

VERJETNOSTNI RAČUN IN STATISTIKA

Aleksandar Jurišić, FRI

1. junij 2010

Kazalo

1	Uvod	1
I	VERJETNOST	7
2	Poskusi, dogodki in definicija verjetnosti	11
2.1	Poskusi in dogodki	12
2.2	Računanje z dogodki	13
2.3	Definicija verjetnosti	15
2.4	Osnovne lastnosti verjetnosti	17
2.5	Aksiomi Kolmogorova	20
3	Pogojna verjetnost	23
3.1	Intriga (po Kvardabri)	23
3.2	Definicija pogojne verjetnosti	26
3.3	Obrazec za popolno verjetnost in večstopenjski poskusi	29
4	Bernoullijevo zaporedje neodvisnih poskusov	33
4.1	Računanje $P_n(k)$	34
5	Slučajne spremenljivke in porazdelitve	45
5.1	Diskretne slučajne spremenljivke	46
5.1.1	Enakomerna diskretna porazdelitev	47
5.1.2	Binomska porazdelitev	47
5.1.3	Poissonova porazdelitev $P(\lambda)$	48
5.1.4	Pascalova porazdelitev $P(m, p)$	49
5.1.5	Hipergeometrijska porazdelitev $H(n; M, N)$	50
5.2	Ponovitev: integrali	50

5.3	Zvezne slučajne spremenljivke	51
5.3.1	Enakomerna porazdelitev zvezne slučajne spremenljivke	52
5.3.2	Normalna ali Gaussova porazdelitev	52
5.3.3	Porazdelitev Poissonovega toka, eksponentna	55
5.3.4	Porazdelitev Gama	56
5.3.5	Porazdelitev hi-kvadrat	57
5.3.6	Cauchyeva porazdelitev	57
6	Slučajni vektorji	63
6.1	Diskretne večrazsežne porazdelitve – polinomska	66
6.2	Ponovitev: dvojni integral	66
6.3	Zvezne večrazsežne porazdelitve	68
6.4	Neodvisnost slučajnih spremenljivk	73
7	Funkcije slučajnih spremenljivk in vektorjev	77
7.1	Funkcije slučajnih spremenljivk	77
7.2	Funkcije in neodvisnost	80
7.3	Funkcije slučajnih vektorjev	81
7.4	Pogojne porazdelitve	83
8	Momenti in kovarianca	87
8.1	Matematično upanje	87
8.2	Disperzija	90
8.3	Standardizirane spremenljivke	91
8.4	Kovarianca	91
8.5	Pogojno matematično upanje	93
8.6	Višji momenti	95
9	Karakteristične funkcije in limitni izreki	97
9.1	Karakteristična funkcija	97
9.2	Limitni izreki	98
9.3	Centralni limitni izrek (CLI)	101
10	Nekaj primerov uporabe	103
10.1	Zamenjalna šifra	103
10.2	Kakšno naključje!!! Mar res?	106

10.3 Ramseyjeva teorija	108
10.4 Teorije kodiranja	114
II STATISTIKA	115
11 Opisna statistika	125
11.1 Vrste spremenljivk oziroma podatkov	125
11.2 Grafična predstavitev kvantitativnih podatkov	127
11.3 Mere za lokacijo in razpršenost	133
11.4 Standardizacija	139
12 Vzorčenje	141
12.1 Osnovni izrek statistike	144
12.2 Vzorčne ocene	145
12.3 Porazdelitve vzorčnih povprečij	146
12.4 Vzorčna statistika	150
12.4.1 (A) Vzorčno povprečje	150
12.4.2 (B) Vzorčna disperzija	153
12.5 Nove porazdelitve	155
12.5.1 Studentova porazdelitev	156
12.5.2 Fisherjeva porazdelitev	157
13 Cenilke	159
13.1 Osnovni pojmi	159
13.2 Rao-Cramérjeva ocena	162
13.3 Učinkovitost cenilk	163
13.4 Metoda momentov	165
12.4 Vzorčna statistika (nadaljevanje)	167
12.4.3 (C) Vzorčne aritmetične sredine	167
12.4.4 (D) Vzorčni deleži	168
12.4.5 (E) Razlika vzorčnih aritmetičnih sredin	170
12.4.6 (F) Razlika vzorčnih deležev	171
14 Intervali zaupanja	173
14.1 Pomen stopnje tveganja	174

14.2	Intervalsko ocenjevanje parametrov	175
14.2.1	Povprečje μ s poznanim σ	177
14.2.2	Velik vzorec za povprečje μ	177
14.2.3	Majhen vzorec za povprečje μ	178
14.2.4	Razlika povprečij $\mu_1 - \mu_2$ s poznanima σ_1 in σ_2	178
14.2.5	Veliki vzorec za razliko povprečij $\mu_1 - \mu_2$ z neznanima σ_1 in σ_2	178
14.2.6	Majhen vzorec za razliko povprečij $\mu_1 - \mu_2$ z neznanima $\sigma_1 = \sigma_2$	179
14.2.7	Majhen vzorec za razliko povprečij $\mu_1 - \mu_2$ z neznanima σ_1 in σ_2	179
14.2.8	Velik vzorec za razliko $\mu_d = \mu_1 - \mu_2$	180
14.2.9	Majhen vzorec za razliko $\mu_d = \mu_1 - \mu_2$	180
14.2.10	Delež π s poznanim $\sigma_{\hat{\pi}}$	181
14.2.11	Veliki vzorec za delež populacije	181
14.2.12	Razlika deležev $\pi_1 - \pi_2$ s poznanim $\sigma_{\hat{p}_1 - \hat{p}_2}$	181
14.2.13	Veliki vzorec za razliko deležev $\pi_1 - \pi_2$	182
14.2.14	Veliki vzorec za varianco σ^2	182
14.2.15	Kvocient varianc σ_1^2/σ_2^2	182
14.3	Izbira velikosti vzorca	183
15	Preverjanje domnev	185
15.1	Ilustrativni primer (ameriški sodni sistem)	187
15.2	Alternativna domneva in definicije napak	188
15.3	P -vrednost	190
15.4	Statistična domneva	191
15.5	Preverjanje predznaka	194
15.6	Wilcoxonov predznačen-rang test	194
15.7	Naloga	197
15.8	Formalen postopek za preverjanje domnev	199
15.8.1	$\mu = \mu_0$ z znanim σ	199
15.8.2	$\mu = \mu_0$ z neznanim σ , $n \geq 30$	201
15.8.3	$\mu = \mu_0$, neznan σ , populacija normalna in $n < 30$	201
15.9	Razlika povprečij $\mu_1 - \mu_2 = D_0$	205
15.9.1	Znana σ_1 in σ_2	205
15.9.2	Neznana σ_1 in/ali σ_2 , $n_1 \geq 30$ in/ali $n_2 \geq 30$	206
15.9.3	Neznana σ_1 in/ali σ_2 , pop. norm., $\sigma_1 = \sigma_2$, $n_1 < 30$ ali $n_2 < 30$	206
15.9.4	Neznana σ_1 in/ali σ_2 , pop. norm., $\sigma_1 \neq \sigma_2$, $n_1 < 30$ ali $n_2 < 30$	207

15.9.5	Povprečje $\mu_d = D_0$ in $n \geq 30$	207
15.9.6	Povprečje $\mu_d = D_0$, populacija razlik normalna, in $n \leq 30$	207
15.9.7	Delež $p = p_0$ z dovolj velikim vzorcem	208
15.10	Preverjanje domneve za delež	211
15.11	Razlika deležev dveh populaciji	212
15.11.1	Velik vzorec za testiranje domneve o $p_1 - p_2$, kadar je $D_0 = 0$	212
15.11.2	Velik vzorec za testiranje domneve o $p_1 - p_2$, kadar je $D_0 \neq 0$	214
15.12	Analiza variance	214
15.12.1	Preverjanje domneve $\sigma^2 = \sigma_0^2$	217
15.12.2	Preverjanje domneve za varianco	218
15.12.3	Preverjanje domneve $\sigma_1^2/\sigma_2^2 = 1$	219
15.13	Preverjanje domnev o porazdelitvi spremenljivke	220
15.13.1	Preverjanje domneve o enakomerni porazdelitvi	220
15.13.2	Preverjanje domneve o normalni porazdelitvi	221
16	Bivariatna analiza in regresija	227
16.1	Preverjanje domneve o povezanosti dveh nominalnih spremenljivk	228
16.2	Koeficienti asociacije	230
16.3	Preverjanje domneve o povezanosti dveh ordinalnih spremenljivk	232
16.4	Preverjanje domneve o povezanosti dveh številskih spremenljivk	234
16.5	Parcialna korelacija	237
16.6	Regresijska analiza	238
16.7	Linearni model	241
16.7.1	Statistično sklepanje o regresijskem koeficientu	245
16.7.2	Pojasnjena varianca (ang. ANOVA)	246
17	Časovne vrste	249
17.1	Primerljivost členov v časovni vrsti	250
17.2	Grafični prikaz časovne vrste	251
17.3	Indeksi	251
17.4	Sestavine dinamike v časovnih vrstah	253
18	Uporaba	257
18.1	Načrtovanje eksperimentov	257

III	KAM NAPREJ	261
A	MATEMATIČNE OSNOVE (ponovitev)	271
A.1	Računala nove dobe	271
A.2	Funkcije/preslikave	274
A.3	Permutacije	275
A.4	Kombinacije	281
A.5	Vrsta za e	283
A.6	Stirlingov obrazec	284
A.7	Normalna krivulja v prostoru	285
A.8	Sredine nenegativnih števil a_1, \dots, a_n	287
A.9	Cauchyjeva neenakost	288
B	PROGRAM R (Martin Raič)	293
B.1	Izvajanje programa	293
B.2	Aritmetika	294
B.3	Najosnovnejše o spremenljivkah	294
B.4	Uporabnikove funkcije	294
B.5	Numerično računanje	295
B.6	Podatkovne strukture	295
B.6.1	Vektorji	295
B.6.2	Matrike	296
B.6.3	Tabele	298
B.6.4	Vektorji, matrike in tabele z označenimi indeksi	299
B.6.5	Zapisi	299
B.6.6	Kontingenčne tabele in vektorji s predpisanimi vrednostmi	300
B.6.7	Preglednice	300
B.7	Osnove programiranja	301
B.7.1	Izvajanje programov, shranjenih v datotekah	301
B.7.2	Najosnovnejši programski ukazi	301
B.7.3	Krmilni stavki	302
B.7.4	Nekaj več o funkcijah	302
B.8	Knjižnice z dodatnimi funkcijami	303
B.9	Vhod in izhod	303
B.9.1	Pisanje	303

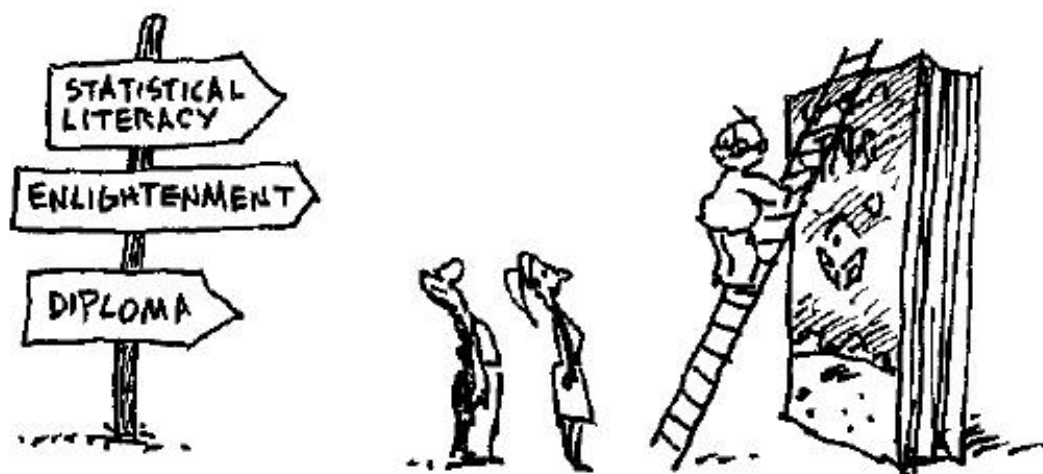
B.9.2 Delo z datotekami	304
B.9.3 Branje	305
B.9.4 Izvoz grafike	305
B.10 Verjetnostne porazdelitve	306
B.11 Simulacije	307
B.12 Statistično sklepanje	307
B.12.1 Verjetnost uspeha poskusa/delež v populaciji	307
B.12.2 Primerjava verjetnosti dveh poskusov/deležev v dveh populacijah	308
B.12.3 Primerjava verjetnosti več poskusov/deležev več populacijah	309
B.12.4 Populacijsko povprečje – T -test	309
B.12.5 Test mediane	309
B.12.6 Primerjava porazdelitev dveh spremenljivk	310
B.12.7 Koreliranost	312

Poglavje 1

Uvod



Naš predmet govori o dogodkih in podatkih. Vsekakor ni samo prepreka na poti do diplome, pač pa se uporablja pri večini drugih predmetov na FRI.



Predstavili bomo osnove teorije verjetnosti in njeno uporabo v statistiki, pa tudi nekaj osnov statistike. Bralec, ki bo predelal večji del naše snovi, bi moral znati opisati napovedljive vzorce, ki na dolgi rok vladajo slučajnim izidom ter doseči osnovno statistično

pismenost, tj. sposobnost sledenja in razumevanja argumentov, ki izhajajo iz podatkov.

V nekem članku i virusih lahko preberemo debato o napadu virusov (črvov), ki so se razširili po Internetu ter upočasnili brskalnike in e-pošto širom po svetu. Koliko računalnikov je bilo okuženih? Strokovnjaki, na katere so se sklicali v članku, pravijo, da je bilo okuženih 39.000 računalnikov, ki so vplivali na stotine tisočev drugih sistemov. Kako so lahko prišli to te številke? Ali ne bi bilo težko priti do take številke? Ali so preverili vsak računalnik na Internetu, da bi se prepričali, če je okužen ali ne? Dejstvo, da je bil članek napisan v manj kot 24 urah od časa napada, sugerira, da je to število samo predpostavka. Vendar pa se lahko vprašamo, zakaj potem 39.000 in ne 40.000?

Statistika je znanost zbiranja, organiziranja in interpretiranja numeričnih dejstev, ki jih imenujemo podatki. Vsakodnevno smo s podatki takorekoč bombardirani. Večina ljudi povezuje "statistiko" z biti podatkov, ki izhajajo v dnevnem časopisju, novicah, reportažah: povprečna temperatura na današni dan, procenti pri košarkaških prostih metih, procent tujih vlaganj na našem trgu, in anketa popularnosti predsednika in premierja. Reklame pogosto trdijo, da podatki kažejo na superiornost njihovega produkta. Vse strani v javnih debatah o ekonomiji, izobraževanju in socialni politiki izhajajo iz podatkov. Kljub temu pa uporabnost statistike presega te vsakodnevne primere.

Podatki so pomembni pri delu mnogih, zato je izobraževanje na področju statistike izredno pomembno pri številnih poklicih. Ekonomisti, finančni svetovalci, vodstveni kader v politiki in gospodarstvu preučujejo najnovejše podatke o nezaposlenosti in inflaciji. Zdravniki morajo razumeti izvor in zanesljivost podatkov, ki so objavljeni v medicinskih revijah. Poslovne odločitve so običajno zasnovane na raziskavah tržišč, ki razkrijejo želje kupcev in njihovo obnašanje. Večina akademskih raziskav uporablja številke in tako hočeš nočes izkorišča statistične metode.

Nič lažje se ni pobegniti podatkom kot se izogniti uporabi besed. Tako kot so besede na papirju brez pomena za nepismenega ali slabo izobraženega človeka, tako so lahko tudi podatki privlačni, zavajajoči ali enostavno nesmiselni. Statistična pismenost, tj. sposobnost sledenja in razumevanja argumentov, ki izhajajo iz podatkov, je pomembna za sleherno osebo.

Na statistiko in njene matematične temelje (verjetnost) lahko gledamo kot na učinkovito orodje, pa ne samo pri teoretičnem računalništvu (teoriji kompleksnosti, randomiziranih algoritmih, teoriji podatkovnih baz), pač pa tudi na praktičnih področjih. V vsakdanjem življenju ni pomembno da Vaš sistem obvlada čisto vse vhodne podatke, učinkovito pa naj opravi vsaj s tistimi, ki pokrijejo 99.99% primerov iz prakse.

Za konec pa še tole. Matematika je dobra osnova, odličen temelj. Če pri konkretnem računalniškem ustvarjanju (pri tem mislim na razvoj programske opreme in reševanje logističnih problemov, ne pa prodajo in popravilo računalniške opreme) nimamo matematične izobrazbe/znanja, potem je skoraj tako kot če bi postavili hišo na blatu. Lahko izgleda lepo, vendar pa se bo začela pogrezati pri naslednjem nalivu.

Pogled od zunaj

Števila so me pogosto begala,
še posebej, če sem imel pred seboj
neko njihovo razvrstitev,
tako da je tem primeru obveljala misel,
ki so jo pripisali Diaraeliju,
z vso pravico in močjo:

“Obstajajo tri vrste laži:

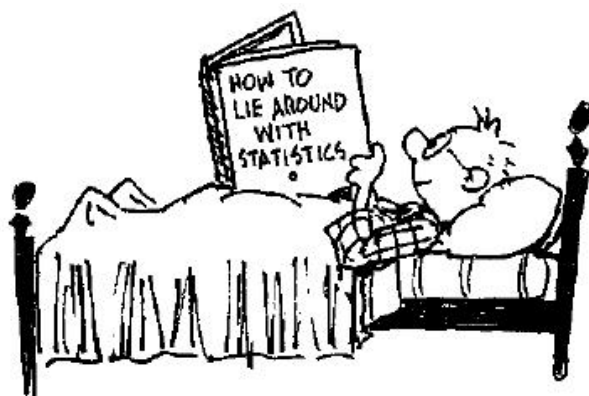
*laži,
preklete laži in
statistika.”*

iz Autobiografije Marka Twaina



*Be good + you will be lonesome.
Mark Twain*

Pa vendar ne misliti, da vas bomo učili laganja – vse prej kot to.



Statistik preučuje podatke, jih zbira, klasificira, povzema, organizira, analizira in interpretira. Glavni veji statistike sta **opisna statistika**, ki se ukvarja z organiziranjem, povzemanjem in opisovanjem zbirk podatkov (reduciranje podatkov na povzetke) ter **analitična statistika** jemlje vzorce podatkov in na osnovi njih naredi zaključke (inferenčnost) o populaciji (ekstrapolacija).

Populacija je podatkovna množica, ki ji je namenjena naša pozornost.

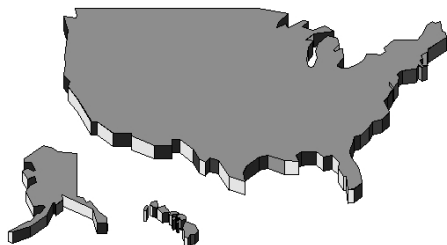


Vzorec je podmnožica podatkov, ki so izbrani iz populacije (po velikosti bistveno manjši od populacije).



- **Populacija**

– vsi objekti, ki jih opazujemo.



Primer:

vsi registrirani glasovalci.

- **Vzorec**

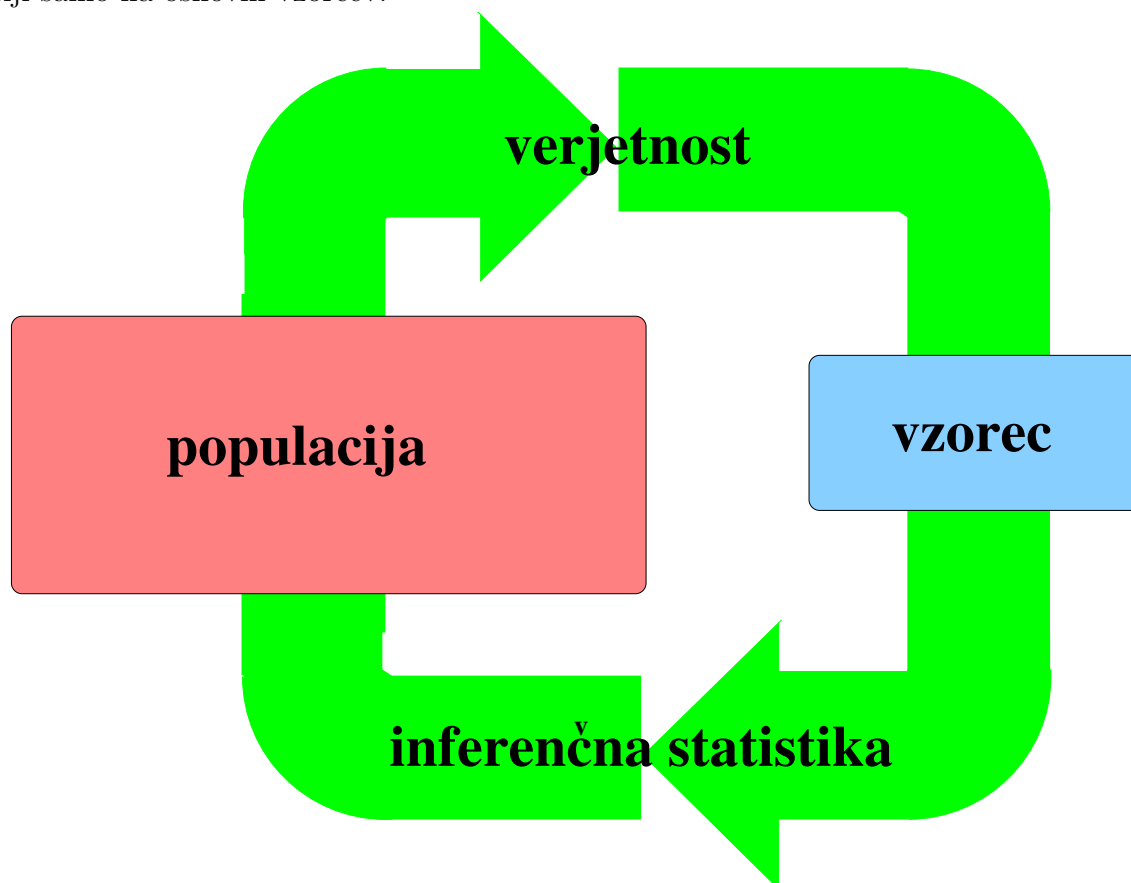
– podmnožica populacije



Primer:

100 registriranih glasovalcev

Pravijo, da je slika vredna 1000 besed, zato si še nazorno poredočimo, kako nam verjetnost pomaga oceniti kakšen bo vzorec, ki ga bomo izbrali iz dane in dobro poznane populacije, medtem ko nam inferenčna statistika pomaga delati zaključke o celotni populaciji samo na osnovni vzorcev.



Del I
VERJETNOST

Ste se kdaj vprašali, zakaj so igre na srečo, ki so za nekatere rekreacija ali pa droga, tako dober posel za igralnice?



Vsak uspešen posel mora iz uslug, ki jih ponuja, kovati napovedljive dobičke. To velja tudi v primeru, ko so te usluge igre na srečo. Posamezni hazarderji lahko zmagajo ali pa izgubijo. Nikoli ne morejo vedeti, če se bo njihov obisk igralnice končal z dobičkom ali z izgubo.



Igralnica pa ne kocka, pač pa dosledno dobiva in država lepo služi na račun loterij ter drugih oblik iger na srečo. Presenetljivo je, da lahko skupni rezultat več tisoč naključnih izidov poznamo s skoraj popolno gotovostjo. Igralnici ni potrebno obtežiti kock, označiti kart ali spremeniti kolesa rulete. Ve, da ji bo na dolgi rok vsak stavljeni euro prinesel približno pet centov dobička.



Splača se ji torej osredotočiti na brezplačne predstave ali poceni avtobusne vozovnice, da bi privabili več gostov in tako povečali število stavljenega denarja. Posledica bo večji dobiček. Splača se ji torej osredotočiti na brezplačne predstave ali poceni avtobusne vozovnice, da bi privabili več gostov in tako povečali število stavljenega denarja. Posledica bo večji dobiček. Igralnice niso edine, ki se okoriščajo z dejstvom, da so velikokratne ponovitve slučajnih izidov napovedljive. Na primer, čeprav zavarovalnica ne ve, kateri od njenih zavarovancev bodo umrli v prihodnjem letu, lahko precej natančno napove, koliko jih bo umrlo. Premije življenjskih zavarovanj postavi v skladu s tem znanjem, ravno tako kot igralnica določi glavne dobitke.

Igralnice niso edine, ki se okoriščajo z dejstvom, da so velikokratne ponovitve slučajnih izidov napovedljive.



Na primer, čeprav zavarovalnica ne ve, kateri od njenih zavarovancev bodo umrli v prihodnjem letu, lahko precej natančno napove, koliko jih bo umrlo. Premije življenjskih zavarovanj postavi v skladu s tem znanjem, ravno tako kot igralnica določi glavne dobitke.

Poglavje 2

Poskusi, dogodki in definicija verjetnosti



Ahil in Ajaks kockata, Amfora, okrog 530 pr.n.š., Eksekias, Vatikan



Cardano

Naključnost so poznale že stare kulture: Egipčani, Grki, ... a je niso poskušale razumeti – razlagale so jo kot voljo bogov.

Leta 1662 je plemič Chevalier de Mere zastavil matematiku Blaise Pascalu vprašanje:

[zakaj določene stave prinašajo dobiček druge pa ne.](#)

Le-ta si je o tem začel dopisovati s Fermatom in iz tega so nastali začetki verjetnostnega računa.

Prvo tovrstno razpravo je napisal že leta 1545 italijanski kockar in matematik Cardano, a ni bila širše znana. Tudi leta 1662 je anglež John Graunt sestavil na osnovi podatkov prve zavarovalniške tabele.

Leta 1713 je Jakob Bernoulli objavil svojo knjigo *Umetnost ugibanja* s katero je verjetnostni račun postal resna in splošno uporabna veda. Njegov pomen je še utrdil Laplace, ko je pokazal njegov pomen pri analizi astronomskih podatkov (1812).

Leta 1865 je avstrijski menih Gregor Mendel uporabil verjetnostno analizo pri razlagi dednosti v genetiki. V 20. stoletju se je uporaba verjetnostnih pristopov razširila skoraj na vsa področja.

2.1 Poskusi in dogodki



Verjetnostni račun obravnava zakonitosti, ki se pokažejo v velikih množicah enakih ali vsaj zelo podobnih pojavov. Predmet verjetnostnega računa je torej empirične narave in njegovi osnovni pojmi so povzeti iz izkušnje. Osnovni pojmi v verjetnostnem računu so: poskus, dogodek in verjetnost dogodka. **Poskus** je realizacija neke množice skupaj nastopajočih dejstev (kompleksa pogojev). Poskus je torej vsako dejanje, ki ga opravimo v natanko določenih pogojih.

Primeri: (1) met igralne kocke, (2) iz kupa 32 igralnih kart izberemo 5 kart, (3) met pikada v tarčo. \diamond

Pojav, ki v množico skupaj nastopajočih dejstev ne spada in se lahko v posameznem poskusu zgodi ali pa ne, imenujemo **dogodek**. Za poskuse bomo privzeli, da jih lahko neomejeno velikokrat ponovimo. Dogodki se bodo nanašali na isti poskus. Poskuse označujemo z velikimi črkami iz konca abecede, npr. X, Y, X_1 . Dogodke pa označujemo z velikimi črkami iz začetka abecede, npr. A, C, E_1 .

Primeri: (1) v poskusu meta igralne kocke je na primer dogodek, da vržemo 6 pik; (2) v poskusu, da vlečemo igralno karto iz kupa 20 kart, je dogodek, da izvlečemo rdečo barvo. (3) v poskusu meta pikada zadanemo center (polje, ki označuje center). \diamond

Izpostavimo nekatere posebne dogodke:

(a) **gotov** dogodek – G : ob vsaki ponovitvi poskusa se zgodi.

Primer: dogodek, da vržemo 1, 2, 3, 4, 5, ali 6 pik pri metu igralne kocke; \diamond

(b) **nemogoč** dogodek – N : nikoli se ne zgodi.

Primer: dogodek, da vržemo 7 pik pri metu igralne kocke;

◇

(c) **slučajen** dogodek: včasih se zgodi, včasih ne.

Primer: dogodek, da vržemo 6 pik pri metu igralne kocke.

◇

2.2 Računanje z dogodki

Dogodek A je **poddogodek** ali **način** dogodka B , kar zapišemo $A \subset B$, če se vsakič, ko se zgodi dogodek A , zagotovo zgodi tudi dogodek B .

Primer: Pri metu kocke je dogodek A , da pade šest pik, način dogodka B , da pade sodo število pik.

◇

Če je dogodek A način dogodka B in sočasno dogodek B način dogodka A , sta dogodka **enaka**: $(A \subset B) \wedge (B \subset A) \iff A = B$. **Vsota** dogodkov A in B , označimo jo z $A \cup B$ ali $A + B$, se zgodi, če se zgodi **vsaj** eden od dogodkov A in B .

Primer: Vsota dogodka A , da vržemo sodo število pik, in dogodka B , da vržemo liho število pik, je gotov dogodek.

◇

V naslednji trditvi zberemo nekatere osnovne lastnosti operacij nad dogodki, ki smo jih vpeljali doslej (gotovo pa ne bi škodilo, če bi prej prebrali še Presekova članka Računala nove dobe 1. in 2., glej lkrv.fri.uni-lj.si).

Trditev 2.1. *Za poljubna dogodka A in B velja:*

- $A \cup B = B \cup A,$
- $A \cup N = A,$
- $B \subset A \iff A \cup B = A,$
- $A \cup A = A,$
- $A \cup G = G,$
- $A \cup (B \cup C) = (A \cup B) \cup C.$ □

Produkt dogodkov A in B , označimo ga z $A \cap B$ ali AB , se zgodi, če se zgodita A in B *hkrati*.

Primer: Produkt dogodka A , da vržemo sodo število pik, in dogodka B , da vržemo liho število pik, je nemogoč dogodek.

◇

Trditev 2.2. *Za poljubna dogodka A in B velja:*

- $A \cup B = B \cup A$, • $A \cap N = N$, • $A \cap (B \cap C) = (A \cap B) \cap C$,
- $A \cap B = B \cap A$, • $A \cap G = A$, • $A \cap (B \cup C) = (A \cap B) \cup (A \cap C)$,
- $A \cap A = A$. • $B \subset A \iff A \cap B = B$, • $A \cup (B \cap C) = (A \cup B) \cap (A \cup C)$. \square

Dogodku A **nasproten** dogodek \bar{A} imenujemo negacijo dogodka A .

Primer: Nasproten dogodek dogodku, da vržemo sodo število pik, je dogodek, da vržemo liho število pik. \diamond

Trditev 2.3. *Za poljubna dogodka A in B velja:*

- $A \cap \bar{A} = N$, • $\bar{N} = G$, • $\overline{A \cup B} = \bar{A} \cap \bar{B}$,
- $A \cup \bar{A} = G$, • $\overline{\bar{A}} = A$, • $\overline{A \cap B} = \bar{A} \cup \bar{B}$. \square

Dogodka A in B sta **nezdružljiva**, če se ne moreta zgoditi hkrati, njun produkt je torej nemogoč dogodek, $A \cap B = N$.

Primer: Dogodka, A – da pri metu kocke pade sodo število pik in B – da pade liho število pik, sta nezdružljiva. \diamond

Poljuben dogodek in njegov nasprotni dogodek sta vedno nezdružljiva. Ob vsaki ponovitvi poskusa se zagotovo zgodi eden od njiju, zato je njuna vsota gotov dogodek:

$$(A \cap \bar{A} = N) \wedge (A \cup \bar{A} = G).$$

Če lahko dogodek A izrazimo kot vsoto nezdružljivih in mogočih dogodkov, rečemo, da je A **sestavljen** dogodek. Dogodek, ki ni sestavljen, imenujemo **osnoven** ali **elementaren** dogodek.

Primer: Pri metu kocke je šest osnovnih dogodkov: E_1 , da pade 1 pika, E_2 , da padeta 2 piki, ..., E_6 , da pade 6 pik. Dogodek, da pade sodo število pik je sestavljen dogodek iz treh osnovnih dogodkov (E_2 , E_4 in E_6). \diamond

Množico dogodkov $S = \{A_1, A_2, \dots, A_n\}$ imenujemo **popoln sistem dogodkov**, če se v vsaki ponovitvi poskusa zgodi natanko eden od dogodkov iz množice S . To pomeni, da ni noben med njimi nemogoč:

$$A_i \neq N, \quad \text{za } i = 1, 2, \dots, n,$$

paroma nezdružljivi

$$A_i \cap A_j = \emptyset \quad i \neq j$$

in njihova vsota je gotov dogodek

$$A_1 \cup A_2 \cup \dots \cup A_n = \Omega.$$

Primer: Popoln sistem dogodkov pri metu kocke sestavljajo na primer osnovni dogodki ali pa tudi dva dogodka: dogodek, da vržem sodo število pik, in dogodek, da vržem liho število pik. \diamond

2.3 Definicija verjetnosti



Opišimo najpreprostejšo verjetnostno zakonitost. Denimo, da smo n -krat ponovili dan poskus in da se je k -krat zgodil dogodek A . Ponovitve poskusa, v katerih se A zgodi, imenujemo ugodne za dogodek A , število

$$f(A) = \frac{k}{n}$$

pa je **relativna frekvenca** (pogostost) dogodka A v opravljenih poskusih. Statistični zakon, ki ga kaže izkušnja, je:

Če poskus X dolgo ponavljamo, se relativna frekvenca slučajnega dogodka ustali in sicer skoraj zmeraj toliko bolj, kolikor več ponovitev poskusa napravimo.

To temeljno zakonitost so empirično preverjali na več načinov. Najbolj znan je poskus s kovanci, kjer so določali relativno frekvenco grba ($f(A)$):

- Buffon je v 4040 metih dobil $f(A) = 0,5069$,
- Pearson je v 12000 metih dobil $f(A) = 0,5016$,
- Pearson je v 24000 metih dobil $f(A) = 0,5005$.

Ti poskusi kažejo, da se relativna frekvenca grba pri metih kovanca običajno ustali blizu 0,5. Ker tudi drugi poskusi kažejo, da je ustalitev relativne frekvence v dovolj velikem številu ponovitev poskusa splošna zakonitost, je smiselna naslednja *statistična definicija verjetnosti*:

Verjetnost dogodka A v danem poskusu je število $P(A)$, pri katerem se navadno ustali relativna frekvenca dogodka A v velikem številu ponovitev tega poskusa.

Ker je relativna frekvenca vedno nenegativna in kvečjemu enaka številu opravljenih poskusov ni težko narediti naslednje zaključke.

Trditev 2.4. *Za poljubna dogodka A in B velja:*

1. $P(A) \geq 0$.
2. $P(G) = 1, P(N) = 0$ in $A \subset B \Rightarrow P(A) \leq P(B)$.
3. Če sta dogodka A in B nezdružljiva, potem velja $P(A \cup B) = P(A) + P(B)$.

Klasični pristop k verjetnosti

Pri določitvi verjetnosti si pri nekaterih poskusih in dogodkih lahko pomagamo s *klasično definicijo verjetnosti*:

Vzemimo, da so dogodki iz popolnega sistema dogodkov $\{E_1, E_2, \dots, E_s\}$ enako verjetni: $P(E_1) = P(E_2) = \dots = P(E_s) = p$. Tedaj je $P(E_i) = 1/s$ za $i = 1, \dots, s$. Če je nek dogodek A sestavljen iz r dogodkov iz tega popolnega sistema dogodkov, potem je njegova verjetnost $P(A) = r/s$.

Primer: Izračunajmo verjetnost dogodka A , da pri metu kocke padejo manj kot 3 pike. Popolni sistem enako verjetnih dogodkov sestavlja 6 dogodkov. Od teh sta le dva ugodna za dogodek A (1 in 2 piki). Zato je verjetnost dogodka A enaka $2/6 = 1/3$. \diamond

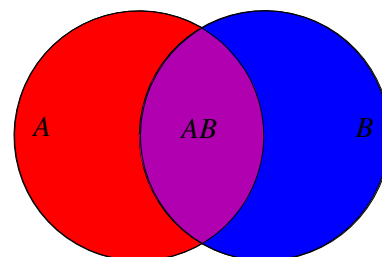
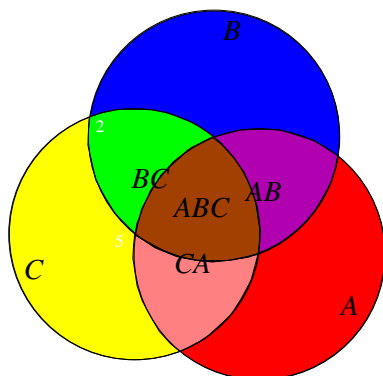
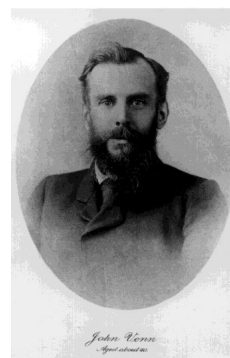
Geometrijska verjetnost

V primerih, ko lahko osnovne dogodke predstavimo kot 'enakovredne' točke na delu premice (ravnine ali prostora), določimo verjetnost sestavljenega dogodka kot razmerje dolžin (ploščin, prostornin) dela, ki ustreza ugodnim izidom, in dela, ki ustreza vsem možnim izidom.

2.4 Osnovne lastnosti verjetnosti

Trditev 2.5. Za poljubna dogodka A in B velja:

$$P(A \cup B) = P(A) + P(B) - P(A \cap B). \quad \square$$



Primer: Denimo, da je verjetnost, da študent naredi izpit iz Sociologije $P(S) = 2/3$. Verjetnost, da naredi izpit iz Politologije je $P(P) = 5/9$. Če je verjetnost, da naredi vsaj enega od obeh izpitov $P(S \cup P) = 4/5$, kolikšna je verjetnost, da naredi oba izpita?

$$P(S \cap P) = P(S) + P(P) - P(S \cup P) = \frac{2}{3} + \frac{5}{9} - \frac{4}{5} = 0,42. \quad \diamond$$

Posledica 2.6. $P(\bar{A}) = 1 - P(A)$. □

Primer: Iz kupa 32 kart slučajno povlečemo 3 karte. Kolikšna je verjetnost, da je med tremi kartami vsaj en as (dogodek A)? Pomagamo si z nasprotnim dogodkom \bar{A} , da med tremi kartami ni asa. Njegova verjetnost po klasični definiciji verjetnosti je določena s

kvocientom števila vseh ugodnih dogodkov v popolnem sistemu dogodkov s številom vseh dogodkov v tem sistemu dogodkov. Vseh dogodkov v popolnem sistemu dogodkov je $\binom{32}{3}$, ugodni pa so tisti, kjer zbiramo med ne-asi, tj. $\binom{28}{3}$. Torej je

$$P(\bar{A}) = \frac{\binom{28}{3}}{\binom{32}{3}} = 0,66; \quad P(A) = 1 - P(\bar{A}) = 1 - 0,66 = 0,34. \quad \diamond$$

Primer: Za n različnih dopisov, namenjenih n osebam, so pripravljene že naslovljene ovojnice. Dopise bomo na slepo razdelili v ovojnice. Kolika je pri tem verjetnost, da niti en dopis ne bo prišel na pravi naslov?

Negacija dogodka A , ki mu iščemo verjetnost, je da pride vsaj eno pismo na pravi naslov. Pisma uredimo po nekem vrstnem redu in naj bo A_i ($1 \leq i \leq n$) dogodek, da pride i -to pismo na pravi naslov. Potem je

$$\bar{A} = A_1 \cup \dots \cup A_n,$$

dogodki na levi strani pa niso nezdružljivi. Torej lahko uporabimo pravilo o vključitvi in izključitvi, pri čemer označimo verjetnost i -te vsote z S_i ($1 \leq i \leq n$). Potem ima vsota S_1 n členov, ki so vsi med seboj enaki, saj gre za izbiranje na slepo:

$$S_1 = n P(A_1).$$

Poskus je v našem primeru razdeljevanje n dopisov po ovojnicah, torej je možnih izidov $n!$ in so vsi med seboj enako verjetni. Med izidi so za A_1 ugodni tisti, pri katerih pride prvi dopis v prvo ovojnico, tj.

$$S_1 = \frac{n(n-1)!}{n!} = 1.$$

Nadalje je

$$S_2 = \binom{n}{2} P(A_1 A_2) = \binom{n}{2} \frac{(n-2)!}{n!} = \frac{1}{2!}.$$

in v splošnem $S_k = 1/k!$ ($1 \leq k \leq n$). Torej je

$$P(\bar{A}) = 1 - \frac{1}{2!} + \frac{1}{3!} - \frac{1}{4!} + \dots + \frac{(-1)^{n-1}}{n!} \quad \text{oziroma} \quad P(A) = \frac{1}{2!} - \frac{1}{3!} + \dots + \frac{(-1)^n}{n!}$$

in končno, če upoštevamo še Trditev A.5,

$$\lim_{n \rightarrow \infty} P(A) = \frac{1}{e},$$

tako da imamo že pri razmeroma majhnih n

$$P(A) \approx \frac{1}{e} \approx 0,369. \quad \diamond$$

Naloga. V družbi je n ljudi. Vsak od njih da svojo vizitko v posodo, nato pa iz nje vsak na slepo izbere po eno vizitko. Dokaži, da je verjetnost, da bo natanko r ljudi ($0 \leq r \leq n$) dobilo svojo vizitko, enaka

$$p_r(n) = \frac{1}{r!} \left(1 - \frac{1}{1!} + \frac{1}{2!} - \frac{1}{3!} + \dots + \frac{(-1)^{n-r}}{(n-r)!} \right).$$

Omenimo še dve posledici, ki prideta pogosto prav. Naslednjo trditev lahko dokažemo na enak način kot Trditev 2.5 ali pa kar z uporabo tega izreka.

Posledica 2.7. *Za dogodke A , B in C velja:*

$$P(A \cup B \cup C) = P(A) + P(B) + P(C) - P(A \cap B) - P(A \cap C) - P(B \cap C) + P(A \cap B \cap C). \quad \square$$

Kako lahko to pravilo posplošimo še na več dogodkov?

Namig: Pravilo o vključitvi in izključitvi za množice A_1, A_2, \dots, A_n :

$$\begin{aligned} |A_1 \cup A_2 \cup \dots \cup A_n| &= \sum_{i=1}^n |A_i| - \sum_{1 \leq i_1 < i_2 \leq n} |A_{i_1} \cap A_{i_2}| \\ &+ \sum_{1 \leq i_1 < i_2 < i_3 \leq n} |A_{i_1} \cap A_{i_2} \cap A_{i_3}| - \dots + (-1)^{n-1} |A_1 \cap A_2 \cap \dots \cap A_n|. \end{aligned}$$

Posledica 2.8. *Če so dogodki A_i , $i \in I$ paroma nezdružljivi, velja*

$$P\left(\bigcup_{i \in I} A_i\right) = \sum_{i \in I} P(A_i).$$

Velja tudi za števno neskončne množice dogodkov. □

Zgornjo relacijo si lažje zapomnimo, če jo na glas “odrecitiramo”:

“**verjetnost vsote** paroma nezdružljivih dogodkov je enaka **vsoti verjetnosti** teh dogodkov.”

Torej gre za nekakšno pravilo o zamenjavi (med vrstnim redom računanja vsote in verjetnosti), ki velja za paroma nezdružljive dogodke.

2.5 Aksiomi Kolmogorova

Dogodek predstavimo z množico zanj ugodnih izidov; gotov dogodek G ustreza univerzalni množici; nemogoč dogodek pa prazni množici. Neprazna družina dogodkov \mathcal{D} je **algebra dogodkov**, če velja:

- $A \in \mathcal{D} \Rightarrow \bar{A} \in \mathcal{D}$,
- $A, B \in \mathcal{D} \Rightarrow A \cup B \in \mathcal{D}$.

Pri neskončnih množicah dogodkov moramo drugo zahtevo posplošiti

- $A_i \in \mathcal{D}, i \in I \Rightarrow \bigcup_{i \in I} A_i \in \mathcal{D}$.

Dobljeni strukturi rečemo **σ -algebra**.



Naj bo \mathcal{D} σ -algebra v G . **Verjetnost na G** je preslikava $P : \mathcal{D} \rightarrow \mathbb{R}$ z lastnostmi:

1. $P(A) \geq 0$.
2. $P(G) = 1$.
3. Če so dogodki $A_i, i \in I$ paroma nezdružljivi, je

$$P\left(\bigcup_{i \in I} A_i\right) = \sum_{i \in I} P(A_i).$$

Trojica (G, \mathcal{D}, P) določa **verjetnostni prostor**.

Andrei Kolmogorov (1903-1987)

<http://www.exploratorium.edu/complexity/CompLexicon/kolmogorov.html>

Kolmogorov was one of the broadest of this century's mathematicians. He laid the mathematical foundations of probability theory and the algorithmic theory of randomness and made crucial contributions to the foundations of statistical mechanics, stochastic processes, information theory, fluid mechanics, and nonlinear dynamics. All of these areas, and their interrelationships, underlie complex systems, as they are studied today. Kolmogorov graduated from Moscow State University in 1925 and then became a professor there in 1931. In 1939 he was elected to the Soviet Academy of Sciences, receiving the Lenin Prize in 1965 and the Order of Lenin on seven separate occasions.

His work on reformulating probability started with a 1933 paper in which he built up probability theory in a rigorous way from fundamental axioms, similar to Euclid's treatment of geometry. Kolmogorov went on to study the motion of the planets and turbulent fluid flows, later publishing two papers in 1941 on turbulence that even today are of fundamental importance. In 1954 he developed his work on dynamical systems in relation to planetary motion, thus demonstrating the vital role of probability theory in physics and re-opening the study of apparent randomness in deterministic systems, much along the lines originally conceived by Henri Poincare.

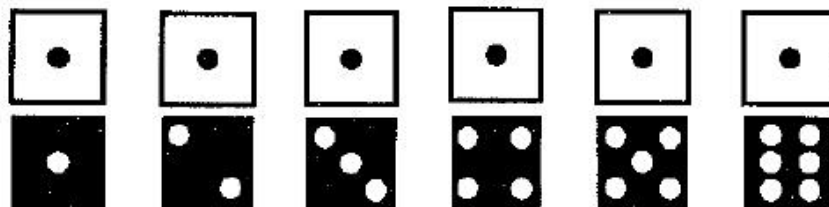
In 1965 he introduced the algorithmic theory of randomness via a measure of complexity, now referred to Kolmogorov Complexity. According to Kolmogorov, the complexity of an object is the length of the shortest computer program that can reproduce the object. Random objects, in his view, were their own shortest description. Whereas, periodic sequences have low Kolmogorov complexity, given by the length of the smallest repeating "template" sequence they contain. Kolmogorov's notion of complexity is a measure of randomness, one that is closely related to Claude Shannon's entropy rate of an information source.

Kolmogorov had many interests outside mathematics research, notable examples being the quantitative analysis of structure in the poetry of the Russian author Pushkin, studies of agrarian development in 16th and 17th century Novgorod, and mathematics education.

Nas bo kasneje zanimal tudi Kolmogorov-Smirnov test
(<http://www.physics.csbsju.edu/stats/KS-test.html>).

Poglavje 3

Pogojna verjetnost



3.1 Intriga (po Kvarkadabri)

Ne podcenjujmo vpliva najrazličnejših rubrik v popularnih časopisnih prilogah, kjer nas domnevni strokovnjaki zasipajo z nasveti vseh vrst, - rubrike krojijo mnenja ljudi in spreminjajo navade celotnih nacij, - sprožajo obsežne polemike tako med širšimi množicami kot tudi v ozki strokovni javnosti. Na področju zdravja in prehrane tako burne odzive seveda pričakujemo, povsem nekaj drugega pa je, če jih sproži preprosto *matematično vprašanje*.

Revija *Parade* - kot prilogo jo vsako nedeljo dodajo več kot 400 ameriškim časopisom in doseže okoli 70 milijonov bralcev, že dolgo izhaja rubrika z imenom "Vprašajte Marilyn." Ureja jo **Marilyn vos Savant**. Sredi 80ih jo je *Guinnessova knjiga rekordov* razglasila za rekorderko z najvišjim inteligenčnim količnikom na planetu. V svoji rubriki zdaj že več kot 20 let odgovarja na najrazličnejša vprašanja bralcev in rešuje njihove težave. Med vsemi vprašanji, ki jih je kdaj obravnavala, ima prav posebno mesto na prvi pogled zelo preprost problem, ki ji ga je 9. septembra 1990 zastavil gospod Craig F. Whitaker:



Dve kozi in avtomobil

“Vzemimo, da sodelujete v nagradni igri, kjer vam ponudijo na izbiro troje vrat. Za enimi se skriva avto, za drugima dvema pa koza. Recimo, da izberete vrata številka 3, voditelj igre, ki ve, kaj se nahaja za posameznimi vrati, pa nato odpre vrata številka 1, za katerimi se pokaže koza. Nato vas vpraša: ‘Bi se sedaj raje odločili za vrata številka 2?’

Zanima me, ali se tekmovalcu spleča zamenjati izbor vrat?”

Poudariti je potrebno, da mora gostitelj nagradne igre vsakič postopati enako. Ne more enkrat ponuditi zamenjavo (npr. takrat, ko vidi, da nastopajoči kaže na vrata za katerimi se skriva avto), drugič pa ne (npr. takrat, ko nastopajoči kaže na vrata za katerimi je koza).

