{\Sigma} = \begin{pmatrix} \underline{\Sigma}_1 & & & \\ & \square & & \\ & & \square & \dots \\ & & & \dots & \square \end{pmatrix}$$

1. $\bar{\lambda}$ katlici mi oblike $\rho I + \underline{1} \cdot \underline{1}^T$.

Potrebno bomo inverze $\bar{\lambda}$ katlic po diagonali. Predpostavimo, da je inverz oblike $(\frac{1}{\rho} I + c \underline{1} \cdot \underline{1}^T)$.

Zmnožimo in sledi

$$\begin{aligned} & (\rho I + \underline{1} \cdot \underline{1}^T) \left(\frac{1}{\rho} I + c \underline{1} \cdot \underline{1}^T \right) \\ &= I + \frac{1}{\rho} \underline{1} \cdot \underline{1}^T + c \rho \underline{1} \cdot \underline{1}^T + c \cdot n \underline{1} \cdot \underline{1}^T \end{aligned}$$

Če želimo, da bo rezultat 1, mora biti

$$\frac{1}{\rho} + c\rho + c \cdot n = 0 \Rightarrow$$

$$c = - \frac{1}{\rho(n + \rho)}$$

S tem lahko izračunamo

$$\underline{x}^T \underline{\Sigma}^{-1} \underline{x} \quad \text{in} \quad \underline{x}^T \underline{\Sigma}^{-1} \underline{y}$$

Označimo $\underline{x}_k = \begin{pmatrix} x_{k1} \\ x_{k2} \\ \vdots \\ x_{kL_k} \end{pmatrix}$, $\underline{y}_k = \begin{pmatrix} y_{k1} \\ \vdots \\ y_{kL_k} \end{pmatrix}$

Označimo $c_k = - \frac{1}{\rho(L_k + \rho)}$

$$\underline{x}^T \underline{\Sigma}^{-1} \underline{x} = \begin{pmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{pmatrix}$$

Večja

$$a_{11} = \underline{1}^T \underline{\Sigma}^{-1} \underline{1}$$

$$= \text{vsota vseh elementov } \underline{\Sigma}^{-1}$$

$$= \frac{1}{\rho} \sum_{k=1}^K L_k + \sum_{k=1}^K c_k L_k^2$$

$$\begin{aligned}
a_{12} = a_{21} &= \underline{\underline{1}}^T \underline{\underline{\Sigma}}^{-1} \underline{\underline{x}} \\
&= \sum_{k=1}^K \underline{\underline{1}}^T \underline{\underline{\Sigma}}_k^{-1} \underline{\underline{x}}_k \\
&= \sum_{k=1}^K \underline{\underline{1}}^T \left(\frac{1}{\rho} \mathbf{I} + c_k \underline{\underline{1}} \cdot \underline{\underline{1}}^T \right) \underline{\underline{x}}_k \\
&= \sum_{k=1}^K \left(\frac{1}{\rho} \sum_{\ell=1}^{L_k} x_{k\ell} \right) \\
&\quad + \sum_{k=1}^K c_k L_k \cdot \sum_{\ell=1}^{L_k} x_{k\ell}
\end{aligned}$$

$$\begin{aligned}
a_{22} &= \sum_{k=1}^K \underline{\underline{x}}_k^T \underline{\underline{\Sigma}}_k^{-1} \underline{\underline{x}}_k \\
&= \sum_{k=1}^K \underline{\underline{x}}_k^T \left(\frac{1}{\rho} \mathbf{I} + c_k \underline{\underline{1}} \cdot \underline{\underline{1}}^T \right) \underline{\underline{x}}_k \\
&= \sum_{k=1}^K \frac{1}{\rho} \sum_{\ell=1}^{L_k} x_{k\ell}^2 \\
&\quad + \sum_{k=1}^K c_k \left(\sum_{\ell=1}^{L_k} x_{k\ell} \right)^2
\end{aligned}$$

Özütüm

$$\underline{x}^T \underline{\Sigma}^{-1} \underline{y} = \begin{pmatrix} b_1 \\ b_2 \end{pmatrix}$$

Veriler

$$b_1 = \underline{1}^T \underline{\Sigma}^{-1} \underline{y}$$

$$= \sum_{k=1}^K \underline{1}^T \left(\frac{1}{\rho} \mathbf{I} + c_k \underline{1} \underline{1}^T \right) \underline{y}_k$$

$$= \frac{1}{\rho} \sum_{k=1}^K \sum_{l=1}^{L_k} y_{kl} + \sum_{k=1}^K c_k L_k \cdot \sum_{l=1}^{L_k} y_{kl}$$

$$b_2 = \underline{x}^T \underline{\Sigma}^{-1} \underline{y}$$

$$= \sum_{k=1}^K \underline{x}_k^T \left(\frac{1}{\rho} \mathbf{I} + c_k \underline{1} \underline{1}^T \right) \underline{y}_k$$

$$= \frac{1}{\rho} \sum_{k=1}^K \sum_{l=1}^{L_k} x_{kl} y_{kl}$$

$$+ \sum_{k=1}^K c_k \left(\sum_{l=1}^{L_k} x_{kl} \right) \left(\sum_{l=1}^{L_k} y_{kl} \right)$$

Özetim

$$\begin{pmatrix} \hat{\alpha} \\ \hat{\beta} \end{pmatrix} = \begin{pmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{pmatrix}^{-1} \begin{pmatrix} b_1 \\ b_2 \end{pmatrix}$$

4.3 Normalni $\underline{\varepsilon}$

V standardnem modelu linearne regresije

$$\underline{y} = \underline{X}\beta + \underline{\varepsilon}$$

lahko predpostavimo $\underline{\varepsilon} \sim N(\underline{0}, \sigma^2 \underline{I})$.

Definicija: Matriko $\underline{H} = \underline{X}(\underline{X}^T \underline{X})^{-1} \underline{X}^T$ imenujemo projekcijska matrika.

Vemo, da je \underline{H} idempotentna. Potem je tudi $\underline{I} - \underline{H}$ idempotentna.

Vemo, da je vektor

$$\begin{pmatrix} \hat{\beta} \\ \hat{\varepsilon} \end{pmatrix} = \begin{pmatrix} (\underline{X}^T \underline{X})^{-1} \underline{X}^T \underline{y} \\ (\underline{I} - \underline{H}) \underline{y} \end{pmatrix}$$

večkrat sežen normalen. Za neodvisnost preverimo (vemo $(\underline{I} - \underline{H}) \underline{X} = \underline{0}$)

$$\text{cov}(\hat{\beta}, \hat{\varepsilon}) = \sigma^2 (\underline{X}^T \underline{X})^{-1} \underline{X}^T (\underline{I} - \underline{H}) = \underline{0}.$$

Posledično je $\hat{\beta}$ neodvisen od $\hat{\varepsilon}$. Poleg tega je

$$\begin{aligned} \hat{\sigma}^2 &= \frac{1}{n-m} \sum_{k=1}^n \varepsilon_k^2 \\ &= \frac{1}{n-m} \underline{\varepsilon}^T \underline{\varepsilon} \\ &= \frac{1}{n-m} \underline{\varepsilon}^T (\underline{I} - \underline{H}) \underline{\varepsilon} \quad (\text{dokazati} \\ &\quad \text{že prej}) \\ &\sim \frac{\sigma^2}{n-m} \cdot \chi^2(n-m) \end{aligned}$$

Sledi, da je

$$T_i = \frac{\hat{\beta}_i - \beta_i}{\sqrt{c_{ii}} \hat{\sigma}} \sim t_{n-m}$$

Pri tem je $\underline{c} = (\underline{X}^T \underline{X})^{-1}$ in je

c_{ii} i-ti diagonalni element \underline{c} .

Zakaj?

$$\frac{\hat{\beta}_i - \beta_i}{\sqrt{c_{ii}}} \sim N(0, \sigma^2)$$

Na podlagi tega lahko testiramo, recimo, ničelno domnevo

$$H_0: \beta_i = 0 \quad \text{proti} \quad H_1: \beta_i \neq 0.$$

Testna statistika $T_i = \frac{\hat{\beta}_i}{\sqrt{c_{ii} \hat{\sigma}^2}}$

ima v primeru, ko H_0 drži, t_{n-k} porazdelitev.

Interpretacija: Če H_0 ne zavrnemo, lahko rečemo, da i -ti stolpec v \underline{X} ne prispeva k $E(\underline{Y})$, torej i -ta neodvisna spremenljivka, nima vpliva na \underline{Y} .

Kako pa bi v splošnem preizkusili

$$H_0: \beta_{p+1} = \beta_{p+2} = \dots = \beta_n = 0.$$

Idea: Uporabimo kvocient
verjetij.

$$L(\beta, \sigma^2 | \underline{x}, \underline{y})$$

$$= \frac{1}{(2\pi)^{n/2} \sigma^n} \exp\left\{-\frac{1}{2} \sum_{k=1}^n (y_k - (\underline{x}\beta)_k)^2 / \sigma^2\right\}$$

ali

$$L(\beta, \sigma^2 | \underline{x}, \underline{y})$$

$$= \frac{n}{2} \log 2\pi - n \log \sigma^2$$

$$- \frac{1}{2\sigma^2} (\underline{y} - \underline{x}\beta)^T (\underline{y} - \underline{x}\beta).$$

Maximum bo očitno dosežen za

$$\hat{\beta} = (\underline{x}^T \underline{x})^{-1} \underline{x}^T \underline{y} \quad \text{in}$$

$$\hat{\sigma}^2 = \frac{1}{n} (\underline{y} - \underline{x}\hat{\beta})^T (\underline{y} - \underline{x}\hat{\beta}).$$

Označimo $\beta = \begin{pmatrix} \beta_1 \\ \vdots \\ \beta_p \\ \beta_{p+1} \\ \vdots \\ \beta_m \end{pmatrix}$

in podobno $X = \begin{pmatrix} X_1 & X_2 \end{pmatrix}$
 $\underbrace{\quad}_p \quad \underbrace{\quad}_{m-p}$