Vprašanja se je prijelo ime “*Monty Hall problem*”, po imenu voditelja popularne ameriške televizijske oddaje *Pogodimo se* (Let’s Make a Deal), v kateri je voditelj Monty Hall goste izzival, da so sprejemali ali zavračali najrazličnejše ponudbe, ki jim jih je zastavljal. Marilyn je bralcu v svoji rubriki odgovorila, da se nam vrata vsekakor spleča zamenjati, saj se tako verjetnost, da bomo zadeli avto, poveča za dvakrat. Tole je njen odgovor: **Seveda se spleča zamenjati vrata**. Prva vrata imajo le $1/3$ verjetnosti za zmago, medtem ko imajo druga verjetnost $2/3$.

Najlažje si vse skupaj predstavljate takole. Predpostavimo, da je na voljo milijon vrat in vi izberete prva. Nato voditelj, ki ve, kaj se nahaja za posameznimi vrati, odpre vsa vrata razen prvih vrat in vrat številka 7777777. V tem primeru bi zelo hitro zamenjali svoj izbor, kajne? Se najinteligentnejša ženska na planetu moti?

Sledila je ploha kritik (več kot 10.000 pisem jeznih bralcev, med katerimi je bilo ogromno učiteljev matematike). Skoraj 1000 pisem je bilo podpisanih z imeni (dr. nazivi, napisana na papirju z glavo katere od ameriških univerz - www.marilynvossavant.com). Marilyn bralce zavraža, **saj se verjetnost za zadetek nikakor ne more spremeniti, če vmes zamenjamo**

izbor vrat. Neki profesor matematike je bil zelo neposreden: “Udarili ste mimo! ... Kot profesionalni matematik sem zelo zaskrbljen nad pomanjkanjem matematičnih veščin v širši javnosti. Prosim, da se opravičite in ste v prihodnosti bolj pazljivi.” Drugi je Marylin celo obtožil, da je ona sama koza.

Polemika je pristala celo na naslovnici New York Timesa, v razpravo so se vključila tudi nekatera znana imena iz sveta matematike. O odgovoru vos Savantove, da naj tekmovalec zamenja vrata, so razpravljali tako na hodnikih Cie kot v oporiščih vojaških pilotov ob Perzijskem zalivu. Analizirali so ga matematiki z MIT in računalniški programerji laboratorijev Los Alamos v Novi Mehiki. Poleg žaljivih pisem, ki so njen odgovor kritizirala, je Marilyn vseeno prejela tudi nekaj pohval. Profesor s prestižnega MIT: “Seveda imate prav. S kolegi v službi smo se poigrali s problemom in moram priznati, da je bila večina, med njimi sem bil tudi sam, sprva prepričana, da se motite!”

Eksperimentalna ugotovitev

Marilyn se kritik ni ustrašila - navsezadnje je objektivno izmerljivo po inteligenčnem količniku pametnejša od vseh svojih kritikov, zato je v eni od svojih naslednjih kolumn vsem učiteljem v državi zadala nalogo, da to preprosto igrico igrajo s svojimi učenci v razredu (seveda ne s pravimi kozami in avtomobilom) in ji pošljejo svoje rezultate. Te je nato tudi objavila in seveda so se povsem skladali z njenim nasvetom, da se v tem konkretnem primeru bistveno bolj spleča spremeniti izbiro vrat. Kdo ima prav?

Razprava o Monty Hall problemu spada na področje, ki mu matematiki pravijo **po-gojna verjetnost**. Najbolj preprosto rečeno je to veda, ki se ukvarja s tem, kako prilagoditi verjetnost za posamezne dogodke, ko se pojavijo novi podatki. Bistvo zapleta, ki je izzval tako obsežno in čustveno nabito reakcijo bralcev, je v tem, da so bralci večinoma spregledali ključni podatek. Zelo pomembno je namreč dejstvo, da **voditelj igre vnaprej ve**, za katerimi vrati je avtomobil.

Ko v drugem delu odpre vrata, za katerimi se pokaže koza, vnaprej ve, da za temi vrati ni avtomobila. Če voditelj te informacije ne bi imel in bi vrata odpiral povsem naključno tako kot igralec, se verjetnost za zadetek ob spremembi vrat res ne bi povečala. Potem bi držale ugotovitve več 1000 bralcev, ki so poslali jezna pisma na uredništvo revije, da Marilyn ne pozna osnov matematike. Matematična intuicija nam namreč pravi, da je verjetnost, da bo avto za enimi ali za drugimi vrati, ko so dvojca še zaprta, enaka. To je seveda res, če zraven ne bi bilo še voditelja, ki ve več kot mi.

Najlažje nejasnost pojasnimo, če analiziramo dogajanje **izza kulis**, od koder ves čas vidimo, za katerimi vrati je avto in kje sta kozi. Če tekmovalec že v prvo izbere vrata,

za katerimi je avto, bo voditelj odprl katera koli od preostalih dveh vrat in zamenjava bo tekmovalcu v tem primeru le škodila. Ampak to velja le za primer, če v prvo izbere vrata, za katerimi je avto, verjetnost za to pa je $1/3$. Če pa v prvo tekmoalec izbere vrata, za katerimi je koza, bo voditelj moral odpreti edina preostala vrata, za katerimi se nahaja koza. V tem primeru se bo tekmovalcu zamenjava vrat v vsakem primeru obrestovala in bo tako z gotovostjo zadel avto.

Če v prvo tekmoalec izbere kozo, se mu vedno splača zamenjati, če pa v prvo izbere avto, se mu zamenjava ne izplača. Verjetnost, da v prvo izbere kozo, je $2/3$, medtem ko je verjetnost, da izbere avto, le $1/3$. Če se tekmoalec odloči za strategijo zamenjave, je zato verjetnost, da zadane avtomobil, $2/3$, če zamenjavo zavrne, pa je verjetnost pol manjša, tj. $1/3$. Če se torej drži strategije zamenjave vrat, ko mu jo voditelj ponudi, bo tako vedno, ko v prvo izbere kozo, ob zamenjavi vrat dobil avto, kar ga do dobitka pripelje v $2 \times$ večjem številu primerov, kot sicer. **Verjetnost za zadetek se mu tako s 33% poveča na 66%**. Če vam ni takoj jasno, se ne sekirajte preveč. Tudi mnogi matematiki so potrebovali kar nekaj časa, da so si razjasnili ta problem.

3.2 Definicija pogojne verjetnosti

Opazujemo dogodek A ob poskusu X , ki je realizacija kompleksa pogojev K . Verjetnost dogodka A je tedaj $P(A)$. Kompleksu pogojev K pridružimo mogoč dogodek B , tj. $P(B) > 0$. Realizacija tega kompleksa pogojev $K' = K \cap B$ je poskus X' in verjetnost dogodka A v tem poskusu je $P_B(A)$, ki se z verjetnostjo $P(A)$ ujema ali pa ne. Pravimo, da je poskus X' poskus X s pogojem B in verjetnost $P_B(A)$ **pogojna verjetnost** dogodka A glede na dogodek B , kar zapišemo takole:

$$P_B(A) = P(A/B).$$

Pogojna verjetnost $P(A/B)$ v poskusu X' je verjetnost dogodka A v poskusu X s pogojem B . Je torej tudi verjetnost (podobno kot npr. je podgrupa grupe zopet grupa), le obravnavani kompleks pogojev, ki mora biti izpolnjen se je spremenil. Pogosto pogojno verjetnost pišejo tudi $P(A/B)$ (oziroma $P(A|B)$). Denimo, da smo n -krat ponovili poskus X in da se je ob tem k_B -krat zgodil dogodek B . To pomeni, da smo v n ponovitvah poskusa X napravili k_B -krat poskus X' . Dogodek A se je zgodil ob poskusu X' le, če se je zgodil tudi B , tj. $A \cap B$. Denimo, da se je dogodek $A \cap B$ zgodil ob ponovitvi poskusa $k_{A \cap B}$ -krat. Potem je

relativna frekvenca dogodka A v opravljenih ponovitvah poskusa X' :

$$f_B(A) = f(A/B) = \frac{k_{A \cap B}}{k_B} = \frac{k_{A \cap B}/n}{k_B/n} = \frac{f(A \cap B)}{f(B)}$$

in smo prišli do naslednje trditve.

Trditev 3.1. *Za dogodka A in B , kjer je $P(B) \neq 0$, velja*

$$P(A/B) = \frac{P(A \cap B)}{P(B)}. \quad \square$$

Pogojna verjetnost P_B ima prav take lastnosti kot brezpogojna. Trojica (B, \mathcal{D}_B, P_B) , $\mathcal{D}_B = \{A \cap B \mid A \in \mathcal{D}\}$ je zopet verjetnostni prostor.

Primer: Denimo, da je v nekem naselju 900 polnoletnih prebivalcev. Zanima nas struktura prebivalcev po spolu (M – moški, Ž – ženski spol) in po zaposlenosti (Z – zaposlen(a), N – nezaposlen(a)). Podatke po obeh spremenljivkah uredimo v dvorazsežno frekvenčno porazdelitev, ki jo imenujemo tudi **kontingenčna tabela**:

<i>spol \ zap.</i>	Z	N	
M	460	40	500
Ž	240	160	400
	700	200	900

Poglejmo, kolikšna je verjetnost, da bo slučajno izbrana oseba moški pri pogoju, da je zaposlena.

$$P(Z) = \frac{700}{900}, \quad P(M \cap Z) = \frac{460}{900}, \quad P(M/Z) = \frac{P(M \cap Z)}{P(Z)} = \frac{460 \cdot 900}{900 \cdot 700} = \frac{460}{700}$$

ali neposredno iz kontingenčne tabele

$$P(M/Z) = \frac{460}{700}. \quad \diamond$$

Iz formule za pogojno verjetnost sledita naslednji zvezi:

$$P(A \cap B) = P(B) P(A/B) \quad \text{in} \quad P(A \cap B) = P(A) P(B/A).$$

Torej velja:

$$P(A) P(B/A) = P(B) P(A/B).$$

Dogodka A in B sta **neodvisna**, če velja

$$P(A/B) = P(A).$$

Zato za neodvisna dogodka A in B velja $P(A \cap B) = P(A) \cdot P(B)$. Za par nezdržljivih dogodkov A in B pa velja $P(A/B) = 0$.

Primer: Iz posode, v kateri imamo 8 belih in 2 rdeči krogli, $2 \times$ na slepo izberemo po eno kroglo. Kolikšna je verjetnost dogodka, da je prva krogla bela (B_1) in druga rdeča (R_2)?

1. Če po prvem izbiranju izvlečeno kroglo ne vrnemo v posodo (odvisnost), je:

$$P(B_1 \cap R_2) = P(B_1) \cdot P(R_2/B_1) = \frac{8}{10} \cdot \frac{2}{9} = 0,18.$$

2. Če po prvem izbiranju izvlečeno kroglo vrnemo v posodo (neodvisnost), je:

$$P(B_1 \cap R_2) = P(B_1) \cdot P(R_2/B_1) = P(B_1) \cdot P(R_2) = \frac{8}{10} \cdot \frac{2}{10} = 0,16. \quad \diamond$$

Trditev 3.2. Dogodka A in B sta neodvisna, če je $P(A/B) = P(A/\bar{B})$.

Nadalje velja

$$P(A \cap B \cap C) = P(A) \cdot P(B/A) \cdot P(C/(A \cap B)). \quad \square$$

(Tudi slednje pravilo lahko posplošimo naprej.) Dogodki A_i , $i \in I$ so **neodvisni**, če je $P(A_j) = P(A_j / \bigcap_{i=1}^{j-1} A_i)$, $j \in I$.

Za neodvisne dogodke A_i , $i \in I$ velja

$$P\left(\bigcap_{i \in I} A_i\right) = \prod_{i \in I} P(A_i).$$

Pomembno je poudariti razliko med *nezdržljivostjo* in *neodvisnostjo*, zato je zopet na vrsti "recitiranje":

“**verjetnost produkta** paroma neodvisnih dogodkov je enaka **produktu verjetnosti** teh dogodkov.”

Torej gre za nekakšno pravilo o zamenjavi (med vrstnim redom računanja produkta in verjetnosti), ki velja za paroma neodvisne dogodke.

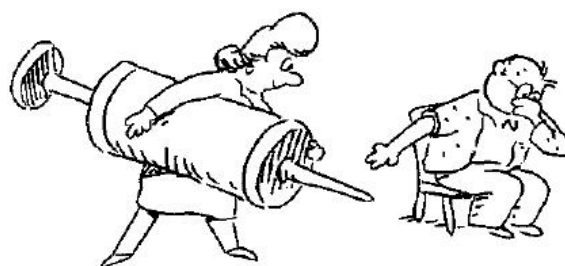
Primer: Redko nalezljivo bolezen dobi ena oseba na 1000. Imamo dober, a ne popoln test za to bolezen:

*če ima neka oseba to bolezen, potem test to pokaže v 99% primerih,
vendar pa test napačno označi tudi 2% zdravih pacientov za bolane.*

V tvojem primeru je bil test pravkar **pozitiven**. Kakšna je verjetnost, da si zares dobili nalezljivo bolezen? Delamo z naslednjimi dogodki:

A: pacient je dobil nalezljivo bolezen,

B: pacientov test je bil pozitiven.



Izrazimo informacijo o učinkovitosti testov:

$P(A) = 0,001$ (en pacient na 1000 se naleze),

$P(B/A) = 0,99$ (test pravilno označi okuženega),

$P(B/\bar{A}) = 0,02$ (test napačno označi zdravega).

Zanima nas $P(A/B)$ (verjetnost, da smo se nalezli, če je test pozitiven). ◇

3.3 Obrazec za popolno verjetnost in večstopenjski poskusi

Naj bo H_i , $i \in I$ popoln sistem dogodkov, tj. **razbitje** gotovega dogodka: $\bigcup_{i \in I} H_i = G$, na paroma nezdružljive dogodke: $H_i \cap H_j = \emptyset$, $i \neq j$. Gotov dogodek smo torej kot hlebec narezali s pomočjo hipotez na posamezne kose, da jih bomo lažje obvladali. Zanima nas verjetnost dogodka A , če poznamo verjetnost $P(H_i)$, in pogojno verjetnost $P(A/H_i)$ za $i \in I$:

$$A = A \cap (H_1 \cup H_2 \cdots H_n) = (A \cap H_1) \cup \cdots \cup (A \cap H_n).$$

Sedaj si bomo ogledali verjetnost dogodka A na posameznih 'kosih', skoraj tako kot pri marmornem kolaču. Ker so tudi dogodki $A \cap H_i$ paroma nezdružljivi, velja:

Trditev 3.3. *Za popoln sistem dogodkov $H_i, i \in I$ in poljuben dogodek A velja*

$$P(A) = \sum_{i \in I} P(A \cap H_i) = \sum_{i \in I} P(H_i)P(A/H_i). \quad \square$$

Zgornji trditvi pravimo tudi *izrek o popolni verjetnosti* ali pa obrazec za razbitje.

Na stvar lahko pogledamo tudi kot na večstopenjski poskus: v prvem koraku se zgodi natanko eden od dogodkov H_i , ki ga imenujemo domneva (hipoteza) (domneve sestavljajo popoln sistem dogodkov). Šele izidi na prejšnjih stopnjah določajo, kako bo potekal poskus na naslednji stopnji.

Omejimo se na poskus z dvema stopnjama. Naj bo A eden izmed mogočih dogodkov na drugi stopnji. Včasih nas zanima po uspešnem izhodu tudi druge stopnje, verjetnost tega, da se je na prvi stopnji zgodil dogodek H_i . Odgovor nam da naslednja trditev.

Trditev 3.4. (Bayesov obrazec) *Za popoln sistem dogodkov $H_i, i \in I$ in poljuben dogodek A velja*

$$P(H_k/A) = \frac{P(H_k) \cdot P(A/H_k)}{\sum_{i \in I} P(H_i) \cdot P(A/H_i)}. \quad \square$$

Primer: Trije lovci so hkrati ustrelili na divjega prašiča in ga ubili. Ko so prišli do njega, so našli v njem eno samo kroglo. Kolikšne so verjetnosti, da je vepra ubil (a) prvi, (b) drugi, (b) tretji. lovec, če poznamo njihove verjetnosti, da zadanejo: 0, 2; 0, 4 in 0, 6? Na ta način jim namreč lahko pomagamo pri pošteni delitvi plena (kajti ne smemo pozabiti, da imajo vsi v rokah nevarno orožje). Sestavimo popoln sistem dogodkov in uporabimo dejstvo, da so lovci med seboj neodvisni, torej $P(A * B * C) = P(A) * P(B) * P(C)$. To nam zna pomagati pri računanju verjetnosti hipotez.

	.2	.4	.6				
	prvi	drugi	tretji	P(H_i)	st.kr.	P(E/H_i)	P(E*H_i)
H1	1	1	1	,2*,4*,6 =0,048	3	0	0
H2	0	1	1	,8*,4*,6 =0,192	2	0	0
H3	1	0	1	,2*,6*,6 =0,072	2	0	0
H4	1	1	0	,2*,4*,4 =0,032	2	0	0
H5	1	0	0	,2*,6*,4 =0,048	1	1	0,048
H6	0	1	0	,8*,4*,4 =0,128	1	1	0,128
H7	0	0	1	,8*,6*,6 =0,288	1	1	0,288
H8	0	0	0	,8*,6*,4 =0,192	0	0	0
vsota				=1,000			0,464

$$P(\text{ena krogla je zadela}) = 0,048 + 0,128 + 0,288 = 0,464 = P(E).$$

Ostale verjetnosti računamo za preiskus:

$$P(\text{nobena krogla ni zadela}) = 0,192 = P(N'),$$

$$P(\text{dve krogli sta zadeli}) = 0,192 + 0,072 + 0,032 = 0,296 = P(D),$$

$$P(\text{tri krogle so zadele}) = 0,048 = P(T).$$

Vsota teh verjetnosti je seveda enaka 1. Končno uporabimo Bayesov obrazec:

$$P(H_5/E) = \frac{P(H_5 * E)}{P(E)} = \frac{0,048}{0,464} = 0,103 = P(\text{prvi je zadel}/E),$$

$$P(H_6/E) = \frac{P(H_6 * E)}{P(E)} = \frac{0,128}{0,464} = 0,276 = P(\text{drugi je zadel}/E),$$

$$P(H_7/E) = \frac{P(H_7 * E)}{P(E)} = \frac{0,288}{0,464} = 0,621 = P(\text{tretji je zadel}/E).$$

Tudi vsota teh verjetnosti pa je enaka 1. Delitev plena se opravi v razmerju 10,3 : 27,6 : 62,1 = 3 : 8 : 18 (in ne 2 : 4 : 6 oziroma 16,6 : 33,3 : 50, kot bi kdo utegnil na hitro pomisliti). \diamond

Bonus vprašanje: Kako bi si razdelili plen, če bi v divjim prašiču našli dve krogli?

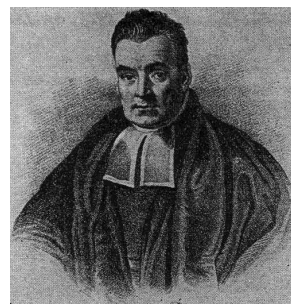
Za konec podajmo še presenetljivo kratko rešitev Monty Hall problema z uporabo pogojne verjetnosti:

Primer: (Rešitev problema Monty Hall) Z A označimo dogodek, da dobimo avto, z V pa da na začetku izberemo prava vrata. Potem obravnavamo naslednji možnosti:

1) Si ne premislimo: $P(A) = P(V) = 1/3$

2) Si premislimo: $P(A) = P(A/V) \cdot P(V) + P(A/\bar{V}) \cdot P(\bar{V}) = 0 \cdot 1/3 + 1 \cdot 2/3 = 2/3$, kjer smo pri drugi možnosti uporabili obrazec za razbitje. Je pa tudi res, da se morata ti dve verjetnosti sešteti v 1 in da nam zato druge v resnici sploh ni treba računati s pomočjo obrazca za razbitje. Vendar smo tako vseeno bolj prepričani o pravilnosti našega sklepanja, kajti kot smo spoznali v sami zgodbi, se je kar veliko ljudi motilo pri tem problemu. \diamond

Leta 2001 je bila na vrsti že 300 letnica rojstva angleškega matematika Bayesa.



REV. T. BAYES

Thomas Bayes (1702-1761) was an English clergyman who set out his theory of probability in 1764. His conclusions were accepted by Laplace in 1781, rediscovered by Condorcet, and remained unchallenged until Boole questioned them. Since then Bayes' techniques have been subject to controversy.

<http://www.york.ac.uk/depts/maths/histstat/bayesbiog.pdf>

<http://www.britannica.com/EBchecked/topic/56807/Thomas-Bayes>

English Nonconformist theologian and mathematician who was the first to use probability inductively and who established a mathematical basis for probability inference (a means of calculating, from the frequency with which an event has occurred in prior trials, the probability that it will occur in future trials. See probability theory: Bayes's theorem.

Bayes set down his findings on probability in "Essay Towards Solving a Problem in the Doctrine of Chances" (1763), published posthumously in the Philosophical Transactions of the Royal Society. That work became the basis of a statistical technique, now called Bayesian estimation, for calculating the probability of the validity of a proposition on the basis of a prior estimate of its probability and new relevant evidence. Disadvantages of the method—pointed out by later statisticians—include the different ways of assigning prior distributions of parameters and the possible sensitivity of conclusions to the choice of distributions.

The only works that Bayes is known to have published in his lifetime are *Divine Benevolence; or, An Attempt to Prove That the Principal End of the Divine Providence and Government Is the Happiness of His Creatures* (1731) and *An Introduction to the Doctrine of Fluxions, and a Defence of the Mathematicians Against the Objections of the Author of The Analyst* (1736), which was published anonymously and which countered the attacks by Bishop George Berkeley on the logical foundations of Sir Isaac Newton's calculus.

Bayes was elected a fellow of the Royal Society in 1742.

Poglavje 4

Bernoullijevo zaporedje neodvisnih poskusov



O zaporedju neodvisnih poskusov $X_1, X_2, \dots, X_n, \dots$ govorimo tedaj, ko so verjetnosti izidov v enem poskusu neodvisne od tega, kaj se zgodi v drugih poskusih.

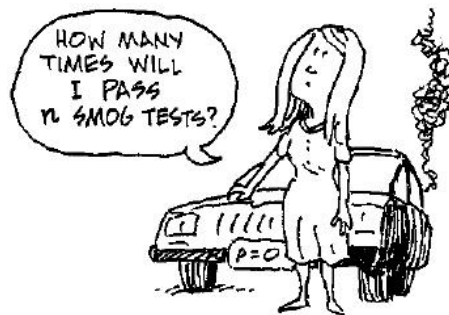
Zaporedje neodvisnih poskusov se imenuje **Bernoullijevo zaporedje**, če se more zgoditi v vsakem poskusu iz zaporedja neodvisnih poskusov le dogodek A z verjetnostjo $P(A) = p$ ali dogodek \bar{A} z verjetnostjo $P(\bar{A}) = 1 - P(A) = 1 - p = q$.



JAKOB BERNOULLI um 1687

Primer: Primer Bernoullijevega zaporedja poskusov je met kocke, kjer ob vsaki ponovitvi poskusa pade šestica (dogodek A) z verjetnostjo $P(A) = p = 1/6$ ali ne pade šestica (dogodek \bar{A}) z verjetnostjo $P(\bar{A}) = 1 - p = q = 5/6$. \diamond

V Bernoullijevem zaporedju neodvisnih poskusov nas zanima, kolikšna je verjetnost, da se v n zaporednih poskusih zgodi dogodek A natanko k -krat. To se lahko zgodi na primer tako, da se najprej zgodi k -krat dogodek A in nato v preostalih $(n - k)$ poskusih zgodi nasprotni dogodek \bar{A} :



$$P\left(\bigcap_{i=1}^k (X_i = A) \cap \bigcap_{i=k+1}^n (X_i = \bar{A})\right) = \prod_{i=1}^k P(A) \cdot \prod_{i=k+1}^n P(\bar{A}) = p^k \cdot q^{n-k}.$$

Dogodek $P_n(k)$, da se dogodek A v n zaporednih poskusih zgodi natanko k -krat, se lahko zgodi tudi na druge načine in sicer je teh toliko, na kolikor načinov lahko izberemo k poskusov iz n poskusov. Teh je $\binom{n}{k}$. Ker so ti načini nezdružljivi med seboj, je verjetnost dogodka $P_n(k)$ enaka

$$P_n(k) = \binom{n}{k} p^k (1-p)^{n-k}.$$

Tej zvezi pravimo **Bernoullijev obrazec**.

Primer: Iz posode, v kateri imamo 8 belih in 2 rdeči krogli, na slepo izberemo po eno kroglo in po izbiranju izvlečeno kroglo vrnemo v posodo. Kolikšna je verjetnost, da v petih poskusih izberemo 3-krat belo kroglo? Dogodek A je, da izvlečem belo kroglo. Potem je

$$p = P(A) = \frac{8}{10} = 0,8 \quad \text{in} \quad q = 1 - p = 1 - 0,8 = 0,2$$

Verjetnost, da v petih poskusih izberemo 3-krat belo kroglo, je:

$$P_5(3) = \binom{5}{3} 0,8^3 (1 - 0,8)^{5-3} = 0,205. \quad \diamond$$

4.1 Računanje $P_n(k)$

Ena možnost je uporaba rekurzije:

Trditev 4.1.

$$P_n(0) = q^n, \quad P_n(k) = \frac{(n-k+1)p}{kq} P_n(k-1), \quad \text{za } k = 1, \dots, n. \quad (4.1)$$

Dokaz.

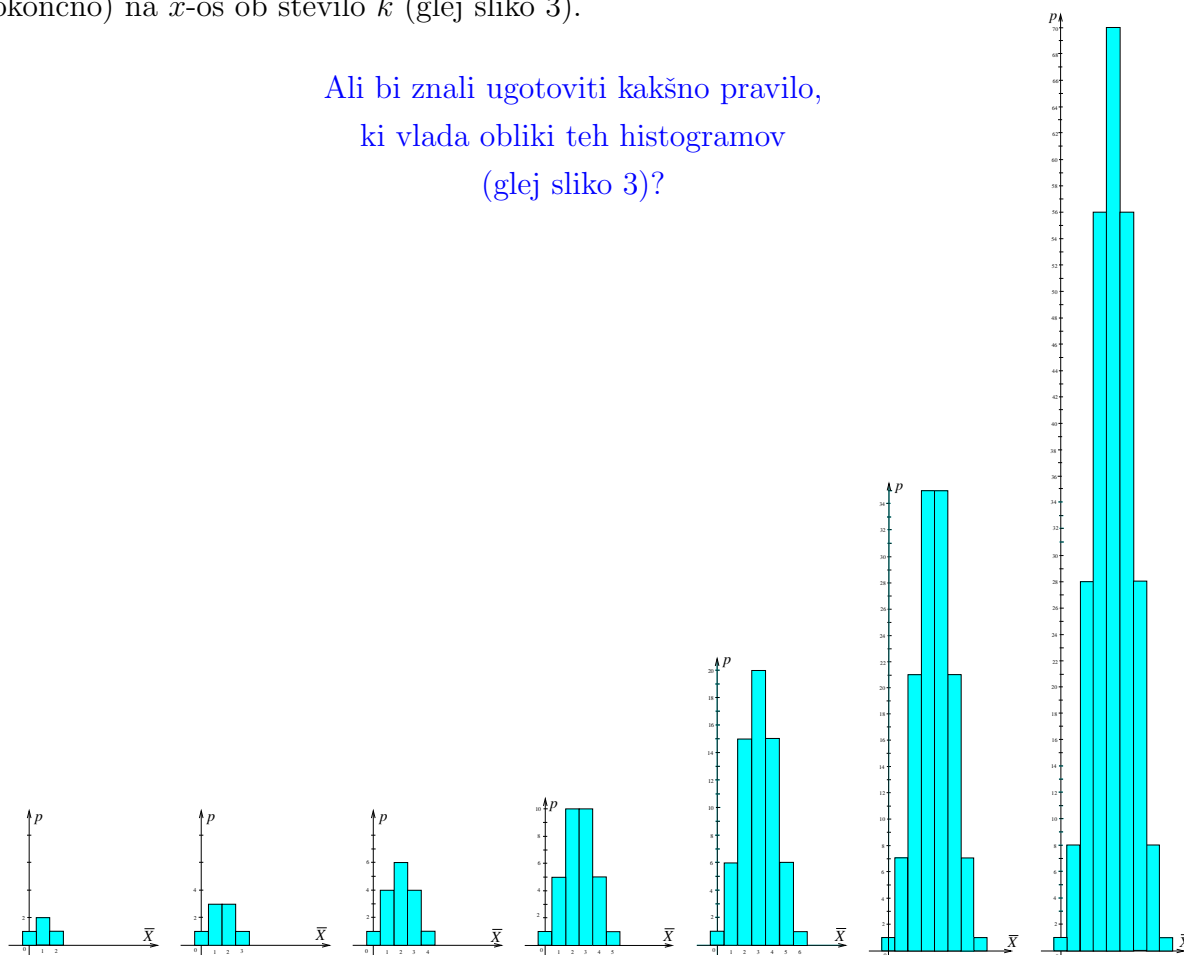
$$\frac{P_n(k)}{P_n(k-1)} = \frac{\binom{n}{k} p^k q^{n-k}}{\binom{n}{k-1} p^{k-1} q^{n-k+1}} = \frac{n! (k-1)! (n-k+1)! p}{k! (n-k)! n! q} = \frac{(n-k+1)p}{kq}. \quad \square$$

Vendar pa je takšno računanje izredno zamudno (eksponentno), tako kot rekurzivno računanje faktorjela, glej razdelek A.6.

Usojena krivulja v Pascalovem trikotniku

Iz vsake vrstice Pascalovega trikotnika lahko narišemo histogram, tj. vsakemu binomskemu simbolu $\binom{n}{k}$ za $k = 0, 1, \dots, n$, priredimo stolpec velikosti $\binom{n}{k} \times 1$ in ga postavimo (pokončno) na x -os ob število k (glej sliko 3).

Ali bi znali ugotoviti kakšno pravilo,
ki vlada obliki teh histogramov
(glej sliko 3)?



Slika 3: Predstavitev 2., 3., ..., 7. in 8. vrstice Pascalovega trikotnika.

Eno je gotovo, višina histograma vrtoglavo raste in zato smo se ustavili že pri $n = 8$. Vsekakor gre za 'enogrbo kamelo', simetrično glede na $n/2$, tj. števila najprej rastejo,

nato se pri lihih n največje število enkrat ponovi in nato (simetrično) padajo. Vsota vseh stolpcev je seveda enaka 2^n , saj po binomskem obrazcu velja:

$$2^n = (1 + 1)^n = \binom{n}{0} + \binom{n}{1} + \binom{n}{2} + \cdots + \binom{n}{n-2} + \binom{n}{n-1} + \binom{n}{n}.$$

Uvedemo nekaj dopolnil, ki bodo izboljšala naše slike, in bomo lahko z malo sreče v sliko spravili še kakšno vrstico Pascalovega trikotnika. Lahko se odločimo, da bo

- višina najvišjega stolpca v vsakem histogramu enaka,
- vsota ploščin vseh stolpcev v vsakem histogramu pa prav tako.

Drugo zahtevo hitro dosežemo tako, da za višino k -tega stolpca v histogramu uporabimo

$$\text{namesto } \binom{n}{k} \quad \text{raje } \binom{n}{k}/2^n,$$

Slednje število je ravno $P_n(k)$ iz Bernoullijevega obrazca za $p = q = 1/2$. Vendar pa sedaj opazimo, da se višina najvišjega stolpca (ki ga dobimo za $k = \lfloor n/2 \rfloor$) z n počasi niža:

$$0.500, 0.500, 0.375, 0.375, 0.313, 0.313, 0.273, 0.273, 0.246, 0.246, \dots$$

Prepričaj se, da je ponavljanje zgornjih vrednosti posledica lastnosti binomskega simbola. Zaradi tega opazujmo le še vrednosti, ki smo jih dobili za lihe n . Kako hitro pada višina najvišjega stolpca? Vsekakor ne linearno, saj se razlike med zaporednimi členi manjšajo. Potrebno bo torej najti funkcijo, ki pada počasneje. Poskusimo s funkcijo \sqrt{n} , ki pa jo zaradi njenega naraščanja obrnemo v padajočo funkcijo $1/\sqrt{n}$. Izkaže se, da smo imeli srečo, vsaj sodeč po zaporedju, ki smo ga dobimo iz prejšnjega z množenjem s \sqrt{n} :

$$0.707, 0.750, 0.765, 0.773, 0.778, 0.781, 0.784, 0.786, 0.787, 0.788, 0.789, 0.790, 0.790, 0.791, 0.791, \dots \quad (4.2)$$

Abraham de Moivre bi nam seveda pritrdil, da smo na pravi poti (saj je nekaj podobnega počel že leta 1733 – **The Doctrine of Chance**).



Limitno vrednost števila iz (4.2) poznamo samo na kakšno decimalno natančno, (če limita sploh obstaja), vendar pa bi ga radi natančno določili (če limita sploh obstaja seveda). Gre v resnici morda za $0.8 = 4/5$? Predno se dodobra zamislimo, nam neutruden računalnik zaupa, da se po nekaj tisoč korakov pride do 0.7979 ali, če smo še bolj potrpežljivi do 0.7978646... Pa nadomestimo število z njegovim kvadratom oziroma s kvadratom recipročne vrednosti, da bo število večje od ena:

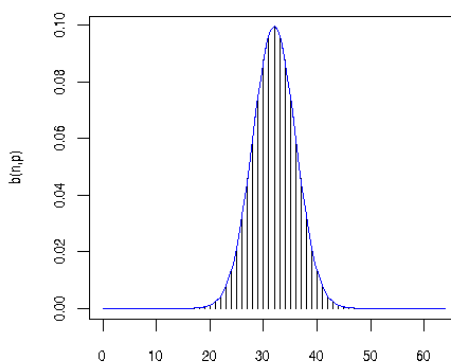
$$1.57158.$$

To število je sumljivo blizu $\pi/2$. Tako blizu, da smo pripravljeni tvegati in postaviti domnevo:

$$\lim_{n \rightarrow \infty} \binom{n}{k} \frac{\sqrt{n\pi/2}}{2^n} = 1, \quad \text{za } k = \lfloor n/2 \rfloor.$$

Sedaj pa bi nam pritrdil celo Stirling, kajti s pomočjo njegove (študentom dobro znane) ocene za faktoriel $n! \approx \sqrt{2\pi n} (n/e)^n$, kjer je e Eulerjeva konstanta ($e = 2.71\dots$), bi prišli do enakega zaključka (ki nam pravzaprav ponuja tudi hitro matematično utemeljitev zgoraj domneve).

Toliko smo se ukvarjali z višino v histogramih (ki se sedaj v limiti približuje neki konstanti), da smo skoraj pozabili na ploščino, da o širini niti ne govorimo. V resnici lahko vse stolpce postavimo enega nad drugega in v tem primeru dobimo en sam stolpec, tj. pravokotnik. Želeli smo, da bi bila njegova ploščina enaka 1, vendar sedaj delimo vse višine stolpcev v histogramu namesto z 2^n le še s številom $2^n/(c\sqrt{n})$, kjer je c poljubna pozitivna konstanta, tj. po deljenju z 2^n še *množimo* s $c\sqrt{n}$. Zato je potrebno širino stolpcev *deliti* s $c\sqrt{n}$. Ker višine stolpcev v zadnjem koraku nismo spreminjali, smo končno izpolnili obe zahtevi (konstantna ploščina in konstantna maksimalna višina histograma).



Slika 5: Histogram za $n = 64$. Stolpce smo nadomestili kar z intervali.

Čeprav še nismo razmišljali o širini histograma, (le-ta vsebuje $n + 1$ stolpcev in je sedaj enaka $(n + 1)/(c\sqrt{n})$, torej gre z rastočim n proti neskončno), zgornja slika kaže, da z njo ne bo težav, višine stolpcev so namreč na robu že zelo majhne (blizu 0) in ne bodo kaj dosti prispevale k obliki. Če želimo še bolje spoznati iskano obliko, je na vrsti študij *konveksnosti*, tj. kdaj moramo zavijati z avtom v desno ozirna v levo, če se vozimo po

krivulji z leve proti desni. Pri tem si pomagamo z opazovanjem spreminjanja predznaka razlike dveh zaporednih binomskih simbolov. Naj bo

$$d_h := \binom{n}{h+1} - \binom{n}{h} \quad \text{za } h = 0, 1, \dots, n-1.$$

Za $h = 0$ (oziroma $h = n-1$) in $h = 1$ (oziroma $h = n-2$) velja

$$d_0 = n = -d_{n-1} \quad \text{in} \quad d_1 = n(n-3)/2 = -d_{n-2}.$$

Zaradi simetrije binomskih simbolov, glede na $n/2$, tj. $\binom{n}{k} = \binom{n}{n-k}$, velja tudi

$$d_i = -d_{n-1-i}, \quad 0 \leq i \leq n-1,$$

zato se bomo omejili le na območje za katerega velja $h \leq \lfloor (n-1)/2 \rfloor$. Naj bo $m \in \mathbb{N}$ in $n = 2m$ oziroma $n = 2m-1$ glede na to ali je n sodo oziroma liho število. Potem je $\lfloor (n+1)/2 \rfloor = m$ in $d_m = 0$, če je n lih in $d_{m-1} = -d_m$, če je n sod. S slike oziroma iz pravkar navedenih lastnosti je očitno, da je potrebno na začetku zaviti v desno, potem pa bomo gotovo morali začeti zavijati še v levo, vsaj tik pred vrhom. Naslednja trditev nam zagotovi, da se na poti do vrha spremeni smer (zavijanja) le enkrat (iz leve v desno), zaradi simetrije pa potem enako velja tudi pri poti 'navzdol' (tokrat iz desne v levo).

Trditev 4.2. *Za $h \in \{1, \dots, \lfloor (n-1)/2 \rfloor\}$ velja*

$$d_{h-1} \leq d_h \quad \text{natanko tedaj, ko velja } h \leq \frac{n}{2} - \frac{\sqrt{n+2}}{2}$$

oziroma ekvivalentno

$$d_0 < d_1 < \dots < d_{k-1} \leq d_k > d_{k+1} > d_{k+2} > \dots > d_{\lfloor (n-1)/2 \rfloor},$$

kjer je $k := \lfloor (n - \sqrt{n+2})/2 \rfloor$. Enačaja velja natanko tedaj, ko je $n+2$ popoln kvadrat in je $h = k$.

Dokaz. Iz

$$d_h = \frac{n!}{(h+1)!(n-h-1)!} - \frac{n!}{h!(n-h)!} = \frac{n(n-2h-1)}{(h+1)!(n-h)!},$$

lahko zaključimo, da je predznak razlike

$$d_h - d_{h-1} = \frac{n!(n-2h-1)}{(h+1)!(n-h)!} - \frac{n!(n-2h+1)}{h!(n-h+1)!} = \frac{n!}{(h+1)!(n-h+1)!} \Delta$$

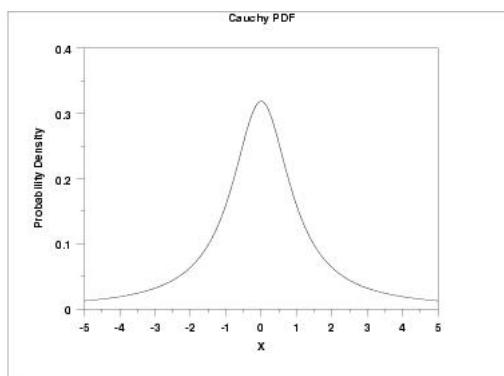
odvisna samo od predznaka razlike

$$\Delta := (n - 2h - 1)(n - h + 1) - (n - 2h + 1)(h + 1) = (n - 2h)^2 - (n + 2). \quad \square$$

Zopet smo pred novim presenečenjem.

Vse popravljene krivulje so si sedaj na moč podobne!

Za katero krivuljo pa gre? Če smo prej našli konstanto c , potem se moramo sedaj nekako podvzati in odkriti še skrivnostno krivuljo. Morda bi kdo najprej pomislil na **parabolo**, a ne bo prava, saj naša krivulja nikoli ne preseka x osi, temveč se ji samo asimptotično približuje. Iz zgornje trditve tudi sledi, da se oblika histogramov vsekakor bistveno razlikuje od oblike parabole, saj se na slednji sploh ne spremeni smer zavijanja. Kako pa je z obliko funkcije **cos x** (na intervalu $[-\pi/2, \pi/2]$)? Na prvi pogled izgleda prava: tudi na njej moramo zaviti najprej v levo, nato v desno. V tem primeru pa je ključnega pomena to, da preidemo iz desnega zavijanja v levo zavijanje pri kosinusu ravno na sredi (od vznožja do vrha), v primeru binomskih simbolov pa se razdalja z večanjem n približuje $n/2$. Tretji poskus s funkcijo $1/(1 + x^2)$ prepustimo bralcu (tudi ta ni prava, a je matematična utemeljitev nekoliko bolj zapletena).



Slika 6: Graf funkcije $1/(\pi(1 + x^2))$, kjer smo pod ulomkovo črto dodali še π , da bo ploščina pod krivuljo enaka 1.

V resnici naša funkcija ni niti racionalna (kvocient dveh polinomov). Morda pa je usoda hotela, da na ta način spoznamo povsem novo funkcijo. Zaenkrat jo poimenujmo kar **usojena** funkcija (da ne bomo pretiravali s tujkami - **fatalka**).

Sledimo Abrahamu de Moivreu (The Doctrine of Chance, 1733), ki predlaga, da vpeljemo

$$x = \frac{k - (n/2)}{\sqrt{n}} \quad \text{oziroma} \quad k = x\sqrt{n} + (n/2)$$

(le ta zamenjava postavi vrh krivulje na y -os in normalizira razdaljo do prevoja) in izračunajmo

$$f(x) = \lim_{n \rightarrow \infty} \frac{\sqrt{n}}{2^n} \binom{n}{x\sqrt{n} + (n/2)} = \lim_{n \rightarrow \infty} \frac{\sqrt{n}}{2^n} \frac{n!}{((n/2) + x\sqrt{n})! ((n/2) - x\sqrt{n})!}.$$

Vpeljimo zamenjavo $n = 4m^2$ (da se znebimo korenov in ulomkov) in si поблиže oglejmo naslednji kvocient

$$\frac{f(x)}{f(0)} = \lim_{n \rightarrow \infty} \frac{(n/2)!(n/2)!}{((n/2) + x\sqrt{n})! ((n/2) - x\sqrt{n})!} = \lim_{m \rightarrow \infty} \frac{(2m^2)! (2m^2)!}{(2m^2 + 2mx)! (2m^2 - 2mx)!}.$$

Okrajšamo, kar se okrajšati da, in dobimo

$$\frac{f(x)}{f(0)} = \lim_{m \rightarrow \infty} \frac{(2m^2)(2m^2 - 1) \cdots (2m^2 - 2mx + 2)(2m^2 - 2mx + 1)}{(2m^2 + 1)(2m^2 + 2) \cdots (2m^2 + 2mx - 1)(2m^2 + 2mx)}.$$

Opazimo, da se istoležeči faktorji, ki ležijo en nad drugim seštejejo v $4m^2 + 1$ in preoblikujemo dobljeni kvocient v naslednjo obliko:

$$\frac{f(x)}{f(0)} = \lim_{m \rightarrow \infty} \frac{\left(2m^2 + \frac{1}{2} - \frac{1}{2}\right) \left(2m^2 + \frac{1}{2} - \frac{3}{2}\right) \left(2m^2 + \frac{1}{2} - \frac{5}{2}\right) \cdots \left(2m^2 + \frac{1}{2} - \frac{4mx - 1}{2}\right)}{\left(2m^2 + \frac{1}{2} + \frac{1}{2}\right) \left(2m^2 + \frac{1}{2} + \frac{3}{2}\right) \left(2m^2 + \frac{1}{2} + \frac{5}{2}\right) \cdots \left(2m^2 + \frac{1}{2} + \frac{4mx - 1}{2}\right)}.$$

Sedaj pa delimo vsak faktor (nad in pod ulomkovo črto) z $2m^2 + 1/2$ in upoštevajmo, da lahko $1/2$ v limiti zanemarimo) in dobimo

$$\frac{f(x)}{f(0)} = \lim_{m \rightarrow \infty} \frac{\left(1 - \frac{1}{4m^2}\right) \left(1 - \frac{3}{4m^2}\right) \left(1 - \frac{5}{4m^2}\right) \cdots \left(1 - \frac{4mx - 1}{4m^2}\right)}{\left(1 + \frac{1}{4m^2}\right) \left(1 + \frac{3}{4m^2}\right) \left(1 + \frac{5}{4m^2}\right) \cdots \left(1 + \frac{4mx - 1}{4m^2}\right)}$$

oziroma

$$\begin{aligned} \ln\left(\frac{f(x)}{f(0)}\right) &= \lim_{m \rightarrow \infty} \ln\left(1 - \frac{1}{4m^2}\right) + \ln\left(1 - \frac{3}{4m^2}\right) + \ln\left(1 - \frac{5}{4m^2}\right) + \cdots + \ln\left(1 - \frac{4mx - 1}{4m^2}\right) \\ &\quad - \ln\left(1 + \frac{1}{4m^2}\right) - \ln\left(1 + \frac{3}{4m^2}\right) - \ln\left(1 + \frac{5}{4m^2}\right) - \cdots - \ln\left(1 + \frac{4mx - 1}{4m^2}\right). \end{aligned}$$

Uporabimo še vrsto $\ln(1 + x) = x - x^2/2 + \cdots$ in opazimo, da gredo z izjemo prvega, vsi členi v tej vrsti s $m \rightarrow \infty$ proti nič:

$$\ln\left(\frac{f(x)}{f(0)}\right) = \lim_{m \rightarrow \infty} -2 \frac{1 + 3 + 5 + \cdots + (4mx - 1)}{4m^2} = \lim_{m \rightarrow \infty} -\frac{(2mx)^2}{2m^2} = -2x^2.$$

Pri zadnjem enačaju smo uporabili še dejstvo, da je vsota prvih N lihih števil enaka kvadratu števila $(N + 1)/2$. Od tod zaključimo

$$f(x) = f(0) e^{-2x^2} = \sqrt{\frac{2}{\pi}} e^{-2x^2} \quad \text{oziroma} \quad N(t) = \frac{e^{\frac{1}{2}x^2}}{\sqrt{2\pi}}, \quad \text{za} \quad N(t) = f(t/2)/2.$$

Izpeljali smo

Izrek 4.3. (De Moivreov točkovni obrazec)

$$P_n(k) \approx \frac{1}{\sqrt{\pi n/2}} e^{-\frac{(k-n/2)^2}{n/2}}$$

Je poseben primer **Laplaceovega točkovnega obrazca**. Slednjega smemo uporabljati, ko je p blizu $1/2$:

$$P_n(k) \approx \frac{1}{\sqrt{2\pi npq}} e^{-\frac{(k-np)^2}{2npq}}.$$



Računanje verjetnosti $P_n(k)$ v R-ju

Program R: Vrednost $P_n(k)$ dobimo z ukazom `dbinom(k, size=n, prob=p)`

```
> dbinom(50, size=1000, prob=0.05)
[1] 0.05778798
```

Izrek 4.4. (Poissonov obrazec) Za majhne verjetnosti,

tj. p blizu 0 velja

$$P_n(k) \approx \frac{(np)^k e^{-np}}{k!}. \quad (4.3)$$



Dokaz. Iz analize se spomnimo naslednje zveze

$$\lim_{n \rightarrow \infty} \left(1 - \frac{\lambda}{n}\right)^n = e^{-\lambda}.$$

Če lahko definiramo binomsko porazdelitev tako, da je $p = \lambda/n$, potem lahko izračunamo limito verjetnosti P za velike n na naslednji način:

$$\begin{aligned} \lim_{n \rightarrow \infty} P_n(k) &= \lim_{n \rightarrow \infty} \binom{n}{k} p^k (1-p)^{n-k} \\ &= \lim_{n \rightarrow \infty} \frac{n!}{(n-k)!k!} \left(\frac{\lambda}{n}\right)^k \left(1 - \frac{\lambda}{n}\right)^{n-k} \\ &= \lim_{n \rightarrow \infty} \underbrace{\left[\frac{n!}{n^k (n-k)!}\right]}_F \underbrace{\left(\frac{\lambda^k}{k!}\right)}_{\rightarrow e^{-\lambda}} \underbrace{\left(1 - \frac{\lambda}{n}\right)^n}_{\rightarrow 1} \underbrace{\left(1 - \frac{\lambda}{n}\right)^{-k}}_{\rightarrow 1}. \end{aligned}$$

Za oceno člena F si najprej oglejmo njegov logaritem:

$$\lim_{n \rightarrow \infty} \ln(F) = \ln(n!) - k \ln(n) - \ln[(n-k)!].$$

Z uporabo Stirlingove aproksimacije je

$$\lim_{n \rightarrow \infty} \ln(n!) \rightarrow n \ln(n) - n,$$

izraz za $\ln(F)$ pa lahko naprej poenostavimo v

$$\begin{aligned} \lim_{n \rightarrow \infty} \ln(F) &= [n \ln(n) - n] - [k \ln(n)] - [(n-k) \ln(n-k) - (n-k)] \\ &= (n-k) \ln\left(\frac{n}{n-k}\right) - k = \underbrace{-\left(1 - \frac{k}{n}\right)}_{\rightarrow -1} \underbrace{\ln\left(1 - \frac{k}{n}\right)^n}_{\rightarrow -k} - k = k - k = 0. \end{aligned}$$

Torej je $\lim_{n \rightarrow \infty} F = e^0 = 1$. Od tod pa sledi, da je porazdelitev v limiti enaka

$$\frac{\lambda^k e^{-\lambda}}{k!},$$

ki ima sedaj obliko Poissonove porazdelitve. □

Jacob Bernoulli (1654–1705)

Jacob Bernoulli (also referred to as Jacques or James) was born in Basel Switzerland on December 27, 1654, the first of the famous Bernoulli family of Swiss mathematicians. He first studied theology, but then against his father's wishes, turned to mathematics, astronomy, and physics. Bernoulli was appointed professor of mathematics at the University of Basel in 1687.

Bernoulli's mathematical work mostly involved the new theory of calculus. He worked on infinite series, conics, cycloids, transcendental curves, isoperimetry, isochronous curves, catenaries, and the logarithmic spiral, publishing his results widely in *Acta Eruditorum*. Bernoulli also wrote what is considered the second book devoted to probability, *Ars Conjectandi*, which was published after his death in 1713. Bernoulli formulated the version of the law of large numbers for independent trials, now called Bernoulli trials in his honor, and studied the binomial distribution.

Jacques Bernoulli died on August 16, 1705 in Basel. The logarithmic spiral is carved into his tombstone.

<http://www.math.uah.edu/STAT/biographies/Bernoulli.xhtml>

Pierre-Simon Laplace (1749–1827)

(b. Beaumont-en-Auge, France; d. Paris, France) French mathematician. Laplace was born into the French middle class but died a Marquis having prospered as well during the French Revolution as during the monarchy. His name is well known to mathematicians both for his work on transformations and for his work on planetary motion. He had a deterministic view of the universe. Reputedly, his reply, when asked by Napoleon where God fitted in, was 'I have no need of that hypothesis'. Napoleon appointed him as Minister of the Interior — but removed him six weeks later for trying 'to carry the spirit of the infinitesimal into administration'. In Statistics he worked on many probability problems including the Laplace distribution: he is credited with independently discovering Bayes's Theorem. He was elected FRS in 1789 and FRSE in 1813. A street in Paris is named after him, as is the Promontorium Laplace on the Moon.

<http://www.answers.com/topic/pierre-simon-laplace>

Abraham de Moivre (1667–1754)

A French Huguenot, de Moivre was jailed as a Protestant upon the revocation of the Edict of Nantes in 1685. When he was released shortly thereafter, he fled to England. In London he became a close friend of Sir Isaac Newton and the astronomer Edmond Halley. De Moivre was elected to the Royal Society of London in 1697 and later to the Berlin and Paris academies. Despite his distinction as a mathematician, he never succeeded in securing a permanent position but eked out a precarious living by working as a tutor and a consultant on gambling and insurance.

De Moivre expanded his paper “De mensura sortis” (written in 1711), which appeared in *Philosophical Transactions*, into *The Doctrine of Chances* (1718). Although the modern theory of probability had begun with the unpublished correspondence (1654) between Blaise Pascal and Pierre de Fermat and the treatise *De Ratiociniis in Ludo Aleae* (1657; “On Ratiocination in Dice Games”) by Christiaan Huygens of Holland, de Moivre’s book greatly advanced probability study. The definition of statistical independence—namely, that the probability of a compound event composed of the intersection of statistically independent events is the product of the probabilities of its components—was first stated in de Moivre’s *Doctrine*. Many problems in dice and other games were included, some of which appeared in the Swiss mathematician Jakob (Jacques) Bernoulli’s *Ars conjectandi* (1713; “The Conjectural Arts”), which was published before de Moivre’s *Doctrine* but after his “De mensura.” He derived the principles of probability from the mathematical expectation of events, just the reverse of present-day practice.

De Moivre’s second important work on probability was *Miscellanea Analytica* (1730; “Analytical Miscellany”). He was the first to use the probability integral in which the integrand is the exponential of a negative quadratic. He originated Stirling’s formula. In 1733 he used Stirling’s formula to derive the normal frequency curve as an approximation of the binomial law.

De Moivre was one of the first mathematicians to use complex numbers in trigonometry. The formula known by his name,

$$(\cos x + i \sin x)^n = \cos nx + i \sin nx,$$

was instrumental in bringing trigonometry out of the realm of geometry and into that of analysis.

<http://www.britannica.com/EBchecked/topic/387796/Abraham-de-Moivre>

Poglavje 5

Slučajne spremenljivke in porazdelitve



Denimo, da imamo poskus, katerega izidi so števila (npr. pri metu kocke so izidi števila pik). Se pravi, da je poskusom prirejena neka količina, ki more imeti različne vrednosti. Torej je spremenljivka. Katero od mogočih vrednosti zavzame v določeni ponovitvi poskusa, je odvisno od slučaja. Zato ji rečemo **slučajna spremenljivka**. Da je slučajna spremenljivka znana, je potrebno vedeti

1. kakšne vrednosti more imeti (*zaloga vrednosti*) in
2. kolikšna je verjetnost vsake izmed možnih vrednosti ali intervala vrednosti.

Predpis, ki določa te verjetnosti, imenujemo **porazdelitveni zakon**.

Slučajne spremenljivke označujemo z velikimi tiskanimi črkami iz konca abecede, vrednosti spremenljivke pa z enakimi malimi črkami. Tako je npr. $(X = x_i)$ dogodek, da slučajna spremenljivka X zavzame vrednost x_i . Porazdelitveni zakon slučajne spremenljivke X je poznan, če je mogoče za vsako realno število x določiti verjetnost

$$F(x) = P(X < x).$$

$F(x)$ imenujemo **porazdelitvena funkcija**. Najpogosteje uporabljamo naslednji vrsti slučajnih spremenljivk:

1. **diskretna** slučajna spremenljivka, pri kateri je zaloga vrednosti neka števna (diskretna) množica
2. **zvezna** slučajna spremenljivka, ki lahko zavzame vsako realno število znotraj določenega intervala.

Lastnosti porazdelitvene funkcije

1. Funkcija F je definirana na vsem \mathbb{R} in velja $0 \leq F(x) \leq 1, x \in \mathbb{R}$.
2. Funkcija F je nepadajoča $x_1 < x_2 \implies F(x_1) \leq F(x_2)$.
3. $F(-\infty) = 0$ in $F(\infty) = 1$.
4. Funkcija je v vsaki točki zvezna od leve $F(x-) = F(x)$.
5. Funkcija ima lahko v nekaterih točkah skok.
Vseh skokov je največ števno mnogo.
6. $P(x_1 \leq X < x_2) = F(x_2) - F(x_1)$.
7. $P(x_1 < X < x_2) = F(x_2) - F(x_1+)$.
8. $P(X \geq x) = 1 - F(x)$.
9. $P(X = x) = F(x+) - F(x)$.

5.1 Diskretne slučajne spremenljivke

Zaloga vrednosti diskretne slučajne spremenljivke X je števna množica

$$\{x_1, x_2, \dots, x_m, \dots\}.$$

Torej je lahko tudi števno neskončna, kot npr. množici naravnih ali celih števil: \mathbb{N}, \mathbb{Z} .
Dogodki

$$X = x_k \quad k = 1, 2, \dots$$

sestavljajo popoln sistem dogodkov. Označimo verjetnost posameznega dogodka s

$$P(X = x_i) = p_i.$$

Vsota verjetnosti vseh dogodkov je enaka 1:

$$p_1 + p_2 + \dots + p_m + \dots = 1.$$

Verjetnostna tabela prikazuje diskretno slučajno spremenljivko s tabelo tako, da so v prvi vrstici zapisane vse vrednosti x_i , pod njimi pa so pripisane pripadajoče verjetnosti:

$$X : \begin{pmatrix} x_1 & x_2 & \cdots & x_m & \cdots \\ p_1 & p_2 & \cdots & p_m & \cdots \end{pmatrix}.$$

Porazdelitvena funkcija je v diskretnem primeru

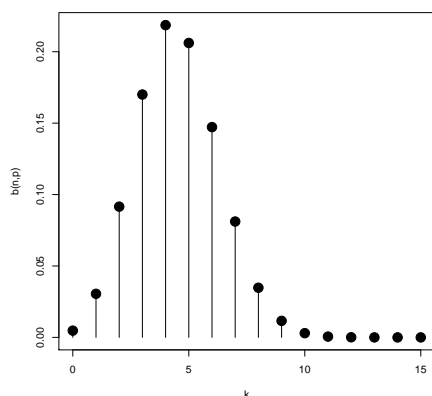
$$F(x_k) = P(X < x_k) = \sum_{i=1}^{k-1} p_i.$$

5.1.1 Enakomerna diskretna porazdelitev

Končna diskretna slučajna spremenljivka se porazdeljuje **enakomerno**, če so vse njene vrednosti enako verjetne. Primer take slučajne spremenljivke je število pik pri metu kocke

$$X : \begin{pmatrix} 1 & 2 & 3 & 4 & 5 & 6 \\ 1/6 & 1/6 & 1/6 & 1/6 & 1/6 & 1/6 \end{pmatrix}.$$

5.1.2 Binomska porazdelitev



```
> h <- dbinom(0:15,size=15,prob=0.3)
> plot(0:15,h,type="h",xlab="k",ylab="b(n,p)")
> points(0:15,h,pch=16,cex=2)
```

Binomska porazdelitev ima zalogo vrednosti $\{0, 1, 2, \dots, n\}$ in verjetnosti, ki jih računamo po Bernoullijevem obrazcu:

$$P(X = k) = \binom{n}{k} p^k (1-p)^{n-k},$$

$k = 0, 1, 2, \dots, n$. Binomska porazdelitev je natančno določena z dvema podatkom – parametroma: n in p . Če se slučajna spremenljivka X porazdeljuje binomsko s parametroma n in p , zapišemo:

$$X : B(n, p).$$

Primer: Naj bo slučajna spremenljivka X določena s številom fantkov v družini s 4 otroki. Denimo, da je enako verjetno, da se v družini rodi fantek ali deklica:

$$P(F) = p = \frac{1}{2}, \quad P(D) = q = \frac{1}{2}.$$

Spremenljivka X se tedaj porazdeljuje binomsko $B(4, \frac{1}{2})$ in njena verjetnostna shema je:

$$X : \begin{pmatrix} 0 & 1 & 2 & 3 & 4 \\ 1/16 & 4/16 & 6/16 & 4/16 & 1/16 \end{pmatrix}.$$

Na primer

$$P(X = 2) = P_4(2) = \binom{4}{2} \left(\frac{1}{2}\right)^2 \left(1 - \frac{1}{2}\right)^{4-2} = \frac{6}{16}$$

Porazdelitev obravnavane slučajne spremenljivke je simetrična. Pokazati se da, da je binomska porazdelitev simetrična, le če je $p = 0,5$. \diamond

5.1.3 Poissonova porazdelitev $P(\lambda)$

Poissonova porazdelitev izraža verjetnost števila dogodkov, ki se zgodijo v danem časovnem intervalu, če vemo, da se ti dogodki pojavijo s poznano povprečno frekvenco in neodvisno od časa, ko se je zgodil zadnji dogodek. Poissonovo porazdelitev lahko uporabimo tudi za število dogodkov v drugih intervalih, npr. razdalja, prostornina,... Ima zalogo vrednosti $\{0, 1, 2, \dots\}$, njena verjetnostna funkcija pa je

$$p_k = P(\text{\#dogodkov} = k) = \lambda^k \frac{e^{-\lambda}}{k!},$$

kjer je $\lambda > 0$ dani parameter – in predstavlja pričakovano pogostost nekega dogodka. Konstanta e je osnova naravnega logaritma, tj. $e = 2.71828 \dots$

$$p_{k+1} = \frac{\lambda}{k+1} p_k, \quad p_0 = e^{-\lambda}.$$



Vidimo, da zaloga vrednosti te slučajne spremenljivke ni omejena, saj je verjetnost, da se v nekem časovnem obdobju zgodi mnogo uspehov različna od nič. To je bistvena razlika v primerjavi z binomsko porazdelitvijo, kjer število uspehov seveda ne more presežati števila Bernoullijevih poskusov n .

Primer: Posebno pomembna je ta porazdelitev v teoriji množične strežbe. Če se dogodek pojavi v povprečju 3-krat na minuto in nas zanima kolikokrat se bo zgodil v četrt ure, potem uporabimo za model Poissonovo porazdelitev z $\lambda = 15 \cdot 3 = 45$. \diamond

Naštejmo še nekaj primerov, ki jih dobro opišemo (modeliramo) s Poissonovo porazdelitvijo:

- število dostopov do omrežnega strežnika na minuto (pod predpostavko homogenosti),
- število telefonskih klicev na bazni postaji na minuto,
- število mutacij v danem intervalu RNK po določeni količini sprejete radiacije,
- število vojakov, ki so umrli vsako leto za posledicami konjske brce v vsaki diviziji Pruske konjenice, (iz knjige Ladislausa Josephovicha Bortkiewicza, 1868–1931).

5.1.4 Pascalova porazdelitev $P(m, p)$

Pascalova porazdelitev ima zalogo vrednosti $\{m, m+1, m+2, \dots\}$, verjetnostna funkcija pa je

$$p_k = \binom{k-1}{m-1} p^m q^{k-m},$$

kjer je $0 < p < 1$ dani parameter za verjetnost dogodka A v posameznem poskusu.

Opisuje porazdelitev števila poskusov potrebnih, da se dogodek A zgodi m -krat.

Za $m = 1$, porazdelitvi $G(p) = P(1, p)$ pravimo **geometrijska porazdelitev**. Opisuje porazdelitev števila poskusov potrebnih, da se dogodek A zgodi prvič.



Primer: Če mečemo kovanec toliko časa, da pade grb in z X označimo število potrebnih metov, vključno z zadnjim, potem je slučajna spremenljivka X geometrijsko porazdeljena.

◇

Če spremenljivka X označuje število metov, vključno z zadnjim, do m -tega grba, potem dobimo **negativno binomsko** slučajno spremenljivko X : $\text{NegBin}(m, p)$ in

$$P(X = k) = \binom{k-1}{m-1} p^m (1-p)^{k-m} \quad \text{za } k \geq m.$$

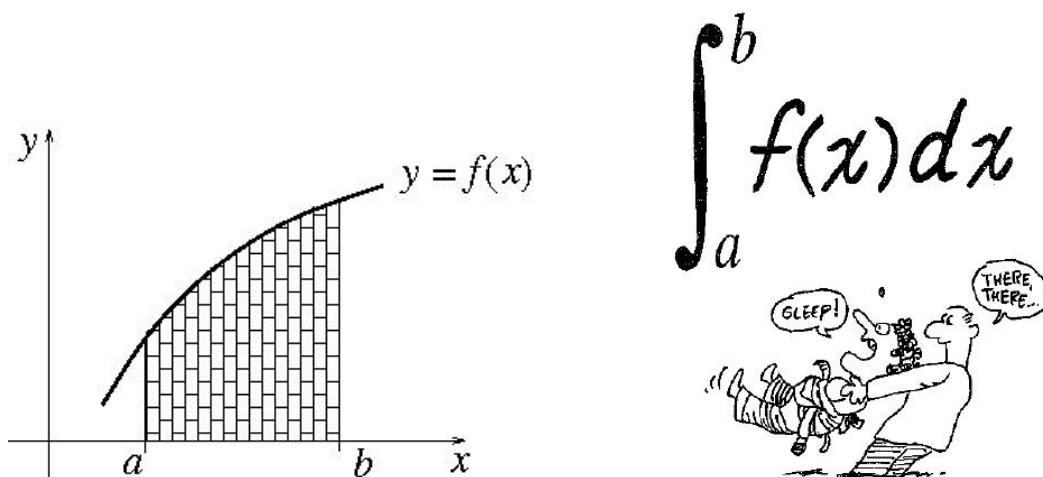
5.1.5 Hipergeometrijska porazdelitev $H(n; M, N)$

Hipergeometrijska porazdelitev ima zalogo vrednosti $\{0, 1, 2, \dots\}$, verjetnostna funkcija pa je

$$p_k = \frac{\binom{M}{k} \binom{N-M}{n-k}}{\binom{N}{n}},$$

kjer so $k \leq n \leq \min(M, N - M)$ dani parametri. Opisuje verjetnost dogodka, da je med n izbranimi kroglicami natanko k belih, če je v posodi M belih in $N - M$ črnih kroglic in izbiramo n -krat brez vračanja.

5.2 Ponovitev: integrali



Določeni integral predstavlja ploščino pod krivuljo. Naj bo funkcija $y = f(x)$ zvezna na $[a, b]$ in nenegativna. Ploščina lika med krivuljo $f(x) \geq 0$, in abscisno osjo na intervalu $[a, b]$ je enaka določenemu integralu

$$\int_a^b f(x) dx.$$

Lastnosti določenega integrala:

Trditev 5.1. (1) Za $a, b \in \mathbb{R}$ velja

$$\int_a^b f(x) dx = - \int_b^a f(x) dx.$$

(2) Če je $f(x) \leq 0 \quad \forall x \in [a, b]$, je vrednost integrala negativna.

(3) Za $c \in [a, b]$ velja

$$\int_a^b f(x) dx = \int_a^c f(x) dx + \int_c^b f(x) dx.$$

(4) Naj bo $f(x) \geq g(x)$, $x \in [a, b]$, potem velja

$$\int_a^b f(x) dx \geq \int_a^b g(x) dx. \quad \square$$

Saj vas razumem!

Potem pa uporabimo še ∞ za mejo pri integriranju.

Brez preplaha!



Iščemo le celotno ploščino pod krivuljo, od enega konca do drugega, le da konca pravzaprav sploh ni.



5.3 Zvezne slučajne spremenljivke

Slučajna spremenljivka X je **zvezno porazdeljena**, če obstaja taka integrabilna funkcija p , imenovana **gostota verjetnosti**, da za vsak $x \in \mathbb{R}$ velja:

$$F(x) = P(X < x) = \int_{-\infty}^x p(t) dt,$$

kjer $p(x) \geq 0$. To verjetnost si lahko predstavimo tudi grafično v koordinatnem sistemu, kjer na abscisno os nanašamo vrednosti slučajne spremenljivke, na ordinatno pa gostoto verjetnosti $p(x)$. Verjetnost je tedaj predstavljena kot ploščina pod krivuljo, ki jo določa $p(x)$. *Velja*

$$\int_{-\infty}^{\infty} p(x) dx = 1 \quad \text{in} \quad P(x_1 \leq X < x_2) = \int_{x_1}^{x_2} p(t) dt$$

ter $p(x) = F'(x)$.

5.3.1 Enakomerna porazdelitev zvezne slučajne spremenljivke

Verjetnostna gostota **enakomerno porazdeljene zvezne** slučajne spremenljivke je:

$$p(x) = \begin{cases} \frac{1}{b-a} & \text{za } a \leq X \leq b \\ 0 & \text{drugod.} \end{cases}$$

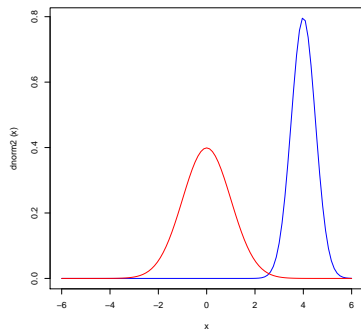
Grafično si jo predstavljamo kot pravokotnik nad intervalom (a, b) višine $\frac{1}{b-a}$.

5.3.2 Normalna ali Gaussova porazdelitev



Leta 1738 je Abraham De Moivre (1667-1754) objavil aproksimacijo binomske porazdelitve, ki je normalna krivulja.

Leta 1809 je Karl Frederic Gauss (1777-1855) raziskoval matematično ozadje planetarnih orbit, ko je prišel do normalne porazdelitvene funkcije.



```
> d2 <-
function(x){dnorm(x,mean=4,sd=0.5)}
> curve(d2,-6,6,col="blue")
> curve(dnorm,-6,6,col="red",add=TRUE)
```

Zaloga vrednosti **normalno porazdeljene** slučajne spremenljivke so vsa realna števila, gostota verjetnosti pa je:

$$p(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2}.$$

Normalna porazdelitev je natanko določena z parametroma: μ in σ . Če se slučajna spremenljivka X porazdeljuje normalno s parametroma μ in σ zapišemo:

$$X : N(\mu, \sigma).$$

Laplaceov intervalski obrazec

Zanima nas, kolikšna je verjetnost $P_n(k_1, k_2)$, da se v Bernoullijevem zaporedju neodvisnih poskusov v n zaporednih poskusih zgodi dogodek A vsaj k_1 -krat in manj kot k_2 -krat. Označimo

$$x_k = \frac{k - np}{\sqrt{npq}} \quad \text{in} \quad \Delta x_k = x_{k+1} - x_k = \frac{1}{\sqrt{npq}}.$$

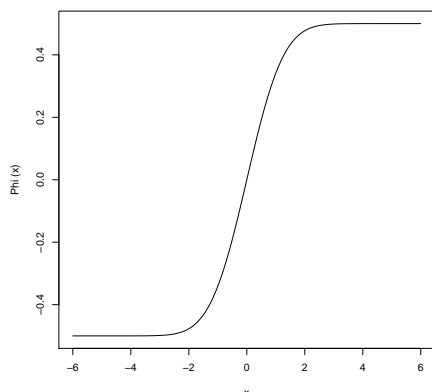
Tedaj je, če upoštevamo Laplaceov točkovni obrazec,

$$P_n(k_1, k_2) = \sum_{k=k_1}^{k_2-1} P_n(k) = \frac{1}{\sqrt{2\pi}} \sum_{k=k_1}^{k_2-1} e^{-\frac{1}{2}x_k^2} \Delta x_k.$$

Za (zelo) velike n lahko vsoto zamenjamo z integralom

$$P_n(k_1, k_2) \approx \frac{1}{\sqrt{2\pi}} \int_{x_{k_1}}^{x_{k_2}} e^{-\frac{1}{2}x^2} dx.$$

Funkcija napake $\Phi(x)$



Funkcija napake imenujemo funkcijo

$$\Phi(x) = \frac{1}{\sqrt{2\pi}} \int_0^x e^{-\frac{1}{2}t^2} dt.$$

Funkcija napake je liha, zvezno odvedljiva, strogo naraščajoča funkcija, za katero velja $\Phi(0) = 0$, $P_n(k_1, k_2) \approx \Phi(x_{k_2}) - \Phi(x_{k_1})$, po izreku A.6 pa še $\Phi(\infty) = 1/2$ in $\Phi(-\infty) = -1/2$.

Vrednosti funkcije napake najdemo v tabelah ali pa je vgrajena v statističnih programih.

```
> Phi <- function(x){pnorm(x)-0.5}
> curve(Phi,-6.6)
```

```
> x2 <- (50 - 1000*0.05)/sqrt(1000*0.05*0.95)
> x1 <- (0 - 1000*0.05)/sqrt(1000*0.05*0.95)
> pnorm(x2)-pnorm(x1)
[1] 0.5
```

$$\begin{aligned} F(x) &= \frac{1}{\sigma\sqrt{2\pi}} \int_{-\infty}^x e^{-\frac{1}{2}\left(\frac{t-\mu}{\sigma}\right)^2} dt \\ &= \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\frac{x-\mu}{\sigma}} e^{-\frac{1}{2}s^2} ds \\ &= \frac{1}{2} + \Phi\left(\frac{x-\mu}{\sigma}\right). \end{aligned}$$



$$P(x_1 \leq X < x_2) = \Phi\left(\frac{x_2 - \mu}{\sigma}\right) - \Phi\left(\frac{x_1 - \mu}{\sigma}\right).$$

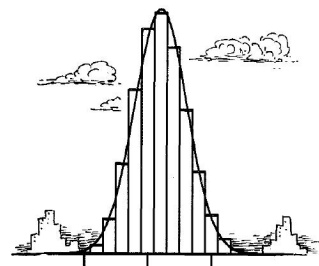
Porazdelitev $N(0, 1)$ je **standardizirana normalna porazdelitev**.

Spremenljivko $X : N(\mu, \sigma)$ pretvorimo z

$$z = \frac{x - \mu}{\sigma}$$

v standardizirano spremenljivko $Z : N(0, 1)$.

Iz Laplaceovega obrazca izhaja $B(n, p) \approx N(np, \sqrt{npq})$.



Izrek 5.2 (Bernoullijev zakon velikih števil (1713)). Naj bo k frekvenca dogodka A v n neodvisnih ponovitvah danega poskusa, v katerem ima dogodek A verjetnost p . Tedaj za vsak $\varepsilon > 0$ velja

$$\lim_{n \rightarrow \infty} P\left(\left|\frac{k}{n} - p\right| < \varepsilon\right) = 1.$$

Dokaz. ŠE PRIDE - ZAENKRAT GLEJ HLADNIK, STR. 21

□

Ta izrek opravičuje statistično definicijo verjetnosti.

Posledica 5.3.

$$P\left(\left|\frac{k}{n} - p\right| < \varepsilon\right) \approx 2\Phi\left(\varepsilon\sqrt{\frac{n}{pq}}\right).$$

Primer: (a) Kolikšna je verjetnost, da se pri metu kovanca relativna frekvenca grba v 3600 metih ne razlikuje od 0,5 za več kot 0,01, se pravi, da grb pade med 1764 in 1836-krat? V tem primeru je

$$p = 1/2, \quad n = 3600, \quad \varepsilon = 0,01,$$

tako, da iz zgornje formule dobimo

$$2\Phi\left(0,01 * \sqrt{\frac{3600}{0,25}}\right) = 2\Phi(1,2) = 2 \cdot 0,385 = 0,77,$$

kar je presenetljivo veliko, kaj ne?

(b) Kolikokrat moramo vreči pošten kovanec, da bo verjetnost dogodka, da se relativna frekvenca grba razlikuje od 0,5 za manj kot 0,05, večja od 0,997? Iz tabele za Φ vidimo, da je $2\Phi(x) > 0,997$ za $x = 3$, zato moramo poiskati tak n , da bo

$$3 < \varepsilon\sqrt{\frac{n}{pq}} = 0,05\sqrt{\frac{n}{0,25}} \quad \text{ozioroma} \quad n > \left(\frac{3}{0,05}\right)^2 \cdot 0,25 = 900. \quad \diamond$$

5.3.3 Porazdelitev Poissonovega toka, eksponentna

Gostota **eksponentne porazdelitve** je enaka

$$p(x) = \lambda e^{-\lambda x}, \quad x \geq 0,$$

porazdelitvena funkcija pa

$$F(x) = \int_0^x \lambda e^{-\lambda t} dt = 1 - e^{-\lambda x}.$$

5.3.4 Porazdelitev Gama

Prejšnjo porazdelitev lahko še precej posplošimo. Naj bosta $b, c > 0$. Tedaj ima **porazdelitev Gama** $\Gamma(b, c)$ gostoto:

$$p(x) = \frac{c^b}{\Gamma(b)} x^{b-1} e^{-cx}, \quad 0 < x$$

in $p(x) = 0$ za $x \leq 0$. Za $b = 1$ seveda dobimo eksponentno porazdelitev. **Funkcijo Gama** lahko definiramo z določenim integralom za $\Re[z] > 0$ (Eulerjeva integralna forma)

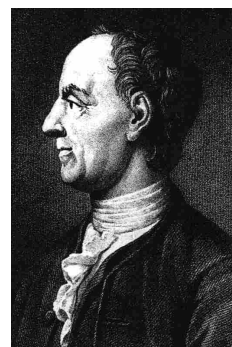
$$\Gamma(z) = \int_0^\infty t^{z-1} e^{-t} dt = 2 \int_0^\infty e^{-t^2} t^{2z-1} dt.$$

Torej je $\Gamma(1) = 1$ in

$$\Gamma(z) = \int_0^1 \left[\ln \frac{1}{t} \right]^{z-1} dt.$$

(Je povsod analitična z izjemo $z = 0, -1, -2, \dots$ in nima ničel.

Glej npr. http://en.wikipedia.org/wiki/Gamma_function in http://en.wikipedia.org/wiki/Gamma_distribution.)



Leonhard Euler (1707-1783) ¹

Integral poskusimo rešiti z integracijo po delih (po realnem argumentu). Za $v = t^x$ in $du = e^{-t} dt$ velja $dv = (x-1)t^{x-2}$ in $u = -e^{-t}$, od koder sledi

$$\begin{aligned} \Gamma(x) &= \int_0^\infty t^{x-1} e^{-t} dt = \left[-t^{x-1} e^{-t} \right]_0^\infty + \int_0^\infty (x-1)t^{x-2} e^{-t} dt \\ &= (x-1) \int_0^\infty t^{x-2} e^{-t} dt = (x-1)\Gamma(x-1). \end{aligned}$$

Za naravno število $x = n \in \{1, 2, \dots\}$ torej dobimo

$$\Gamma(n) = (n-1)\Gamma(n-1) = (n-1)(n-2)\Gamma(n-2) = (n-1)(n-2)\dots 1 = (n-1)!,$$

kar pomeni, da se v tem primeru Γ funkcija zreducira v 'faktorijel'.

¹ Eulerju se lahko zahvalimo za zapis $f(x)$ za funkcijo (1734), e za osnovo naravnega logaritma (1727), i za kvadratni koren števila -1 (1777), π for pi, Σ za vsoto (1755), in mnoge druge oznake/koncepte, ki smo jih danes sprejeli za same po sebi umevne. Glej <http://www.gap-system.org/history/Biographies/Euler.html>

5.3.5 Porazdelitev hi-kvadrat

Porazdelitev **hi-kvadrat** je poseben primer porazdelitve Gama:

$$\chi^2(n) = \Gamma\left(\frac{n}{2}, \frac{1}{2}\right)$$

($n \in \mathbb{N}$ je število prostostnih stopenj) in ima gostoto

$$p(x) = \frac{1}{2^{\frac{n}{2}} \Gamma(\frac{n}{2})} x^{\frac{n}{2}-1} e^{-\frac{x}{2}}, \quad \text{kjer je } x > 0.$$



Ernst Abbe
(1840-1905)

Leta 1863 jo je prvi izpeljal nemški fizik Ernst Abbe, ko je preučeval porazdelitev vsote kvadratov napak.

Leta 1878 je Ludwig Boltzmann izpeljal hi-kvadrat porazdelitev z dvema in tremi prostostnimi stopnjami, ko je študiral kinetično energijo molekul.

Karl Pearson (1875-1937) je demonstriral uporabnost hi-kvadrat porazdelitve statistikom.



5.3.6 Cauchyeva porazdelitev

je porazdelitev z gostoto

$$p(x) = \frac{a}{\pi} \frac{1}{1 + a^2(x - b)^2}, \quad \text{za } -\infty < x < \infty, \quad a > 0$$

ima porazdelitveno funkcijo

$$F(x) = \frac{a}{\pi} \int_{-\infty}^x \frac{1}{1 + a^2(x - b)^2} dx = \frac{1}{\pi} \operatorname{arctg}(a(x - b)) + \frac{1}{2}.$$

Porazdelitve v R-ju

V R-ju so za delo s pomembnejšimi porazdelitvami na voljo funkcije:

`dime` – gostota porazdelitve *ime* $p_{ime}(x)$

`pime` – porazdelitvena funkcija *ime* $F_{ime}(q)$

`qime` – obratna funkcija: $q = F_{ime}(p)$

`rim` – slučajno zaporedje iz dane porazdelitve

Za *ime* lahko postavimo:

- `unif` – zvezna enakomerna,
- `binom` – binomska,
- `norm` – normalna,
- `exp` – eksponentna,
- `lnorm` – logaritmičnonormalna,
- `chisq` – porazdelitev χ^2 ,
- ...

Opis posamezne funkcije in njenih parametrov dobimo z ukazom `help`. Na primer `help(rnorm)`.

Simeon Poisson (1781–1840)

“Life is good for two things, learning mathematics and teaching mathematics.”

(b. Pithviers, d. Paris). Simeon Poisson developed many novel applications of mathematics for statistics and physics. His father had been a private soldier, and on his retirement was given a small administrative post in his native village. When the French revolution broke out, his father assumed the government of the village, and soon became a local dignitary.

He was educated by his father who prodded him to be a doctor. His uncle offered to teach him medicine, and began by making him prick the veins of cabbage-leaves with a lancet. When he had perfected this, he was allowed to practice on humans, but in the first case that he did this by himself, the patient died within a few hours. Although the other physicians assured him that this was not an uncommon occurrence, he vowed he would have nothing more to do with the medical profession.

Upon returning home, he discovered a copy of a question set from the Polytechnic school among the official papers sent to his father. This chance event determined his career. At the age of seventeen he entered the Polytechnic. A memoir on finite differences which he wrote when only eighteen was so impressive that it was rapidly published in a prestigious journal. As soon as he had finished his studies he was appointed as a lecturer. Throughout his life he held various scientific posts and professorships. He made the study of mathematics his hobby as well as his business.

Over his life Simeon Poisson wrote between 300-400 manuscripts and books on a variety of mathematical topics, including pure mathematics, the application of mathematics to physical problems, the probability of random events, the theory of electrostatics and magnetism (which led the forefront of the new field of quantum mechanics), physical astronomy, and wave theory.

One of Simeon Poisson’s contributions was the development of equations to analyze random events, later dubbed the Poisson Distribution. The fame of this distribution is often attributed to the following story. Many soldiers in the Prussian Army died due to kicks from horses. To determine whether this was due to a random occurrence or the wrath of god, the Czar commissioned the Russian mathematician Ladislaus Bortkiewicz to determine the statistical significance of the events. Fourteen corps were examined, each for twenty years. For over half the corps-year combinations there were no deaths from horse kicks; for the other combinations the number of deaths ranged up to four. Presumably the risk of lethal horse kicks varied over years and corps, yet the over-all distribution fit remarkably well to a Poisson distribution.

Johann Carl Friedrich Gauss (1777–1855)

German mathematician who is sometimes called the “prince of mathematics.” He was a prodigious child, at the age of three informing his father of an arithmetical error in a complicated payroll calculation and stating the correct answer. In school, when his teacher gave the problem of summing the integers from 1 to 100 (an arithmetic series) to his students to keep them busy, Gauss immediately wrote down the correct answer 5050 on his slate. At age 19, Gauss demonstrated a method for constructing a heptadecagon using only a straightedge and compass which had eluded the Greeks. (The explicit construction of the heptadecagon was accomplished around 1800 by Erchinger.) Gauss also showed that only regular polygons of a certain number of sides could be in that manner (a heptagon, for example, could not be constructed.)

Gauss proved the fundamental theorem of algebra, which states that every polynomial has a root of the form $a+bi$. In fact, he gave four different proofs, the first of which appeared in his dissertation. In 1801, he proved the fundamental theorem of arithmetic, which states that every natural number can be represented as the product of primes in only one way.

At age 24, Gauss published one of the most brilliant achievements in mathematics, *Disquisitiones Arithmeticae* (1801). In it, Gauss systematized the study of number theory (properties of the integers). Gauss proved that every number is the sum of at most three triangular numbers and developed the algebra of congruences.

In 1801, Gauss developed the method of least squares fitting, 10 years before Legendre, but did not publish it. The method enabled him to calculate the orbit of the asteroid Ceres, which had been discovered by Piazzi from only three observations. However, after his independent discovery, Legendre accused Gauss of plagiarism. Gauss published his monumental treatise on celestial mechanics *Theoria Motus* in 1806. He became interested in the compass through surveying and developed the magnetometer and, with Wilhelm Weber measured the intensity of magnetic forces. With Weber, he also built the first successful telegraph.

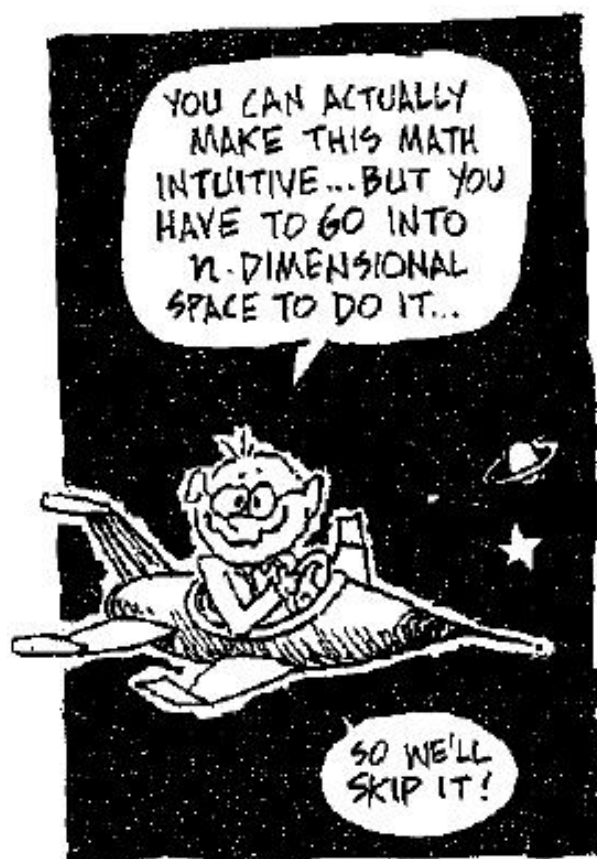
Gauss is reported to have said “There have been only three epoch-making mathematicians: Archimedes, Newton and Eisenstein” (Boyer 1968, p. 553). Most historians are puzzled by the inclusion of Eisenstein in the same class as the other two. There is also a story that in 1807 he was interrupted in the middle of a problem and told that his wife was dying. He is purported to have said, “Tell her to wait a moment ’til I’m through” (Asimov 1972, p. 280).

Gauss arrived at important results on the parallel postulate, but failed to publish them. Credit for the discovery of non-Euclidean geometry therefore went to Janos Bolyai and Lobachevsky. However, he did publish his seminal work on differential geometry in *Disquisitiones circa superticies curvas*. The Gaussian curvature (or “second” curvature) is named for him. He also discovered the Cauchy integral theorem for analytic functions, but did not publish it. Gauss solved the general problem of making a conformal map of one surface onto another.

Unfortunately for mathematics, Gauss reworked and improved papers incessantly, therefore publishing only a fraction of his work, in keeping with his motto “*pauca sed matura*” (few but ripe). Many of his results were subsequently repeated by others, since his terse diary remained unpublished for years after his death. This diary was only 19 pages long, but later confirmed his priority on many results he had not published. Gauss wanted a heptadecagon placed on his gravestone, but the carver refused, saying it would be indistinguishable from a circle. The heptadecagon appears, however, as the shape of a pedestal with a statue erected in his honor in his home town of Braunschweig.

Poglavje 6

Slučajni vektorji



Slučajni vektor je n -terica slučajnih spremenljivk $X = (X_1, \dots, X_n)$. Opišemo ga s porazdelitveno funkcijo ($x_i \in \mathbb{R}$)

$$F(x_1, \dots, x_n) = P(X_1 < x_1, \dots, X_n < x_n),$$

(pri čemer slednja oznaka pomeni $P(\{X_1 < x_1\} \cap \dots \cap \{X_n < x_n\})$) in za katero velja: $0 \leq F(x_1, \dots, x_n) \leq 1$. Funkcija F je za vsako spremenljivko naraščajoča in od leve zvezna,

veljati pa mora tudi

$$F(-\infty, \dots, -\infty) = 0 \quad \text{in} \quad F(\infty, \dots, \infty) = 1.$$

Funkciji $F_i(x_i) = F(\infty, \dots, \infty, x_i, \infty, \dots, \infty)$ pravimo **robna porazdelitvena funkcija** spremenljivke X_i .

Primer: Katere od naslednjih funkcij so lahko porazdelitvene funkcije nekega slučajnega vektorja (X, Y) :

- (a) $F(x, y)$ je enaka $1 - e^{-x-y} - e^{-x} - e^{-y}$, ko sta x in y oba nenegativna, sicer pa je 0,
- (b) $F(x, y)$ je enaka $1 + e^{-x-y} - e^{-x} - e^{-y}$, ko sta x in y oba nenegativna, sicer pa je 0,
- (c) $F(x, y) = x^2$,
- (d) $F(x, y) = x^2 - y^2$.
- (e) $F(x, y) = 0$.

Funkcija iz (a) nima vseh vrednosti na intervalu $[0, 1]$, npr. $F(0, 0) = 1 - 1 - 1 - 1 = -2 < 0$, zato ne more biti porazdelitvena funkcija. Podobno je tudi v primerih (c): $F(2, 0) = 4 \notin [0, 1]$ in (d): $F(0, 1) = -1 \notin [0, 1]$. V primeru (e) pa velja $F(\infty, \infty) = 0 \neq 1$, kar pomeni, da nam ostane le še možnost (b). V tem primeru lahko zapišemo $F(x, y) = (1 - e^{-x})(1 - e^{-y})$ od koder vidimo, da za $x \geq 0$ in $y \geq 0$ velja $F(x, y) \in [0, 1]$. Preverimo še $F(0, 0) = 0$ in $F(\infty, \infty) = 1$. \diamond

Slučajni vektorji – primer

Naj bo

$$A(x, y) = \{(u, v) \in \mathbb{R}^2 : u < x \wedge v < y\}$$

(levi spodnji kvadrant glede na (x, y)). Naj porazdelitvena funkcija opisuje verjetnost, da je slučajna točka (X, Y) v množici $A(x, y)$

$$F(x, y) = P(X < x, Y < y) = P((X, Y) \in A(x, y)).$$

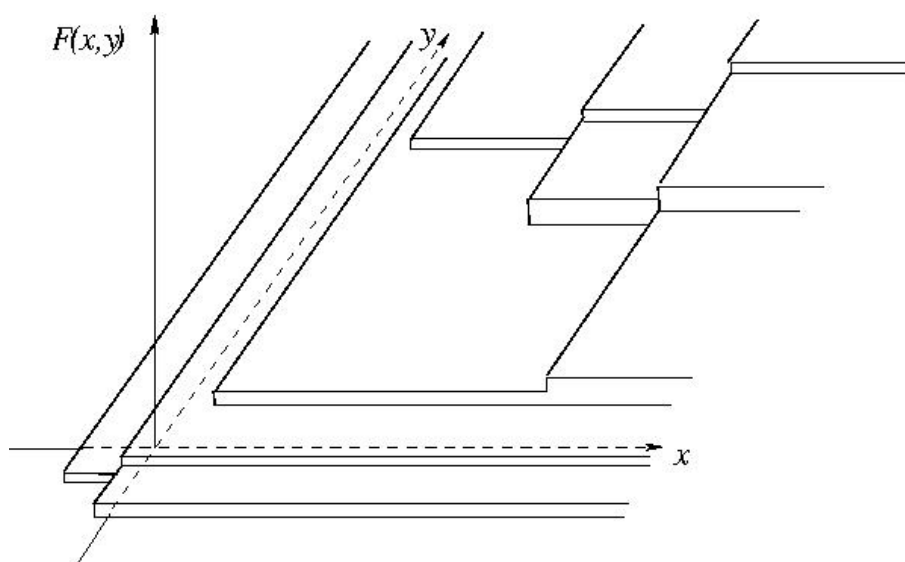
Tedaj je verjetnost, da je slučajna točka (X, Y) v pravokotniku $[a, b] \times [c, d]$ enaka

$$P\left((X, Y) \in [a, b] \times [c, d]\right) = F(b, d) - F(a, d) - F(b, c) + F(a, c) \quad (6.1)$$

Zaloga vrednosti je kvečjemu števna množica. Opišemo jo z **verjetnostno funkcijo** $p_{k_1, \dots, k_n} = P(X_1 = x_{k_1}, \dots, X_n = x_{k_n})$. Za $n = 2$, $X : \{x_1, x_2, \dots, x_k\}$, $Y : \{y_1, \dots, y_m\}$ in $P(X = x_i, Y = y_j)$, sestavimo **verjetnostno tabelo**:

$X \setminus Y$	y_1	y_2	\dots	y_m	X
x_1	p_{11}	p_{12}	\dots	p_{1m}	p_1
x_2	p_{21}	p_{22}	\dots	p_{2m}	p_2
\dots	\dots	\dots	\dots	\dots	\dots
x_k	p_{k1}	p_{k2}	\dots	p_{km}	p_k
Y	q_1	q_2	\dots	q_m	1

$$p_i = P(X = x_i) = \sum_{j=1}^m p_{ij} \quad \text{in} \quad q_j = P(Y = y_j) = \sum_{i=1}^k p_{ij}$$



Slika: Porazdelitvena funkcija $F(x, y)$, v primeru, ko sta spremenljivki X in Y diskretni.

Primer: Naj bosta X in Y diskretni slučajni spremenljivki z zalogami vrednosti $X : \{1, 2\}$ in $Y : \{0, 1\}$. Naj bo

$$P(X = x, Y = y) = p(x, y) = \frac{x - y + a}{5},$$

za neko konstanto a . **Določi a !**

Imamo štiri možne pare za (x, y) : $(1, 0)$, $(1, 1)$, $(2, 0)$, $(2, 1)$, po zgornji formuli pa velja: $P((x, y) = (1, 0)) = (1 + a)/5$, $P((x, y) = (1, 1)) = a/5$, $P((x, y) = (2, 0)) = (2 + a)/5$, $P((x, y) = (2, 1)) = (1 + a)/5$. Vsota vseh verjetnosti je enaka 1, torej je $4 + 4a = 5$ oziroma $a = (5 - 4)/4 = 1/4$. \diamond

6.1 Diskretne večrazsežne porazdelitve – polinomska

Polinomska porazdelitev $P(n; p_1, p_2, \dots, p_r)$, $\sum p_i = 1$, $\sum k_i = n$ je določena s predpisom

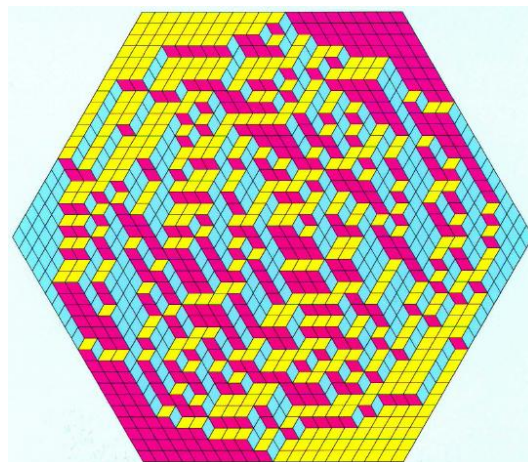
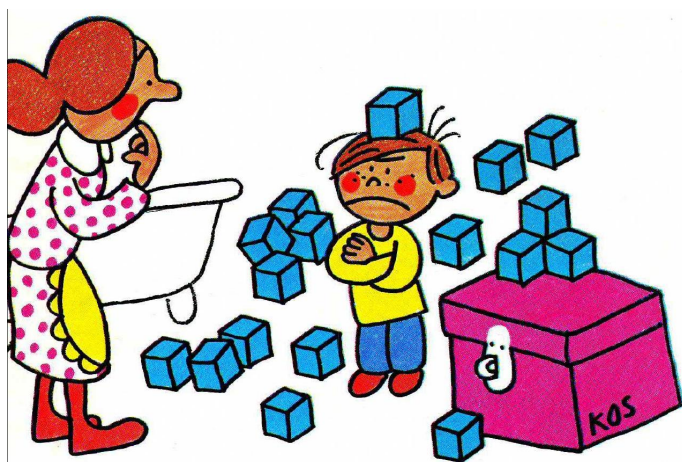
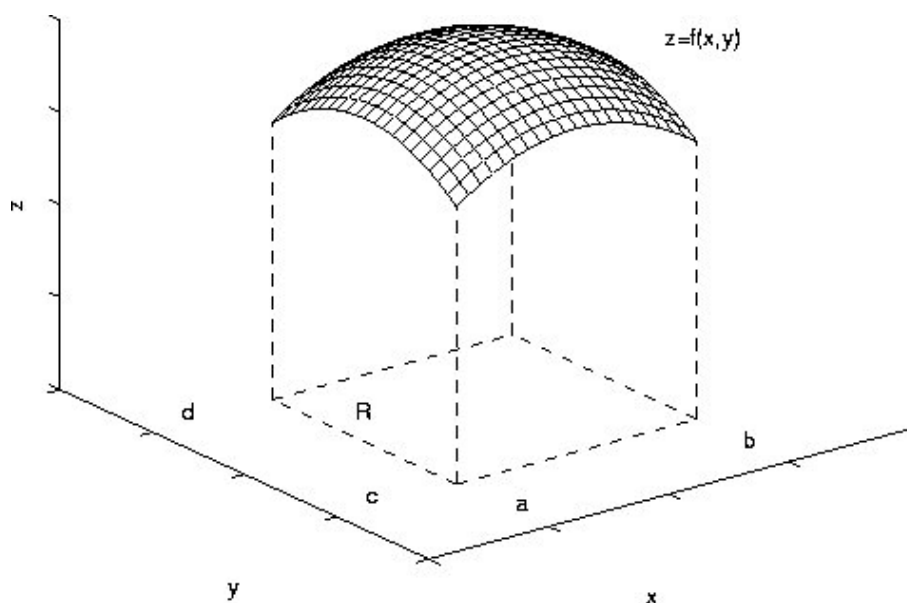
$$P(X_1 = k_1, \dots, X_r = k_r) = \frac{n!}{k_1! \dots k_r!} p_1^{k_1} \dots p_r^{k_r}.$$

Keficient šteje permutacije s ponavljanjem.

http://en.wikipedia.org/wiki/Multinomial_distribution

Za $r = 2$ dobimo binomsko porazdelitev, tj. $B(n, p) = P(n; p, q)$.