Če se omejimo na $\beta_{p+1} = \dots = \beta_m = 0$

dobimo oceni

$$\hat{\beta} = (X_1^T X_1)^{-1} X_1^T Y$$

$$\hat{\sigma}^2 = \frac{1}{n} (Y - X_1 \hat{\beta})^T (Y - X_1 \hat{\beta})$$

Notavimo in dobimo

$$\lambda = -2n \log \hat{\sigma}^2 + 2n \log \tilde{\sigma}^2$$

$$= n \log \left(\frac{\hat{\sigma}^2}{\tilde{\sigma}^2} \right)$$

Ničelno domnevo zavzememo, če je λ prevelik. Ampak to je isto kot izjava, da je kvocient $\frac{\hat{\sigma}^2}{\hat{\sigma}^2}$ prevelik. Označimo

$$\underline{H} = \underline{X} (\underline{X}^T \underline{X})^{-1} \underline{X}^T, \quad \underline{H}_1 = \underline{X}_1 (\underline{X}_1^T \underline{X}_1)^{-1} \underline{X}_1^T.$$

Velja, če H_0 drži,

$$\begin{aligned} \frac{\hat{\sigma}^2}{\hat{\sigma}^2} &= \frac{\underline{\varepsilon}^T (\underline{I} - \underline{H}_1) \underline{\varepsilon}}{\underline{\varepsilon}^T (\underline{I} - \underline{H}) \underline{\varepsilon}} \\ &= \frac{\underline{\varepsilon}^T (\underline{I} - \underline{H}) \underline{\varepsilon} + \underline{\varepsilon}^T (\underline{H} - \underline{H}_1) \underline{\varepsilon}}{\underline{\varepsilon}^T (\underline{I} - \underline{H}) \underline{\varepsilon}} \end{aligned}$$

Velja $\underline{H}_1 \underline{H} = \underline{H} \underline{H}_1 = \underline{H}_1$, torej je

$$\begin{aligned} (\underline{H} - \underline{H}_1)^2 &= \underline{H}^2 - 2 \underline{H} \underline{H}_1 + \underline{H}_1^2 \\ &= \underline{H} - \underline{H}_1. \end{aligned}$$

Matrica $\underline{H} - \underline{H}_1$ je idempotentna.

Poleg tega je

$$\begin{aligned} & \text{cov}((\underline{I} - \underline{H})\underline{Y}, (\underline{H} - \underline{H}_1)\underline{Y}) \\ &= (\underline{I} - \underline{H}) \sigma^2 \underline{I} (\underline{H} - \underline{H}_1) \\ &= \sigma^2 (\underline{H} - \underline{H}^2 - \underline{H}_1 + \underline{H} \cdot \underline{H}_1) \\ &= 0 \end{aligned}$$

Torej je

$$\frac{\hat{\sigma}^2}{\hat{\sigma}^2} = 1 + \frac{\underline{\varepsilon}^T (\underline{H} - \underline{H}_1) \underline{\varepsilon}}{\underline{\varepsilon}^T (\underline{I} - \underline{H}) \underline{\varepsilon}}$$

prevelik, če je kvocient na desni prevelik. Če je $\underline{x}_p \in \mathcal{R}(\underline{H}_1)$,

$$\text{je } \underline{\varepsilon}^T (\underline{H} - \underline{H}_1) \underline{\varepsilon} \sim \sigma^2 \chi^2(m-p)$$

$$\text{in } \underline{\varepsilon}^T (\underline{I} - \underline{H}) \underline{\varepsilon} \sim \sigma^2 \chi^2(n-m),$$

poleg tega pa sta števec in imenovalec neodvisna.

Kvocijent lahko še množimo s konstanto. Definiramo

$$F = \frac{\underline{y}^T (\underline{H} - \underline{H}_1) \underline{y} / (m-p)}{\underline{y}^T (\underline{I} - \underline{H}) \underline{y} / (n-m)}$$

$$\sim F_{m-p, n-m}$$

H_0 zavrnemo, če je F prevelik, kritično vrednost pa vzberemo iz F -porazdelitve.

Opomba: Zgorajemu večemo preizkus sklopne linearne domneve.

Opomba: Če je prvi stolpec v \underline{X} enak $\underline{1}$ in večemo $\underline{X} = (\underline{1}; \underline{X}_1)$ potem količini

$$\dots$$

$$R^2 = \frac{\sum_{k=1}^n (\hat{y}_k - \bar{y})^2}{\sum_{k=1}^n (y_k - \bar{y})^2}$$

$$= \frac{y^T \left(\mathbb{H} - \frac{1}{n} \mathbb{1}\mathbb{1}^T \right) y}{y^T \left(\mathbb{I} - \frac{1}{n} \mathbb{1}\mathbb{1}^T \right) y}$$

imevamo R^2 koeficient in ga interpretiramo kot mero napovedne moči modela. Večji R^2 pomeni večjo napovedno moč modela.