6.2 Ponovitev: dvojni integral



Dvojni integral predstavlja prostornino pod neko ploskvijo. Naj bo funkcija $z = f(x, y) \geq 0$ zvezna na nekem območju R v ravnini \mathbb{R}^2 (npr. kar $[a, b] \times [c, d]$). Ploščina telesa med ploskvijo, ki je podana z $z = f(x, y)$, in ravnino $z = 0$ je enaka dvojnemu integralu

$$\iint_R f(x, y) \, dx dy,$$

ki ga v primeru $R = [a, b] \times [c, d]$ izračunamo s pomočjo dvakratnega integrala

$$\int_c^d \left(\int_a^b f(x, y) \, dx \right) dy = \int_a^b \left(\int_c^d f(x, y) \, dy \right) dx.$$

Lastnosti dvojnega integrala

Trditev 6.1. 1) Če je $f(x, y) \leq 0 \forall (x, y) \in R$, je vrednost dvojnega integrala negativna.

2) Naj bo območje $R = R_1 \cup R_2$, kjer je $R_1 \cap R_2 = \emptyset$. Potem velja

$$\iint_R f(x, y) \, dx dy = \iint_{R_1} f(x, y) \, dx dy + \iint_{R_2} f(x, y) \, dx dy.$$

3) Naj bo $f(x, y) \leq g(x, y)$, za vse točke $(x, y) \in R$, potem velja

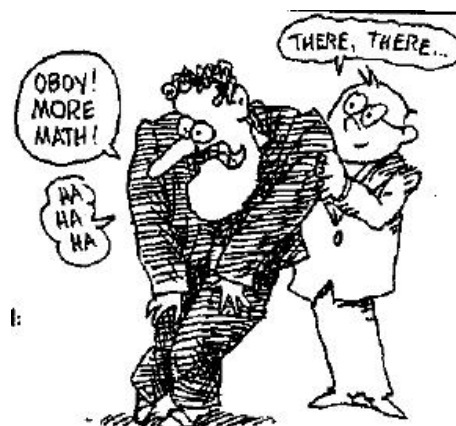
$$\iint_R f(x, y) \, dx dy \leq \iint_R g(x, y) \, dx dy. \quad \square$$

Več o dvojnih integralih najdete npr. na:

<http://www.math.oregonstate.edu/home/programs/undergrad/CalculusQuestStudyGuides/vcalc/255doub/255doub.html>

Računanje dvojnih integralov na pravokotnem območju se prevede na dva običajna (enkratna) integrala.

Kot bomo videli kasneje na primerih, pa je težje izračunati dvojni integral na območju, ki ni pravokotno, ampak je omejeno s poljubnimi krivuljami.



6.3 Zvezne večrazsežne porazdelitve

Slučajni vektor $X = (X_1, X_2, \dots, X_n)$ je **zvezno porazdeljen**, če obstaja integrabilna funkcija (**gostota verjetnosti**) $p(x_1, x_2, \dots, x_n) \geq 0$ z lastnostjo

$$F(x_1, x_2, x_3, \dots, x_n) = \int_{-\infty}^{x_1} \left(\int_{-\infty}^{x_2} \left(\dots \left(\int_{-\infty}^{x_n} p(t_1, t_2, \dots, t_n) dt_n \right) \dots \right) dt_2 \right) dt_1$$

in

$$F(\infty, \infty, \infty, \dots, \infty) = 1.$$

Zvezne dvorazsežne porazdelitve

$$F(x, y) = \int_{-\infty}^x \left(\int_{-\infty}^y p(u, v) dv \right) du,$$

$$P((X, Y) \in [a, b] \times [c, d]) = \int_a^b \left(\int_c^d p(u, v) dv \right) du.$$

Velja

$$\frac{\partial F}{\partial x} = \int_{-\infty}^y p(x, v) dv \quad \text{in} \quad \frac{\partial^2 F}{\partial x \partial y} = p(x, y).$$

Robni verjetnostni gostoti sta

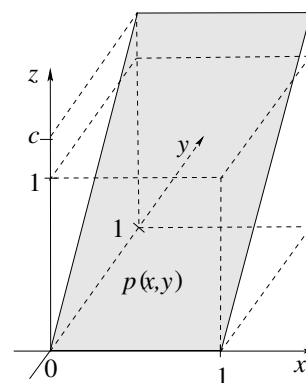
$$p_X(x) = F'_X(x) = \int_{-\infty}^{\infty} p(x, y) dy,$$

in

$$p_Y(y) = F'_Y(y) = \int_{-\infty}^{\infty} p(x, y) dx.$$

Primer: Naj bo gostota porazdelitve vektorja (X, Y) podana s

$$p(x, y) = \begin{cases} cy & \text{če je } 0 \leq x \leq 1 \text{ in } 0 \leq y \leq 1 \\ 0 & \text{sicer.} \end{cases}$$



Določi vrednost konstante c ter robni gostoti za slučajni spremenljivki X in Y !

Dvojni integral gostote verjetnosti je po eni strani enak 1, po drugi pa prostornini telesa, ki je pod osenčenim delom in nad ravnino xy , se pravi, da gre za polovico kvadra in znaša

$1 \times 1 \times c \times 1/2$, od koder dobimo $c = 2$. Slučajna spremenljivka X je na intervalu $[0, 1]$ porazdeljena enakomerno, se pravi, da je $p(x) = 1$. To vidimo tako s slike (prečni prerez), kakor tudi iz definicije:

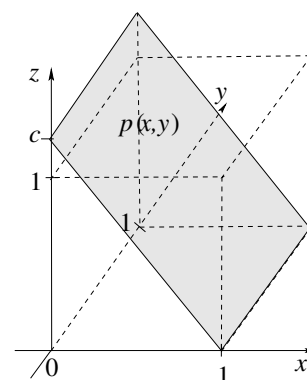
$$p_X(x) = \int_0^1 2y \, dy = y^2 \Big|_{y=0}^1 = 1.$$

Za gostoto verjetnosti slučajne spremenljivke Y pa na intervalu $[0, 1]$ velja

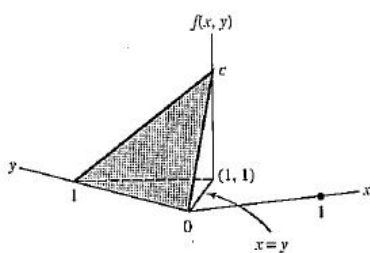
$$p_Y(y) = \int_0^1 2y \, dx = 2xy \Big|_{x=0}^1 = 2y. \quad \diamond$$

Za vajo poskusi odgovoriti na ista vprašanja še za

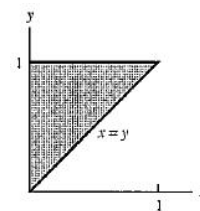
$$p(x, y) = \begin{cases} cx & \text{če je } 0 \leq x \leq 1 \text{ in } 0 \leq y \leq 1 \\ 0 & \text{sicer.} \end{cases}$$



Primer: Naj bo gostota porazdelitve slučajnega vektorja (X, Y) podana s $p(x, y) = f(x, y)$, kjer je



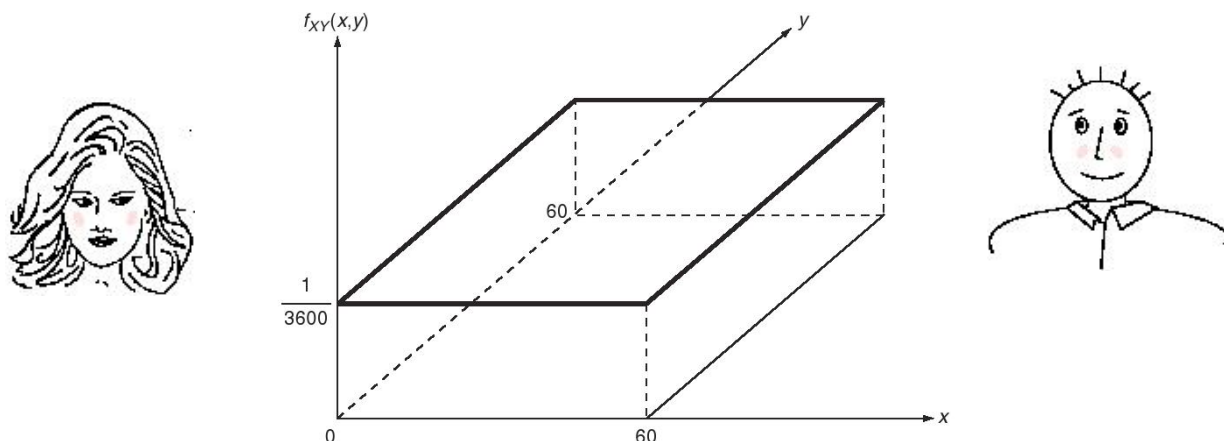
$$f(x, y) = \begin{cases} cx & \text{če je } 0 \leq x \leq y \text{ in } 0 \leq y \leq 1 \\ 0 & \text{sicer.} \end{cases}$$



Izberi vrednost konstante c med:

- (a) 6, (b) 2, (c) 3, (d) 1. \diamond

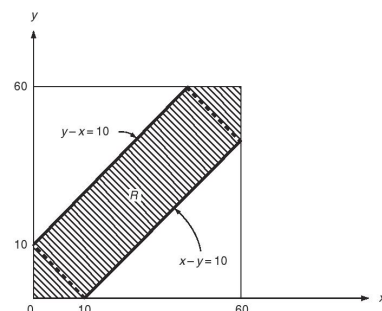
Primer: Dekle in fant se želita srečati na določenem mestu med 9-o in 10-o uro, pri čemer noben od njiju ne bo čakal drugega dlje od 10-ih minut. Če je vsak čas med 9-o in 10-o za vsakega od njiju enako verjeten, in sta njuna časa prihodov neodvisna, poišči verjetnost, da se bosta srečala. Naj bo čas prihoda fanta X minut po 9-i, pravtako pa naj bo čas prihoda dekleta Y minut po 9-i.



Ploskev, ki jo določa gostota porazdelitve, je ravnina, ker pa je prostornina pod njo enaka 1, je oddaljena od ravnine $z = 0$ za $1/3600$.

Prostornina, ki jo iščemo, se nahaja nad področjem R , ki je določeno z $|X - Y| \leq 10$, torej je verjetnost srečanja enaka:

$$P(|X - Y| \leq 10) = \frac{(2 \times 5 \times 10 + 10\sqrt{2} \cdot 50\sqrt{2})}{3600} = \frac{11}{36}.$$



Pri bolj zapletenih gostotah verjetnosti, moramo dejansko izračunati integral

$$F(x, y) = \iint_R p(x, y) dy dx.$$

Za vajo izračunajmo obe robni verjetnostni gostoti. Očitno velja:

$$F(x, y) = 0 \text{ za } (x, y) < (0, 0) \quad \text{in} \quad F(x, y) = 1 \text{ za } (x, y) > (60, 60).$$

Sedaj pa za $(0, 0) \leq (x, y) \leq (60, 60)$ velja

$$F(x, y) = \int_0^y \int_0^x \left(\frac{1}{3600} \right) dy dx = \frac{xy}{3600}.$$

in

$$p_X(x) = F'_X(x) = \int_0^{60} \left(\frac{1}{3600} \right) dy = \frac{1}{60} \quad \text{za } 0 \leq y \leq 60,$$

$$p_Y(y) = F'_Y(y) = \int_0^{60} \left(\frac{1}{3600} \right) dx = \frac{1}{60} \quad \text{za } 0 \leq x \leq 60,$$

za vse ostale x in y pa je $p_X(x) = 0$ ter $p_Y(x) = 0$, torej sta X in Y obe enakomerno porazdeljeni slučajni spremenljivki na intervalu $[0, 60]$. \diamond

Večrazsežna normalna porazdelitev

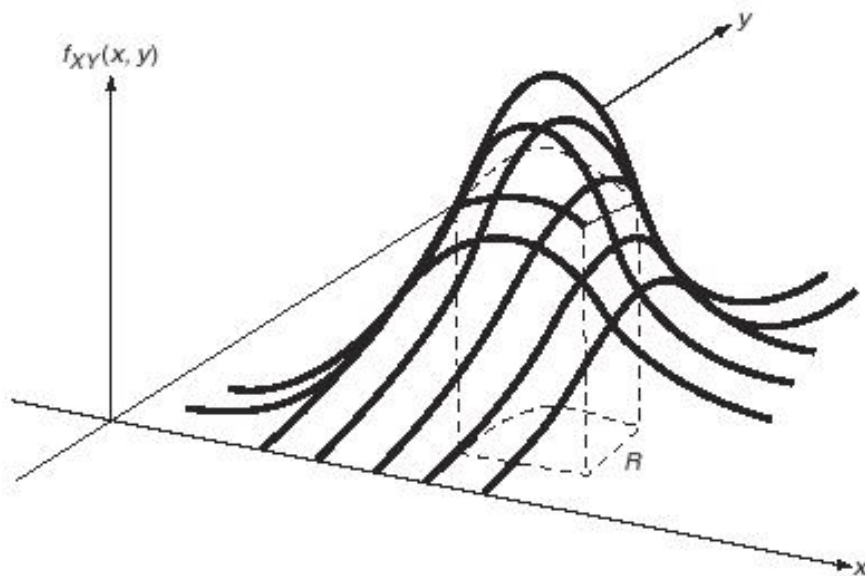
V dveh razsežnostih označimo normalno porazdelitev z $N(\mu_x, \mu_y, \sigma_x, \sigma_y, \rho)$ in ima gostoto

$$p(x, y) = \frac{1}{2\pi\sigma_x\sigma_y\sqrt{1-\rho^2}} e^{-\frac{1}{2(1-\rho^2)} \left(\left(\frac{x-\mu_x}{\sigma_x} \right)^2 - 2\rho \frac{x-\mu_x}{\sigma_x} \frac{y-\mu_y}{\sigma_y} + \left(\frac{y-\mu_y}{\sigma_y} \right)^2 \right)}.$$

V splošnem pa jo zapišemo v matrični obliki

$$p(\mathbf{x}) = \sqrt{\frac{\det A}{(2\pi)^n}} e^{-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T A(\mathbf{x} - \boldsymbol{\mu})},$$

kjer je A simetrična pozitivno definitna matrika, $\boldsymbol{\mu} = (\mu_1, \dots, \mu_n)^T$, T pa pomeni transpiranje. Obe robni porazdelitvi sta normalni.



Primer: Pri študiju upora Y strukturnega elementa in sile X , ki deluje nanj, smatramo za slučajni spremenljivki. Verjetnost napake n_f je definirana z $P(Y \leq X)$. Predpostavimo, da je

$$p(x, y) = abe^{-(ax+by)} \quad \text{za } (x, y) > 0$$

in $p(x, y) = 0$ sicer, pri čemer sta a in b poznani pozitivni števili. Želimo izračunati n_f , tj.

$$F(x, y) = \iint_R p(x, y) dy dx,$$

kjer je območje R določeno s pogojem $Y \leq X$. Ker slučajni spremenljivki X in Y zavzameta samo pozitivne vrednosti, velja

$$n_f = \int_0^\infty \int_y^\infty ab e^{-(ax+by)} dx dy = \int_0^\infty \int_0^x ab e^{-(ax+by)} dy dx.$$

Tu izračunajmo prvi integral, ki smo ga pričeli računati že na prejšnjih predavanjih (bodite pozorni na to, da so sedaj meje popravljene). Upoštevamo $a dx = d(ax) = -d(-ax - by)$:

$$\begin{aligned} \int_0^\infty \int_y^\infty ab e^{-(ax+by)} dx dy &= -b \int_0^\infty \left(\int_y^\infty e^{-(ax+by)} d(-ax - by) \right) dy \\ &= -b \int_0^\infty \left(e^{-(ax+by)} \Big|_{x=y}^\infty \right) dy = b \int_0^\infty e^{-y(a+b)} dy \\ &= \frac{-b}{a+b} \int_0^\infty e^{-y(a+b)} d(-y(a+b)) = \frac{-b}{a+b} \left(e^{-y(a+b)} \Big|_{y=0}^\infty \right) = \frac{b}{a+b}. \quad \diamond \end{aligned}$$

Za vajo izračunajte tudi drugi dvojni integral.

6.4 Neodvisnost slučajnih spremenljivk

Podobno kot pri dogodkih pravimo za slučajne spremenljivke X_1, X_2, \dots, X_n , da so med seboj **neodvisne**, če za poljubne vrednosti $x_1, x_2, \dots, x_n \in \mathbb{R}$ velja

$$F(x_1, x_2, \dots, x_n) = F_1(x_1) \cdot F_2(x_2) \cdots F_n(x_n),$$

kjer je F porazdelitvena funkcija vektorja, F_i pa so porazdelitvene funkcije njegovih komponent.

Trditev 6.2. *Diskretni slučajni spremenljivki X in Y z verjetnostnima tabelama*

$$X : \begin{pmatrix} x_1 & x_2 & \cdots \\ p_1 & p_2 & \cdots \end{pmatrix} \quad \text{in} \quad Y : \begin{pmatrix} y_1 & y_2 & \cdots \\ q_1 & q_2 & \cdots \end{pmatrix}$$

ter verjetnostno funkcijo p_{ij} slučajnega vektorja (X, Y) sta X in Y neodvisni natanko takrat, ko je $p_{ij} = p_i q_j$ za vsak par naravnih števil i, j .

Dokaz. Zgornji pogoj za neodvisnost lahko zapišemo z verjetnostjo v naslednji obliki:

$$P(X < x_k, Y < y_\ell) = P(X < x_k) \cdot P(Y < y_\ell), \quad (6.2)$$

kjer sta k in ℓ poljubni naravni števili. Predpostavimo najprej, da je $p_{ij} = p_i q_j$ za vsak par naravnih števil i, j . Dokaz relacije (6.2) je sedaj precej direkten:

$$\begin{aligned} P(X < x_k, Y < y_\ell) &= \sum_{i < k} \sum_{j < \ell} p_{ij} = \sum_{i < k} \sum_{j < \ell} p_i \cdot q_j \\ &= \sum_{i < k} p_i \cdot \sum_{j < \ell} q_j = P(X < x_k) \cdot P(Y < y_\ell). \end{aligned}$$

Sedaj pa privzemimo pogoj (6.2). Zapišimo diskretno varianto relacije (6.1):

$$\begin{aligned} P(X = x_i, Y = y_j) &= P(X < x_{i+1}, Y < y_{j+1}) - P(X < x_{i+1}, Y < y_j) \\ &\quad - P(X < x_i, Y < y_{j+1}) + P(X < x_i, Y < y_j). \end{aligned}$$

Nariši si sliko, s katero se prepričaš o veljavnosti te relacije, nato pa uporabiš (6.2) na vsaki izmed verjetnosti na desni strani zgornje relacije:

$$\begin{aligned} p_{ij} &= P(X < x_{i+1}) \cdot P(Y < y_{j+1}) - P(X < x_i) \cdot P(Y < y_{j+1}) \\ &\quad - P(X < x_{i+1}) \cdot P(Y < y_j) + P(X < x_i) \cdot P(Y < y_j) \\ &= (P(X < x_{i+1}) - P(X < x_i)) \cdot (P(Y < y_{j+1}) - P(Y < y_j)) = p_i \cdot q_j. \end{aligned}$$

Sklicevanju na sliko pa se lahko izognemo na naslednji način. Najprej se ukvarjamo s spremenljivko X , tako da odštejemo naslednji enakosti:

$$\begin{aligned} P(X < x_i, Y < y_j) &= P(X < x_i) \cdot P(Y < y_j) \\ P(X < x_{i+1}, Y < y_j) &= P(X < x_{i+1}) \cdot P(Y < y_j), \end{aligned}$$

kar nam da

$$P(X = x_i, Y < y_j) = p_i \cdot P(Y < y_j).$$

Potem seveda velja tudi

$$P(X = x_i, Y < y_{j+1}) = p_i \cdot P(Y < y_{j+1}),$$

se pravi, da se lahko posvetimo še spremenljivki Y . Razlika zadnjih dveh relacij nam sedaj da $p_{ij} = p_i \cdot q_j$, kar smo želeli dokazati. \square

Podobno dokažemo tudi naslednjo trditev za zvezno porazdeljeni slučajni spremenljivki (le da v prvemu delu namesto seštevanja integriramo, v drugem delu pa namesto odštevanja parcialno odvajamo).

Trditev 6.3. Če sta X in Y zvezno porazdeljeni slučajni spremenljivki z gostotama p_X in p_Y ter je $p(x, y)$ gostota zvezno porazdeljenega slučajnega vektorja (X, Y) , potem sta X in Y neodvisni natanko takrat, ko za vsak par x, y velja $p(x, y) = p_X(x) \cdot p_Y(y)$. \square

Primer: Naj bo dvorazsežni slučajni vektor (X, Y) z normalno porazdelitvijo $N(\mu_x, \mu_y, \sigma_x, \sigma_y, \rho)$. Če je $\rho = 0$, je

$$p(x, y) = \frac{1}{2\pi\sigma_x\sigma_y} e^{-\frac{1}{2}\left(\left(\frac{x-\mu_x}{\sigma_x}\right)^2 + \left(\frac{y-\mu_y}{\sigma_y}\right)^2\right)} = p_X(x) \cdot p_Y(y).$$

Torej sta komponenti X in Y neodvisni. \diamond

Brez dokaza pa omenimo še močnejšo trditev.

Izrek 6.4. Zvezno porazdeljeni slučajni spremenljivki X in Y sta neodvisni natanko takrat, ko lahko gostoto $p(x, y)$ verjetnosti slučajnega vektorja (X, Y) zapišemo v obliki

$$p(x, y) = f(x) \cdot g(y). \quad \square$$

Naj bosta zvezno porazdeljeni slučajni spremenljivki X in Y tudi neodvisni ter A in B poljubni (Borelovi) podmnožici v \mathbb{R} . Potem sta neodvisna tudi dogodka $X \in A$ in $Y \in B$.

Trditev velja tudi za diskretni slučajni spremenljivki X in Y .

Pogosto pokažemo odvisnost spremenljivk X in Y tako, da najdemo množici A in B , za kateri je

$$P(X \in A, Y \in B) \neq P(X \in A) \cdot P(Y \in B).$$

Primer: Naj slučajna spremenljivka X predstavlja število naprav, ki so na voljo, slučajna spremenljivka Y pa število zaporednih operacij, ki jih moramo opraviti za procesiranje kosa materiala. Verjetnostna funkcija $P(X = x, Y = y) = p(x, y)$ je definirana z naslednjo tabelo:

$Y \setminus X$	1	2	3	4
0	0	0,10	0,20	0,10
1	0,03	0,07	0,10	0,05
2	0,05	0,10	0,05	0
3	0	0,10	0,05	0

Poišči verjetnostno tabelo spremenljivke X !

Seštejemo verjetnosti po stolpcih:

$Y \setminus X$	1	2	3	4	Y
0	0	0,10	0,20	0,10	0,40
1	0,03	0,07	0,10	0,05	0,25
2	0,05	0,10	0,05	0	0,20
3	0	0,10	0,05	0	0,15
X	0,08	0,37	0,40	0,15	1

in dobimo

$$X : \begin{pmatrix} 1 & 2 & 3 & 4 \\ 0,08 & 0,37 & 0,40 & 0,15 \end{pmatrix}.$$

Enako lahko storimo tudi za Y . *Ali sta slučajni spremenljivki X in Y neodvisni?*

Ne nista, saj velja npr.: $P(x = 4, y = 3) = 0 \neq 0,15 \cdot 0,15 = P(x = 4) \cdot P(y = 3)$. \diamond

Poglavje 7

Funkcije slučajnih spremenljivk in vektorjev



7.1 Funkcije slučajnih spremenljivk

Naj bo $X : G \rightarrow \mathbb{R}$ slučajna spremenljivka in $f : \mathbb{R} \rightarrow \mathbb{R}$ neka realna funkcija. Tedaj je njun kompozitum $Y = f \circ X$ določen s predpisom $Y(e) = f(X(e))$, za vsak $e \in G$, določa novo preslikavo $Y : G \rightarrow \mathbb{R}$. *Kdaj je tudi Y slučajna spremenljivka na (G, \mathcal{D}, P) ?* V ta namen mora biti za vsak $y \in \mathbb{R}$ množica

$$(Y < y) = \{e \in G : Y(e) < y\} = \{e \in G : X(e) \in f^{-1}(-\infty, y)\}$$

dogodek – torej v \mathcal{D} . Če je ta pogoj izpolnjen, imenujemo Y **funkcija!slučajne spremenljivke** X in jo zapišemo kar $Y = f(X)$. Njena porazdelitvena funkcija je v tem

primeru

$$F_Y(y) = P(Y < y).$$

Če je funkcija f linearna, potem se porazdelitev verjetnosti ne spremeni, sicer pa se lahko, kot bomo videli v naslednjem primeru.

Primer: Naj bo diskretna slučajna spremenljivka X podana z

$$\begin{pmatrix} -1 & 0 & 1 \\ 1/2 & 1/3 & 1/6 \end{pmatrix}.$$

Potem je porazdelitev za slučajno spremenljivko $2X$ enaka

$$\begin{pmatrix} -2 & 0 & 2 \\ 1/2 & 1/3 & 1/6 \end{pmatrix}.$$

Sedaj pa izberimo še porazdelitev za slučajno spremenljivko X^2 med naslednjimi možnostmi:

$$(a) \begin{pmatrix} -1 & 0 & 1 \\ 1/2 & 1/3 & 2/3 \end{pmatrix}, (b) \begin{pmatrix} -1 & 0 & 1 \\ 1/4 & 1/9 & 1/36 \end{pmatrix}, (c) \begin{pmatrix} -1 & 0 & 1 \\ 0 & 1/3 & 2/3 \end{pmatrix}, (d) \begin{pmatrix} -1 & 0 & 1 \\ 0 & 1/4 & 3/4 \end{pmatrix}.$$

Hitro opazimo, da se v primeru (a) in (b) verjetnosti sploh ne seštejejo v 1, v primeru (d) pa se je spremenila verjetnost za vrednost 0, kar pa tudi ni mogoče. Ostane nam samo še možnost (c), ki pa je seveda prava, saj se verjetnost pri -1 spremenila v 0, verjetnost pri 0 pa je ostala nespremenjena. \diamond

Borelove množice

Vprašanje: kakšna mora biti množica A , da je množica

$$X^{-1}(A) = \{e \in G : X(e) \in A\} \text{ v } \mathcal{D}?$$

Zadoščajo množice A , ki so ali intervali, ali števne unije intervalov, ali števniki preseki števnih unij intervalov – **Borelove množice**.

Kdaj je $f^{-1}(-\infty, y)$ Borelova množica?

Vsekakor je to res, ko je f zvezna funkcija.

V nadaljevanju nas bodo zanimali samo taki primeri.



Emile Borel

Primer: zvezne strogo naraščajoče funkcije

Naj bo $f : \mathbb{R} \rightarrow \mathbb{R}$ zvezna in strogo naraščajoča funkcija. Tedaj je taka tudi funkcija f^{-1} in velja

$$f^{-1}(-\infty, y) = \{x \in \mathbb{R} : f(x) < y\} = \{x \in \mathbb{R} : x < f^{-1}(y)\} = (-\infty, f^{-1}(y))$$

in potemtakem tudi $F_Y = F_X \circ f^{-1}$, o čemer se prepričamo takole

$$F_Y(y) = P(Y < y) = P(f(X) < y) = P(X < f^{-1}(y)) = F_X(f^{-1}(y))$$

Če je X porazdeljena zvezno z gostoto $p(x)$, je

$$F_Y(y) = \int_{-\infty}^{f^{-1}(y)} p(x) dx$$

in, če je f odvedljiva, še

$$p_Y(y) = p(f^{-1}(y))f^{-1}(y)'$$

Če funkcija ni monotona, lahko njeno definicijsko območje razdelimo na intervale monotoni in obravnavamo vsak interval ločeno.

Primer: Obravnavajmo kvadrat normalno porazdeljene spremenljivke, tj. naj bo $X : N(0, 1)$ in $Y = X^2$. Tedaj je $F_Y(y) = P(Y < y) = P(X^2 < y) = 0$ za $y \leq 0$, za $y > 0$ pa velja

$$F_Y(y) = P(|X| < \sqrt{y}) = F_X(\sqrt{y}) - F_X(-\sqrt{y})$$

in ker je $p_X(x)$ soda funkcija

$$p_Y(y) = p_X(\sqrt{y})\frac{1}{2\sqrt{y}} + p_X(-\sqrt{y})\frac{1}{2\sqrt{y}} = \frac{1}{\sqrt{y}}p_X(\sqrt{y})$$

Vstavimo še standardizirano normalno porazdelitev

$$p_Y(y) = \begin{cases} \frac{1}{\sqrt{2\pi}}y^{-\frac{1}{2}}e^{-\frac{y}{2}} & y > 0 \\ 0 & y \leq 0 \end{cases}$$

pa dobimo porazdelitev $\chi^2(1)$. ◇

7.2 Funkcije in neodvisnost

Trditev 7.1. Če sta X in Y neodvisni slučajni spremenljivki ter f in g zvezni funkciji na \mathbb{R} , sta tudi $U = f(X)$ in $V = g(Y)$ neodvisni slučajni spremenljivki.

Dokaz. Za poljubna $u, v \in \mathbb{R}$ velja

$$\begin{aligned}
 P(U < u, V < v) &= P(f(X) < u, g(Y) < v) \\
 &= P(X \in f^{-1}(-\infty, u), Y \in g^{-1}(-\infty, v)) \\
 &\quad (\text{ker sta } X \text{ in } Y \text{ sta neodvisni uporabimo Trditev 6.3}) \\
 &= P(X \in f^{-1}(-\infty, u)) \cdot P(Y \in g^{-1}(-\infty, v)) \\
 &\quad (\text{in naprej}) \\
 &= P(f(X) < u) \cdot P(g(Y) < v) \\
 &= P(U < u) \cdot P(V < v). \quad \square
 \end{aligned}$$

Funkcije slučajnih vektorjev

Imejmo slučajni vektor $\mathbf{X} = (X_1, X_2, \dots, X_n) : G \rightarrow \mathbb{R}^n$ in zvezno vektorsko preslikavo $f = (f_1, f_2, \dots, f_m) : \mathbb{R}^n \rightarrow \mathbb{R}^m$. Tedaj so $Y_j = f_j(X_1, X_2, \dots, X_n)$, $j = 1, \dots, m$ slučajne spremenljivke – komponente slučajnega vektorja $\mathbf{Y} = (Y_1, Y_2, \dots, Y_m)$.

Pravimo tudi, da je \mathbf{Y} **funkcija slučajnega vektorja** \mathbf{X} , tj. $\mathbf{Y} = f(\mathbf{X})$.

Porazdelitve komponent dobimo na običajen način

$$F_{Y_j}(y) = P(Y_j < y) = P(f_j(\mathbf{X}) < y) = P(\mathbf{X} \in f_j^{-1}(-\infty, y))$$

in, če je \mathbf{X} zvezno porazdeljen z gostoto $p(x_1, x_2, \dots, x_n)$, potem je

$$F_{Y_j}(y) = \int \int \dots \int_{f_j^{-1}(-\infty, y)} p(x_1, x_2, \dots, x_n) dx_1 dx_2 \dots dx_n.$$

Vsota slučajnih spremenljivk

Oglejmo si en enostaven primer. Naj bo $Z = X + Y$, kjer je (X, Y) zvezno porazdeljen slučajni vektor z gostoto $p(x, y)$. Tedaj je

$$\begin{aligned}
 F_Z(z) &= P(Z < z) = P(X + Y < z) = \\
 &= \int \int_{x+y < z} p(x, y) dx dy = \int_{-\infty}^{\infty} dx \int_{-\infty}^{z-x} p(x, y) dy
 \end{aligned}$$

in

$$p_Z(z) = F'_Z(z) = \int_{-\infty}^{\infty} p(x, z-x) dx = \int_{-\infty}^{\infty} p(z-y, y) dy.$$

Če sta spremenljivki X in Y neodvisni dobimo naprej zvezo

$$p_Z(z) = \int_{-\infty}^{\infty} p_X(x) p_Y(z-x) dx.$$

Gostota $p_Z = p_X * p_Y$ je **konvolucija** funkcij p_X in p_Y .

Primer: Če je $(X, Y) : N(\mu_x, \mu_y, \sigma_x, \sigma_y, \rho)$, je vsota $Z = X + Y$ zopet normalno porazdeljena $Z : N(\mu_x + \mu_y, \sqrt{\sigma_x^2 + 2\rho\sigma_x\sigma_y + \sigma_y^2})$.

Če sta $X : \chi^2(n)$ in $Y : \chi^2(m)$ neodvisni slučajni spremenljivki, je tudi njuna vsota $Z = X + Y$ porazdeljena po tej porazdelitvi $Z : \chi^2(n+m)$. \diamond

Dosedanje ugotovitve lahko združimo v naslednjo trditev.

Trditev 7.2. Če so X_1, X_2, \dots, X_n neodvisne standardizirano normalne slučajne spremenljivke, je slučajna spremenljivka $Y = X_1^2 + X_2^2 + \dots + X_n^2$ porazdeljena po $\chi^2(n)$.

7.3 Funkcije slučajnih vektorjev

Naj bo sedaj $f : (x, y) \mapsto (u, v)$ transformacija slučajnega vektorja (X, Y) v slučajni vektor (U, V) določena z zvezama $u = u(x, y)$ in $v = v(x, y)$, torej je $U = u(X, Y)$ in $V = v(X, Y)$. Porazdelitveni zakon za nov slučajni vektor (U, V) je

$$\begin{aligned} F_{U,V}(u, v) &= P(U < u, V < v) = P((U, V) \in A(u, v)) \\ &= P((X, Y) \in f^{-1}(A(u, v))). \end{aligned}$$

Pri zvezno porazdeljenem slučajnem vektorju (X, Y) z gostoto $p(x, y)$ je

$$F_{U,V}(u, v) = \iint_{f^{-1}(A(u, v))} p(x, y) dx dy.$$

Če je f bijektivna z zveznimi parcialnimi odvodi, lahko nadaljujemo

$$F_{U,V}(u, v) = \iint_{A(u, v)} p(x(u, v), y(u, v)) |J(u, v)| du dv,$$

kjer je (glej učbenik <http://rkb.home.cern.ch/rkb/titleA.html>)

$$J(u, v) = \frac{\partial(u, v)}{\partial(x, y)} = \det \begin{pmatrix} \frac{\partial u}{\partial x} & \frac{\partial u}{\partial y} \\ \frac{\partial v}{\partial x} & \frac{\partial v}{\partial y} \end{pmatrix}.$$

Jacobijeva determinanta (glej <http://en.wikipedia.org/wiki/Jacobian> za kakšen primer).

Za gostoto $q(u, v)$ vektorja (U, V) dobimo od tu

$$q(u, v) = p(x(u, v), y(u, v)) |J(u, v)|.$$



Primer:

$$\Omega = \{(x, y) \mid 0 < x \leq 1, 0 < y \leq 1\}.$$

Naj bo

$$\begin{aligned} r &= \sqrt{-2 \log(x)}, & \varphi &= 2\pi y, \\ u &= r \cos \varphi, & v &= r \sin \varphi. \end{aligned}$$

Potem po pravilu za odvajanje posrednih funkcij in definiciji Jacobijeve matrike velja

$$\frac{\partial(u, v)}{\partial(x, y)} = \begin{pmatrix} \frac{\partial(u, v)}{\partial(r, \varphi)} \end{pmatrix} \begin{pmatrix} \frac{\partial(r, \varphi)}{\partial(x, y)} \end{pmatrix} = \begin{pmatrix} \cos \varphi & -r \sin \varphi \\ \sin \varphi & r \cos \varphi \end{pmatrix} \begin{pmatrix} \frac{-1}{rx} & 0 \\ 0 & 2\pi \end{pmatrix}.$$

Jacobijeva determinanta je

$$\det \left(\frac{d\mathbf{u}}{d\mathbf{x}} \right) = \det \left(\frac{\partial(u, v)}{\partial(x, y)} \right) = \det \left(\frac{\partial(u, v)}{\partial(r, \varphi)} \right) \det \left(\frac{\partial(r, \varphi)}{\partial(x, y)} \right) = r \frac{-2\pi}{rx} = \frac{-2\pi}{x}$$

in

$$d^2\mathbf{x} = \left| \det \left(\frac{d\mathbf{x}}{d\mathbf{u}} \right) \right| d^2\mathbf{u} = \left| \det \left(\frac{d\mathbf{u}}{d\mathbf{x}} \right) \right|^{-1} d^2\mathbf{u} = \frac{x}{2\pi} d^2\mathbf{u} = \frac{e^{\frac{u^2 + v^2}{2}}}{2\pi} d^2\mathbf{u}.$$

Od tod zaključimo, da za neodvisni slučajni spremenljivki x in y , ki sta enakomerno porazdeljeni med 0 in 1, zgoraj definirani slučajni spremenljivki u in v pravtako neodvisni in porazdeljeni normalno.



7.4 Pogojne porazdelitve

Naj bo B nek mogoč dogodek, tj. $P(B) > 0$. Potem lahko vpeljemo **pogojno porazdelitveno funkcijo**

$$F(x|B) = P(X < x | B) = \frac{P(X < x, B)}{P(B)}$$

V diskretnem primeru je:

$$p_{ik} = P(X = x_i, Y = y_k), \quad B = (Y = y_k) \quad \text{in} \quad P(B) = P(Y = y_k) = q_k.$$

Tedaj je pogojna porazdelitvena funkcija

$$F_X(x|y_k) = F_X(x|Y = y_k) = P(X < x | Y = y_k) = \frac{P(X < x, Y = y_k)}{P(Y = y_k)} = \frac{1}{q_k} \sum_{x_i < x} p_{ik}$$

Vpeljimo **pogojno verjetnostno funkcijo** s $p_{i|k} = \frac{p_{ik}}{q_k}$. Tedaj je $F_X(x|y_k) = \sum_{x_i < x} p_{i|k}$.

Primer: Nadaljujmo primer s katerim smo končali razdelek 6.4. Zapiši pogojno verjetnostno porazdelitev slučajne spremenljivke X glede na pogoj $y = 2$!

$Y \setminus X$	1	2	3	4	Y
0	0	0,10	0,20	0,10	0,40
1	0,03	0,07	0,10	0,05	0,25
2	0,05	0,10	0,05	0	0,20
3	0	0,10	0,05	0	0,15
X	0,08	0,37	0,40	0,15	1

Verjetnosti v vrstici pri $y = 2$ moramo deliti s $P(Y = 2)$, ki je enaka 0,2:

$$X|y = 2 : \begin{pmatrix} 1 & 2 & 3 & 4 \\ 0,25 & 0,50 & 0,25 & 0 \end{pmatrix}.$$

◇

Primer: Dostavni tovornjak potuje od A do B in nazaj vsak dan. Na poti ima tri semaforje. Naj bo X število rdečih semaforjev na katere naleti tovornjak na poti do dostavne točke B , in y število rdečih luči nazaj na poti do točke A . Inženir za promet je določil naslednjo verjetnostno porazdelitev:

$Y \setminus X$	0	1	2	3	Y
0	0,01	0,02	0,07	0,01	0,11
1	0,03	0,06	0,10	0,06	0,25
2	0,05	0,12	0,15	0,08	0,40
3	0,02	0,09	0,08	0,05	0,24
X	0,11	0,29	0,40	0,20	1

Poišči robno porazdelitev za Y .

$$Y : \begin{pmatrix} 0 & 1 & 2 & 3 \\ 0,11 & 0,25 & 0,40 & 0,24 \end{pmatrix}.$$

Če vemo, da je tovornjak naletel na $x = 2$ luči do točke B , potem določi porazdelitev za Y .

$$Y|x=2 : \begin{pmatrix} 0 & 1 & 2 & 3 \\ 7/40 & 1/4 & 3/8 & 1/5 \end{pmatrix}. \quad \diamond$$

Zvezne pogojne porazdelitve

Postavimo $B = (y \leq Y < y + h)$ za $h > 0$ in zahtevajmo $P(B) > 0$.

$$F_X(x|B) = P(X < x | B) = \frac{P(X < x, y \leq Y < y + h)}{P(y \leq Y < y + h)} = \frac{F(x, y + h) - F(x, y)}{F_Y(y + h) - F_Y(y)}.$$

Če obstaja limita (za $h \rightarrow 0$)

$$F_X(x|y) = F_X(x | Y = y) = \lim_{h \rightarrow 0} \frac{F(x, y + h) - F(x, y)}{F_Y(y + h) - F_Y(y)},$$

jo imenujemo **pogojna porazdelitvena funkcija** slučajne spremenljivke X glede na dogodek ($Y = y$).

Gostota zvezne pogojne porazdelitve

Naj bosta gostoti $p(x, y)$ in $p_Y(y)$ zvezni ter $p_Y(y) > 0$. Tedaj je

$$F_X(x|y) = \lim_{h \rightarrow 0} \frac{\frac{F(x, y + h) - F(x, y)}{h}}{\frac{F_Y(y + h) - F_Y(y)}{h}} = \frac{\frac{\partial F}{\partial y}(x, y)}{F'_Y(y)} = \frac{1}{p_Y(y)} \int_{-\infty}^x p(u, y) du$$

oziroma, če vpeljemo **pogojno gostoto**

$$p_X(x|y) = \frac{p(x, y)}{p_Y(y)},$$

tudi $F_X(x|y) = \int_{-\infty}^x p_X(u|y) du$.

Primer: Za dvorazsežno normalno porazdelitev dobimo

$$p_X(x|y) : N\left(\mu_x + \rho \frac{\sigma_x}{\sigma_y}(y - \mu_y), \sigma_x \sqrt{1 - \rho^2}\right). \quad \diamond$$

Emile Borel (1871–1956)

French mathematician educated at the École Normale Supérieure (1889-1892) who received his Doctor es sciences (1894). He was appointed “Maitre de conférence” at the University of Lille (1893), and subsequently at the École Normale Supérieure in 1897. He was professor of the theory of functions at the Sorbonne (1909-1920) and professor of Calcul des Probabilités et de Physique mathématiques (1920-1941). He was scientific director of the École Normale Supérieure (1911), and became a member of the Académie des sciences in 1921. He was also a member of French Chamber of Deputies (1924-1936), and Minister of the Navy for a few months in 1925 (not the fifteen years claimed in many biographies).

Borel worked on divergent series, the theory of functions, probability, game theory, and was the first to define games of strategy. In particular, he found an elementary proof of Picard’s theorem in 1896. Borel founded measure theory, which is the application of the theory of sets to the theory of functions, thus becoming founding with Lebesgue and Baire of modern theory of functions of real variables. Borel’s work culminated in the Heine-Borel theorem. He showed that a “sum” could be defined for some divergent series. Borel wrote more thirty books and near 300 papers, including a number of popular works of high scientific quality, and devoted more fifty papers to history and philosophy of sciences. <http://scienceworld.wolfram.com/biography/Borel.html>

Poglavje 8

Momenti in kovarianca



8.1 Matematično upanje

Matematično upanje EX (pričakovana vrednost) je posplošitev povprečne vrednosti diskretne spremenljivke $X : \begin{pmatrix} x_1 & x_2 & \cdots & x_n \\ p_1 & p_2 & \cdots & p_n \end{pmatrix}$, tj.

$$\bar{X} = \frac{1}{N} \sum_{i=1}^n x_i k_i = \sum_{i=1}^n x_i f_i,$$

od koder izhaja

$$EX = \sum_{i=1}^n x_i p_i.$$

Diskretna slučajna spremenljivka X z verjetnostno funkcijo p_k ima matematično upanje $EX = \sum_{i=1}^{\infty} x_i p_i$, če je

$$\sum_{i=1}^{\infty} |x_i| p_i < \infty.$$

Zvezna slučajna spremenljivka X z gostoto $p(x)$ ima matematično upanje $EX = \int_{-\infty}^{\infty} xp(x) dx$, če je

$$\int_{-\infty}^{\infty} |x| p(x) dx < \infty.$$

Primer: Omenimo dve slučajni spremenljivki, za katere matematično upanje ne obstaja:

Diskretna: $x_k = (-1)^{k+1}2^k/k$ in $p_k = 2^{-k}$.

Zvezna: $X : p(x) = \frac{1}{\pi(1+x^2)}$ – Cauchyeva porazdelitev.

V prvem primeru bi morala biti končna naslednja vsota:

$$S = \sum_{k=1}^{\infty} \frac{1}{k}.$$

Opazimo naslednje:

$$\frac{1}{3} + \frac{1}{4} > \frac{1}{4} + \frac{1}{4} = \frac{1}{2} \quad \text{in} \quad \frac{1}{5} + \frac{1}{6} + \frac{1}{7} + \frac{1}{8} > \frac{1}{8} + \frac{1}{8} + \frac{1}{8} + \frac{1}{8} = \frac{1}{2}$$

ter v splošnem

$$\frac{1}{2^n + 1} + \cdots + \frac{1}{2^{n+1}} > \frac{2^n}{2^{n+1}} = \frac{1}{2}.$$

Na ta način pridemo do spodnje meje

$$\sum_{k=1}^{\infty} \frac{1}{k} > \sum_{n=1}^{\infty} \frac{1}{2},$$

od koder je očitno, da je vsota S neskončna. V drugem primeru pa

$$\int_{-\infty}^{\infty} \frac{|x| dx}{\pi(1+x^2)} = 2 \int_0^{\infty} \frac{x dx}{\pi(1+x^2)} > \int_0^{\infty} \frac{d(x^2+1)}{\pi(1+x^2)} = \frac{1}{\pi} \int_1^{\infty} \frac{dz}{z} = \frac{1}{\pi} \ln z \Big|_1^{\infty} = \infty.$$

Velja omeniti, da je zadnji integral večji od S , tako da nam sploh ne bi bilo potrebno integrirati, da bi se prepričali, da integral ni končen. \diamond

Lastnosti matematičnega upanja

Naj bo a realna konstanta. Če je $P(X = a) = 1$, velja $\mathbf{E}X = a$.

Slučajna spremenljivka X ima matematično upanje natanko takrat, ko ga ima slučajna spremenljivka $|X|$. Očitno velja $|\mathbf{E}X| \leq \mathbf{E}|X|$. Za diskretno slučajno spremenljivko je $\mathbf{E}|X| = \sum_{i=1}^{\infty} |x_i|p_i$, za zvezno pa $\mathbf{E}|X| = \int_{-\infty}^{\infty} |x|p(x) dx$.

Velja splošno: matematično upanje funkcije $f(X)$ obstaja in je enako za diskretno slučajno spremenljivko $\mathbf{E}f(X) = \sum_{i=1}^{\infty} f(x_i)p_i$, za zvezno pa $\mathbf{E}f(X) = \int_{-\infty}^{\infty} f(x)p(x) dx$, če ustrezni izraz absolutno konvergira.

Naj bo a realna konstanta. Če ima slučajna spremenljivka X matematično upanje, potem ga ima tudi spremenljivka aX in velja $\mathbf{E}(aX) = a\mathbf{E}X$. Če imata slučajni spremenljivki X in Y matematično upanje, ga ima tudi njuna vsota $X + Y$ in velja $\mathbf{E}(X + Y) = \mathbf{E}X + \mathbf{E}Y$. Za primer dokažimo zadnjo lastnost za zvezne slučajne spremenljivke. Naj bo p gostota slučajnega vektorja (X, Y) in $Z = X + Y$. Kot vemo, je $p_Z(z) = \int_{-\infty}^{\infty} p(x, z - x) dx$. Pokažimo najprej, da Z ima matematično upanje.

$$\begin{aligned} \mathbf{E}|X + Y| &= \mathbf{E}|Z| = \int_{-\infty}^{\infty} |z| p_Z(z) dz \\ &= \int_{-\infty}^{\infty} |z| \left(\int_{-\infty}^{\infty} p(x, z - x) dx \right) dz = \int_{-\infty}^{\infty} \left(\int_{-\infty}^{\infty} |x + y| p(x, y) dx \right) dy \\ &\leq \int_{-\infty}^{\infty} \left(\int_{-\infty}^{\infty} |x| p(x, y) dx \right) dy + \int_{-\infty}^{\infty} \left(\int_{-\infty}^{\infty} |y| p(x, y) dx \right) dy \\ &= \int_{-\infty}^{\infty} |x| p_X(x) dx + \int_{-\infty}^{\infty} |y| p_Y(y) dy = \mathbf{E}|X| + \mathbf{E}|Y| < \infty. \end{aligned}$$

Sedaj pa še zvezo

$$\begin{aligned} \mathbf{E}(X + Y) &= \mathbf{E}Z = \int_{-\infty}^{\infty} z p_Z(z) dz \\ &= \int_{-\infty}^{\infty} z \left(\int_{-\infty}^{\infty} p(x, z - x) dx \right) dz = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} (x + y) p(x, y) dx dy \\ &= \int_{-\infty}^{\infty} \left(\int_{-\infty}^{\infty} x p(x, y) dx \right) dy + \int_{-\infty}^{\infty} \left(\int_{-\infty}^{\infty} y p(x, y) dx \right) dy \\ &= \int_{-\infty}^{\infty} x p_X(x) dx + \int_{-\infty}^{\infty} y p_Y(y) dy = \mathbf{E}X + \mathbf{E}Y. \end{aligned}$$

Torej je matematično upanje \mathbf{E} **linearen funkcional**, tj.

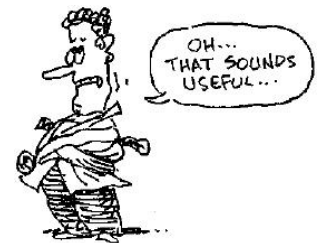
$$\mathbf{E}(aX + bY) = a\mathbf{E}X + b\mathbf{E}Y.$$

Z indukcijo posplošimo to na poljubno končno število členov

$$\mathbf{E}(a_1X_1 + a_2X_2 + \cdots + a_nX_n) = a_1\mathbf{E}X_1 + a_2\mathbf{E}X_2 + \cdots + a_n\mathbf{E}X_n.$$

Trditev 8.1. Če obstajata matematični upanji $\mathbf{E}X^2$ in $\mathbf{E}Y^2$, obstaja tudi matematično upanje produkta $\mathbf{E}XY$ in velja ocena $\mathbf{E}|XY| \leq \sqrt{\mathbf{E}X^2\mathbf{E}Y^2}$.

Enakost velja natanko takrat, ko velja $Y = \pm\sqrt{\mathbf{E}Y^2/\mathbf{E}X^2}X$ z verjetnostjo 1.



Trditev 8.2. Če sta slučajni spremenljivki, ki imata matematično upanje, neodvisni, obstaja tudi matematično upanje njenega produkta in velja $\mathbf{EXY} = \mathbf{EX} \cdot \mathbf{EY}$. \square

Obstajajo tudi odvisne spremenljivke, za katere velja gornja zveza.

Primer: Naj bo slučajna spremenljivka X porazdeljena standardizirano normalno. Potem je $Y = X^2$ porazdeljena po $\chi^2(1)$. Velja tudi $\mathbf{EX} = 0$, $\mathbf{EXY} = \mathbf{E}(X^3) = 0$ in zato $\mathbf{EXY} = 0 = \mathbf{EX} \cdot \mathbf{EY}$. Po drugi strani pa je $P(0 \leq X < 1, Y \geq 1) = 0$, $P(0 \leq X < 1) = \Phi(1) > 0$ in $P(Y \geq 1) = 1 - P(Y < 1) = 1 - P(-1 < X < 1) = 1 - 2\Phi(1) > 0$. \diamond

Spremenljivki, za kateri velja $\mathbf{EXY} \neq \mathbf{EX} \cdot \mathbf{EY}$ imenujemo **korelirani**.

Primer: Življenska doba varovalk (merjena v stotinah ur), ki jih uporabljamo pri računalniških monitorjih ima eksponentno porazdelitev s parametrom $\lambda = 5$. Vsak monitor ima dve varovalki, pri čemer ena deluje kot 'backup' in prične delovati šele ko prva odpove.

- Če imata dve taki varovalki neodvisni življenski dobi X in Y , potem poišči gostoto porazdelitve $p(x, y)$.
- Efektivna skupna življenska doba dveh varovalk je $(X + Y)$. Poišči pričakovano skupno efektivno življensko dobo para dveh varovalk za monitor.

Odgovor: (a)

$$p(x, y) = \begin{cases} 25e^{-5(x+y)} & x, y > 0 \\ 0 & \text{sicer} \end{cases}$$

(b) $2/5$. \diamond

Primer: Koliko trčenj (rojstni dan na isti dan) lahko pričakujemo v skupini 100ih ljudi? \diamond

8.2 Disperzija

Disperzija ali **varianca** \mathbf{DX} slučajne spremenljivke, ki ima matematično upanje, je določena z izrazom

$$\mathbf{DX} = \mathbf{E}(X - \mathbf{EX})^2.$$

Disperzija je vedno nenegativna, $\mathbf{DX} \geq 0$, je pa lahko tudi neskončna. Velja zveza

$$\mathbf{DX} = \mathbf{EX}^2 - (\mathbf{EX})^2.$$

Naj bo a realna konstanta. Če je $P(X = a) = 1$, je $DX = 0$. Iz linearnosti matematičnega upanja sledi tudi $D(aX) = a^2DX$ in $D(X + a) = DX$.

Trditev 8.3. Če obstaja DX in je a realna konstanta, obstaja tudi $E(X - a)^2$ in velja $E(X - a)^2 \geq DX$. Enakost velja natanko za $a = EX$.

Količino $\sigma X = \sqrt{DX}$ imenujemo **standardna deviacija** ali **standardni odklon**.

8.3 Standardizirane spremenljivke

Slučajno spremenljivko X **standardiziramo** s transformacijo

$$X_S = \frac{X - \mu}{\sigma},$$

kjer sta $\mu = EX$ in $\sigma = \sqrt{DX}$. Za X_S velja $EX_S = 0$ in $DX_S = 1$, saj je

$$EX_S = E\frac{X - \mu}{\sigma} = \frac{E(X - \mu)}{\sigma} = \frac{\mu - \mu}{\sigma} = 0,$$

kjer smo upoštevali linearnost matematičnega upanja, ter

$$DX_S = D\frac{X - \mu}{\sigma} = \frac{D(X - \mu)}{\sigma^2} = \frac{\sigma^2 - 0}{\sigma^2} = 1.$$

Matematična upanja in disperzije nekaterih porazdelitev

porazdelitev	EX	DX
binomska $B(n, p)$	np	npq
Poissonova $P(\lambda)$	λ	λ
Pascalova $P(m, p)$	m/p	mq/p^2
geometrijska $G(p)$	$1/p$	q/p^2
enakomerna zv. $E(a, b)$	$(a + b)/2$	$(b - a)^2/12$
normalna $N(\mu, \sigma)$	μ	σ^2
Gama $\Gamma(b, c)$	b/c	b/c^2
hi-kvadrat $\chi^2(n)$	n	$2n$

8.4 Kovarianca

Kovarianca $\text{Cov}(X, Y)$ slučajnih spremenljivk X in Y je definirana z izrazom

$$\text{Cov}(X, Y) = E((X - EX)(Y - EY)).$$

Zgornji izraz pa lahko poenostavimo na enak način kot pri varianci:

$$\text{Cov}(X, Y) = \text{EXY} - \text{EXEY}.$$

Velja tudi: $\text{Cov}(X, Y) = \text{Cov}(Y, X)$ (simetričnost) in

$\text{Cov}(aX + bY, Z) = a \text{Cov}(X, Z) + b \text{Cov}(Y, Z)$ (bilinearnost).

Trditev 8.4. Če obstajata DX in DY , obstaja tudi $\text{Cov}(X, Y)$ in velja

$$|\text{Cov}(X, Y)| \leq \sqrt{\text{DXDY}} = \sigma_X \sigma_Y.$$

Enakost velja natanko takrat, ko je

$$Y - \text{EY} = \pm \frac{\sigma_Y}{\sigma_X} (X - \text{EX})$$

z verjetnostjo 1.

Spremenljivki X in Y sta nekorelirani natanko takrat, ko je $\text{Cov}(X, Y) = 0$. Če imata spremenljivki X in Y končni disperziji, jo ima tudi njuna vsota $X + Y$ in velja

$$\text{D}(X + Y) = \text{DX} + \text{DY} + 2\text{Cov}(X, Y).$$

Če pa sta spremenljivki nekorelirani, je enostavno

$$\text{D}(X + Y) = \text{DX} + \text{DY}.$$

Zvezo lahko posplošimo na

$$\text{D}\left(\sum_{i=1}^n X_i\right) = \sum_{i=1}^n \text{DX}_i + \sum_{i \neq j} \text{Cov}(X_i, X_j)$$

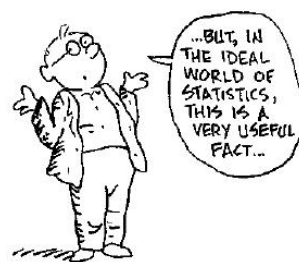
in za paroma nekorelirane spremenljivke

$$\text{D}\left(\sum_{i=1}^n X_i\right) = \sum_{i=1}^n \text{DX}_i.$$

Korelacijski koeficient

Korelacijski koeficient slučajnih spremenljivk X in Y je definiran z izrazom

$$r(X, Y) = \frac{\text{Cov}(X, Y)}{\sigma_X \sigma_Y} = \frac{\text{E}((X - \text{EX})(Y - \text{EY}))}{\sigma_X \sigma_Y}.$$



Za $(X, Y) : N(\mu_x, \mu_y, \sigma_x, \sigma_y, \rho)$ je $r(X, Y) = \rho$.

Primer: Torej sta normalno porazdeljeni slučajni spremenljivki X in Y neodvisni natanko takrat, ko sta nekorelirani. \diamond

Velja še:

$$-1 \leq r(X, Y) \leq 1.$$

$r(X, Y) = 0$ natanko takrat, ko sta X in Y nekorelirani.

$r(X, Y) = 1$ natanko takrat, ko je $Y = \frac{\sigma_Y}{\sigma_X}(X - EX) + EY$ z verjetnostjo 1;

$r(X, Y) = -1$ natanko takrat, ko je $Y = -\frac{\sigma_Y}{\sigma_X}(X - EX) + EY$ z verjetnostjo 1.

Torej, če je $|r(X, Y)| = 1$, obstaja med X in Y linearna zveza z verjetnostjo 1.

8.5 Pogojno matematično upanje

Pogojno matematično upanje je matematično upanje pogojne porazdelitve:

Diskretna slučajna spremenljivka X ima pri pogoju $Y = y_k$ pogojno verjetnostno funkcijo $p_{i|k} = p_{ik}/q_k$, $i = 1, 2, \dots$ in potemtakem pogojno matematično upanje

$$E(X|y_k) = \sum_{i=1}^{\infty} x_i p_{i|k} = \frac{1}{q_k} \sum_{i=1}^{\infty} x_i p_{ik}.$$

Slučajna spremenljivka

$$E(X|Y) : \begin{pmatrix} E(X|y_1) & E(X|y_2) & \cdots \\ q_1 & q_2 & \cdots \end{pmatrix}$$

ima enako matematično upanje kot spremenljivka X :

$$E(E(X|Y)) = \sum_{k=1}^{\infty} q_k E(X|y_k) = \sum_{k=1}^{\infty} \sum_{i=1}^{\infty} x_i p_{ik} = \sum_{i=1}^{\infty} x_i \sum_{k=1}^{\infty} p_{ik} = \sum_{i=1}^{\infty} x_i p_i = EX.$$

Pogojno matematično upanje zvezne spremenljivke

Zvezna slučajna spremenljivka X ima pri pogoju $Y = y$ ima pogojno verjetnostno gostoto $p(x|y) = p(x, y)/p_Y(y)$, $x \in \mathbb{R}$ in potemtakem pogojno matematično upanje

$$E(X|y) = \int_{-\infty}^{\infty} xp(x|y) dx = \frac{1}{p_Y(y)} \int_{-\infty}^{\infty} xp(x, y) dx.$$

Slučajna spremenljivka $E(X|Y)$ z gostoto $p_Y(y)$ ima enako matematično upanje kot spremenljivka X

$$E(E(X|Y)) = \int_{-\infty}^{\infty} E(X|y)p_Y(y) dy = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} xp(x,y) dx dy = \int_{-\infty}^{\infty} xp_X(x)dx = EX.$$

Regresijska funkcija

Preslikavo $x \mapsto E(Y|x)$ imenujemo **regresija** slučajne spremenljivke Y glede na slučajno spremenljivko X .

Primer: Naj bo $(X, Y) : N(\mu_x, \mu_y, \sigma_x, \sigma_y, \rho)$. Tedaj je, kot vemo

$$p_X(x|y) : N\left(\mu_x + \rho \frac{\sigma_x}{\sigma_y}(y - \mu_y), \sigma_x \sqrt{1 - \rho^2}\right).$$

Torej je pogojno matematično upanje

$$E(X|y) = \mu_x + \rho \frac{\sigma_x}{\sigma_y}(y - \mu_y)$$

in prirejena spremenljivka

$$E(X|Y) = \mu_x + \rho \frac{\sigma_x}{\sigma_y}(Y - \mu_y).$$

Na podoben način vpeljemo regresijo slučajne spremenljivke X glede na slučajno spremenljivko Y . Za dvorazsežno normalno porazdelitev dobimo

$$E(Y|X) = \mu_y + \rho \frac{\sigma_y}{\sigma_x}(X - \mu_x).$$

Obe regresijski funkciji sta **linearni**.

Kovariančna matrika

Matematično upanje slučajnega vektorja $\mathbf{X} = (X_1, X_2, \dots, X_n)$ je vektor $E\mathbf{X} = (EX_1, EX_2, \dots, EX_n)$.

Primer: Za $(X, Y) : N(\mu_x, \mu_y, \sigma_x, \sigma_y, \rho)$ je $E(X, Y) = (\mu_x, \mu_y)$. ◇

Matematično upanje slučajne spremenljivke Y , ki je linearna kombinacija spremenljivk X_1, X_2, \dots, X_n , je potem

$$EY = E\left(\sum_{i=1}^n a_i X_i\right) = \sum_{i=1}^n a_i EX_i.$$

Za disperzijo spremenljivke Y pa dobimo $DY = E(Y - EY)^2 =$

$$E\left(\sum_{i=1}^n \sum_{j=1}^n a_i a_j (X_i - EX_i)(X_j - EX_j)\right) = \sum_{i=1}^n \sum_{j=1}^n a_i a_j \text{Cov}(X_i, X_j) = \mathbf{a}^T \mathbf{K} \mathbf{a},$$

kjer je $\text{Cov}(X_i, X_j) = E((X_i - EX_i)(X_j - EX_j))$ kovarianca spremenljivk X_i in X_j , $\mathbf{K} = [\text{Cov}(X_i, X_j)]$ **kovariančna matrika** vektorja X , ter $\mathbf{a} = (a_1, a_2, \dots, a_n)^T$.

Lastnosti kovariančne matrike

Kovariančna matrika $\mathbf{K} = [K_{ij}]$ je *simetrična*: $K_{ij} = K_{ji}$. Diagonalne vrednosti so disperzije spremenljivk: $K_{ii} = DX_i$. Ker je $\mathbf{a}^T \mathbf{K} \mathbf{a} = DY \geq 0$, je pozitivno semidefinitna matrika. Naj bo \mathbf{a} , $\|\mathbf{a}\| = 1$ lastni vektor, ki pripada lastni vrednosti λ kovariančne matrike \mathbf{K} , tj. $\mathbf{K} \mathbf{a} = \lambda \mathbf{a}$. Tedaj je $0 \leq DY = \mathbf{a}^T \mathbf{K} \mathbf{a} = \lambda$, kar pomeni, da so vse lastne vrednosti kovariančne matrike nenegativne. Če je kaka lastna vrednost enaka 0, je vsa verjetnost skoncentrirana na neki hiperravnini – porazdelitev je *izrojena*. To se zgodi natanko takrat, ko kovariančna matrika \mathbf{K} ni obrnljiva, oziroma ko je $\det \mathbf{K} = 0$.

Primer: Za $(X, Y) : N(\mu_x, \mu_y, \sigma_x, \sigma_y, \rho)$ je $\mathbf{K} = \begin{bmatrix} \sigma_x^2 & \rho\sigma_x\sigma_y \\ \rho\sigma_x\sigma_y & \sigma_y^2 \end{bmatrix}$.

Ker je $|\rho| < 1$, je $\det \mathbf{K} = \sigma_x^2 \sigma_y^2 (1 - \rho^2) > 0$ in je potemtakem porazdelitev vedno neizrojena. Za $N(\boldsymbol{\mu}, \mathbf{A})$ je $\mathbf{K} = \mathbf{A}^{-1}$. \diamond

Poglejmo še, kako se spremeni kovariančna matrika pri linearni transformaciji vektorja $X' = \mathbf{A}X$, kjer je \mathbf{A} poljubna matrika reda $n \times n$. Vemo, da je $D(\mathbf{a}^T X) = \mathbf{a}^T \mathbf{K} \mathbf{a}$. Tedaj je, če označimo kovariančno matriko vektorja X' s \mathbf{K}' ,

$$\mathbf{a}^T \mathbf{K}' \mathbf{a} = D(\mathbf{a}^T X') = D(\mathbf{a}^T \mathbf{A}X) = D((\mathbf{A}^T \mathbf{a})^T X) = (\mathbf{A}^T \mathbf{a})^T \mathbf{K} (\mathbf{A}^T \mathbf{a}) = \mathbf{a}^T \mathbf{A} \mathbf{K} \mathbf{A}^T \mathbf{a}$$

in potemtakem

$$\mathbf{K}' = \mathbf{A} \mathbf{K} \mathbf{A}^T.$$

8.6 Višji momenti

Višji momenti so posplošitev pojmov matematičnega upanja in disperzije. **Moment reda** $k \in \mathbb{N}$ *glede na točko* $a \in \mathbb{R}$ imenujemo količino

$$m_k(a) = E((X - a)^k).$$

Moment obstaja, če obstaja matematično upanje $E(|X - a|^k) < \infty$. Za $a = 0$ dobimo **začetni moment** $z_k = m_k(0)$; za $a = EX$ pa **centralni moment** $m_k = m_k(EX)$.

Primer: $EX = z_1$ in $DX = m_2$. ◇

Če obstaja moment $m_n(a)$, potem obstajajo tudi vsi momenti $m_k(a)$ za $k < n$. Če obstaja moment z_n , obstaja tudi moment $m_n(a)$ za vse $a \in \mathbb{R}$.

$$m_n(a) = E((X - a)^n) = \sum_{k=0}^n \binom{n}{k} (-a)^{n-k} z_k.$$

Posebej za centralni moment velja

$$m_n = m_n(z_1) = \sum_{k=0}^n \binom{n}{k} (-z_1)^k z_{n-k}$$

$$m_0 = 1, m_1 = 0, m_2 = z_2 - z_1^2, m_3 = z_3 - 3z_2z_1 + 2z_1^3, \dots$$

Asimetrija spremenljivke X imenujemo količino $A(X) = \frac{m_3}{\sigma^3}$.

Sploščenost spremenljivke X imenujemo količino $K(X) = \frac{m_4}{\sigma^4} - 3$, kjer je $\sigma = \sqrt{m_2}$.

Za simetrično glede na $z_1 = EX$ porazdeljene spremenljivke so vsi lihi centralni momenti enaki 0.

Primer: Za $X : N(\mu, \sigma)$ so $m_{2k+1} = 0$ in $m_{2k} = (2k - 1)!!\sigma^{2k}$. Zato sta tudi $A(X) = 0$ in $K(X) = 0$. ◇

Če sta spremenljivki X in Y neodvisni, je $m_3(X + Y) = m_3(X) + m_3(Y)$.

Primer: Za binomsko porazdeljeno spremenljivko $X : B(n, p)$ pa je

$$m_3(X) = npq(q - p) \quad \text{in dalje} \quad A(X) = \frac{q - p}{\sqrt{npq}}. \quad \diamond$$

Kadar spremenljivka nima momentov, uporabljamo kvantile. **Kvantil reda** $p \in (0, 1)$ je vsaka vrednost $x \in \mathbb{R}$, za katero velja $P(X \leq x) \geq p$ in $P(X \geq x) \geq 1 - p$ oziroma $F(x) \leq p \leq F(x+)$. Kvantil reda p označimo z x_p . Za zvezno spremenljivko je $F(x_p) = p$. Kvantil $x_{\frac{1}{2}}$ imenujemo **mediana**; $x_{\frac{i}{4}}$, $i = 0, 1, 2, 3, 4$ so **kvartilni**. Kot nadomestek za standardni odklon uporabljamo **kvartilni razmik**

$$\frac{1}{2}(x_{\frac{3}{4}} - x_{\frac{1}{4}}).$$

Poglavje 9

Karakteristične funkcije in limitni izreki



9.1 Karakteristična funkcija

Naj bo Z kompleksna slučajna spremenljivka, tj. $Z = X + iY$ za slučajni spremenljivki X in Y . Njeno upanje izračunamo z

$$E(Z) = E(X) + iE(Y),$$

disperzijo pa z

$$D(Z) = E(|Z - E(Z)|^2) = D(X) + D(Y),$$

Kompleksna funkcija realne slučajne spremenljivke je kompleksna slučajna spremenljivka, npr. e^{iX} .

Karakteristična funkcija realne slučajne spremenljivke X je kompleksna funkcija $\varphi_X(t)$ realne spremenljivke t določena z zvezo $\varphi_X(t) = Ee^{itX}$. Karakteristične funkcije

vedno obstajajo in so močno računsko orodje. Posebej pomembni lastnosti sta:

Če obstaja začetni moment z_n , je karakteristična funkcija n -krat odvedljiva v vsaki točki in velja $\varphi_X^{(k)}(0) = i^k z_k$. Za neodvisni spremenljivki X in Y je $\varphi_{X+Y}(t) = \varphi_X(t)\varphi_Y(t)$.

Pojem karakteristične funkcije lahko posplošimo tudi na slučajne vektorje.

Reprodukcijska lastnost normalne porazdelitve

Vsaka linearna kombinacija *neodvisnih in normalno porazdeljenih slučajnih spremenljivk* je tudi sama **normalno** porazdeljena.

Če so slučajne spremenljivke X_1, \dots, X_n neodvisne in normalno porazdeljene $N(\mu_i, \sigma_i)$, potem je njihova vsota tudi normalno porazdeljena:

$$N\left(\sum \mu_i, \sqrt{\sum \sigma_i^2}\right).$$

Da ne bi vsota povprečij rastle z n , nadomestimo vsoto spremenljivk X_i z njihovim povprečjem \bar{X} in dobimo

$$N\left(\bar{\mu}, \sqrt{\sum \left(\frac{\sigma_i}{n}\right)^2}\right).$$

Če privzamemo $\mu_i = \mu$ in $\sigma_i = \sigma$, dobimo $N(\mu, \sigma/\sqrt{n})$.

9.2 Limitni izreki

Zaporedje slučajnih spremenljivk X_n **verjetnostno konvergira** k slučajni spremenljivki X , če za vsak $\varepsilon > 0$ velja

$$\lim_{n \rightarrow \infty} P(|X_n - X| \geq \varepsilon) = 0$$

ali enakovredno

$$\lim_{n \rightarrow \infty} P(|X_n - X| < \varepsilon) = 1.$$

Zaporedje slučajnih spremenljivk X_n **skoraj gotovo konvergira** k slučajni spremenljivki X , če velja

$$P(\lim_{n \rightarrow \infty} X_n = X) = 1.$$

Če zaporedje slučajnih spremenljivk X_n skoraj gotovo konvergira k slučajni spremenljivki X , potem za vsak $\varepsilon > 0$ velja

$$\lim_{m \rightarrow \infty} P(|X_n - X| < \varepsilon \quad \text{za vsak } n \geq m) = 1.$$

Od tu izhaja:

če konvergira skoraj gotovo $X_n \rightarrow X$,
potem konvergira tudi verjetnostno $X_n \rightarrow X$.

Šibki in krepki zakon velikih števil

Naj bo X_1, \dots, X_n zaporedje spremenljivk, ki imajo matematično upanje.

Označimo $S_n = \sum_{k=1}^n X_k$ in

$$Y_n = \frac{S_n - \mathbf{E}S_n}{n} = \frac{1}{n} \sum_{k=1}^n (X_k - \mathbf{E}X_k) = \frac{1}{n} \sum_{k=1}^n X_k - \frac{1}{n} \sum_{k=1}^n \mathbf{E}X_k.$$

Pravimo, da za zaporedje slučajnih spremenljivk X_k velja:

- **šibki zakon velikih števil**, če gre verjetnostno $Y_n \rightarrow 0$, tj., če $\forall \varepsilon > 0$ velja

$$\lim_{n \rightarrow \infty} P\left(\left|\frac{S_n - \mathbf{E}S_n}{n}\right| < \varepsilon\right) = 1;$$

- **krepki zakon velikih števil**, če gre skoraj gotovo $Y_n \rightarrow 0$, tj., če velja

$$P\left(\lim_{n \rightarrow \infty} \frac{S_n - \mathbf{E}S_n}{n} = 0\right) = 1.$$

Če za zaporedje X_1, \dots, X_n velja krepki zakon, velja tudi šibki.

Neenakost Čebiševa

Če ima slučajna spremenljivka X končno disperzijo, tj. $\mathbf{D}X < \infty$,

velja za vsak $\varepsilon > 0$ **neenakost Čebiševa**

$$P(|X - \mathbf{E}X| \geq \varepsilon) \leq \frac{\mathbf{D}X}{\varepsilon^2}.$$



Dokaz: Pokažimo jo za zvezne spremenljivke

$$\begin{aligned} P(|X - \mathbf{E}X| \geq \varepsilon) &= \int_{|x - \mathbf{E}X| \geq \varepsilon} p(x) dx = \frac{1}{\varepsilon^2} \int_{|x - \mathbf{E}X| \geq \varepsilon} \varepsilon^2 p(x) dx \\ &\leq \frac{1}{\varepsilon^2} \int_{-\infty}^{\infty} (x - \mathbf{E}X)^2 p(x) dx = \frac{\mathbf{D}X}{\varepsilon^2}. \quad \square \end{aligned}$$

Neenakost Čebiševa – posledice

Izrek 9.1. (*Markov*) Če gre za zaporedje slučajnih spremenljivk X_i izraz

$$\frac{DS_n}{n^2} \rightarrow 0,$$

ko gre $n \rightarrow \infty$, velja za zaporedje šibki zakon velikih števil.

Izrek 9.2. (*Čebišev*) Če so slučajne spremenljivke X_i paroma nekorelirane in so vse njihove disperzije omejene z isto konstanto C , tj.

$$DX_i < C \quad \text{za vsak } i,$$

velja za zaporedje šibki zakon velikih števil.

Dokaz Bernoullijevega izreka

Za Bernoullijevo zaporedje X_i so spremenljivke paroma neodvisne, $DX_i = pq$, $S_n = k$. Pogoji izreka Čebiševa so izpolnjeni in dobimo:

Izrek 9.3. (*Bernoulli 1713*) Za vsak $\varepsilon > 0$ velja

$$\lim_{n \rightarrow \infty} P\left(\left|\frac{k}{n} - p\right| < \varepsilon\right) = 1.$$

Še nekaj izrekov

Izrek 9.4. (*Hinčin*) Če so neodvisne slučajne spremenljivke X_i enako porazdeljene in imajo matematično upanje $EX_i = a$ za vsak i , potem velja zanje šibki zakon velikih števil, tj. za vsak $\varepsilon > 0$ je

$$\lim_{n \rightarrow \infty} P\left(\left|\frac{S_n}{n} - a\right| < \varepsilon\right) = 1.$$

Izrek 9.5. (*Kolmogorov*) Če so slučajne spremenljivke X_i neodvisne, imajo končno disperzijo in velja $\sum_{n=1}^{\infty} \frac{DS_n}{n^2} < \infty$, potem velja krepki zakon velikih števil:

$$P\left(\lim_{n \rightarrow \infty} \frac{S_n - ES_n}{n} = 0\right) = 1.$$

Izrek 9.6. (Kolmogorov) Če so slučajne spremenljivke X_i neodvisne, enako porazdeljene in imajo matematično upanje $\mathbf{E}X_i = \mu$, potem velja krepki zakon velikih števil

$$P\left(\lim_{n \rightarrow \infty} \frac{1}{n} \sum_{i=1}^n X_i = \mu\right) = 1.$$

Izrek 9.7. (Borel 1909) Za Bernoullijevo zaporedje velja

$$P\left(\lim_{n \rightarrow \infty} \frac{k}{n} = p\right) = 1.$$

9.3 Centralni limitni izrek (CLI)



Leta 1810 je Pierre Laplace (1749-1827) študiral anomalije orbit Jupitra in Saturna, ko je izpeljal razširitev De Moivrevega limitnega izreka,

“Vsaka vsota ali povprečje, če je število členov dovolj veliko, je približno normalno porazdeljena.”

Centralni limitni zakon

Opazujmo sedaj zaporedje standardiziranih spremenljivk

$$Z_n = \frac{S_n - \mathbf{E}S_n}{\sigma(S_n)}.$$

Za zaporedje slučajnih spremenljivk X_i velja **centralni limitni zakon**, če porazdelitvene funkcije za Z_n gredo proti porazdelitveni funkciji standardizirane normalne porazdelitve, to je, če za vsak $x \in \mathbb{R}$ velja

$$\lim_{n \rightarrow \infty} P\left(\frac{S_n - \mathbb{E}S_n}{\sigma(S_n)} < x\right) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^x e^{-\frac{t^2}{2}} dt.$$

(Osnovni CLI) Če so slučajne spremenljivke X_i neodvisne, enako porazdeljene s končnim matematičnim upanjem in končno disperzijo, potem zanje velja centralni limitni zakon.

Skica dokaz centralnega limitnega izreka

Naj bo $Z_i = \frac{X_i - \mu}{\sigma}$. Potem je

$$M_Z(t) = 1 - \frac{t^2}{2!} + \frac{t^3}{3!}E(Z_i^3) + \dots$$

Za $Y_n = \frac{\bar{X} - \mu}{\sigma/\sqrt{n}} = \frac{1}{\sqrt{n}} \frac{\sum_{i=1}^n X_i - n\mu}{\sigma} = \frac{1}{\sqrt{n}} \sum_{i=1}^n Z_i$ velja

$$M_n(t) = \left[M_Z\left(\frac{t}{\sqrt{n}}\right) \right]^n = \left(1 - \frac{t^2}{2n} + \frac{t^3}{3!n^{3/2}}k + \dots \right)^n,$$

kjer je $k = E(Z_i^3)$.

$$\log M_n(t) = n \log \left(1 - \frac{t^2}{2n} + \frac{t^3}{3!n^{3/2}}k + \dots \right)$$

Za $x = \left(-\frac{t^2}{2n} + \frac{t^3}{3!n^{3/2}} + \dots \right)$ velja

$$\begin{aligned} \log M_n(t) &= n \log(1+x) = n \left(x - \frac{x^2}{2} + \dots \right) = \\ &= n \left[\left(-\frac{t^2}{2n} + \frac{t^3}{3!n^{3/2}} + \dots \right) - \frac{1}{2} \left(-\frac{t^2}{2n} + \frac{t^3}{3!n^{3/2}} + \dots \right)^2 + \dots \right] \end{aligned}$$

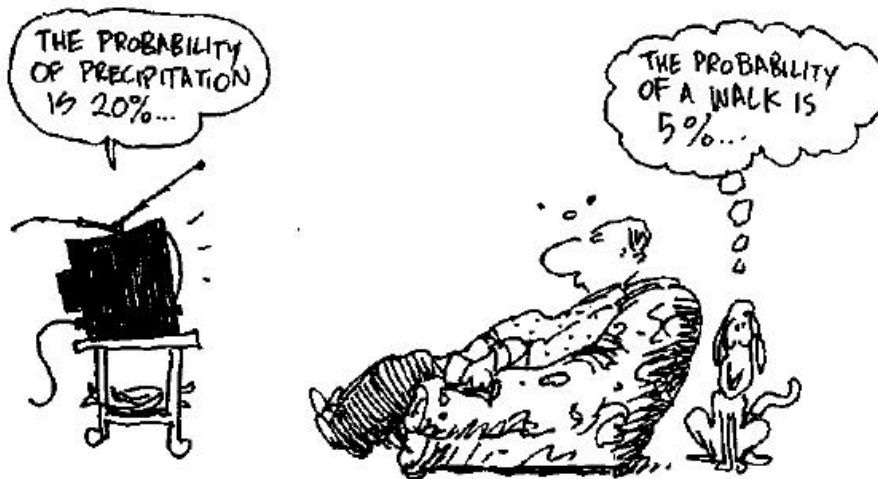
in od tod končno še

$$\lim_{n \rightarrow \infty} \log M_n(t) = -\frac{t^2}{2} \quad \text{oziroma} \quad \lim_{n \rightarrow \infty} M_n(t) = e^{-t^2/2}.$$

Iz konvergence karakterističnih funkcij φ_{Y_n} proti karakteristični funkciji standardizirane normalne porazdelitve lahko sklepamo po obratnem konvergenčnem izreku, da tudi porazdelitvene funkcije za Y_n konvergirajo proti porazdelitveni funkciji standardizirane normalne porazdelitve. Torej velja centralni limitni zakon. \square

Poglavje 10

Nekaj primerov uporabe



10.1 Zamenjalna šifra

Tomaž Pisanski, Skrivnostno sporočilo, *Presek* V/1, 1977/78, str. 40-42.

YHW?HD+CVODHVTHVO-!JV G: CDCYJ (JV/-V?HV (-T?HVW-4YC4 (?-DJV/- (?S-V03CWC%J (-V4-DC
V!CW-?CVNJDJVD-?+-V03CWC%J (-VQW-DQ-VJ+V?HVDWHN-V3C: CODCV!H+?-DJVD-?+CV3JO-YC

(črko Č smo zamenjali s C, črko Ć pa z D). Imamo $26! = 40329146112665635584000000$ možnosti z direktnim preizkušanjem, zato v članku dobimo naslednje nasvete:

(0) Relativna frekvenca črk in presledkov v slovenščini: presledek 173,

E	A	I	O	N	R	S	L	J	T	V	D	K	M	P	U	Z	B	G	Č	H	Š	C	Ž	F
89	84	74	73	57	44	43	39	37	37	33	30	29	27	26	18	17	15	12	12	9	9	6	6	1

- (1) Na začetku besed so najpogostejše črke N, S, K, T, J, L.
- (2) Najpogostejše končnice pa so E, A, I, O, U, R, N.
- (3) Ugotovi, kateri znaki zagotovo predstavljajo samoglasnike in kateri soglasnike.
- (4) V vsaki besedi je vsaj en samoglasnik ali samoglasniški R.
- (5) V vsaki besedi z dvema črkama je ena črka samoglasnik, druga pa soglasnik.
- (6) detektivska sreča

Pa začnimo z reševanjem (oziroma kakor pravijo kriptografi: z razbijanjem):

(0) V - C D J ? H W O (+ 3 Y 4 ! / Q : % T N S G
 23 19 16 12 11 10 9 7 6 6 5 4 4 3 3 2 2 2 2 2 2 2 1 1

Zaključek V --> ' ' (drugi znaki z visoko frekvenco ne morejo biti). Dve besedi se ponovita: 03CWC%J(-, opazimo pa tudi eno sklanjatev: D-?+- ter D-?+C. Torej nadaljujemo z naslednjim tekstom:

YHW?HD+C ODH TH O-!J G:CDYJ(J /- ?H (-T?H W-4YD4(?-DJ /-(?S- 03CWC%J(- 4-DC
 !CW-?C NJDJ D-?+- 03CWC%J(- QW-DQ- J+ ?H DWHN- 3C:CODC !H+?-DJ D-?+C 3JO-YC

- (3) Kandidati za samoglasnike e,a,i,o so znaki z visokimi frekvencami. Vzamemo:

$$\{e,a,i,o\} = \{-,C,J,H\}$$

(saj D izključi -,H,J,C in ? izključi -,H,C, znaki -,C,J,H pa se ne izključujejo)

Razporeditev teh znakov kot samoglasnikov izgleda prav verjetna. To potrdi tudi gostota končnic, gostota parov je namreč:

AV CV HV JV VO ?H -D DC JM W- DJ UC CW -? VD
 7 5 5 5 4 4 4 3 3 3 3 3 3 3 3

- (5) Preučimo besede z dvema črkama:

Samoglasnik na koncu

Samoglasnik na začetku

- | | | |
|---------------------------------------|----------------------|------------|
| 1) da ga na pa ta za (ha ja la) | 1) ar as | (ah aj au) |
| 2) če je le me ne se še te ve že (he) | 2) en ep | (ej eh) |
| 3) bi ji ki mi ni si ti vi | 3) in iz ig | |
| 4) bo do (ho) jo ko no po so to | 4) on ob od os on | (oh oj) |
| 5) ju mu tu (bu) | 5) uk up uš ud um ur | (uh ut) |
| 6) rž rt | | |

in opazujemo besedi: /- ?H ter besedi: J+ ?H. J+ ima najmanj možnosti, + pa verjetno ni črka n, zato nam ostane samo še:

J+ ?H	DWHN-
/- ?H	
iz te	(ne gre zaradi: D-?+C)
ob ta(e,o)	(ne gre zaradi: D-?+C)
od te	(ne gre zaradi: D-?+C)

tako da bo potrebno nekaj spremeniti in preizkusiti še naslednje: on bo; on jo; in so; in se; in je; in ta; en je; od tu ...

(6) Če nam po dolgem premisleku ne uspe najti rdeče niti, bo morda potrebno iskati napako s prijatelji (tudi računalniški program z metodo lokalne optimizacije ni zmožal problema zaradi premajhne dolžine tajnopisa, vsekakor pa bi bilo problem mogoče rešiti s pomočjo elektronskega slovarja). Tudi psihološki pristop pomaga, je svetoval Martin Juvan in naloga je bila rešena (poskusite sami!).

Kaj pa tuji jeziki

Podobna naloga je v angleščini dosti lažja, saj je v tem jeziku veliko členov THE, A in AN, vendar pa zato običajno najprej izpustimo presledke iz teksta, ki ga želimo spraviti v tajnopis. V angleščini imajo seveda črke drugačno gostoto kot v slovenščini. Razdelimo jih v naslednjih pet skupin:

1. E, z verjetnostjo okoli 0,120,
2. T, A, O, I, N, S, H, R, vse z verjetnostjo med 0,06 in 0,09,
3. D, L, obe z verjetnostjo okoli 0,04,
4. C, U, M, W, F, G, Y, P, B, vse z verjetnostjo med 0,015 in 0,028,
5. V, K, J, X, Q, Z, vse z verjetnostjo manjšo od 0,01.

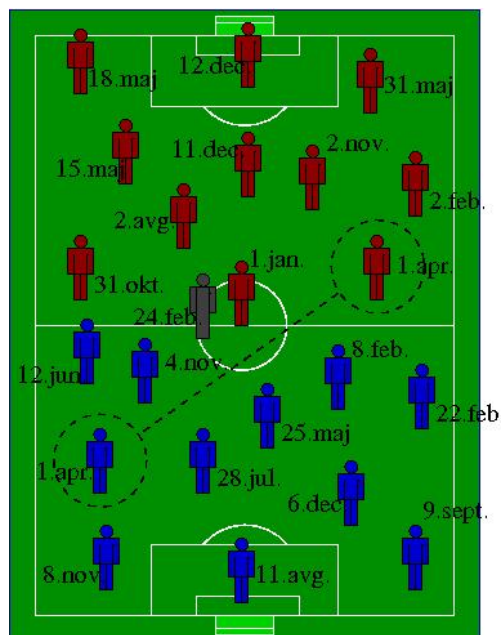
Najbolj pogosti pari so (v padajočem zaporedju): TH, HE, IN, ER, AN, RE, ED, ON, ES, ST, EN, AT, TO, NT, HA, ND, OU, EA, NG, AS, OR, TI, IS, ET, IT, AR, TE, SE, HI in OF. Najbolj pogoste trojice pa so (v padajočem zaporedju): THE, ING, AND, HER, ERE, ENT, THA, NTH, WAS, ETH, FOR in DTH.

10.2 Kakšno naključje!!! Mar res?

Na nogometni tekmi sta
na igrišču dve enajsterici
in sodnik, skupaj
23 oseb.

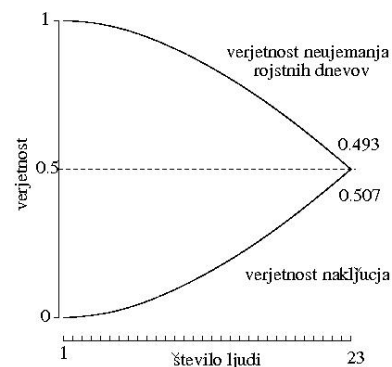
Kakšna je verjetnost,
da imata **dve osebi**
isti rojstni dan?

Ali je ta verjetnost lahko večja od **0,5**?



Ko vstopi v sobo k -ta oseba, je verjetnost, da je vseh k rojstnih dnevov različnih enaka:

$$\frac{365}{365} \times \frac{364}{365} \times \frac{363}{365} \times \dots \times \frac{365 - k + 1}{365} = \begin{cases} 0,493; & \text{če je } k=22 \\ 0,507; & \text{če je } k=23 \end{cases}$$



**V poljubni skupini 23-ih ljudi je verjetnost,
da imata vsaj dva skupni rojstni dan $> 1/2$.**

Čeprav je 23 majhno število, je med 23 osebami 253 različnih parov. To število je veliko bolj povezano z iskano verjetnostjo. Testirajte to na zabavah z več kot 23 osebami. Organizirajte stave in dolgoročno boste gotovo na boljšem, na velikih zabavah pa boste zlahka zmagovali.

Napad s pomočjo paradoksa rojstnih dnevov (angl. *Birthday Attack*)

To seveda ni paradoks, a vseeno ponavadi zavede naš občutek.

Ocenimo še splošno verjetnost. Mečemo k žogic v n posod in gledamo, ali sta v kakšni posodi vsaj dve žogici. Poiščimo spodnjo mejo za verjetnost zgoraj opisanega dogodka:

$$\left(1 - \frac{1}{n}\right) \left(1 - \frac{2}{n}\right) \cdots \left(1 - \frac{k-1}{n}\right) = \prod_{i=1}^{k-1} \left(1 - \frac{i}{n}\right)$$

Iz Taylorjeve vrste

$$e^{-x} = 1 - x + \frac{x^2}{2!} - \frac{x^3}{3!} + \cdots$$

ocenimo $1 - x \approx e^{-x}$ in dobimo

$$\prod_{i=1}^{k-1} \left(1 - \frac{i}{n}\right) \approx \prod_{i=1}^{k-1} e^{-\frac{i}{n}} = e^{-\frac{k(k-1)}{2n}}.$$

Torej je verjetnost trčenja

$$1 - e^{-\frac{k(k-1)}{2n}}.$$

Potem velja

$$e^{-\frac{k(k-1)}{2n}} \approx 1 - \varepsilon$$

oziroma

$$\frac{-k(k-1)}{2n} \approx \log(1 - \varepsilon), \quad \text{tj.} \quad k^2 - k \approx 2n \log \frac{1}{1 - \varepsilon}$$

in če ignoriramo $-k$, dobimo končno

$$k \approx \sqrt{2n \log \frac{1}{1 - \varepsilon}}.$$

Za $\varepsilon = 0,5$ je

$$k \approx 1,17\sqrt{n},$$

kar pomeni, da, če zgostimo nekaj več kot \sqrt{n} elementov, je bolj verjetno, da pride do trčenja kot da ne pride do trčenja. **V splošnem je k proporcionalen s \sqrt{n} .**

Raba v kriptografiji

Napad s pomočjo paradoksa rojstnih dnevov s tem določi spodnjo mejo za velikost zaloge vrednosti zgoščevalnih funkcij, ki jih uporabljamo v kriptografiji in računalniški varnosti. 40-bitna zgostitev ne bi bila varna, saj bi prišli do trčenja z nekaj več kot 2^{20} (se pravi milijon) naključnimi zgostitvami z verjetnostjo vsaj $1/2$. V praksi je priporočena najmanj 128-bitna zgostitev in standard za shema digitalnega podpisa (160 bitov) to vsekakor upošteva. Podobno si lahko pomagamo tudi pri napadih na DLP in še kje.

10.3 Ramseyjeva teorija

- intuitivna ideja
- Ramseyjev izrek
- Erdösev izrek
- primeri uporabe



Po 3,500 let starem zapisu je antični sumerski učenjak pogledal v nebo in zagledal leva, bika in škorpijona. **Ali gre za kozmične sile?** Astronom bi rekel: kolekcija zvezd, tj. začasna konfiguracija zvezd, ki jo gledamo z roba navadne galaksije.

1928 Frank Plumpton Ramsey (26 let, angleški matematik, filozof in ekonomist)

Popoln nered je nemogoč.

Ramseyjeva teorija: Vsaka dovolj velika struktura vsebuje urejeno podstrukturo. Konkretna naloga: **Koliko objektov nam zagotavlja željeno podstrukturo?**

Izrek (SIM). *V družbi šestih ljudi obstaja trojica v kateri se vsaka dva poznata ali pa vsaka dva ne poznata.*

- naivni prestop: preverimo $2^{15} = 32.768$ možnosti,
- barvanje povezav polnega grafa K_6 in Dirichletov princip.

Nekaj težja naloga: **V družbi 17ih znanstvenikov se vsaka dva dopisujeta o eni izmed treh tem. Dokaži, da obstajajo trije, ki se dopisujejo o isti temi!**

Ramseyjevo število $r(k, \ell)$ je najmanjše število za katerega vsak graf na $r(k, \ell)$ vozliščih vsebuje bodisi k -kliko bodisi ℓ -antikliko. Prepričaj se, da je $r(k, \ell) = r(\ell, k)$.

Primeri: $r(k, 1) = 1 = r(1, \ell)$, $r(2, \ell) = \ell$, $r(k, 2) = k$, SIM: $r(3, 3) \leq 6$.

Ramseyjev izrek. $\forall k, \ell \in \mathbb{N}$

$$r(k, \ell) \leq r(k, \ell - 1) + r(k - 1, \ell).$$

Če sta obe števili na desni strani neenakosti sodi, potem velja stroga neenakost.

Zgled uporabe: $r(3, 3) \leq r(3, 2) + r(2, 3) = 3 + 3 = 6$.

Dokaz: (1935 Erdős & Szekeres, 1955 Greenwood & Gleason) Naj bo G graf na $r(k, \ell - 1) + r(k - 1, \ell)$ vozliščih. Potem velja ena izmed naslednjih možnosti:

(a) Vozlišče v ni sosednje množici S z vsaj $r(k, \ell - 1)$ vozlišči.

kar pomeni, da $G[S]$ vsebuje ali k -kliko ali $(\ell - 1)$ -antikliko.

(b) Vozlišče v je sosednje množici T z vsaj $r(k - 1, \ell)$ vozlišči.

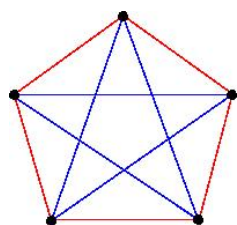
kar pomeni, da $G[T]$ vsebuje ali $(k - 1)$ -kliko ali ℓ -antikliko.

Od tod sledi, da G vsebuje bodisi k -kliko bodisi ℓ -antikliko. Naj bosta $r(k, \ell - 1)$ in $r(k - 1, \ell)$ sodi števili in $|G| = r(k, \ell - 1) + r(k - 1, \ell) - 1$. Potem obstaja vozlišče $v \in V(G)$, katerega stopnja je sodo število. Torej v ni soseden točno $r(k - 1, \ell) - 1$ vozliščem in velja bodisi (a) bodisi (b). \square

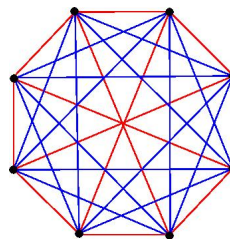
Pokaži:

$$r(3, 4) \leq 9, \quad r(3, 5) \leq 14, \quad r(4, 4) \leq 18, \quad r(k, \ell) \leq \binom{k + \ell - 2}{k - 1}.$$

To je bila zgornja meja. Kaj pa spodnja meja?



$$5 < r(3, 3) = 6$$



$$8 < r(3, 4) = 9$$

Podobno dobimo tudi $13 < r(3, 5) = 14$, $17 < r(3, 6) = 18$,
 $22 < r(3, 7) = 23$, $27 < r(3, 8) \leq 29$ in $35 < r(3, 9) = 36$.

Erdösev Izrek. $\forall k \in \mathbb{N} \quad r(k, k) \geq 2^{k/2}$.

Zgled uporabe: $r(3, 3) \geq 3$ and $r(4, 4) \geq 4$.

Če Marsovci napadejo Zemljo nam morda uspe izračunati $r(5, 5) \in [43, 49]$ (Exoo 1989, McKay and Radziszowski 1995), nikakor pa ne moremo izračunati $r(6, 6) \in [102, 165]$ (Kalbfleisch 1965, Mackey 1994).

Znana Ramseyeva števila:

$k \setminus \ell$	3	4	5	6	7	8	9	10		
3	6	9	14	18	23	28	36	?		
4	9	18	25	?	?	?	?	?	?	?
6	18	?	?	?	?	?	?	?		

[Ester Klein](#) je leta 1933 predstavil naslednjo geometrijsko nalogo:

Med petimi točkami v ravnini, od katerih nobene tri niso kolinearne (ležijo na premici), lahko vedno izberemo štiri, ki določajo konveksen četverkotnik.

Rešitev: Vpeljemo pojem **konveksne ogrinjače** ...

Če je konveksna ogrinjača teh petih točk

- (a) **petkotnik**, potem vsake 4 točke med njimi sestavljajo konveksen četverkotnik,
- (b) **štirikotnik**, potem so njegovi vrhovi tiste 4 točke, ki smo jih iskali,
- (c) **trikotnik**, potem ga lahko označimo z A , B in C , preostali točki pa z D in E , tako da sta točki A in B na isti s strani premice DE .

V tem primeru je četverkotnik $ABCD$ konveksen. □

Nalogo lahko posplošimo na 9 točk in iskanje konveksnega petkotnika ter počasi pridemo do Erdöseve domneve, da za konveksen k -kotnik potrebujemo v ravnini vsaj

$$n = 1 + 2^{k-2}$$

točk od katerih nobene 3 niso kolinearne. Pravzaprav se je najprej Szekeres prepričal, da za dovolj velik n vedno obstaja konveksen k -kotnik, potem pa je Erdös postavil svojo domnevo.

Erdőseva probabilistična metoda (1947)

34 točk določa 561 premic. Da se to zgodi v eni barvi, je verjetnost

$$2^{-561} \approx 2,6 \cdot 10^{-169}.$$

Velja tudi $\binom{1.000.000}{34} = 3,4 \cdot 10^{165}$. Torej lahko pričakujemo $\binom{10^6}{34} = 3,4 \cdot 10^{165} \approx 0,01$ oziroma 0,01% enobarvnih. To pomeni, da v 99,9% ne dobimo enobarvnega K_{34} .

Slednjo idejo pretvorimo v Erdősev dokaz.

Dokaz Erdősevega izreka: Probabilistična metoda (ni konstruktivna) in štetje. Naj bo \mathcal{G}_n množica grafov z vozlišči v_1, v_2, \dots, v_n . Naj bo \mathcal{G}_n^k množica grafov iz \mathcal{G}_n , ki vsebujejo k -kliko. Potem je $|\mathcal{G}_n| = 2^{\binom{n}{2}}$, $|\mathcal{G}_n^k| = 2^{\binom{n}{2} - \binom{k}{2}} \binom{n}{k}$ in

$$q = |\mathcal{G}_n^k|/|\mathcal{G}_n| \leq \frac{n^k 2^{-\binom{k}{2}}}{k!}.$$

Če je $n < 2^{k/2}$, velja $q \leq \frac{2^{\frac{k^2}{2} - \binom{k}{2}}}{k!} < \frac{1}{2}$.

Se pravi, da manj kot polovica grafov iz \mathcal{G}_n vsebuje k -klike.

Iz $\mathcal{G}_n = \{G \mid \bar{G} \in \mathcal{G}_n\}$ pa sledi, da manj kot polovica grafov iz \mathcal{G}_n vsebuje k -antiklike. \square

Posledica. Za $m := \min(k, \ell)$ velja $r(k, \ell) \geq 2^{m/2}$.

Uporaba: Pobarvaj z modro in rdečo števila 1 2 3 4 5 6 7 8 9.

Posledica Ramseyjevega izreka (Waerden 1926):

3 rdeča ali 3 modra števila tvorijo aritmetično zaporedje.

PODVOJ(x) = $2x$, EKSPONENT(x) = 2^x , STOLP(x) = $2^{2^{2^{\dots^2}}}$ (x dvojk)

UAU(1) = STOLP(1)=2. UAU(2) = STOLP(2)=4. UAU(3) = STOLP(4)=65,536

UAU(4) = prevelik za vse knjige, za vse računalnike ..., UAU(x) = ...

Zaporedje 1, 2, ..., ACKERMANN(k) pobarvamo z dvema barvama. Potem obstaja monokromatično (enobarvno) aritmetično podzaporedje s k členi.

Paul Erdős (1913 – 1996)

http://www.maa.org/mathland/mathland_10_7.html

<http://www.britannica.com/EBchecked/topic/191138/Paul-Erdos>: Paul Hoffman

Hungarian “freelance” mathematician (known for his work in number theory and combinatorics) and legendary eccentric who was arguably the most prolific mathematician of the 20th century, in terms of both the number of problems he solved and the number of problems he convinced others to tackle. Erdős did mathematics with a missionary zeal, often 20 hours a day, turning out some 1,500 papers, an order of magnitude higher than his most prolific colleagues produced. He was active to the last days of his life. At least 50 papers on which he is listed as a coauthor are yet to appear, representing the results of various recent collaborative efforts. His interests were mainly in number theory and combinatorics, though they ranged into topology and other areas of mathematics. He was fascinated by relationships among numbers, and numbers served as the raw materials for many of his conjectures, questions, and proofs.

In 1930, at age 17, Erdős entered the Péter Pázmány University in Budapest, where in four years he completed his undergraduate work and earned a Ph.D. in mathematics. As a college freshman (18), he made a name for himself in mathematical circles with a stunningly simple proof of Chebyshev’s theorem, which says that a prime can always be found between any integer n (greater than 1) and its double $2n$. A little later, he proved his own theorem that there is always a prime of the form $4k + 1$ and $4k + 3$ between n and $2n$. For example, the interval between 100 and 200 contains the prime-number pair 101 and 103 ($k = 25$). Paul Erdős has the theory that God has a book containing all the theorems of mathematics with their absolutely most beautiful proofs, and when [Erdős] wants to express particular appreciation of a proof, he exclaims, ‘This is one from the book!’ During his university years he and other young Jewish mathematicians, who called themselves the Anonymous group, championed a fledgling branch of mathematics called Ramsey theory.

In 1934 Erdős, disturbed by the rise of anti-Semitism in Hungary, left the country for a four-year postdoctoral fellowship at the University of Manchester in England. In September 1938 he emigrated to the United States, accepting a one-year appointment at the Institute for Advanced Study in Princeton, New Jersey, where he cofounded the field of probabilistic number theory. During the 1940s he wandered around the United States from one university to the next—Purdue, Stanford, Notre Dame, Johns Hopkins—spurning full-time job offers so that he would have the freedom to work with anyone at any time on any problem of his choice. Thus began half a century of nomadic existence that would make him a legend in the mathematics community. With no home, no wife, and no job to tie him down, his wanderlust took him to Israel, China, Australia, and 22 other countries (although sometimes he was turned away at the border—during the Cold War, Hungary feared he was an American spy, and the United States feared he was a communist spy). Erdős would show up—often unannounced—on the doorstep of a fellow mathematician, declare “My brain is open!” and stay as long as his colleague served up interesting mathematical challenges.

In 1949 Erdős had his most satisfying victory over the prime numbers when he and Atle Selberg gave The Book proof of the prime number theorem (which is a statement about the frequency of primes at larger and larger numbers). In 1951 John von Neumann presented the Cole Prize to Erdős for his work in prime number theory. In 1959 Erdős attended the first International Conference on Graph Theory, a field he helped found. During the next three decades he continued to do important work in combinatorics, partition theory, set theory, number theory, and geometry—the diversity of the fields he worked in was unusual. In 1984 he won the most lucrative award in mathematics, the Wolf Prize, and used all but \$720 of the \$50,000 prize money to establish a scholarship in his parents’ memory in Israel. He was elected to many of the world’s most prestigious scientific societies, including the Hungarian Academy of Science (1956), the U.S. National Academy of Sciences (1979), and the British Royal Society (1989). Defying the conventional wisdom that mathematics was a young man’s game, Erdős went on proving and conjecturing until the age of 83, succumbing to a heart attack only hours after disposing of a nettlesome problem in geometry at a conference in Warsaw.

Erdős had once remarked that mathematics is eternal because it has an infinity of problems. In the same spirit, his own contributions have enriched mathematics. Erdős problems – solved and unsolved – abound in the mathematical literature, lying in wait to provoke thought and elicit surprise.

Erdős loved problems that people could understand without learning a mass of definitions. His hallmark was the deceptively simple, precisely stated problem and the succinct and ingenious argument to settle the issue. Though simply stated, however, his problems were often notoriously difficult to solve. Here's a sample, not-so-difficult Erdős problem that concerns sequences of +1's and -1's. Suppose there are equal numbers of +1's and -1's lined up in a row. If there are two +1's and two -1's, for example, a row could consist of +1 +1 -1 -1. Because these terms can be listed in any order, there are in fact six different ways to write such a row. Of course, the sum of all the numbers in a row is zero. However, it's interesting to look at the partial sums in each row. In the example above, the partial sums are +1 (after one term), +2 (after two terms), +1 (after three terms), and 0 (after four terms). The problem is to determine how many rows out of all the possibilities yield no partial sum that is negative. Of the six different rows for $n = 2$, only two escape a negative partial sum. Of the 20 rows for $n = 3$, just five have exclusively nonnegative partial sums; for $n = 4$, 14 out of 70 rows have this particular characteristic; and so on. The answer turns out to be a sequence called the Catalan numbers: $1/(n + 1)$ times the number of different rows for $n + 1$'s and $n - 1$'s. One can liken these rows to patrons lined up at a theater box office. The price of admission is 50 cents, and half the people have the exact change while the other half have one-dollar bills. Thus, each person provides one unit of change for the cashier's later use or uses up one unit of change. In how many ways can the patrons be lined up so that a cashier, who begins with no money of her own, is never stuck for change?

He turned mathematics into a social activity, encouraging his most hermetic colleagues to work together. The collective goal, he said, was to reveal the pages in the Book. Erdős himself published papers with 507 coauthors. In the mathematics community those 507 people gained the coveted distinction of having an "Erdős number of 1," meaning that they wrote a paper with Erdős himself. Someone who published a paper with one of Erdős's coauthors was said to have an Erdős number of 2, and an Erdős number of 3 meant that someone wrote a paper with someone who wrote a paper with someone who worked with Erdős. Albert Einstein's Erdős number, for instance, was 2. The highest known Erdős number is 15; this excludes nonmathematicians, who all have an Erdős number of infinity.

Erdős enjoyed offering monetary rewards for solving particular problems, ranging from \$10,000 for what he called "a hopeless problem" in number theory to \$25 for something that he considered not particularly difficult but still tricky, proposed in the middle of a lecture. One problem worth a \$3,000 reward concerns an infinite sequence of integers, the sum of whose reciprocals diverges. The conjecture is that such a sequence contains arbitrarily long arithmetic progressions. "This would imply that the primes contain arbitrarily long arithmetic progressions," Erdős remarked. "This would be really nice. And I don't expect to have to pay this money, but I should leave some money for it in case I leave."

Frank Plumpton Ramsey (1903–1930)

The Cambridge philosopher Ramsey, a wunderkind of first order, wrote three important contributions to economics.

The first, "Truth and Probability" (written in 1926, published 1931), was the first paper to lay out the theory of subjective probability and begin to axiomatize choice under (subjective) uncertainty, a task completed decades later by Bruno de Finetti and Leonard Savage. (This was written in opposition to John Maynard Keynes's own information-theoretic *Treatise on Probability*.) Ramsey's second contribution was his theory of taxation (1927), generating the famous "Boiteux-Ramsey" pricing rule. Ramsey's third contribution was his exercise in determining optimal savings (1928), the famous "optimal growth" model - what has since become known as the "Ramsey model" - one of the earliest applications of the calculus of variations to economics.

Frank Ramsey died on January 27, 1930, just before his 27th birthday. In his tragically short life he produced an extraordinary amount of profound and original work in economics, mathematics and logic as well as in philosophy: work which in all these fields is still extremely influential.

10.4 Teorije kodiranja

Claude Shannon je postavil teoretične osnove **teorije informacij** in zanesljivega prenosa digitalnih podatkov kmalu po koncu druge svetovne vojne.



Glavni mejniki teorija kodiranja

- 1947-48:** začetki teorije informacij: znamenita izreka o “**Source Coding**” in pa “**Channel Capacity**” (C. Shannon)
- 1949-50:** odkritje *prvih kod* za odpravljanje napak (M. Golay, R. Hamming).
- 1959-60:** odkritje **BCH-kod** (R. Bose, D. Ray-Chaudhuri, A. Hochquenghem).
- 1967:** Viterby algoritm za odkodiranje **konvolucijskih kod**, (ki sta jih predlagala Elias 1955, Hagelbarger 1959).
- 1993:** razvoj **turbo kod** (C. Berrou, A. Glavieux, P. Titimajshima).

Del II
STATISTIKA



Skozi življenje se prebijamo z odločitvami,
ki jih naredimo na osnovi nepopolnih informacij ...

Pridobivanje podatkov

Novice so polne številčk. Televizijski napovedovalec pove, da se je stopnja nezaposlenosti zmanjšala na 4,7%. Raziskava trdi, da je 45% Američanov zaradi kriminala strah ponoči zapustiti domove. Od kod pridejo te številke? Ne vprašamo vseh ljudi, če so zaposleni ali ne. Raziskovalne agencije vprašajo le nekaj posameznikov, če zaradi strahu pred ropi ostajajo ponoči doma. Vsak dan se v novicah pojavi nov naslov. Eden od teh trdi: Aspirin preprečuje srčne infarkte. Nadaljnje branje razkrije, da je raziskava obravnavala 22 tisoč zdravnikov srednjih let. Polovica zdravnikov je vsak drugi dan vzela aspirin, druga polovica pa je dobila neaktivno tableto. V skupini, ki je jemala aspirin, je 139 zdravnikov doživelo srčni infarkt. V drugi skupini je bilo v enakem časovnem obdobju 239 infarktov. Ali je ta razlika dovolj velika, da lahko trdimo, da aspirin res preprečuje srčne infarkte?

Da bi ubežali neprijetnostim kot sta nezaposlenost in srčni infarkt, prižgimo televizijo. V pogovorni oddaji voditelj povabi gledalce, da sodelujejo v anketi. Tema pogovora je dobrodelnost in voditelja zanima, če gledalci redno prispevajo denar ali oblačila v dobrodelne namene. Med oddajo sprejmejo 50 tisoč klicev in 83% gledalcev trdi, da redno sodelujejo v tovrstnih akcijah. Ali je res, da smo tako zelo humanitarno osveščeni? Zanesljivost teh številčk je v prvi vrsti odvisna od njihovega izvora. Podatkom o nezaposlenosti lahko zaupamo, v tistih 83% iz pogovorne oddaje pa najbrž lahko utemeljeno podvomimo. Naučili se bomo prepoznati dobre in slabe metode pridobivanja podatkov. Razumevanje metod, s katerimi lahko pridobimo zaupanja vredne podatke, je prvi (in najpomembnejši) korak k pridobivanju sposobnosti odločanja o pravilnosti sklepov, ki jih izpeljemo na osnovi danih podatkov. Izpeljava zaupanja vrednih metod za pridobivanje podatkov je področje, kjer vstopimo v svet statistike, znanosti o podatkih.

Obdelava podatkov

Za sodobno družbo je značilna poplava podatkov. Podatki, ali numerična dejstva, so bistveni pri odločanju na skoraj vseh področjih življenja in dela. Kot druge velike poplave nam poplava podatkov grozi, da nas bo pokopala pod sabo. Moramo jo kontrolirati s premišljeno organizacijo in interpretacijo podatkov. Baza podatkov kakšnega podjetja na primer vsebuje velikansko število podatkov: o zaposlenih, prodaji, inventarju, računih strank, opremi, davkih in drugem. Ti podatki so koristni le v primeru, ko jih lahko organiziramo in predstavimo tako, da je njihov pomen jasen. Posledice neupoštevanja podatkov so lahko hude. Veliko bank je izgubilo na milijarde dolarjev pri nedovoljenih špekulacijah njihovih zaposlenih, ki so ostale skrite med goro podatkov, ki jih odgovorni niso dovolj

pozorno pregledali.

Statistično sklepanje

Sklepanje je proces, pri katerem pridemo do zaključkov na podlagi danih dokazov. Dokazi so lahko v mnogo različnih oblikah. V sojenju zaradi umora jih lahko predstavljajo izjave prič, posnetki telefonskih pogovorov, analize DNK iz vzorcev krvi in podobno. Pri statističnem sklepanju nam dokaze priskrbijo podatki. Po domače statistično sklepanje velikokrat temelji na grafični predstavitvi podatkov. Formalno sklepanje, tema tega predmeta, uporablja verjetnost, da pove, do kakšne mere smo lahko prepričani, da so naši zaključki pravilni.

Nekaj statističnih izzivov za začetnike

Trgovec je vašemu podjetju prodal 10.000 sodov rjavega fižola. Cena le-tega je na trgu za 10% višja od sivega (bolj obstojen in večja hranljiva vrednost). Še predno plačamo, odidemo do skladišča in odpremo naključno izban sod, ugotovimo, da je res napolnjen do vrha s fižolom, vendar pa so zrna rjava ali siva. Kako najhitreje ugotovimo, za koliko moramo znižati plačilo, če se odločimo, da bomo fižol vseeno prevzeli?

Dal bi vam toliko "odpustkov", kolikor las imam na glavi. Koliko las pa imamo na glavi?

Napisali smo diplomu, ki je dolga 100 strani, kolega pa ima za 20 strani daljšo diplomu. Če za trenutek pustimo ob strani samo vsebino (kvaliteto), je še vedno vprašanje ali je bil res boljši od nas v kvantiteti. Uporabljal je drugačen font, njegov rob je nekoliko večji,... Kako lahko na hitro ocenimo dejansko stanje (brez da bi primerjali sami datoteki)?

Nadaljujmo z branjem časopisov

Napoved vremena predstavlja naslednje področje statistike za množice, s svojimi napovedmi za dnevne najvišje in najnižje temperature (kako se lahko odločijo za 10 stopinj ne pa za 9 stopinj?). (Kako pridejo do teh števil? Z jemanjem vzorcev? Koliko vzorcev morajo zbrati in kje jih zbirajo? Najdete tudi napovedi za 3 dni naprej, morda celo teden, mesec in leto! Kako natančne so vremenske napovedi v današnjem času? Glede na to kolikokrat nas je ujel dež, ali kolikokrat so napovedali sonce, lahko zaključite, da morajo nadaljevati z raziskovanjem na tem področju. Verjetnost in računalniško modeliranje igra pomembno vlogo pri napovedovanju vremena. Posebej uspešni so pri večjih dogodkih kot

so orkani, potresi in vulkanski izbruhi. Seveda pa so računalniki le tako pametni kot ljudje, ki so napisali programske opreme, ki jih poganja. Raziskovalci bodo imeli še veliko dela, predno bodo uspeli napovedati tornade še pred njihovim začetkom.

Poglejmo tisti del časopisa, ki se je posvečen filmom, ki jih trenutno vrtijo v kinematografih. Vsaka reklama vsebuje citate izbranih kritikov, npr. "Nepozabno!", "Vrhunska predstava našega časa", "Žares osupljivo", ali "En izmed 10 najboljših filmov tega leta!" Ali vam kritike kaj pomenijo? Kako se odločite katere filme si želite ogledati? Strokovnjaki so mnenja, da čeprav lahko vplivamo na popularnost filma s kritikami (dober ali slab) na samem začetku, pa je v celoti najbolj pomembno za film ustno izročilo. Študije so pokazale tudi, da bolj ko je dramatičen film, več kokic je prodanih. Res je, zabavna industrija beleži celo koliko hrustanja opravite med gledanjem. Kako v resnici zberejo vse te informacije in kako to vpliva na zvrsti filmov, ki jih delajo? Tudi to je del statistike: načrtovanje in izdelava študij, ki pomagajo določiti gledalce in ugotoviti kaj imajo radi, ter uporabiti informacijo za pomoč pri vodenju izdelave produkta/izdelka. Če Vas naslednjič nekdo ustavi z anketo in želi nekaj Vašega časa, si ga boste morda res vzeli v upanju, da bo upoštevana tudi Vaša volja.

Loterija in stave. Ko opazujemo zlorabo številke v vsakdanjem življenju, ne moremo mimo športnih stavnic, več milijardno industrijo (letno) ki prevzame tako občasnega stavca, kakor tudi profesionalnega igralca in impulzivnega zasvojenca z igrami na srečo. Na kaj lahko stavimo? Pravzaprav na takorekoč vse kar se konča na dva različna načina.

Številkam se ni mogoče izogniti niti s skokom v sekcijo potovanja. Tam najdemo tudi najbolj pogosto vprašanje naslovljeno na Urad za odzivni center transporta in varnosti, ki prejme tedensko povprečno 2.000 telefonskih klicev, 2.500 e-sporočil in 200 pisem (Bi želeli biti en izmed tistih, ki mora vse to prešteti?): "Ali lahko nesem to-in-to na letalo?", pri čemer se "to-in-to" nanaša na takorekoč karkoli od živali do velikanske konzerve kokic (slednjega ne priporočam, saj je konzervo potrebno shraniti v vodoravni legi, med letom pa se stvari običajno premaknejo, pokrov se odpre in po pristanku vse skupaj pade na Vaše sopotnike - to se je enkrat celo v resnici zgodilo). To nas pripelje do zanimivega statističnega vprašanja: koliko telefonistov je potrebno v različnih časovnih obdobjih tokom dneva, da obdelajo vse klice? Ocena števila klicev je samo prvi korak, in če nismo zadeli prave vrednosti, nas bo to bodisi drago stalo (v primeru, če je bila ocena prevelika) ali pa bomo prišli na slab glas (če je bila ocena prenizka).

Naslednja stvar, ki zbudi našo pozornost, je poročilo o povečanem številu mrtvih na naših cestah. Strokovnjaki nas opozarjajo, da se je število povečalo za več kot 50% od

leta 1997 in nihče ne zna ugotoviti zakaj. Statistika nam pove zanimivo zgodbo. V letu 1997 je umrlo 2,116 motoristov, v letu 2001 pa je statistični urad (National Highway Traffic Safety Administration - NHTSA) poročal o 3,181 žrtvah. V članku je obdelanih več možnih razlogov za povečanje števila žrtev, vključno z dejstvom, da so danes motoristi starejši (povprečna starost ponesrečenih motoristov se je povzpela z 29 let v letu 1990 na 36 let v letu 2001). Velikost dvokolesnikov je opisana kot druga možnost. Prostornina se je v povprečju povečala za skoraj 25% (iz 769 kubičnih centimeterov v letu 1990 na 959 kubičnih centimeters v letu 2001). Naslednja možnost je, da nekatere države ne izvajajo več tako strog nadzor nad zakonom o čeladah. V članku citirajo strokovnjake, da je potrebna veliko natančnejša študija, vendar pa najverjetneje ne bo opravljena, saj bi stala med 2 in 3 milijoni. En aspekt, ki v članku ni omenjen, je število motoristov v letu 2001 v primerjavi s številom v letu 1997. Večje število ljudi na cesti v glavnem pomeni tudi več žrtev, če vsi ostali faktorji ostanejo nespremenjeni. Kljub temu pa je v članku prikazan tudi graf, ki predstavi število smrtnih žrtev na 100 milijonov prepotovanih km od leta 1997 do 2001; ali ta podatek odgovori na vprašanje glede števila ljudi na cesti? Predstavljen je tudi stolpčni graf (diagram), ki primerja število smrtnih žrtev motoristov s številom nezgod s smrtnim izidom, ki so se pripetile z drugimi vozili. Le-ta prikaže 21 ponesrečenih motoristov na 100 milijonov prepotovanih km v primerjavi s samo 1,7 nezgodami s smrtnim izidom pri enakem številu prepotovanih km z avtom. Ta članek vsebuje veliko števil in statistike, toda kaj vse to sploh pomeni?



Statistika je veda, ki proučuje množične pojave.

Ljudje običajno besedo *statistika* povezujejo z zbiranjem in urejanjem podatkov o nekem pojavu, izračunom raznih značilnosti iz teh podatkov, njih predstavitvijo in razlago. To je najstarejši del statistike in ima svoje začetke že v antiki – z nastankom večjih združb (držav) se je pojavila potreba po poznavanju stanja – 'računovodstvo', astronomija, ... Sama beseda *statistika* naj bi izvirala iz latinske besede *status* – v pomenu država. Tej veji statistike pravimo *opisna statistika*. Druga veja, *inferenčna statistika*, poskuša spoznanja iz zbranih podatkov posplošiti (razširiti, podaljšati, napovedati, ...) in oceniti kakovost teh posplošitev. Statistiko lahko razdelimo tudi na *uporabno* in *teoretično* (računalniško in matematično) statistiko.

(Statistična) enota – posamezna proučevana stvar ali pojav (npr. redni študent na Univerzi v Ljubljani v tekočem študijskem letu).

Populacija – množica vseh proučevanih enot; pomembna je natančna opredelitev populacije (npr. časovno in prostorsko). Npr. vsi redni študentje na UL v tekočem študijskem letu.

Vzorec – podmnožica populacije, na osnovi katere ponavadi sklepamo o lastnostih celotne populacije (npr. vzorec 300 slučajno izbranih rednih študentov).

Spremenljivka – lastnost enot; označujemo jih npr. z X , Y , X_1 . Vrednost spremenljivke X na i -ti enoti označimo z x_i (npr. spol, uspeh iz matematike v zadnjem razredu srednje šole, izobrazba matere in višina mesečnih dohodkov staršev študenta itd.).

Posamezne spremenljivke in odnose med njimi opisujejo ustrezne porazdelitve.

Parameter je značilnost populacije, običajno jih označujemo z malimi grškimi črkami.

Statistika je značilnost vzorca; običajno jih označujemo z malimi latinskimi črkami.

Vrednost statistike je lahko za različne vzorce različna.

Eno izmed osnovnih vprašanj statistike je,
kako z uporabo ustreznih statistik oceniti
vrednosti izbranih parametrov.

Poglavje 11

Opisna statistika

11.1 Vrste spremenljivk oziroma podatkov

Glede na vrsto vrednosti jih delimo na:

1. *številске* (ali numerične) spremenljivke – vrednosti lahko izrazimo s števili (npr. starost). Za podatke rečemo, da so *kvantitativni* kadar predstavljajo kvantiteto ali količino nečesa.



2. *opisne* (ali atributivne) spremenljivke – vrednosti lahko opišemo z imeni razredov (npr. poklic, uspeh, spol); Za podatke rečemo, da so *kvalitativni* (kategorični) kadar jih delimo v kategorije in zanje ni kvantitativnih interpretacij.



Če so kategorije brez odgovarjajočega vrstnega reda/urejenosti, rečemo, da so podatki *nominalni*. Kakor hitro imajo kategorije neko urejenost pa rečemo, da so podatki *ordinalni/številski*.

Glede na vrsto merske lestvice:

1. *imenske* (ali nominalne) spremenljivke – vrednosti lahko le razlikujemo med seboj: dve vrednosti sta enaki ali različni (npr. spol);
2. *urejenostne* (ali ordinalne) spremenljivke – vrednosti lahko uredimo od najmanjše do največje (npr. uspeh);
3. *razmične* (ali intervalne) spremenljivke – lahko primerjamo razlike med vrednostima dvojic enot (npr. temperatura);
4. *razmernostne* spremenljivke – lahko primerjamo razmerja med vrednostima dvojic enot (npr. starost).
5. *absolutne* spremenljivke – štetja (npr. število prebivalcev).



<i>dovoljene transformacije</i>	<i>vrsta lestvice</i>	<i>primeri</i>
$f(x) = x$ (identiteta)	absolutna	štetje
$f(x) = a \cdot x, a > 0$ podobnost	razmernostna	masa temperatura (K)
$f(x) = a \cdot x + b, a > 0$	razmična	temperatura (C,F) čas (koledar)
$x \geq y \Leftrightarrow f(x) \geq f(y)$ strogo naraščajoča	urejenostna	šolske ocene, kakovost zraka, trdost kamnin
f je povratno enolična	imenska	barva las, narodnost

Vrste spremenljivk so urejene od tistih z najslabšimi merskimi lastnostmi do tistih z najboljšimi. Urejenostne spremenljivke zadoščajo lastnostim, ki jih imajo imenske spremenljivke; in podobno razmernostne spremenljivke zadoščajo lastnostim, ki jih imajo razmične, urejenostne in imenske spremenljivke.

$$\text{absolutna} \subset \text{razmernostna} \subset \text{razmična} \subset \text{urejenostna} \subset \text{imenska}$$

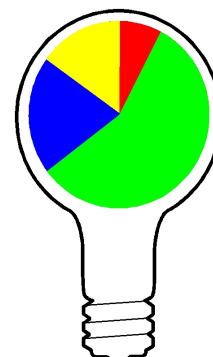
Posamezne statistične metode predpostavljajo določeno vrsto spremenljivk. Največ učinkovitih statističnih metod je razvitih za številske spremenljivke.

V teoriji merjenja pravimo, da je nek stavek *smiselno*, če ohranja resničnost/lažnost pri zamenjavi meritev z enakovrednimi (glede na dovoljene transformacije) meritvami.

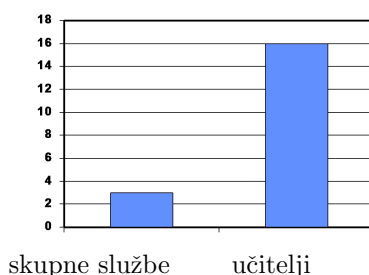
11.2 Grafična predstavitev kvantitativnih podatkov

Oddelek sistemskih inženirjev

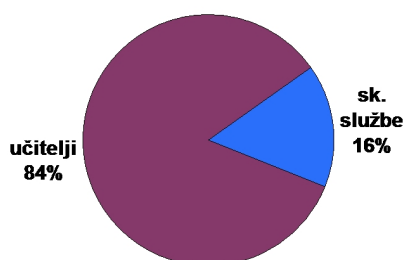
kategorija	frekvenca	relativna frekvenca
vrsta	število	
zaposlenih	zaposlenih	delež
učitelji	16	0,8421
skupne službe	3	0,1579
skupaj	19	1,0000



Stolpčni prikaz (tudi stolpčni graf, poligonski diagram): Na eni osi prikažemo (urejene) razrede. Nad vsakim naredimo stolpec/črto višine sorazmerne frekvenci razreda.

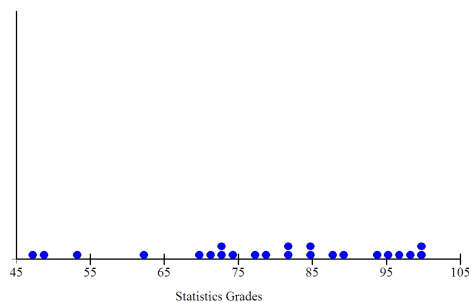
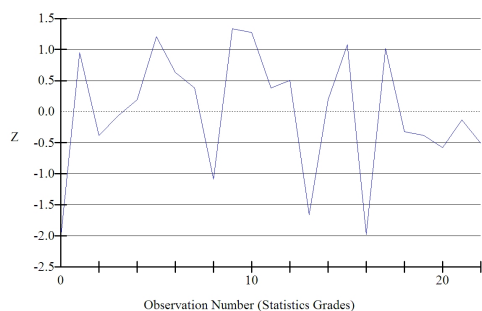


Krožni prikaz (tudi strukturni krog, pogača, kolač): Vsakemu razredu priredimo krožni izsek s kotom $\alpha_i = \frac{f_i}{n} 360$ stopinj.

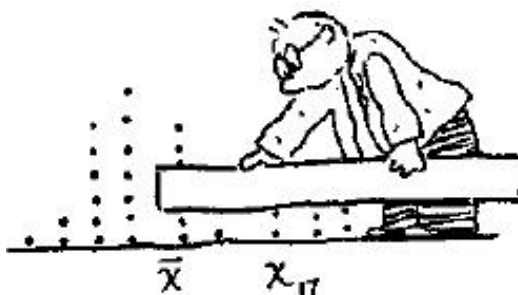


Poligon: v koordinatnem sistemu zaznamujemo točke (x_i, f_i) , kjer je x_i sredina i -tega razreda in f_i njegova frekvenca. K tem točkam dodamo še točki $(x_0, 0)$ in $(x_{k+1}, 0)$, če je v frekvenčni porazdelitvi k razredov. Točke zvežemo z daljicami.

Runs chart/plot in Dot plot



Frekvenčna porazdelitev



Število vseh možnih vrednosti proučevane spremenljivke je lahko preveliko za pregledno prikazovanje podatkov. Zato sorodne vrednosti razvrstimo v skupine. Posamezni skupini priredimo ustrezno reprezentativno vrednost, ki je nova vrednost spremenljivke. Skupine vrednosti morajo biti določene *enolično*: vsaka enota s svojo vrednostjo je lahko uvrščena v natanko eno skupino vrednosti. *Frekvenčna porazdelitev* spremenljivke je *tabela*, ki jo določajo *vrednosti ali skupine vrednosti* in njihove *frekvence*. Če je spremenljivka vsaj urejenostna, vrednosti (ali skupine vrednosti) uredimo od najmanjše do največje. Skupine vrednosti številskih spremenljivk imenujemo *razredi*. Če zapišem podatke v vrsto po njihovi numerični velikosti pravimo, da gre za **urejeno zaporedje** oziroma *ranžirano vrsto*, ustreznemu mestu pa pravimo *rang*.

x_{min} in x_{max} – *najmanjša* in *največja* vrednost spremenljivke X .

$x_{i,min}$ in $x_{i,max}$ – *spodnja* in *zgornja meja* i -tega razreda.

Meje razredov so določene tako, da velja $x_{i,max} = x_{i+1,min}$.

Širina i -tega razreda je $d_i = x_{i,max} - x_{i,min}$.

Če je le mogoče, vrednosti razvrstimo v razrede enake širine.

Sredina i -tega razreda je $x_i = \frac{x_{i,min} + x_{i,max}}{2}$ in je značilna vrednost – predstavnik razreda.

Kumulativa (ali nakopičena frekvenca) je frekvenca do spodnje meje določenega razreda. Velja $F_{i+1} = F_i + f_i$, kjer je F_i kumulativa in f_i frekvenca v i -tem razredu.

Primer zaporedja podatkov (nal. 2.48, str.64)

	88	103	113	122	132
	92	108	114	124	133
	95	109	116	124	133
(a) Konstruiraj urejeno zaporedje.	97	109	116	124	135
	97	111	117	128	136
	97	111	118	128	138
(b) Nariši steblo-list diagram.	98	112	119	128	138
	98	112	120	131	142
	100	112	120	131	146
(c) Naredi histogram.	100	113	122	131	150

Koraki za konstrukcijo steblo-list predstavitev

1. Razdeli vsako opazovanje-podatke na dva dela: **stebila** (angl. stem) in **listi** (angl. leaf).
2. Naštej stebila po vrsti v stolpec, tako da začneš pri najmanjšem in končaš pri največjem.
3. Upoštevaj vse podatke in postavi liste za vsak dogodek/meritev v ustrezno vrstico/steblo.
4. Preštej frekvenca za vsako steblo.

Steblo-list diagram

stebila	listi	rel. ν	ν
08	8	1	2%
09	2 5 7 7 7 8 8	7	14%
10	0 0 3 8 9 9	6	12%
11	1 1 2 2 2 3 3 4 6 6 7 8 9	13	26%
12	0 0 2 2 4 4 4 8 8 8	10	20%
13	1 1 1 2 3 3 5 6 8 8	10	20%
14	2 6	2	4%
15	0	1	2%
		50	100%



Histogrami

Histogram: drug poleg drugega rišemo stolpce – pravokotnike, katerih ploščina je sorazmerna frekvenci v razredu. Če so razredi enako široki, je višina sorazmerna tudi frekvenci.

Ogiva: grafična predstavitev kumulativne frekvenčne porazdelitve s poligonom, kjer v koordinatni sistem nanašamo točke $(x_{i,min}, F_i)$.



(1) Kako zgradimo histogram

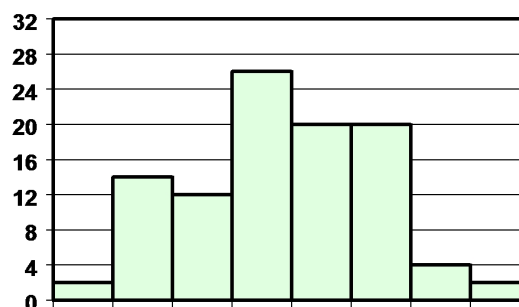
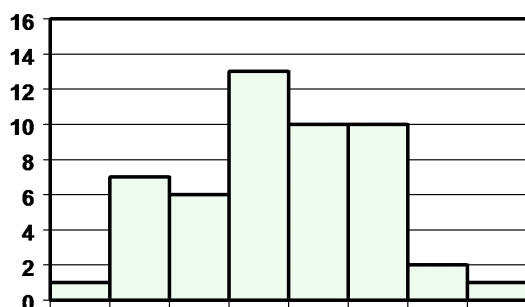
- Izračunaj **razpon** podatkov.
- Razdeli razpon na **5 do 20 razredov** enake širine.
- Za vsak razred preštej število vzorcev, ki spadajo v ta razred. To število imenujemo **frekvenca razreda**.
- Izračunaj vse **relativne frekvence razredov**

(2) Pravilo za določanje števila razredov v histogramu

število vzorcev v množici podatkov	število razredov
manj kot 25	5 ali 6
25 – 50	7 – 14
več kot 50	15 – 20

(3,4) Frekvenčna porazdelitev

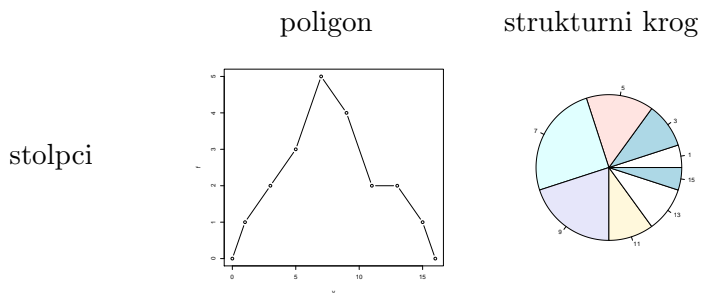
razred	interval razreda	frekvenca	relativna frekvenca
1	80 – 90	1	2%
2	90 – 100	7	14%
3	100 – 110	6	12%
4	110 – 120	13	26%
5	120 – 130	10	20%
6	130 – 140	10	20%
7	140 – 150	2	4%
8	150 – 160	1	2%

Frekvenčni in procentni histogram**Nekaj ukazov v R-ju**

```

> X <- c(5,11,3,7,5,7,15,1,13,11,9,9,3,13,9,7,7,5,9,7)
> n <- length(X)
> t <- tabulate(X)
> t
[1] 1 0 2 0 3 0 5 0 4 0 2 0 2 0 1
> v <- (1:max(X))[t>0]
> f <- t[t>0]
> rbind(v,f)
  [,1] [,2] [,3] [,4] [,5] [,6] [,7] [,8]
v   1   3   5   7   9  11  13  15
f   1   2   3   5   4   2   2   1
> plot(v,f,type="h")
> plot(c(0,v,16),c(0,f,0),type="b",xlab="v",ylab="f")
> pie(f,v)
> plot(c(0,v,16),c(0,cumsum(f)/n,1),col="red",type="š",
  xlab="v",ylab="f")
> x <- sort(rnorm(100,mean=175,sd=30))
> y <- (1:100)/100
> plot(x,y,main="Normalna porazdelitev, n=100",type="š")
> curve(pnorm(x,mean=175,sd=30),add=T,col="red")

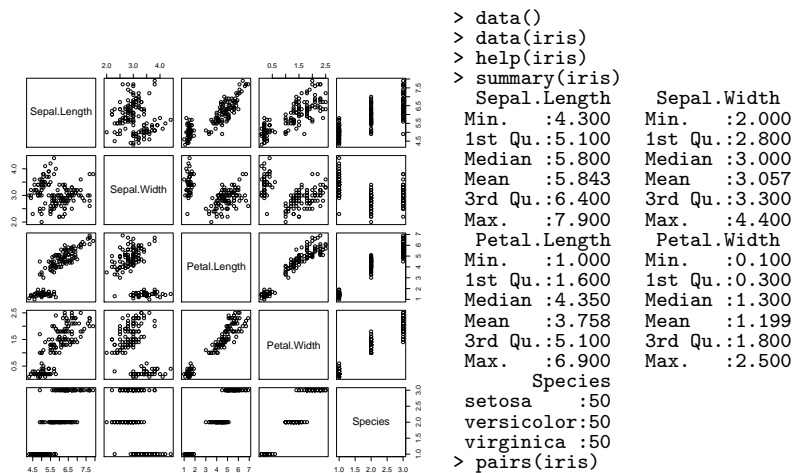
```

```
> x <- rnorm(1000,mean=175,sd=30)
> mean(x)
[1] 175.2683
> sd(x)
[1] 30.78941
> var(x)
[1] 947.9878
> median(x)
[1] 174.4802
> min(x)
[1] 92.09012
> max(x)
[1] 261.3666
> quantile(x,seq(0,1,0.1))
  0%      10%     20%     30%
92.09012 135.83928 148.33908 158.53864
  40%     50%     60%     70%
166.96955 174.48018 182.08577 191.29261
  80%     90%    100%
200.86309 216.94009 261.36656

> summary(x)
   Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
 92.09  154.20  174.50  175.30  195.50  261.40
> hist(x,freq=F)
> curve(dnorm(x,mean=175,sd=30),add=T,col="red")
```

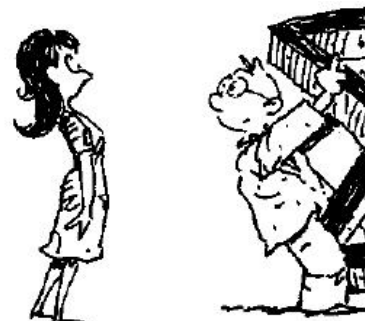
Fisherjeve oziroma Andersonove perunike (Iris data)



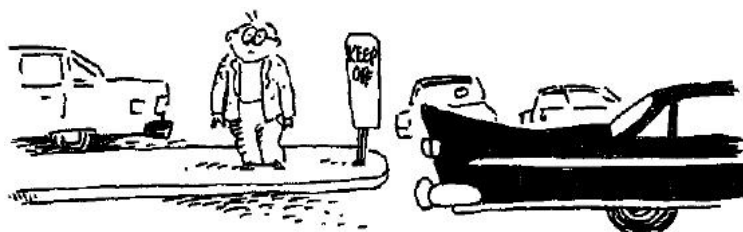
Parni prikaz.

11.3 Mere za lokacijo in razpršenost

- srednje vrednosti
- razpon (min/max)
- centili, kvartili
- varianca
- standardni odklon
- Z-vrednosti



Modus (oznaka M_0) množice podatkov je tista vrednost, ki se pojavi z največjo frekvenco.



Da bi prišli do **mediane** (oznaka M_e) za neko množico podatkov, naredimo naslednje:

1. Podatke uredimo po velikosti v naraščujočem vrstnem redu,
2. Če je število podatkov liho, potem je mediana podatek na sredini,
3. Če je število podatkov sodo, je mediana enaka povprečju dveh podatkov na sredini.

Oznake: mediana populacije: μ mediana vzorca: m



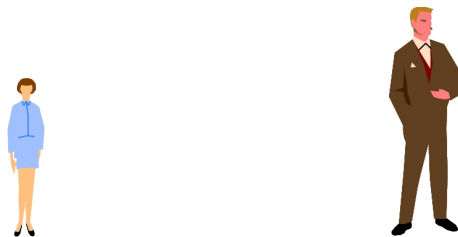
Povrečje populacije:

$$\mu = \frac{1}{N}(y_1 + \cdots + y_N) = \frac{\sum_{i=1}^N y_i}{N}$$

Povrečje vzorca:

$$\bar{y} = \frac{1}{n}(y_1 + \cdots + y_n) = \frac{\sum_{i=1}^n y_i}{n}$$

Razpon je razlika med največjo in najmanjšo meritvijo v množici podatkov.



100p-ti centil ($p \in [0, 1]$) je definiran kot število, od katerega ima 100p % meritev manjšo ali enako numerično vrednost. 100p-ti centil določimo tako, da izračunamo vrednost $p(n+1)$ in jo zaokrožimo na najbližje celo število. Naj bo to število enako i . Izmerjena vrednost z i -tim rangom je 100p-ti centil.

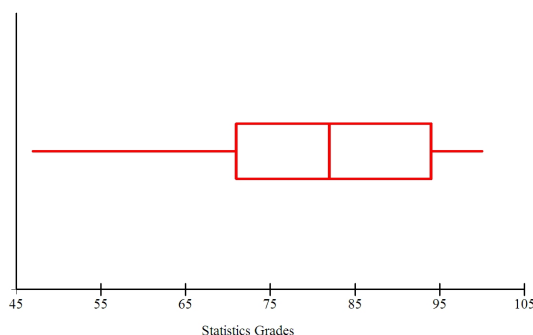
25. centil se imenuje tudi **1. kvartil**.

50. centil se imenuje **2. kvartil** ali **mediana**.

75. centil se imenuje tudi **3. kvartil**.



Škatla z brki (angl. box plot)



Še nekaj ukazov v R-ju: škatle in Q-Q-prikazi

Škatle (box-and-whiskers plot; grafikon kvantilov) **boxplot**:

škatla prikazuje notranja kvartila razdeljena z mediansko črto.

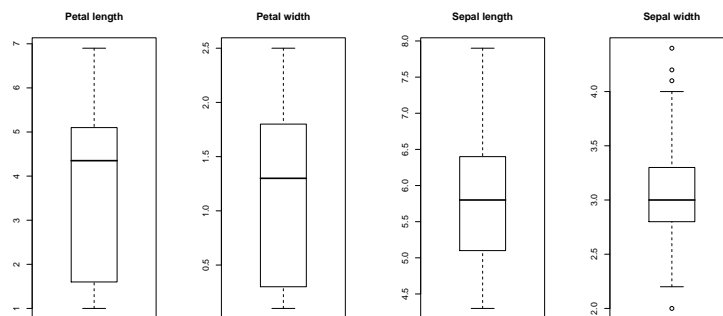
Daljci – brka vodita do robnih podatkov, ki sta največ za 1,5 dolžine škatle oddaljena od nje. Ostali podatki so prikazani posamično.

Q-Q-prikaz `qqnorm` je namenjen prikazu normalnosti porazdelitve danih n podatkov. Podatke uredimo in prikažemo pare točk sestavljene iz vrednosti k -tega podatka in pričakovane

vrednosti k -tega podatka izmed n normalno porazdeljenih podatkov. Če sta obe porazdelitvi normalni, ležijo točke na premici. Premica `qqline` nariše premico skozi prvi in tretji kvartil.

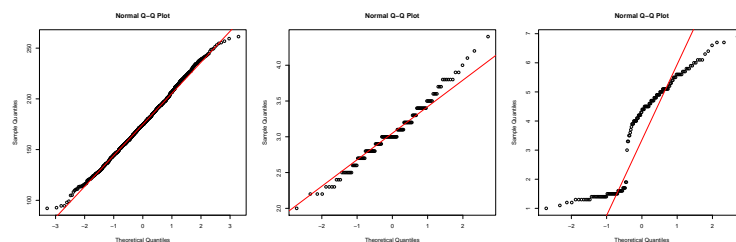
Obstaja tudi splošnejši ukaz `qqplot`, ki omogoča prikaz povezanosti poljubnega para porazdelitev. S parametrom `datax=T` zamenjamo vlogo koordinatnih osi.

Škatle



```
> par(mfrow=c(1,2))
> boxplot(iris$Petal.Length,main="Petal length")
> boxplot(iris$Petal.Width,main="Petal width")
> boxplot(iris$Sepal.Length,main="Sepal length")
> boxplot(iris$Sepal.Width,main="Sepal width")
> par(mfrow=c(1,1))
```

Q-Q-prikaz



```
> qqnorm(x)
> qqline(x,col="red")
> qqnorm(iris$Sepal.Width)
> qqline(iris$Sepal.Width,col="red")
> qqnorm(iris$Petal.Length)
> qqline(iris$Petal.Length,col="red")
```

Mere razpršenosti

varianca

- kvadrat pričakovanega odklona (populacije)
- vsota kvadratov odklonov deljena s stopnjo prostosti (vzorca)

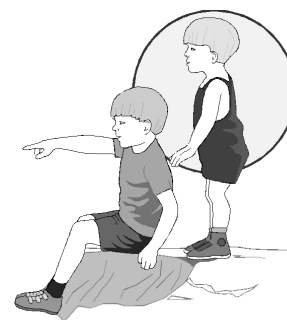
standardni odklon (deviacija)

- pozitivni kvadratni koren variance

koeficient variacije

- standardni odklon deljen s povprečjem

	populacija	vzorec
varianca	σ^2 D, V	S^2, s^2
standardni odklon	σ	S, s



Za vzorec smo vzeli osebje FRI in zabeležili naslednje število otrok: 1, 2, 2, 1, 2, 5, 1, 2.

Varianca in standardni odklon

Varianca populacije (končne populacije z N elementi):

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

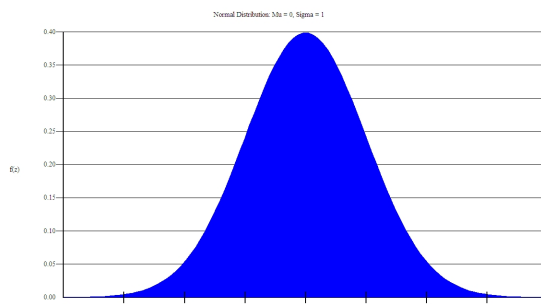
Varianca vzorca (n meritvami):

$$s^2 = \frac{\sum_{i=1}^n (y_i - \mu)^2}{n-1} = \frac{\sum_{i=1}^n y_i^2 - \frac{\left(\sum_{i=1}^n y_i\right)^2}{n}}{n-1}$$

Standardni odklon je pozitivno predznačen kvadratni koren variance.

Normalna porazdelitev

Veliko podatkovnih množic ima porazdelitev približno **zvonaste oblike** (unimodalna oblika - ima en sam vrh):



Empirična pravila

Če ima podatkovna množica porazdelitev približno **zvonaste oblike**, potem veljajo naslednja pravila (angl. rule of thumb), ki jih lahko uporabimo za opis podatkovne množice:

1. Približno **68,3%** vseh meritev leži na razdalji
 $1 \times$ *standardnega odklona* od njihovega povprečja.
2. Približno **95,4%** meritev leži na razdalji do
 $2 \times$ *standardnega odklona* od njihovega povprečja.
3. Skoraj vse meritve (**99,7%**) ležijo na razdalji
 $3 \times$ *standardnega odklona* od njihovega povprečja.

Mere oblike

Če je spremenljivka približno normalno porazdeljena, potem jo statistični karakteristiki **povprečje** in **standardni odklon** zelo dobro opisujeta. V primeru unimodalne porazdelitve spremenljivke, ki pa je bolj asimetrična in bolj ali manj sploščena (koničasta), pa je potrebno izračunati še stopnjo *asimetrije* in *sploščenosti* (koničavosti).

ℓ -ti centralni moment je

$$m_\ell = \frac{(y_1 - \mu)^\ell + \dots + (y_n - \mu)^\ell}{n}.$$

$$m_1 = 0, m_2 = \sigma^2, \dots$$

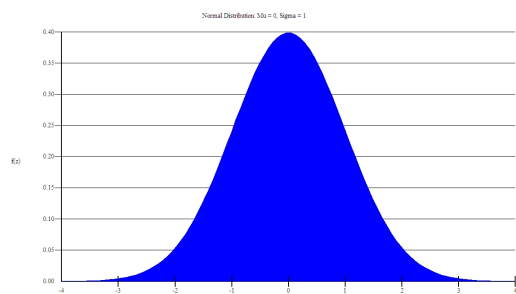
Koeficient asimetrije (s centralnimi momenti): $g_1 = m_3/m_2^{3/2}$. Mere asimetrije dobim tako, da opazujemo razlike med srednjimi vrednostimi. Le-te so tem večje čim bolj je porazdelitev asimetrična:

$$KA_{M_0} = (\mu - M_0)/\sigma, \quad KA_{M_e} = 3(\mu - M_e)/\sigma.$$

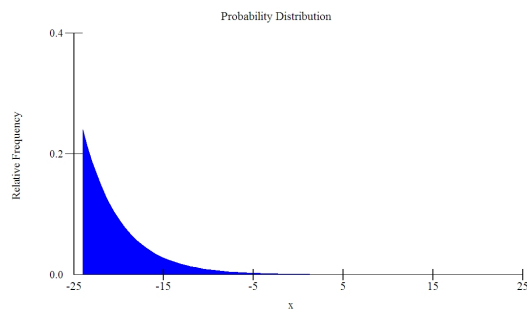
Koeficient sploščenosti (kurtosis) (s centralnimi momenti): $K = g_2 = m_4/m_2^2 - 3$

- $K = 3$ (ali 0) normalna porazdelitev zvonaste-oblike (*mesokurtic*),
- $K < 3$ (ali < 0) bolj kopasta kot normalna porazdelitev, s krajšimi repi (*platykurtic*),
- $K > 3$ (ali > 0) bolj špičasta kot normalna porazdelitev, z daljšimi repi (*leptokurtic*).

Normalna in asimetrična porazdelitev

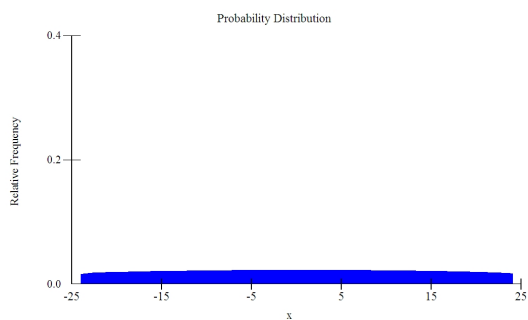


asim.= 0, sploščenost= 3 (mesokurtic).

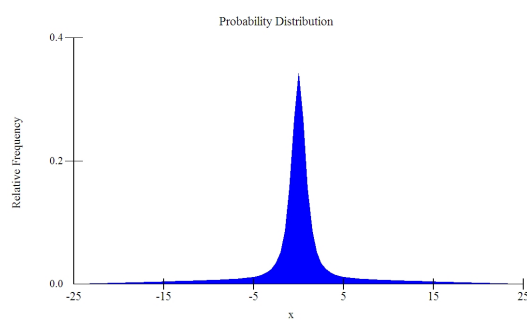


asim.= 1,99, sploščenost= 8,85.

Kopasta in špičasta porazdelitev



asim.= 0, sploščenost= 1,86 (platykurtic)



asim.= -1,99, sploščenost= 8,85 (leptokurtic).

11.4 Standardizacija

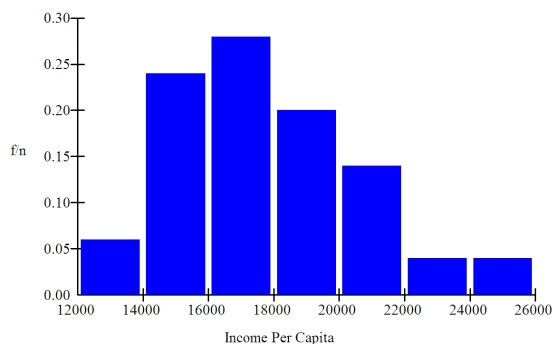
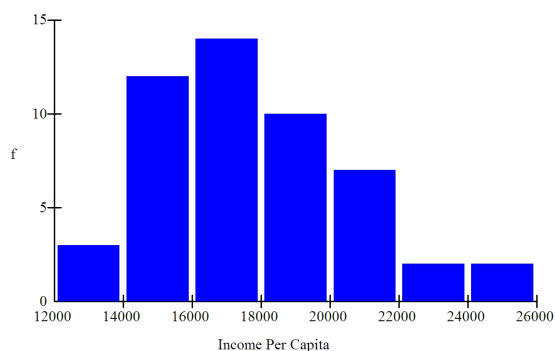
Vsaki vrednosti x_i spremenljivke X odštejemo njeno povprečje μ in delimo z njenim standardnim odklonom σ :

$$z_i = \frac{x_i - \mu}{\sigma}.$$

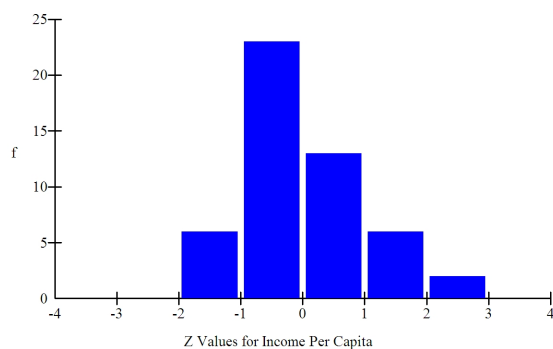
Za novo spremenljivko Z bomo rekli, da je **standardizirana**, z_i pa je **standardizirana vrednost**.

Potem je $\mu(Z) = 0$ in $\sigma(Z) = 1$.

Frekvenčni in relativni frekvenčni histogram

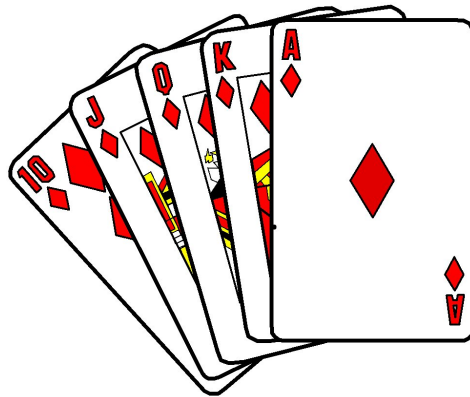


Histogram standardiziranih Z -vrednosti

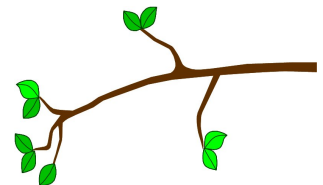


Poglavje 12

Vzorčenje



Analitična statistika je veja statistike, ki se ukvarja z uporabo vzorčnih podatkov, da bi z njimi naredili zaključek (inferenco) o populaciji.



Zakaj vzorčenje?

- cena
- čas
- destruktivno testiranje

Glavno vprašanje statistike je:

kakšen mora biti vzorec, da lahko iz podatkov zbranih na njem veljavno sklepamo o lastnostih celotne populacije.

Kdaj vzorec dobro predstavlja celo populacijo?

Preprost odgovor je:

- vzorec mora biti izbran *nepristransko*,
- vzorec mora biti *dovolj velik*.

Recimo, da merimo spremenljivko X , tako da n -krat naključno izberemo neko enoto in na njej izmerimo vrednost spremenljivke X . Postopku ustreza slučajni vektor

$$(X_1, X_2, \dots, X_n),$$

ki mu rečemo *vzorec*. Število n je *velikost* vzorca.



Ker v vzorcu merimo isto spremenljivko in posamezna meritev ne sme vplivati na ostale, lahko predpostavimo:

1. vsi členi X_i vektorja imajo *isto* porazdelitev, kot spremenljivka X ,
2. členi X_i so med seboj *neodvisni*.

Takemu vzorcu rečemo *enostavni slučajni vzorec*. Večina statistične teorije temelji na predpostavki, da imamo opravka enostavnim slučajnim vzorcem. Če je populacija končna, lahko dobimo enostavni slučajni vzorec, tako da slučajno izbiramo (z vračanjem) enote z enako verjetnostjo. Z vprašanjem, kako sestaviti dobre vzorce v praksi, se ukvarja posebno področje statistike – *teorija vzorčenja*.

Načini vzorčenja

- ocena
 - priročnost
- naključno
 - enostavno: pri enostavnem naključnem vzorčenju je vsak član populacije izbran/vključen z *enako verjetnostjo*.
 - deljeno: razdeljen naključni vzorec dobimo tako, da razdelimo populacijo na disjunktne množice oziroma dele (razrede) in nato izberemo enostavne naključne vzorce za vsak del posebej.
 - grozdno: takšno vzorčenje je enostavno naključno vzorčenje skupin ali klastrov/grozdov elementov.

12.1 Osnovni izrek statistike

Spremenljivka X ima na populaciji G porazdelitev $F(x) = P(X < x)$. Toda tudi vsakemu vzorcu ustreza neka porazdelitev. Za realizacijo vzorca $(x_1, x_2, x_3, \dots, x_n)$ in $x \in \mathbb{R}$ postavimo

$$K(x) = |\{x_i : x_i < x, i = 1, \dots, n\}| \quad \text{in} \quad V_n(x) = K(x)/n.$$

Slučajni spremenljivki $V_n(x)$ pravimo *vzorčna porazdelitvena funkcija*. Ker ima, tako kot tudi $K(x)$, $n + 1$ možnih vrednosti k/n , $k = 0, \dots, n$, je njena verjetnostna funkcija $B(n, F(x))$

$$P(V_n(x) = k/n) = \binom{n}{k} F(x)^k (1 - F(x))^{n-k}.$$

Če vzamemo n neodvisnih Bernoullijevih spremenljivk

$$Y_i(x) : \begin{pmatrix} 1 & 0 \\ F(x) & 1 - F(x) \end{pmatrix},$$

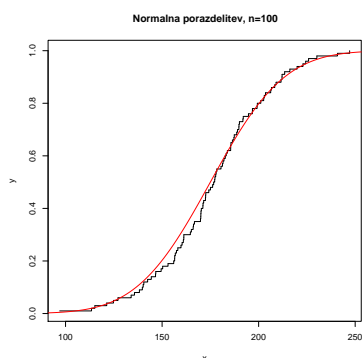
velja

$$V_n(x) = \frac{1}{n} \sum_{i=1}^n Y_i(x).$$

Krepki zakon velikih števil tedaj zagotavlja, da za vsak x velja

$$P\left(\lim_{n \rightarrow \infty} V_n(x) = F(x)\right) = 1.$$

To je v bistvu Borelov zakon, da relativna frekvenca dogodka $(X < x)$ skoraj gotovo konvergira proti verjetnosti tega dogodka. Velja pa še več. $V_n(x)$ je stopničasta funkcija, ki se praviloma dobro prilega funkciji $F(x)$.



Odstopanje med $V_n(x)$ in $F(x)$ lahko izmerimo s slučajno spremenljivko

$$D_n = \sup_{x \in \mathbb{R}} |V_n(x) - F(x)|$$

za $n = 1, 2, \dots$. Zanj lahko

pokažemo *osnovni izrek statistike*

$$P\left(\lim_{n \rightarrow \infty} D_n = 0\right) = 1.$$

Torej se z rastjo velikosti vzorca $V_n(x)$ enakomerno vse bolj prilega funkciji $F(x)$ – vse bolj povzema razmere na celotni populaciji.

12.2 Vzorčne ocene

Najpogostejša parametra, ki bi ju radi ocenili sta: *sredina populacije* μ glede na izbrano lastnost – matematično upanje spremenljivke X na populaciji; in *povprečni odklon* od sredine σ – standardni odklon spremenljivke X na populaciji. Statistike/ocene za te parametre so izračunane iz podatkov vzorca. Zato jim tudi rečemo *vzorčne ocene*.

Sredinske mere

Kot sredinske mere se pogosto uporabljajo:

Vzorčni modus – najpogostejša vrednost (smiselna tudi za imenske).

Vzorčna mediana – srednja vrednost, glede na urejenost, (smiselna tudi za urejenostne).

Vzorčno povprečje – povprečna vrednost (smiselna za vsaj razmične)

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$$

Vzorčna geometrijska sredina – (smiselna za vsaj razmernostne)

$$G(x) = \sqrt[n]{\prod_{i=1}^n x_i}$$

Mere razpršenosti

Za oceno populacijskega odklona uporabljamo *mere razpršenosti*.

Vzorčni razmah = $\max_i x_i - \min_i x_i$.

Vzorčna disperzija $s_0^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2$.

Popravljen vzorčna disperzija $s^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2$.

ter ustrezna *vzorčna odklona* s_0 in s .

12.3 Porazdelitve vzorčnih povprečij

Denimo, da je v populaciji N enot in da iz te populacije slučajno izbiramo n enot v enostavni slučajni vzorec ali na kratko slučajni vzorec (vsaka enota ima enako verjetnost, da bo izbrana v vzorec, tj. $1/N$).

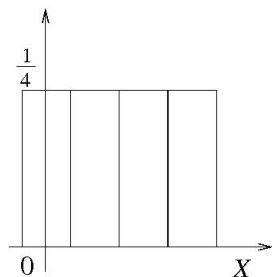
Če hočemo dobiti slučajni vzorec, moramo izbrane enote pred ponovnim izbiranjem vrniti v populacijo (vzorec s ponavljanjem).

Če je velikost vzorca v primerjavi s populacijo majhna, se ne pregrešimo preveč, če imamo za slučajni vzorec tudi vzorec, ki nastane s slučajnim izbiranjem brez vračanja.

Predstavljajmo si, da smo iz populacije izbrali vse možne vzorce. Dobili smo populacijo vseh možnih vzorcev. Teh je v primeru enostavnih slučajnih vzorcev (s ponavljanjem) N^n ; kjer je N število enot v populaciji in n število enot v vzorcu.

Število slučajnih vzorcev brez ponavljanja pa je $\binom{N}{n}$, če ne upoštevamo vrstnega reda izbranih enot v vzorcu, oziroma $\binom{N+n-1}{n}$, če upoštevamo vrstni red.

Primer: Vzemimo populacijo z $N = 4$ enotami, ki imajo naslednje vrednosti spremenljivke X : 0, 1, 2, 3. Grafično si lahko porazdelitev spremenljivke X predstavimo s histogramom:



in izračunamo populacijsko aritmetično sredino in varianco:

$$\mu = \frac{1}{N} \sum_{i=1}^N x_i = \frac{3}{2},$$

$$\sigma^2 = \frac{1}{N} \sum_{i=1}^N (x_i - \mu)^2 = \frac{5}{4}.$$

Sedaj pa tvorimo vse možne vzorce velikosti $n = 2$ s ponavljanjem, in na vsakem izračunajmo vzorčno aritmetično sredino \bar{X} :

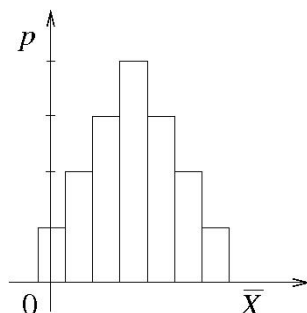
$$\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i = \frac{1}{2}(X_1 + X_2).$$

vzorci	\bar{X}	vzorci	\bar{X}
0, 0	0	2, 0	1
0, 1	0,5	2, 1	1,5
0, 2	1	2, 2	2
0, 3	1,5	2, 3	2,5
1, 0	0,5	3, 0	1,5
1, 1	1	3, 1	2
1, 2	1,5	3, 2	2,5
1, 3	2	3, 3	3

Zapišimo verjetnostno shemo slučajne spremenljivke vzorčno povprečje \bar{X} :

$$\bar{X} : \begin{pmatrix} 0 & 0,5 & 1 & 1,5 & 2 & 2,5 & 3 \\ 1/16 & 2/16 & 3/16 & 4/16 & 3/16 & 2/16 & 1/16 \end{pmatrix}$$

Grafično jo predstavimo s histogramom:



... in izračunajmo matematično upanje ter disperzijo vzorčnega povprečja:

$$E(\bar{X}) = \sum_{i=1}^m \bar{X}_i p_i = \frac{0 + 1 + 3 + 6 + 6 + 5 + 3}{16} = \frac{3}{2},$$

$$D(\bar{X}) = \sum_{i=1}^m \left(\bar{X}_i - E(\bar{X}) \right)^2 p_i = \frac{5}{8}.$$

◇

S tem primerom smo pokazali, da je statistika ‘vzorčna aritmetična sredina’ slučajna spremenljivka s svojo porazdelitvijo. Poglejmo, kaj lahko rečemo v splošnem o porazdelitvi vzorčnih aritmetičnih sredin.

Spomnimo se, da nam **Centralni limitni izrek** v primeru povprečja pove:

Če je naključni vzorec velikosti n izbran iz populacije s končnim povprečjem μ in varianco σ^2 , potem je lahko, če je n dovolj velik, vzorčna porazdelitev povprečja \bar{y} aproksimirana z gostoto normalne porazdelitve.

Naj bo y_1, y_2, \dots, y_n naključni vzorec, ki je sestavljen iz n meritev populacije s končnim povprečjem μ in končnim standardnim odklonom σ . Potem sta povprečje in standardni odklon vzorčne porazdelitve \bar{y} enaka

$$\mu_{\bar{Y}} = \mu, \quad \text{in} \quad \sigma_{\bar{Y}} = \frac{\sigma}{\sqrt{n}}.$$

Hitrost centralne tendence pri CLI

Dokaz CLI je precej tehničen, kljub temu pa nam ne da občutka kako velik mora biti n , da se porazdelitev slučajne spremenljivke

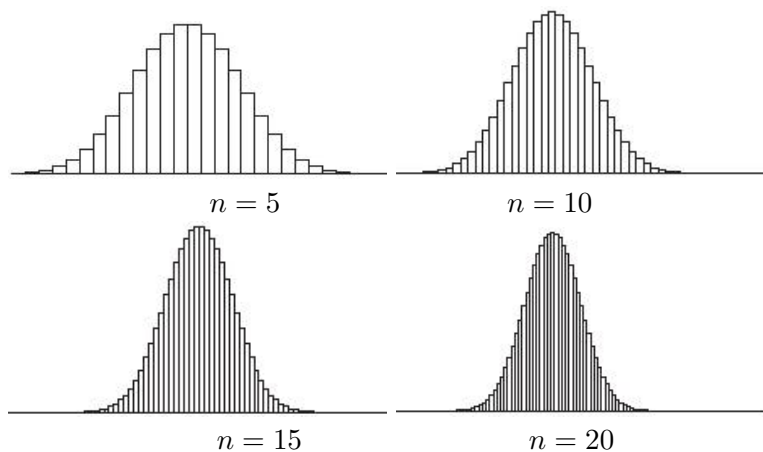
$$X_1 + \dots + X_n$$

približa normalni porazdelitvi. Hitrost približevanja k normalni porazdelitvi je odvisna od tega kako simetrična je porazdelitev. To lahko potrdimo z eksperimentom: mečemo (ne)pošteno kocko, X_k naj bo vrednost, ki jo kocka pokaže pri k -tem metu.

Centralna tendenca za pošteno kocko

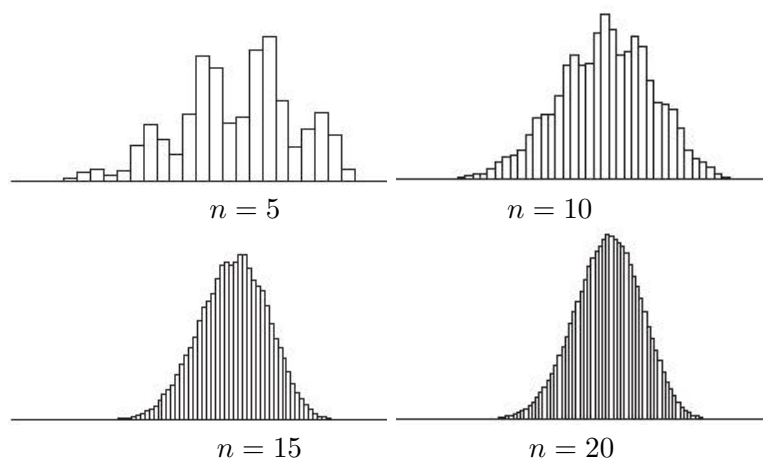
$$p_1 = 1/6, \quad p_2 = 1/6, \quad p_3 = 1/6, \quad p_4 = 1/6, \quad p_5 = 1/6, \quad p_6 = 1/6.$$

in slučajno spremenljivko $X_1 + X_2 + \dots + X_n$:

**Centralna tendenca za goljufivo kocko**

$$p_1 = 0,2, \quad p_2 = 0,1, \quad p_3 = 0, \quad p_4 = 0, \quad p_5 = 0,3, \quad p_6 = 0,4.$$

in slučajno spremenljivko $X_1 + X_2 + \dots + X_n$:



12.4 Vzorčna statistika

Vzorčna statistika je poljubna simetrična funkcija (tj. njena vrednost je neodvisna od permutacije argumentov) vzorca

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

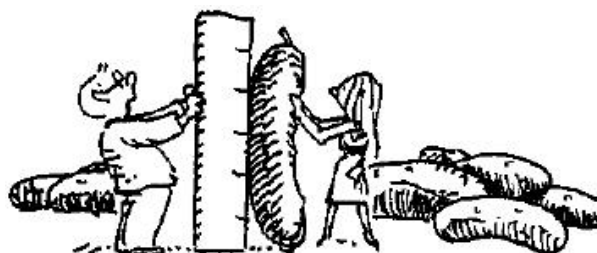
Tudi vzorčna statistika je slučajna spremenljivka, za katero lahko določimo porazdelitev iz porazdelitve vzorca. Najzanimivejši sta značilni vrednosti

- njeno matematično upanje EY ,
- standardni odklon σY , ki mu pravimo tudi *standardna napaka* statistike Y (angl. standard error – zato oznaka $SE(Y)$).



12.4.1 (A) Vzorčno povprečje

Proizvajalec embalaže za kumare bi rad ugotovil **povprečno dolžino** kumarice (da se odloči za velikost embalaže), ne da bi izmeril dolžino čisto vsake.



Zato naključno izbere n kumar in izmeri njihove dolžine X_1, \dots, X_n . Sedaj nam je že blizu ideja, da je vsaka dolžina X_i **slučajna spremenljivka** (numerični rezultat naključnega eksperimenta). Če je μ (iskano/neznano) povprečje dolžin, in je σ standardni odklon porazdelitve dolžin kumar, **potem velja**

$$EX_i = \mu, \quad \text{in} \quad DX_i = \sigma^2,$$

za vsak i , ker bi X_i bila lahko dolžina katerekoli kumare.



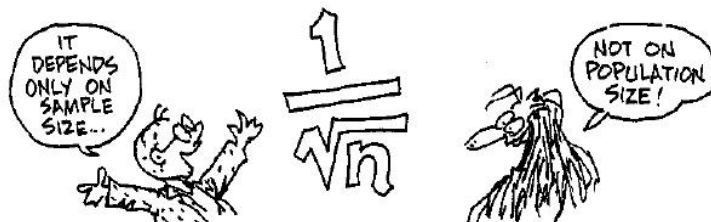
Oglejmo si **vzorčno povprečje**, določeno z zvezo

$$\bar{X} = \frac{X_1 + \dots + X_n}{n} = \frac{1}{n} \sum_{i=1}^n X_i,$$

ki je tudi slučajna spremenljivka. Tedaj je

$$E\bar{X} = \frac{1}{n} \sum_{i=1}^n EX_i = \mu \quad \text{in} \quad D\bar{X} = \frac{1}{n^2} \sum_{i=1}^n DX_i = \frac{1}{n^2} \sum_{i=1}^n \sigma^2 = \frac{\sigma^2}{n}.$$

Iz druge zveze vidimo, da standardna napaka $\sigma\bar{X} = \frac{\sigma}{\sqrt{n}}$ statistike \bar{X} pada z naraščanjem velikosti vzorca, tj. $\bar{X} \rightarrow \mu$; (enako nam zagotavlja tudi krepki zakon velikih števil).



Denimo, da se spremenljivka X na populaciji porazdeljuje normalno $N(\mu, \sigma)$. Na vsakem vzorcu (s ponavljanjem) izračunamo vzorčno aritmetično sredino \bar{X} . Dokazati se da, da je **porazdelitev vzorčnih aritmetičnih sredin** normalna, kjer je

- matematično upanje vzorčnih aritmetičnih sredin enako aritmetični sredini spremenljivke na populaciji

$$E(\bar{X}) = \mu,$$

- standardni odklon vzorčnih aritmetičnih sredin

$$SE(\bar{X}) = \frac{\sigma}{\sqrt{n}}.$$

Če tvorimo vzorce iz končne populacije brez vračanja, je standardni odklon vzorčnih aritmetičnih sredin

$$\text{SE}(\bar{X}) = \frac{\sigma}{\sqrt{n}} \sqrt{\frac{N-n}{N-1}}.$$

Za dovolj velike vzorce ($n > 30$) je porazdelitev vzorčnih aritmetičnih sredin približno normalna, tudi če spremenljivka X ni normalno porazdeljena. Če se statistika X porazdeljuje vsaj približno normalno s standardno napako $\text{SE}(X)$, potem se

$$Z = \frac{X - E(X)}{\text{SE}(X)}$$

porazdeljuje standardizirano normalno.

Vzorčno povprečje in normalna porazdelitev

Naj bo $X : N(\mu, \sigma)$. Tedaj je $\sum_{i=1}^n X_i : N(n\mu, \sigma\sqrt{n})$ in dalje $\bar{X} : N(\mu, \sigma/\sqrt{n})$. Tedaj je vzorčna statistika

$$Z = \frac{\bar{X} - \mu}{\sigma/\sqrt{n}} : N(0, 1)$$

Kaj pa če porazdelitev X ni normalna? Izračun porazdelitve se lahko zelo zaplete. Toda pri večjih vzorcih ($n > 30$), lahko uporabimo centralni limitni izrek, ki zagotavlja, da je spremenljivka Z porazdeljena skoraj standardizirano normalno. Vzorčno povprečje

$$\bar{X} = \frac{\sigma}{\sqrt{n}} Z + \mu$$

ima tedaj porazdelitev približno $N(\mu, \sigma/\sqrt{n})$.

Primer: Kolikšna je verjetnost, da bo pri 36 metih igralne kocke povprečno število pik večje ali enako 4? X je slučajna spremenljivka z vrednostmi 1,2,3,4,5,6 in verjetnostmi 1/6. Zanj je $\mu = 3,5$ in standardni odklon $\sigma = 1,7$. Vseh 36 ponovitev meta lahko obravnavamo kot slučajni vzorec velikost 36. Tedaj je

$$P(\bar{X} \geq 4) = P\left(Z \geq (4 - \mu)\sqrt{n}/\sigma\right) = P(Z \geq 1,75) \approx 0,04.$$

```
> x <- 1:6
> m <- mean(x)
> s <- sd(x)*sqrt(5/6)
> z <- (4-m)*6/s
> p <- 1-pnorm(z)
> cbind(m,s,z,p)
      m      s      z      p
[1,] 3.5 1.707825 1.75662 0.03949129
```

◇

12.4.2 (B) Vzorčna disperzija

Imejmo normalno populacijo $N(\mu, \sigma)$. Kako bi določili porazdelitev za vzorčno disperzijo ali popravljeno vzorčno disperzijo, tj.

$$S_0^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2 \quad \text{ozioroma} \quad S^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2?$$

Raje izračunamo porazdelitev za statistiko

$$\chi^2 = \frac{nS_0^2}{\sigma^2} = \frac{(n-1)S^2}{\sigma^2} = \frac{1}{\sigma^2} \sum_{i=1}^n (X_i - \bar{X})^2.$$

Preoblikujemo jo lahko takole:

$$\begin{aligned} \chi^2 &= \frac{1}{\sigma^2} \sum_{i=1}^n (X_i - \mu + \mu - \bar{X})^2 \\ &= \frac{1}{\sigma^2} \sum_{i=1}^n (X_i - \mu)^2 + \frac{2}{\sigma^2} (\mu - \bar{X}) \sum_{i=1}^n (X_i - \mu) + \frac{n}{\sigma^2} (\mu - \bar{X})^2 \end{aligned}$$

in, ker je $\sum_{i=1}^n (X_i - \mu) = n(\bar{X} - \mu) = -n(\mu - \bar{X})$, dalje

$$\chi^2 = \sum_{i=1}^n \left(\frac{X_i - \mu}{\sigma} \right)^2 - \frac{1}{n} \left(\sum_{i=1}^n \frac{X_i - \mu}{\sigma} \right)^2 = \sum_{i=1}^n Y_i^2 - \frac{1}{n} \left(\sum_{i=1}^n Y_i \right)^2,$$

kjer so Y_1, Y_2, \dots, Y_n paroma neodvisne standardizirano normalno porazdeljene slučajne spremenljivke, $Y_i = (X_i - \mu)/\sigma$. Porazdelitvena funkcija za χ^2 je

$$F_{\chi^2} = P(\chi^2 < z) = \iiint \dots \int_{\sum_{i=1}^n y_i^2 - \frac{1}{n} (\sum_{i=1}^n y_i)^2 < z} e^{-(y_1^2 + y_2^2 + \dots + y_n^2)/2} dy_n \dots dy_1,$$

z ustrezno ortogonalno transformacijo v nove spremenljivke z_1, z_2, \dots, z_n dobimo po nekaj računanja (glej Hladnik)

$$F_{\chi^2} = \frac{1}{(2\pi)^{(n-1)/2}} \iiint \dots \int_{\sum_{i=1}^{n-1} z_i^2 < z} e^{-(z_1^2 + z_2^2 + \dots + z_{n-1}^2)/2} dz_{n-1} \dots dz_1.$$

Pod integralom je gostota vektorja $(Z_1, Z_2, \dots, Z_{n-1})$ z neodvisnimi standardizirano normalnimi členi. Integral sam pa ustreza porazdelitveni funkciji vsote kvadratov

$$Z_1^2 + Z_2^2 + \dots + Z_{n-1}^2.$$

Tako je porazdeljena tudi statistika χ^2 . Kakšna pa je ta porazdelitev? Ker so tudi kvadrati $Z_1^2, Z_2^2, \dots, Z_{n-1}^2$ med seboj neodvisni in porazdeljeni po zakonu $\chi^2(1)$, je njihova vsota

porazdeljena po zakonu $\chi^2(n-1)$. Tako je torej porazdeljena tudi statistika χ^2 . Ker vemo, da je $E\chi^2(n) = n$ in $D\chi^2(n) = 2n$, lahko takoj izračunamo

$$ES_0^2 = E\frac{\sigma^2\chi^2}{n} = \frac{(n-1)\sigma^2}{n}, \quad ES^2 = E\frac{\sigma^2\chi^2}{n-1} = \sigma^2$$

in

$$DS_0^2 = D\frac{\sigma^2\chi^2}{n} = \frac{2(n-1)\sigma^4}{n^2} \quad DS^2 = D\frac{\sigma^2\chi^2}{n-1} = \frac{2\sigma^4}{n-1}$$

Če je n zelo velik, je po centralnem limitnem izreku statistika χ^2 porazdeljena približno normalno in sicer po zakonu

$$N(n-1, \sqrt{2(n-1)}),$$

vzorčna disperzija S_0^2 približno po

$$N\left(\frac{(n-1)\sigma^2}{n}, \frac{\sqrt{2(n-1)}\sigma^2}{n}\right)$$

in popravljena vzorčna disperzija S^2 približno po

$$N\left(\sigma^2, \sqrt{\frac{2}{n-1}}\sigma^2\right).$$

12.5 Nove porazdelitve

Pri normalno porazdeljeni slučajni spremenljivki X je tudi porazdelitev \bar{X} normalna, in sicer $N(\mu, \sigma/\sqrt{n})$. Statistika

$$Z = \frac{\bar{X} - \mu}{\sigma} \sqrt{n}$$

je potem porazdeljena standardizirano normalno.

Pri ocenjevanju parametra μ z vzorčnim povprečjem \bar{X} to lahko uporabimo le, če poznamo σ ; sicer ne moremo oceniti standardne napake – ne vemo, kako dobra je ocena za μ . Kaj lahko naredimo, če σ ne poznamo?

Parameter σ lahko ocenimo s S_0 ali S . *Toda* S je slučajna spremenljivka in porazdelitev statistike

$$\frac{\bar{X} - \mu}{S} \sqrt{n}$$

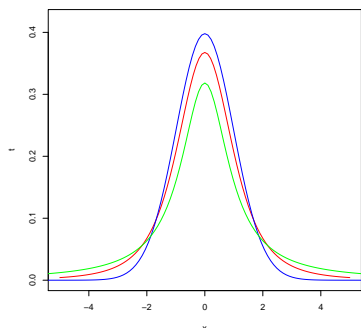
ni več normalna $N(0, 1)$ (razen, če je n zelo velik in S skoraj enak σ). Kakšna je porazdelitev nove vzorčne statistike

$$T = \frac{\bar{X} - \mu}{S} \sqrt{n} \text{ ?}$$



"Student" in 1908

12.5.1 Studentova porazdelitev



Leta 1908 je W.S. Gosset (1876-1937) pod psevdonimom 'Student' objavil članek, v katerem je pokazal, da ima statistika T porazdelitev $S(n-1)$ z gostoto

$$p(t) = \frac{\left(1 + \frac{t^2}{n-1}\right)^{-n/2}}{\sqrt{n-1} B\left(\frac{n-1}{2}, \frac{1}{2}\right)}$$

Tej porazdelitvi pravimo **Studentova porazdelitev** z $n-1$ prostostnimi stopnjami.

```
> plot(function(x) dt(x,df=3),-5,5,ylim=c(0,0.42),ylab="t",
  col="red")
> curve(dt(x,df=100),col="blue",add=T)
> curve(dt(x,df=1),col="green",add=T)
```

Za $S(1)$ dobimo Cauchyovo porazdelitev z gostoto

$$p(t) = \frac{1}{\pi(1+t^2)}$$

Za $n \rightarrow \infty$ pa gre

$$\frac{1}{\sqrt{n-1} B\left(\frac{n-1}{2}, \frac{1}{2}\right)} \rightarrow \sqrt{2\pi} \quad \text{in} \quad \left(1 + \frac{t^2}{n-1}\right)^{-n/2} \rightarrow e^{-t^2/2}.$$

Torej ima limitna porazdelitev gostoto

$$p(t) = \frac{1}{\sqrt{2\pi}} e^{-t^2/2}$$

standardizirane normalne porazdelitve.

Če zadnji sliki dodamo

```
> curve(dnorm(x),col="magenta",add=T)
```

ta pokrije modro krivuljo.



12.5.2 Fisherjeva porazdelitev

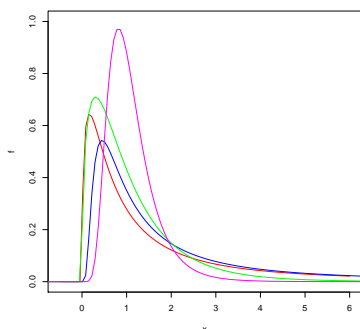
Poskusimo najti še porazdelitev kvocienta $Z = \frac{U}{V}$,

kjer sta $U : \chi^2(m)$ in $V : \chi^2(n)$ ter sta U in V neodvisni.

Z nekaj računanja (glej Hladnik) je mogoče pokazati, da je za $x > 0$ gostota ustrezne porazdelitve $F(m, n)$ enaka

$$p(x) = \frac{m^{m/2} n^{n/2}}{B\left(\frac{m}{2}, \frac{n}{2}\right)} \frac{x^{(m-2)/2}}{(n + mx)^{(m+n)/2}}$$

in je enaka 0 drugje.



Porazdelitvi $F(m, n)$ pravimo **Fisherjeva** ali tudi **Snedecorjeva porazdelitev F z (m, n) prostostnimi stopnjami**.

```
> plot(function(x) df(x,df1=3,df2=2), -0.5,6,ylim=c(0,1),ylab="f",
  col="red")
> curve(df(x,df1=20,df2=2), col="blue", add=T)
> curve(df(x,df1=3,df2=20), col="green", add=T)
> curve(df(x,df1=20,df2=20), col="magenta", add=T)
```

Po zakonu $F(m-1, n-1)$ je na primer porazdeljena statistika

$$F = \frac{S_X^2/\sigma_X^2}{S_Y^2/\sigma_Y^2}$$

saj vemo, da sta spremenljivki

$$U = (m-1)S_X^2/\sigma_X^2 \quad \text{in} \quad V = (n-1)S_Y^2/\sigma_Y^2$$

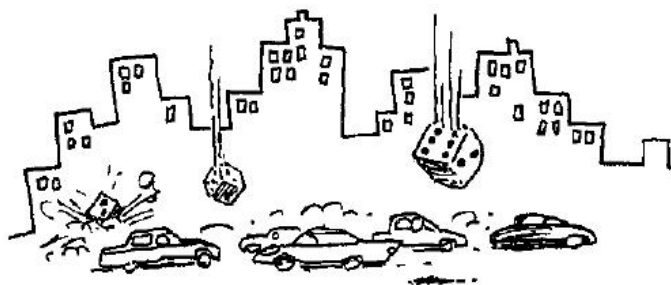
porazdeljeni po χ^2 z $m-1$ oziroma $n-1$ prostostnimi stopnjami in sta neodvisni. Velja še:

če je $U : F(m, n)$, je $1/U : F(n, m)$,

če je $U : S(n)$, je $U^2 : F(1, n)$.

Poglavje 13

Cenilke



13.1 Osnovni pojmi

Točkovna cenilka je pravilo ali formula, ki nam pove, kako izračunati numerično oceno parametra populacije na osnovi merjenj vzorca.

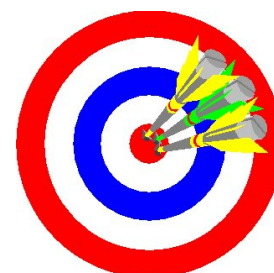
Število, ki je rezultat izračuna, se imenuje **točkovna ocena** (in mu ne moremo zaupati v smislu verjetnosti).

Cenilka parametra ζ je vzorčna statistika $C = C(X_1, \dots, X_n)$, katere porazdelitveni zakon je odvisen le od parametra ζ , njene vrednosti pa ležijo v prostoru parametrov. Seveda je odvisna tudi od velikosti vzorca n .

Primer: Vzorčna mediana \tilde{X} in vzorčno povprečje \bar{X} sta cenilki za populacijsko povprečje μ ; popravljena vzorčna disperzija S^2 pa je cenilka za populacijsko disperzijo σ^2 . \diamond

Cenilka C parametra ζ (grška črka zeta) je **dosledna**, če z rastočim n zaporedje C_n verjetnostno konvergira k parametru ζ , tj. za vsak $\varepsilon > 0$ velja

$$\lim_{n \rightarrow \infty} P(|C_n - \zeta| < \varepsilon) = 1.$$



Primeri: vzorčno povprečje \bar{X} je dosledna cenilka za populacijsko povprečje μ . Tudi vsi **vzorčni začetni momenti**

$$Z_k = \frac{1}{n} \sum_{i=1}^n X_i^k$$

so dosledne cenilke ustreznih začetnih populacijskih momentov $z_k = EX^k$, če le-ti obstajajo.

Primer: Vzorčna mediana \tilde{X} je dosledna cenilka za populacijsko mediano. \diamond

Trditev 13.1. Če pri pogoju $n \rightarrow \infty$ velja $EC_n \rightarrow \zeta$ in $DC_n \rightarrow 0$, je C_n dosledna cenilka parametra ζ .

Dokaz. To sprevidimo takole:

$$1 - P(|C_n - \zeta| < \varepsilon) = P(|C_n - \zeta| \geq \varepsilon) \leq P(|C_n - EC_n| + |EC_n - \zeta| \geq \varepsilon),$$

upoštevajmo še, da za dovolj velike n velja $|EC_n - \zeta| < \varepsilon/2$, in uporabimo neenakost Čebiševa

$$P(|C_n - EC_n| \geq \varepsilon/2) \leq \frac{4DC_n}{\varepsilon^2} \rightarrow 0. \quad \square$$

Primer: Naj bo $X : N(\mu, \sigma)$. Ker za $n \rightarrow \infty$ velja

$$ES_0^2 = \frac{(n-1)\sigma^2}{n} \rightarrow \sigma^2 \quad \text{in} \quad DS_0^2 = \frac{2(n-1)\sigma^4}{n^2} \rightarrow 0,$$

je vzorčna disperzija S_0^2 dosledna cenilka za σ^2 . \diamond

Nepristrana cenilka z najmanjšo varianco

Cenilka C_n parametra ζ je **nepristranska**, če je $EC_n = \zeta$ (za vsak n); in je **asimptotično nepristranska**, če je $\lim_{n \rightarrow \infty} EC_n = \zeta$. Količino $B(C_n) = EC_n - \zeta$ imenujemo **pristranost** (angl. *bias*) cenilke C_n .

Primer: Vzorčno povprečje \bar{X} je nepristranska cenilka za populacijsko povprečje μ ; vzorčna disperzija S_0^2 je samo asimptotično nepristranska cenilka za σ^2 , popravljena vzorčna disperzija S^2 pa je nepristranska cenilka za σ^2 . \diamond

Disperzija nepristranskih cenilk

Izmed nepristranskih cenilk istega parametra ζ je boljša tista, ki ima manjšo disperzijo – v povprečju daje bolj točne ocene. Če je razred cenilk parametra ζ *konveksen* (vsebuje tudi njihove konveksne kombinacije), obstaja v bistvu ena sama cenilka z najmanjšo disperzijo: Naj bo razred nepristranskih cenilk parametra ζ konveksen. Če sta C in C' nepristranski cenilki, obe z najmanjšo disperzijo σ^2 , je $C = C'$ z verjetnostjo 1. Za to poglejmo

$$D((C + C')/2) = \frac{1}{4}(DC + DC' + 2\text{Cov}(C, C')) \leq \left(\frac{1}{2}(\sqrt{DC} + \sqrt{DC'})\right)^2 = \sigma^2.$$

Ker sta cenilki minimalni, mora biti tudi $D((C + C')/2) = \sigma^2$ in dalje $\text{Cov}(C, C') = \sigma^2$ oziroma $r(C, C') = 1$. Torej je $C' = aC + b$, $a > 0$ z verjetnostjo 1. Iz $DC = DC'$ izhaja $a = 1$, iz $EC = EC'$ pa še $b = 0$.

Srednja kvadratična napaka

Včasih je celo bolje vzeti pristransko cenilko z manjšo disperzijo, kot jo ima druga, sicer nepristranska, cenilka z veliko disperzijo. Mera *učinkovitosti* cenilk parametra ζ je *srednja kvadratična napaka*

$$q(C) = E(C - \zeta)^2.$$

Ker velja

$$q(C) = E(C - EC + EC - \zeta)^2 = E(C - EC)^2 + (EC - \zeta)^2,$$

jo lahko zapišemo tudi v obliki

$$q(C) = DC + B(C)^2.$$

Za nepristranske cenilke je $B(C) = 0$ in zato $q(C) = DC$. Če pa je disperzija cenilke skoraj 0, je $q(C) \approx B(C)^2$.

13.2 Rao-Cramérjeva ocena

Naj bo f gostotna ali verjetnostna funkcija slučajne spremenljivke X in naj bo odvisna še od parametra ζ , tako da je $f(x; \zeta)$ njena vrednost v točki x . Združeno gostotno ali verjetnostno funkcijo slučajnega vzorca (X_1, \dots, X_n) označimo z L in ji pravimo *funkcija verjetja* (tudi *zanesljivosti*, angl. *likelihood*)

$$L(x_1, \dots, x_n; \zeta) = f(x_1; \zeta) \cdots f(x_n; \zeta).$$

Velja

$$\int \int \dots \int L(x_1, \dots, x_n; \zeta) dx_1 \dots dx_n = 1. \quad (13.1)$$

$L(X_1, \dots, X_n)$ je funkcija vzorca – torej slučajna spremenljivka. Privzemimo, da je funkcija L vsaj dvakrat zvezno odvedljiva po ζ na nekem intervalu I in naj na tem intervalu tudi integral odvoda L po ζ enakomerno konvergira. Odvajajmo enakost (13.1) po ζ in upoštevajmo $\frac{\partial \ln L}{\partial \zeta} = \frac{1}{L} \frac{\partial L}{\partial \zeta}$ pa dobimo

$$\int \int \dots \int \frac{\partial \ln L}{\partial \zeta} L dx_1 dx_2 \dots dx_n = 0,$$

kar lahko tolmačimo kot

$$\mathbb{E} \frac{\partial \ln L}{\partial \zeta} = 0.$$

Naj bo sedaj C nepristranska cenilka parametra ζ , torej $\mathbb{E}C = \zeta$, oziroma zapisano z integrali $\int \int \dots \int CL dx_1 dx_2 \dots dx_n = \zeta$. Ker C ni odvisna od ζ , dobimo z odvajanjem po parametru ζ :

$$\int \int \dots \int C \frac{\partial \ln L}{\partial \zeta} L dx_1 dx_2 \dots dx_n = 1$$

kar z drugimi besedami pomeni

$$\mathbb{E} \left(C \frac{\partial \ln L}{\partial \zeta} \right) = 1.$$

Če to enakost združimo s prejšnjo (pomnoženo s ζ), dobimo:

$$\mathbb{E} \left((C - \zeta) \frac{\partial \ln L}{\partial \zeta} \right) = 1.$$

Od tu po $(\mathbb{E}XY)^2 \leq \mathbb{E}X^2 \mathbb{E}Y^2$ izhaja naprej

$$1 = \left(\mathbb{E} \left((C - \zeta) \frac{\partial \ln L}{\partial \zeta} \right) \right)^2 \leq \mathbb{E}(C - \zeta)^2 \mathbb{E} \left(\frac{\partial \ln L}{\partial \zeta} \right)^2 = DC \mathbb{E} \left(\frac{\partial \ln L}{\partial \zeta} \right)^2,$$

kar da *Rao-Cramérjevo oceno*

$$DC \geq \left(\mathbb{E} \left(\frac{\partial \ln L}{\partial \zeta} \right)^2 \right)^{-1} = \left(-\mathbb{E} \frac{\partial^2 \ln L}{\partial \zeta^2} \right)^{-1} = \left(n \mathbb{E} \left(\frac{\partial \ln f}{\partial \zeta} \right)^2 \right)^{-1}.$$

13.3 Učinkovitost cenilk

Rao-Cramérjeva ocena da absolutno spodnjo mejo disperzije za vse nepristranske cenilke parametra ζ (v dovolj gladkih porazdelitvah). Ta meja ni nujno dosežena. Cenilka, ki jo doseže, se imenuje *najučinkovitejša cenilka* parametra ζ in je ena sama (z verjetnostjo 1).

Kdaj pa je ta spodnja meja dosežena?

V neenakosti $(EY)^2 \leq EX^2EY^2$, ki je uporabljena v izpeljavi Rao-Cramérjeve ocene, velja enakost natanko takrat, ko je $Y = cX$ z verjetnostjo 1. Torej velja v Rao-Cramérjevi oceni enakost natanko takrat, ko je

$$\frac{\partial \ln L}{\partial \zeta} = A(\zeta)(C - \zeta),$$

kjer je $A(\zeta)$ konstanta, odvisna od ζ in neodvisna od vzorca. Zato je tudi

$$(DC)^{-1} = E\left(\frac{\partial \ln L}{\partial \zeta}\right)^2 = A(\zeta)^2 E(C - \zeta)^2 = A(\zeta)^2 DC$$

oziroma končno

$$DC = |A(\zeta)|^{-1}.$$

Najučinkovitejše cenilke za parametre normalne porazdelitve

Naj bo $X : N(\mu, \sigma)$. Tedaj je

$$L = \frac{1}{(2\pi)^{n/2}\sigma^n} e^{-((\frac{X_1-\mu}{\sigma})^2 + \dots + (\frac{X_n-\mu}{\sigma})^2)/2}$$

in

$$\ln L = \ln \frac{1}{(2\pi)^{n/2}\sigma^n} - \frac{1}{2} \left(\left(\frac{X_1 - \mu}{\sigma}\right)^2 + \dots + \left(\frac{X_n - \mu}{\sigma}\right)^2 \right)$$

ter dalje

$$\frac{\partial \ln L}{\partial \mu} = \frac{X_1 - \mu}{\sigma^2} + \dots + \frac{X_n - \mu}{\sigma^2} = \frac{n}{\sigma^2} (\bar{X} - \mu).$$

Torej je vzorčno povprečje \bar{X} najučinkovitejša cenilka za μ z disperzijo $D\bar{X} = \sigma^2/n$. Prvi člen v izrazu za $\ln L$ lahko zapišemo tudi $-\frac{n}{2}(\ln 2\pi + \ln \sigma^2)$. Tedaj je, če privzamemo, da je μ znano število

$$\frac{\partial \ln L}{\partial \sigma^2} = -\frac{n}{2\sigma^2} + \frac{1}{2\sigma^4} ((X_1 - \mu)^2 + \dots + (X_n - \mu)^2) = \frac{n}{2\sigma^4} (S_\mu^2 - \sigma^2).$$

To pomeni, da je $S_\mu^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \mu)^2$ najučinkovitejša cenilka za parameter σ^2 z disperzijo $DS_\mu^2 = 2\sigma^4/n$.

Primer: Za Poissonovo porazdelitev $P(\lambda)$ s parametrom λ , tj. $p_k = \lambda^k e^{-\lambda}/k!$, je

$$L = e^{-n\lambda} \frac{\lambda^{x_1 + \dots + x_n}}{x_1! \cdots x_n!}$$

in dalje

$$\ln L = -n\lambda + (x_1 + \dots + x_n) \ln \lambda - \ln(x_1! \cdots x_n!)$$

ter končno

$$\frac{\partial \ln L}{\partial \lambda} = -n + \frac{x_1 + \dots + x_n}{\lambda} = \frac{n}{\lambda}(\bar{X} - \lambda). \quad (13.2)$$

Najučinkovitejša cenilka za parameter λ je \bar{X} z disperzijo $D\bar{X} = \lambda/n$. \diamond

Naj bo C_0 najučinkovitejša cenilka parametra ζ in C kaka druga nepristranska cenilka. Tedaj je *učinkovitost* cenilke C določena s predpisom

$$e(C) = \frac{DC_0}{DC}.$$

Učinkovitost najučinkovitejše cenilke je $e(C_0) = 1$. Če najučinkovitejša cenilka ne obstaja, vzamemo za vrednost DC_0 desno stran v Rao-Cramérjevi oceni.

Primer: Naj bo $X : N(\mu, \sigma)$. Pri velikih n -jih je vzorčna mediana \tilde{X} – ocena za μ , porazdeljena približno po $N(\mu, \sigma\sqrt{\pi/2n})$. Torej je

$$e(\tilde{X}) = \frac{D\bar{X}}{D\tilde{X}} = \frac{\frac{\sigma^2}{n}}{\frac{\pi\sigma^2}{2n}} = \frac{2}{\pi} \approx 0,64. \quad \diamond$$

Primer: Naj bo $X : N(\mu, \sigma)$. Če poznamo μ , je najučinkovitejša cenilka za σ^2 statistika S_μ^2 z disperzijo $DS_\mu^2 = 2\sigma^4/n$. Popravljen vzorčna disperzija S^2 pa je nepristranska cenilka istega parametra z disperzijo $DS^2 = 2\sigma^4/(n-1)$. Tokrat je

$$e(S^2) = \frac{DS_\mu^2}{DS^2} = \frac{\frac{2\sigma^4}{n}}{\frac{2\sigma^4}{n-1}} = \frac{n-1}{n}.$$

Iz tega vidimo, da $e(S^2) \rightarrow 1$, ko $n \rightarrow \infty$. Pravimo, da je cenilka S^2 *asimptotično najučinkovitejša cenilka* za σ^2 . \diamond

13.4 Metoda momentov

Recimo, da je za zvezno slučajno spremenljivko X njena gostota f odvisna od m parametrov $f(x; \zeta_1, \dots, \zeta_m)$ in naj obstajajo momenti

$$z_k = z_k(\zeta_1, \dots, \zeta_m) = \int_{-\infty}^{\infty} x^k f(x; \zeta_1, \dots, \zeta_m) dx \quad \text{za } k = 1, \dots, m.$$

Če se dajo iz teh enačb enolično izračunati parametri ζ_1, \dots, ζ_m kot funkcije momentov z_1, \dots, z_m :

$$\zeta_k = \varphi_k(z_1, \dots, z_m),$$

potem so

$$C_k = \varphi_k(Z_1, \dots, Z_m)$$

cenilke parametrov ζ_k po *metodi momentov*. k -ti vzorčni začetni moment $Z_k = \frac{1}{n} \sum_{i=1}^n X_i^k$ je cenilka za ustrežni populacijski moment z_k . Cenilke, ki jih dobimo po metodi momentov so dosledne.

Primer: Naj bo $X : N(\mu, \sigma)$. Tedaj je $z_1 = \mu$ in $z_2 = \sigma^2 + \mu^2$. Od tu dobimo $\mu = z_1$ in $\sigma^2 = z_2 - z_1^2$. Ustrežni cenilki sta

$$Z_1 = \bar{X} \quad \text{za } \mu \quad \text{in} \quad Z_2 - Z_1^2 = \overline{X^2} - \bar{X}^2 = S_0^2 \quad \text{za } \sigma^2,$$

torej vzorčno povprečje in disperzija. ◇

Metoda največjega verjetja

Funkcija verjetja

$$L(x_1, \dots, x_n; \zeta) = f(x_1; \zeta) \cdots f(x_n; \zeta)$$

je pri danih x_1, \dots, x_n odvisna še od parametra ζ . Izberemo tak ζ , da bo funkcija L dosegla največjo vrednost. Če je L vsaj dvakrat zvezno odvedljiva, mora veljati

$$\frac{\partial L}{\partial \zeta} = 0 \quad \text{in} \quad \frac{\partial^2 L}{\partial \zeta^2} < 0.$$

Največja vrednost parametra je še odvisna od x_1, \dots, x_n : $\zeta_{\max} = \varphi(x_1, \dots, x_n)$. Tedaj je cenilka za parameter ζ enaka

$$C = \varphi(X_1, \dots, X_n).$$

Metodo lahko posplošimo na večje število parametrov. Pogosto raje iščemo maksimum funkcije $\ln L$. Če najučinkovitejša cenilka obstaja, jo dobimo s to metodo.

Primer: Naj bo $X : B(1, p)$. Tedaj je $f(x; p) = p^x(1 - p)^{1-x}$, kjer je $x = 0$ ali $x = 1$. Ocenjujemo parameter p . Funkcija verjetja ima obliko $L = p^x(1 - p)^{n-x}$, kjer je sedaj $x \in \{0, \dots, n\}$. Ker je $\ln L = x \ln p + (n - x) \ln(1 - p)$, dobimo

$$\frac{\partial \ln L}{\partial p} = \frac{x}{p} - \frac{n - x}{1 - p},$$

ki je enak 0 pri $p = x/n$. Ker je v tem primeru

$$\frac{\partial^2 \ln L}{\partial p^2} = -\frac{x}{p^2} - \frac{n - x}{(1 - p)^2} < 0,$$

je v tej točki maksimum. Cenilka po metodi največjega verjetja je torej $P = X/n$, kjer je X binomsko porazdeljena spremenljivka – frekvenca v n ponovitvah. Cenilka P je nepristranska, saj je $EP = EX/n = p$. Ker gre $DP = DX/n^2 = p(1 - p)/n \rightarrow 0$ za $n \rightarrow \infty$, je P dosledna cenilka. Je pa tudi najučinkovitejša

$$\frac{\partial \ln L}{\partial p} = \frac{X}{p} - \frac{n - X}{1 - p} = \frac{n}{p(1 - p)} \left(\frac{X}{n} - p \right) = \frac{n}{p(1 - p)} (P - p). \quad \diamond$$

Primer: Nadaljujmo primer Poissonove porazdelitve. Odvod (13.2) je enak 0 za $\lambda = \bar{X}$. Drugi odvod v tej točki je

$$\frac{\partial^2 \ln L}{\partial \lambda^2} = -\frac{x_1 + \dots + x_n}{\lambda^2} < 0,$$

kar pomeni, da je v tej točki maksimum. Cenilka za λ je po metodi največjega verjetja vzorčno povprečje \bar{X} . Je tudi najučinkovitejša cenilka za λ z disperzijo $D\bar{X} = \lambda/n$. \diamond

12.4 Vzorčna statistika (nadaljevanje)

12.4.3 (C) Vzorčne aritmetične sredine

Primer: Denimo, da se spremenljivka inteligenčni kvocient na populaciji porazdeljuje normalno z aritmetično sredino $\mu = 100$ in standardnim odklonom $\sigma = 15$, tj.

$$X : N(100, 15)$$

Denimo, da imamo vzorec velikosti $n = 225$. Tedaj se vzorčne aritmetične sredine porazdeljujejo normalno

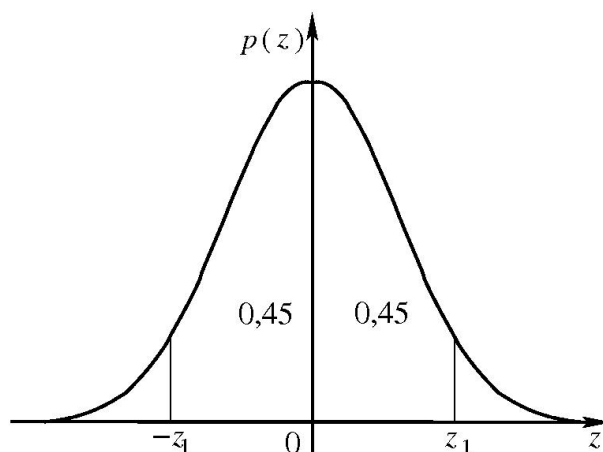
$$\bar{X} : N\left(100, \frac{15}{\sqrt{225}}\right) = N(100, 1)$$

Izračunajmo, kolikšne vzorčne aritmetične sredine ima 90% vzorcev (simetrično na povprečje). 90% vzorčnih aritmetičnih sredin se nahaja na intervalu:

$$P(\bar{X}_1 < \bar{X} < \bar{X}_2) = 0,90$$

$$P(-z_1 < z < z_1) = 0,90 \implies 2\Phi(z_1) = 0,90$$

$$\Phi(z_1) = 0,45 \implies z_1 = 1,65$$



Potem se vzorčne aritmetične sredine nahajajo v intervalu

$$P\left(\mu - z_1 \frac{\sigma}{\sqrt{n}} < \bar{X} < \mu + z_1 \frac{\sigma}{\sqrt{n}}\right) = 0,90$$

oziroma konkretno

$$P\left(100 - 1,65 < \bar{X} < 100 + 1,65\right) = 0,90$$

90% vseh slučajnih vzorcev velikosti 225 enot bo imelo povprečja za inteligenčni kvocient na intervalu

$$(98,35; 101,65).$$

Lahko preverimo, da bi bil ta interval v primeru večjega vzorca ožji. Npr. v primeru vzorcev velikosti $n = 2500$ je ta interval

$$P\left(100 - 1,65 \frac{15}{\sqrt{2500}} < \bar{X} < 100 + 1,65 \frac{15}{\sqrt{2500}}\right) = 0,90$$

oziroma

$$(99,5; 100,5).$$

◇

12.4.4 (D) Vzorčni deleži

Denimo, da želimo na populaciji oceniti delež enot π z določeno lastnostjo.



Zato na vsakem vzorcu poiščemo vzorčni delež p . Pokazati se da, da se za dovolj velike slučajne vzorce s ponavljanjem (za deleže okoli 0,5 je dovolj 20 enot ali več) vzorčni deleži porazdeljujejo približno normalno z

- aritmetično sredino vzorčnih deležev, ki je enaka deležu na populaciji

$$E p = \pi,$$

- standardnim odklonom vzorčnih deležev

$$SE(p) = \sqrt{\frac{\pi(1-\pi)}{n}}.$$

Za manjše vzorce se vzorčni deleži porazdeljujejo binomsko. Cenilka populacijskega deleža je nepristranska cenilka, ker velja $E p = \pi$.

Primer: V izbrani populaciji prebivalcev je polovica žensk $\pi = 0,5$. Če tvorimo vzorce

po $n = 25$ enot, nas zanima, kolikšna je verjetnost, da je v vzorcu več kot 55 % žensk? To pomeni, da iščemo verjetnost $P(p > 0,55)$. Vzorčni deleži p se porazdeljujejo približno normalno:

$$p : N\left(0,5, \sqrt{\frac{\pi(1-\pi)}{n}}\right) = N\left(0,5, \sqrt{\frac{0,5 \cdot 0,5}{25}}\right) = N(0,5, 0,1).$$

Zato je

$$\begin{aligned} P(p > 0,55) &= P\left(Z > \frac{0,55 - 0,5}{0,1}\right) = P(Z > 0,5) \\ &= 0,5 - \Phi(0,5) = 0,5 - 0,1915 = 0,3085. \end{aligned}$$

Rezultat pomeni, da lahko pričakujemo, da bo pri približno 31% vzorcev delež žensk večji od 0,55. Poglejmo, kolikšna je ta verjetnost, če bi tvorili vzorce velikosti $n = 2500$ enot:

$$P(p > 0,55) = P\left(Z > \frac{0,55 - 0,5}{\sqrt{\frac{0,5(1-0,5)}{2500}}}\right) = P(Z > 5) = 0,5 - \Phi(5) = 0,5 - 0,5 = 0.$$

V 10-krat večjih vzorcih kot prej ne moremo pričakovati več kot 55% žensk. \diamond

12.4.5 (E) Razlika vzorčnih aritmetičnih sredin

Denimo, da imamo dve populaciji velikosti N_1 in N_2 in se spremenljivka X na prvi populaciji porazdeljuje normalno $N(\mu_1, \sigma)$, na drugi populaciji pa $N(\mu_2, \sigma)$ (standardna odklona sta na obeh populacijah enaka!). V vsaki od obeh populacij tvorimo neodvisno slučajne vzorce velikosti n_1 in n_2 . Na vsakem vzorcu (s ponavljanjem) prve populacije izračunamo vzorčno aritmetično sredino \bar{X}_1 in podobno na vsakem vzorcu druge populacije \bar{X}_2 . Dokazati se da, da je porazdelitev razlik vzorčnih aritmetičnih sredin normalna, kjer je

- matematično upanje razlik vzorčnih aritmetičnih sredin enako

$$\mathbf{E}(\bar{X}_1 - \bar{X}_2) = \mathbf{E}\bar{X}_1 - \mathbf{E}\bar{X}_2 = \mu_1 - \mu_2,$$

- disperzija razlik vzorčnih aritmetičnih sredin enaka

$$\mathbf{D}(\bar{X}_1 - \bar{X}_2) = \mathbf{D}\bar{X}_1 + \mathbf{D}\bar{X}_2 = \frac{\sigma^2}{n_1} + \frac{\sigma^2}{n_2} = \sigma^2 \cdot \frac{n_1 + n_2}{n_1 n_2}.$$

Primer: Dvema populacijama študentov na neki univerzi (tehnikom in družboslovcem) so izmerili neko sposobnost s povprečjema $\mu_t = 70$ in $\mu_d = 80$ točk in standardnim odklonom, ki je na obeh populacijah enak, $\sigma = 7$ točk.

Kolikšna je verjetnost, da je aritmetična sredina slučajnega vzorca družboslovcev ($n_d = 36$) večja za več kot 12 točk od aritmetične sredine vzorca tehnikov ($n_t = 64$)? Zanima nas torej verjetnost:

$$\begin{aligned} P(\bar{X}_d - \bar{X}_t > 12) &= P\left(Z > \frac{12 - 10}{7\sqrt{\frac{36+64}{36 \cdot 64}}}\right) = P(Z > 1,37) = 0,5 - \Phi(1,37) \\ &= 0,5 - 0,4147 = 0,0853. \end{aligned}$$

Torej, približno 8,5% parov vzorcev je takih, da je povprečje družboslovcev večje od povprečja tehnikov za 12 točk. \diamond

12.4.6 (F) Razlika vzorčnih deležev

Podobno kot pri porazdelitvi razlik vzorčnih aritmetičnih sredin naj bosta dani dve populaciji velikosti N_1 in N_2 z deležema enot z neko lastnostjo π_1 in π_2 . Iz prve populacije tvorimo slučajne vzorce velikosti n_1 in na vsakem izračunamo delež enot s to lastnostjo p_1 .

Podobno naredimo tudi na drugi populaciji: tvorimo slučajne vzorce velikosti n_2 in na njih določimo deleže p_2 . Pokazati se da, da se za dovolj velike vzorce razlike vzorčnih deležev porazdeljujejo približno normalno z

- matematičnim upanjem razlik vzorčnih deležev

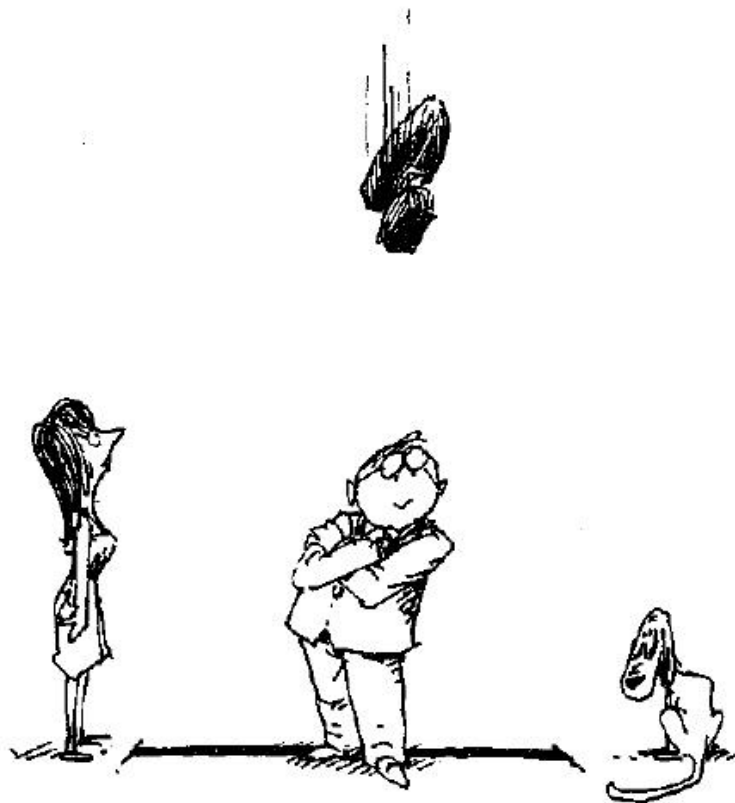
$$E(p_1 - p_2) = Ep_1 - Ep_2 = \pi_1 - \pi_2,$$

- disperzijo razlik vzorčnih deležev

$$D(p_1 - p_2) = Dp_1 + Dp_2 = \frac{\pi_1(1 - \pi_1)}{n_1} + \frac{\pi_2(1 - \pi_2)}{n_2}.$$

Poglavje 14

Intervali zaupanja



Denimo, da s slučajnim vzorcem ocenjujemo parameter γ . Poskušamo najti statistiko g , ki je nepristranska, tj. $Eg = \gamma$ in se na vseh možnih vzorcih vsaj približno normalno porazdeljuje s standardno napako $SE(g)$. Nato poskušamo najti interval, v katerem se bo z dano gotovostjo $(1 - \alpha)$ nahajal ocenjevani parameter:

$$P(a < \gamma < b) = 1 - \alpha,$$

kjer je a je spodnja meja zaupanja, b je zgornja meja zaupanja, α verjetnost tveganja oziroma $1 - \alpha$ verjetnost gotovosti. Ta interval imenujemo **interval zaupanja** in ga interpretiramo takole: z verjetnostjo tveganja α se parameter γ nahaja v tem intervalu.

Konstruirajmo interval zaupanja. Na osnovi omenjenih predpostavk o porazdelitvi statistike g lahko zapišemo, da se statistika

$$Z = \frac{g - \mathbb{E}g}{\text{SE}(g)} = \frac{g - \gamma}{\text{SE}(g)}$$

porazdeljuje standardizirano normalno $N(0, 1)$. Če tveganje α porazdelimo polovico na levo in polovico na desno na konce normalne porazdelitve, lahko zapišemo

$$P\left(-z_{\alpha/2} < \frac{g - \gamma}{\text{SE}(g)} < z_{\alpha/2}\right) = 1 - \alpha.$$

Po ustrezni preureditvi lahko izpeljemo naslednji interval zaupanja za parameter γ

$$P\left(g - z_{\alpha/2} \text{SE}(g) < \gamma < g + z_{\alpha/2} \text{SE}(g)\right) = 1 - \alpha$$

$z_{\alpha/2}$ je določen le s stopnjo tveganja α . Vrednosti $z_{\alpha/2}$ lahko razberemo iz tabele za verjetnosti za standardizirano normalno porazdelitev, ker velja

$$\Phi(z_{\alpha/2}) = 0,5 - \frac{\alpha}{2}$$

Podajmo vrednost $z_{\alpha/2}$ za nekaj najbolj standardnih tveganj:

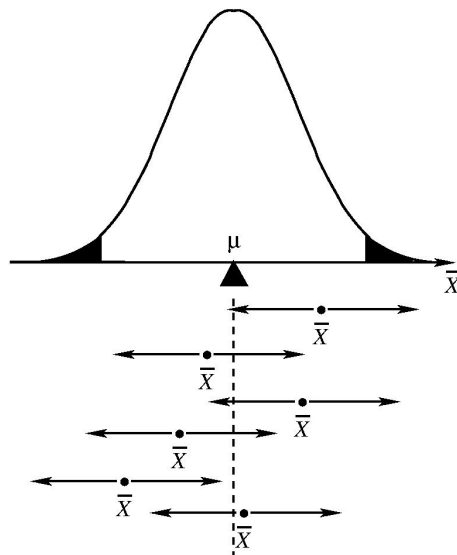
- $\alpha = 0,10$, $z_{\alpha/2} = 1,65$
- $\alpha = 0,05$, $z_{\alpha/2} = 1,96$
- $\alpha = 0,01$, $z_{\alpha/2} = 2,58$

14.1 Pomen stopnje tveganja

Za vsak slučajni vzorec lahko ob omenjenih predpostavkah izračunamo ob izbrani stopnji tveganja α interval zaupanja za parameter γ . Ker se podatki vzorcev razlikujejo, se razlikujejo vzorčne ocene parametrov in zato tudi izračunani intervali zaupanja za parameter γ . To pomeni, da se intervali zaupanja od vzorca do vzorca razlikujejo. Meji intervala sta slučajni spremenljivki.

Primer: Vzemimo stopnjo tveganja $\alpha = 0,05$. Denimo, da smo izbrali 100 slučajnih vzorcev in za vsakega izračunali interval zaupanja za parameter γ . Tedaj lahko pričakujemo,

da 5 intervalov zaupanja od 100 ne bo pokrilo iskanega parametra γ . Povedano je lepo predstavljeno tudi grafično:



V tem primeru ocenjujemo parameter aritmetično sredino inteligenčnega kvocienta. Kot vemo, se vzorčne aritmetične sredine \bar{X} za dovolj velike vzorce porazdeljujejo normalno. Denimo, da v tem primeru poznamo vrednost parametra ($\mu = 100$). Za več slučajnih vzorcev smo izračunali in prikazali interval zaupanja za μ ob stopnji tveganja $\alpha = 0,05$. Predstavitve več intervalov zaupanja za aritmetično sredino μ pri 5% stopnji tveganja: približno 95% intervalov pokrije parameter μ . \diamond

14.2 Intervalsko ocenjevanje parametrov

Naj bo X slučajna spremenljivka na populaciji G z gostoto verjetnosti odvisno od parametra ζ . Slučajna množica $M \subset \mathbb{R}$, ki je odvisna le od slučajnega vzorca, ne pa od parametra ζ , se imenuje *množica zaupanja* za parameter ζ , če obstaja tako število α , $0 < \alpha < 1$, da velja $P(\zeta \in M) = 1 - \alpha$. Število $1 - \alpha$ imenujemo tedaj *stopnja zaupanja*; število α pa *stopnja tveganja*. Stopnja zaupanja je običajno 95% ali 99% – $\alpha = 0,05$ ali $\alpha = 0,01$. Pove nam, kakšna je verjetnost, da M vsebuje vrednost parametra ζ ne glede na to, kakšna je njegova dejanska vrednost. Če je množica M interval $M = [A, B]$, ji rečemo *interval zaupanja* (za parameter ζ). Njegovi krajišči sta funkciji slučajnega vzorca – torej statistiki.

Naj bo $X : N(\mu, \sigma)$ in recimo, da poznamo parameter σ in ocenjujemo parameter μ . Izberimo konstanti a in b , $b > a$, tako da bo $P(a \leq Z \leq b) = 1 - \alpha$, kjer je $Z = \frac{\bar{X} - \mu}{\sigma} \sqrt{n}$.

Tedaj je

$$P\left(\bar{X} - \frac{b\sigma}{\sqrt{n}} \leq \mu \leq \bar{X} - \frac{a\sigma}{\sqrt{n}}\right) = 1 - \alpha$$

Označimo $A = \bar{X} - b\sigma/\sqrt{n}$ in $B = \bar{X} - a\sigma/\sqrt{n}$.

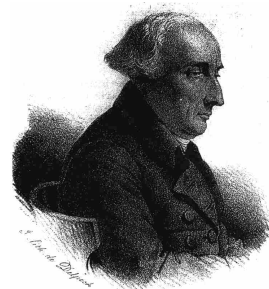
Za katera a in b je interval $[A, B]$ najkrajši?

Pokazati je mogoče (Lagrangeova funkcija), da mora biti $a = -b$ in $\Phi(b) = (1 - \alpha)/2$; oziroma, če označimo $b = z_{\alpha/2}$, velja $P(Z > z_{\alpha/2}) = \alpha/2$. Iskani interval je torej

$$A = \bar{X} - z_{\alpha/2} \frac{\sigma}{\sqrt{n}} \quad \text{in} \quad B = \bar{X} + z_{\alpha/2} \frac{\sigma}{\sqrt{n}}$$

tj., z verjetnostjo $1 - \alpha$ je $|\bar{X} - \mu| < z_{\alpha/2}\sigma/\sqrt{n}$.

Od tu dobimo, da mora za to, da bo napaka manjša od ε z verjetnostjo $1 - \alpha$, veljati $n > (z_{\alpha/2}\sigma/\varepsilon)^2$.



Če pri porazdelitvi $X : N(\mu, \sigma)$ tudi parameter σ ni znan, ga nadomestimo s cenilko S in moramo zato uporabiti Studentovo statistiko $T = \frac{\bar{X} - \mu}{S} \sqrt{n}$. Ustrezni interval je tedaj

$$A = \bar{X} - t_{\alpha/2} \frac{S}{\sqrt{n}}, \quad B = \bar{X} + t_{\alpha/2} \frac{S}{\sqrt{n}}$$

kjer je $P(T > t_{\alpha/2}) = \alpha/2$. Če pa bi ocenjevali parameter σ^2 , uporabimo statistiko $\chi^2 = (n - 1)S^2/\sigma^2$, ki je porazdeljena po $\chi^2(n - 1)$. Tedaj je

$$A = \frac{(n - 1)S^2}{b} \quad \text{in} \quad B = \frac{(n - 1)S^2}{a}$$

Konstanti a in b včasih določimo iz pogojev

$$P(\chi^2 < a) = \alpha/2 \quad \text{in} \quad P(\chi^2 > b) = \alpha/2,$$

najkrajši interval pa dobimo, ko velja zveza $a^2p(a) = b^2p(b)$ in seveda $\int_a^b p(t)dt = 1 - \alpha$.

Teoretična interpretacija koeficienta zaupanja $(1 - \alpha)$

Če zaporedoma izbiramo vzorce velikosti n iz dane populacije in konstruiramo $[(1 - a)100]\%$ interval zaupanja za vsak vzorec, potem lahko pričakujemo, da bo $[(1 - a)100]\%$ intervalov dalo prvo vrednost parametra.

stopnja tveganja = 1 - stopnja zaupanja

14.2.1 Povprečje μ s poznanim σ

Točki

$$\bar{y} \pm z_{\alpha/2} \left(\frac{\sigma}{\sqrt{n}} \right)$$

prestavljata krajišči intervala zaupanja, pri čemer je:

$z_{\alpha/2}$ vrednost spremenljivke, ki zavzame površino $\alpha/2$ na svoji desni;

σ je standardni odklon za populacijo;

n je velikost vzorca;

\bar{y} je vrednost vzorčnega povprečja.

14.2.2 Velik vzorec za povprečje μ

$$\bar{y} \pm z_{\alpha/2} \left(\frac{s}{\sqrt{n}} \right),$$

kjer je s standardni odklon vzorca.

Primer: Na vzorcu velikosti $n = 151$ podjetnikov v majhnih podjetjih v Sloveniji, ki je bil izveden v okviru ankete 'Drobno gospodarstvo v Sloveniji' (Prašnikar, 1993), so izračunali, da je povprečna starost anketiranih podjetnikov $\bar{X} = 40,4$ let in standardni odklon $s = 10,2$ let. Pri 5 % tveganju želimo z intervalom zaupanja oceniti povprečno starost podjetnikov v majhnih podjetjih v Sloveniji.

$$P \left(\bar{y} - z_{\alpha/2} \frac{s}{\sqrt{n}} < \mu < \bar{y} + z_{\alpha/2} \frac{s}{\sqrt{n}} \right) = 1 - \alpha.$$

oziroma

$$40,4 - \frac{1,96 \times 10,2}{\sqrt{151}} < \mu < 40,4 + \frac{1,96 \times 10,2}{\sqrt{151}}$$

in končno

$$40,4 - 1,6 < \mu < 40,4 + 1,6$$

95 % interval zaupanja za povprečno starost podjetnikov v majhnih podjetjih v Sloveniji je med 38,8 in 42,0 leti. \diamond

14.2.3 Majhen vzorec za povprečje μ

$$\bar{y} \pm t_{\alpha/2} \left(\frac{s}{\sqrt{n}} \right),$$

kjer je porazdelitev spremenljivke y vzeta na osnovi $(n - 1)$ prostostnih stopenj.

Privzeli smo: populacija, iz katere smo izbrali vzorec, ima **približno normalno porazdelitev**.

Primer: Vzemimo, da se spremenljivka X - število ur branja dnevnih časopisov na teden - porazdeljuje normalno $N(\mu, \sigma)$. Na osnovi podatkov za 7 slučajno izbranih oseb ocenimo interval zaupanja za aritmetično sredino pri 10% tveganju. Podatki in ustrezni izračuni so:

x_i	$x_i - \bar{X}$	$(x_i - \bar{X})^2$
5	-2	4
7	0	0
9	2	4
7	0	0
6	-1	1
10	3	9
5	-2	4
49	0	22

$$\bar{X} = \frac{\sum_{i=1}^n x_i}{n} = \frac{49}{7} = 7, \quad s^2 = \frac{\sum_{i=1}^n (x_i - \bar{X})^2}{n - 1} = \frac{22}{6} = 3,67.$$

Iz tabele za t -porazdelitev preberemo, da je $t_{\alpha/2}(n - 1) = t_{0,05}(6) = 1,943$ in interval zaupanja je

$$7 - 1,943 \cdot \frac{1,9}{\sqrt{7}} < \mu < 7 + 1,943 \cdot \frac{1,9}{\sqrt{7}} \quad \text{oziroma} \quad 7 - 1,4 < \mu < 7 + 1,4. \quad \diamond$$

14.2.4 Razlika povprečij $\mu_1 - \mu_2$ s poznanima σ_1 in σ_2

$$\bar{y}_1 - \bar{y}_2 \pm z_{\alpha/2} \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}$$

14.2.5 Veliki vzorec za razliko povprečij $\mu_1 - \mu_2$ z neznanima σ_1 in σ_2

$$\bar{y}_1 - \bar{y}_2 \pm z_{\alpha/2} \sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}.$$

14.2.6 Majhen vzorec za razliko povprečij $\mu_1 - \mu_2$ z neznanima $\sigma_1 = \sigma_2$

$$\bar{y}_1 - \bar{y}_2 \pm t_{\alpha/2, n_1+n_2-2} \sqrt{s_p^2 \left(\frac{1}{n_1} + \frac{1}{n_2} \right)}, \quad \text{kjer je} \quad s_p^2 = \frac{(n_1 - 1)s_1^2 + (n_2 - 1)s_2^2}{n_1 + n_2 - 2}.$$

Privzeli smo:

- obe populaciji sta porazdeljeni **približno normalni**,
- varianci sta **enaki**,
- naključni vzorci so izbrani **neodvisno**.

14.2.7 Majhen vzorec za razliko povprečij $\mu_1 - \mu_2$ z neznanima σ_1 in σ_2

$$\bar{y}_1 - \bar{y}_2 \pm t_{\alpha/2, \nu} \sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}, \quad \text{kjer je} \quad \nu = \left\lfloor \frac{(s_1^2/n_1 + s_2^2/n_2)^2}{(s_1^2/n_1)^2/(n_1 - 1) + (s_2^2/n_2)^2/(n_2 - 1)} \right\rfloor.$$

Če ν ni naravno število, zaokroži ν navzdol do najbližjega naravnega števila za uporabo t -tabele.

Primer: Naslednji podatki predstavljajo dolžine filmov, ki sta jih naredila dva filmska studija. Izračunaj 90%-ni interval zaupanja za razliko med povprečnim časom filmov, ki sta jih producirala ta dva studija. Predpostavimo, da so dolžine filmov porazdeljene **približno normalno**. Čas (v minutah)

Studio 1:	103	94	110	87	98		
Studio 2:	97	82	123	92	175	88	118



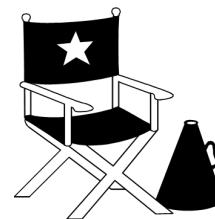
Podatke vnesemo v Minitab (Film.MTV):

Studio 1:	Studio 2:
103	97
94	82
110	123
87	92
98	175
	88
	118

Dva vzorca T -Test in interval zaupanja

Dva vzorca T za C1 : C2				
	N	povpr.	St.odk.	SE povpr.
C1	5	98.40	8.73	3.9
C2	7	110.7	32.2	12

90%-ni interval zaupanja za $\mu_{C1} - \mu_{C2}$: (-36.5, 12)
 T-TEST $\mu_{C1} = \mu_{C2}$ (vs μ_0):
 T = -0.96 P = 0.37 DF = 7



14.2.8 Velik vzorec za razliko $\mu_d = \mu_1 - \mu_2$

$$\bar{d} \pm z_{\alpha/2} \left(\frac{s_d}{\sqrt{n}} \right), \text{ kjer je } n \text{ število parov.}$$

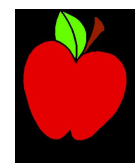
14.2.9 Majhen vzorec za razliko $\mu_d = \mu_1 - \mu_2$

$$\bar{d} \pm t_{\alpha/2, n-1} \left(\frac{s_d}{\sqrt{n}} \right), \text{ kjer je } n \text{ število parov.}$$

Privzeli smo: populacija razlik parov je normalno porazdeljena.

Primer: Špricanje jabolk lahko pozroči kontaminacijo zraka. Zato so v času najbolj intenzivnega špricanja zbrali in analizirali vzorce zraka za vsak od 11ih dni. Raziskovalci želijo vedeti ali se povprečje ostankov škropiv (diazinon) razlikuje med dnevom in nočjo. Analiziraj podatke za 90% interval zaupanja.

Datum	Diazinon Residue		razlika
	dan	noč	dan-noč
Jan. 11	5,4	24,3	-18,9
12	2,7	16,5	-13,8
13	34,2	47,2	-13,0
14	19,9	12,4	7,5
15	2,4	24,0	-21,6
16	7,0	21,6	-14,6
17	6,1	104,3	-98,2
18	7,7	96,9	-89,2
19	18,4	105,3	-86,9
20	27,1	78,7	-51,6
21	16,9	44,6	-27,7



Podatke vnesemo v Minitab (Ex8-39.MTW), pri čemer sta drugi in tretji stolpec zgoraj C1 in C2.

```
MTB > Let C3=C1-C2.
```

```
T interval zaupanja
```

Spremen.	N	povpr.	Stdev	SEpovpr.
C3	11	-38.9	36.6	11.0

Torej je 90,0% interval zaupanja (58,9; 18,9). ◇

Za p – delež populacije, \hat{p} – delež vzorca, kjer je $\hat{p} = y/n$ in je y število uspehov v n poskusih.

14.2.10 Delež π s poznanim $\sigma_{\hat{p}}$

$$\hat{p} \pm z_{\alpha/2} \sigma_{\hat{p}}$$

14.2.11 Veliki vzorec za delež populacije

$$\hat{p} \pm z_{\alpha/2} \sqrt{\frac{\hat{p}\hat{q}}{n}}$$

Privzeli smo: velikost vzorca n je dovolj velika, da je aproksimacija veljavna.

Izkušnje kažejo, da je za izpolnitev pogoja “dovolj velik vzorec” priporočljivo privzeti (angl. rule of thumb):

$$n\hat{p} \geq 4 \quad \text{in} \quad n\hat{q} \geq 4.$$

Primer: Na vzorcu ($n = 151$), ki je bil izveden v okviru ankete ‘Drobno gospodarstvo v Sloveniji’, so izračunali, da je delež obrtnih podjetij $p = 0,50$. Pri 5 % tveganju želimo z intervalom zaupanja oceniti delež obrtnih majhnih podjetij v Sloveniji.

$$0,50 - 1,96 \frac{0,50 \times 0,50}{\sqrt{151}} < \pi < 0,50 + 1,96 \frac{0,50 \times 0,50}{\sqrt{151}}$$

oziroma

$$0,50 - 0,08 < \pi < 0,50 + 0,08.$$

S 5% stopnjo tveganja trdimo, da je delež obrtnih majhnih podjetij v Sloveniji glede na vsa majhna podjetja med 0,42 in 0,58. ◇

14.2.12 Razlika deležev $\pi_1 - \pi_2$ s poznanim $\sigma_{\hat{p}_1 - \hat{p}_2}$

$$(\hat{p}_1 - \hat{p}_2) \pm z_{\alpha/2} \sigma_{\hat{p}_1 - \hat{p}_2}.$$

14.2.13 Veliki vzorec za razliko deležev $\pi_1 - \pi_2$

$$(\hat{p}_1 - \hat{p}_2) \pm z_{\alpha/2} \sqrt{\frac{\hat{p}_1 \hat{q}_1}{n_1} + \frac{\hat{p}_2 \hat{q}_2}{n_2}}$$

Privzeli smo: velikost vzorca n je dovolj velika, da je aproksimacija veljavna.

Izkušnje kažejo, da je za izpolnitev pogoja “dovolj velik vzorec” priporočljivo privzeti:

$$n_1 \hat{p}_1 \geq 4, \quad n_1 \hat{q}_1 \geq 4,$$

$$n_2 \hat{p}_2 \geq 4 \quad \text{in} \quad n_2 \hat{q}_2 \geq 4.$$

14.2.14 Veliki vzorec za varianco σ^2

$$\frac{(n-1)s^2}{\chi_{(\alpha/2, n-1)}^2} \leq \sigma^2 \leq \frac{(n-1)s^2}{\chi_{(1-\alpha/2, n-1)}^2}$$

Privzeli smo: populacija iz katere izbiramo vzorce, ima **približno normalno porazdelitev**.

Primer: Vzemimo prejšnji primer spremenljivke o številu ur branja dnevnih časopisov na teden. Za omenjene podatke iz vzorca ocenimo z intervalom zaupanja varianco pri 10% tveganju. Iz tabele za χ^2 -porazdelitev preberemo, da je

$$\chi_{1-\alpha/2}^2(n-1) = \chi_{0,95}^2(6) = 12,6,$$

$$\chi_{\alpha/2}^2(n-1) = \chi_{0,05}^2(6) = 1,64.$$

90 % interval zaupanja za varianco je tedaj

$$\frac{6 \cdot 3,67}{12,6} < \sigma^2 < \frac{6 \cdot 3,67}{1,64} \quad \text{ozioroma} \quad 1,75 < \sigma^2 < 13,43. \quad \diamond$$

14.2.15 Kvocient varianc σ_1^2/σ_2^2

$$\frac{s_1^2}{s_2^2} \cdot \frac{1}{F_{(\alpha/2, n_1-1, n_2-1)}} \leq \frac{\sigma_1^2}{\sigma_2^2} \leq \frac{s_1^2}{s_2^2} \cdot \frac{1}{F_{(\alpha/2, n_2-1, n_1-1)}}$$

Privzeli smo:

- obe populaciji iz katerih izbiramo vzorce, imata **približno normalni porazdelitvi** relativnih frekvenc.
- naključni vzorci so izbrani **neodvisno** iz obeh populacij.

14.3 Izbira velikosti vzorca

V tem razdelku bomo izbirali velikosti vzorcev za oceno različnih parametrov, ki je pravilna znotraj H enot z verjetnostjo $(1 - \alpha)$.

- Populacijsko povprečje μ :

$$n = \left(\frac{z_{\alpha/2} \sigma}{H} \right)^2$$

Populacijski odklon mora biti običajno aproksimiran.

- Razlika med parom populacijskih povprečij, tj. $\mu_1 - \mu_2$:

$$n_1 = n_2 = \left(\frac{z_{\alpha/2}}{H} \right)^2 (\sigma_1^2 + \sigma_2^2)$$

- Populacijski delež π :

$$n = \left(\frac{z_{\alpha/2}}{H} \right)^2 pq$$

Opozorilo: v tem primeru potrebujemo oceni za p in q . Če nimamo nobene na voljo, potem uporabimo $p = q = 0,5$ za konzervativno izbiro števila n .

- Razlika med parom populacijskih deležev, tj. $p_1 - p_2$

$$n_1 = n_2 = \left(\frac{z_{\alpha/2}}{H} \right)^2 (p_1 q_1 + p_2 q_2)$$



Poglavje 15

Preverjanje domnev



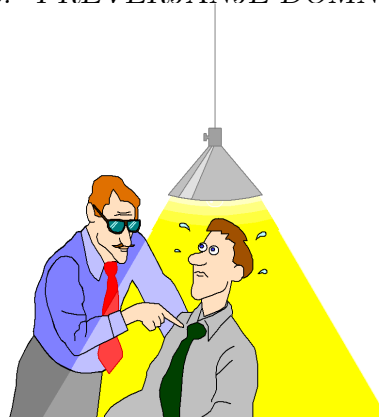
Načrt

- postopek
- elementi
 - napake 1. in 2. vrste
 - značilno razlikovanje
 - moč statističnega testa
- testi
 - centralna tendenca
 - delež
 - varianca



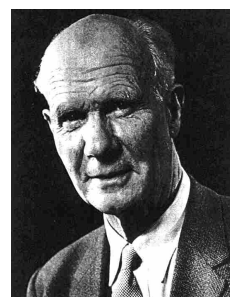
Uvod

- postavimo trditev o populaciji,
- izberemo vzorec, s katerim bomo preverili trditev,
- zavrnemo ali sprejmemo trditev.



Domneva je testirana z določanjem verjetja, da dobimo določen rezultat, kadar jemljemo vzorce iz populacije s predpostavljenimi vrednostimi parametrov.

Teorijo preizkušanja domnev sta v 20. in 30. letih prejšnjega stoletja razvila J. Neyman in E.S. Pearson.



Statistična domneva (ali hipoteza) je vsaka domneva o porazdelitvi slučajne spremenljivke X na populaciji. Če poznamo vrsto (obliko) porazdelitve $f(x; \zeta)$ in postavljamo/raziskujemo domnevo o parametru ζ , govorimo o *parametrični domnevi*. Če pa je vprašljiva tudi sama vrsta porazdelitve, je domneva *neparametrična*. Domneva je *enostavna*, če natančno določa porazdelitev (njeno vrsto in točno vrednost parametra); sicer je *sestavljena*.

Primer: Naj bo $X : N(\mu, \sigma)$. Če poznamo σ , je domneva $H : \mu = 0$ enostavna; Če pa parametra σ ne poznamo, je sestavljena. Primer sestavljene domneve je tudi $H : \mu > 0$. \diamond

Statistična domneva je lahko pravilna ali napačna. Želimo seveda sprejeti pravilno domnevo in zavrniti napačno. Težava je v tem, da o pravilnosti/napačnosti domneve ne moremo biti gotovi, če jo ne preverimo na celotni populaciji. Ponavadi se odločamo le na podlagi vzorca. Če vzorčni podatki preveč odstopajo od domneve, rečemo, da niso *skladni* z domnevo, oziroma, da so *razlike značilne*, in domnevo zavrnemo. Če pa podatki domnevo podpirajo, jo ne zavrnemo – včasih jo celo sprejmemo. To ne pomeni, da je domneva pravilna, temveč da ni zadostnega razloga za zavrnitev.

- **Ničelna domneva (H_0)**
 - je trditev o lastnosti populacije za katero predpostavimo, da drži (oziroma za katero verjamemo, da je resnična),
 - je trditev, ki jo test skuša ovreči.
- **Alternativna (nasprotna) domneva (H_a)**
 - je trditev nasprotna ničelni domnevi,
 - je trditev, ki jo s testiranjem skušamo dokazati.

Okvirni postopek testiranja domneve

- postavi ničelno in alternativno domnevo,
- izberi testno statistiko,
- določi zavrnitveni kriterij,
- izberi naključni vzorec,
- izračunaj vrednost na osnovi testne statistike,
- sprejmi odločitev,
- naredi ustrezen zaključek.

15.1 Ilustrativni primer (ameriški sodni sistem)

- ničelna domneva (H_0): obtoženec je nedolžen,
- alternativna domneva (H_a): obtoženec je kriv.

Odločitev in zaključek

- Porota je spoznala obtoženca za **krivega**. Zaključimo, da je bilo dovolj dokazov, ki nas prepričajo, da je obtoženec storil kaznivo dejanje.
- Porota je spoznala obtoženca za **nedolžnega**. Zaključimo, da je ni bilo dovolj dokazov, ki bi nas prepričali, da je obtoženec storil kaznivo dejanje.

Elementi testiranja domneve



		<i>odločitev</i>	
		nedolžen	kriv
<i>dejansko stanje</i>	nedolžen	pravilna odločitev	napaka 1. vrste (α)
	kriv	napaka 2. vrste (β)	moč ($1 - \beta$)

- Verjetnost napake 1. vrste (α) je verjetnost za obtožbo nedolžnega obtoženca.
- Značilno razlikovanje (signifikantno) oziroma **stopnja značilnosti**.
- Količina dvoma (α), ki ga bo porota še sprejela:
 - kriminalna tožba: Beyond a reasonable doubt...
 - civilna tožba: The preponderance of evidence must suggest...
- Verjetnost napake 2. vrste (β) je verjetnost, da spoznamo krivega obtoženca za nedolžnega.
- Moč testa ($1 - \beta$) je verjetnost, da obtožimo krivega obtoženca.

Sodba

- breme dokazov,
- potrebno je prepričati poroto, da je obtoženi kriv (alternativna domneva) preko določene stopnje značilnosti:
 - kriminalna tožba: Reasonable Doubt,
 - civilna tožba: Preponderance of evidence.

Obramba

- Ni bremena dokazovanja.
- Povzročiti morajo dovolj dvoma pri poroti, če je obtoženi resnično kriv.

15.2 Alternativna domneva in definicije napak

Nekaj konkretnih primerov statističnih ničelnih domnev:

$H_0 : \mu = 9mm$ (Premer 9 milimetrskega kroga),

$H_0 : \mu = 600$ km (Proizvalajec trdi, da je to doseg novih vozil),

$H_0 : \mu = 3$ dnevi

Neusmerjena alternativna domneva

Merjenje 9 milimetrskega kroga:

$H_0 : \mu = 9mm$,

$H_a : \mu \neq 9mm$.



“Manj kot” alternativna domneva

Proizvalajec trdi, da je to doseg novih vozil:

$H_0 : \mu = 600$ km,

$H_a : \mu < 600$ km.

“Več kot” alternativna domneva

Čas odsotnosti določenega artikla pri neposredni podpori:

$H_0 : \mu = 3$ dnevi,

$H_a : \mu > 3$ dnevi.

SASA slika/tabela

Definicije

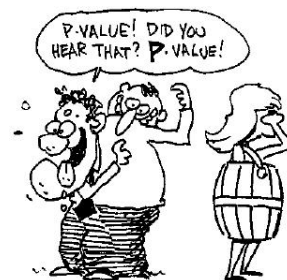
1. **Napaka 1. vrste** je zavrnitev ničelne domneve, če je le-ta pravilna. Verjetnost, da naredimo napako 1. vrste, označimo s simbolom α in ji pravimo **stopnja tveganja**, $(1 - \alpha)$ pa je **stopnja zaupanja**.
2. **Stopnja značilnosti testa (signifikantnosti)** je največji α , ki ga je vodja eksperimenta pripravljen sprejeti (zgornja meja za napako 1. vrste).
3. Če ne zavrnemo ničelno domnevo, v primeru, da je napačna, pravimo, da gre za **napako 2. vrste**. Verjetnost, da naredimo napako 2. vrste, označimo s simbolom β .
4. **Moč statističnega testa**, $(1 - \beta)$, je verjetnost zavrnitve ničelne domneve v primeru, ko je le-ta v resnici napačna.

		<i>odločitev</i>	
		FTR H_0	zavrni H_0
<i>dejansko stanje</i>	H_0 je pravilna		(α)
	H_0 je napačna	(β)	$(1 - \beta)$

velikost vzorca	napaka 1.vrste	napaka 2.vrste	moč
n	α	β	$1 - \beta$
konst.	↑	↓	↑
konst.	↓	↑	↓
povečanje	↓	↓	↑
zamnjšanje	↑	↑	↓

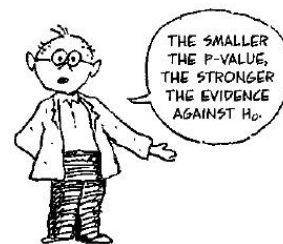
15.3 P -vrednost

P -vrednost (ali ugotovljena bistvena stopnja za določen statistični test) je verjetnost (ob predpostavki, da drži H_0), da ugotovimo vrednost testne statistike, ki je vsaj toliko v protislovju s H_0 in podpira H_a kot tisto, ki je izračunana iz vzorčnih podatkov.



- Sprejemljivost domneve H_0 na osnovi vzorca
 - Verjetnost, da je opazovani vzorec (ali podatki) bolj ekstremni, če je domneva H_0 pravilna.
- Najmanjši α pri katerem zavrnemo domnevo H_0 :

- če je P -vrednost $> \alpha$, potem FTR H_0 ,
- če je P -vrednost $< \alpha$, potem zavrne H_0 .



15.4 Statistična domneva

- ničelna domneva $H_0 : q = q_0$
- alternativna domneva
 - $H_a : q \neq q_0$
 - $H_a : q > q_0$
 - $H_a : q < q_0$

Primer: Predpostavimo, da je dejanska mediana (τ) pH iz določene regije 6,0. Da bi preverili to trditev, bomo izbrali 10 vzorcev zemlje iz te regije, da ugotovimo, če empirični vzorci močno podpirajo, da je dejanska mediana manjša ali enaka 6,0?

Predpostavke

- naključni vzorec
 - neodvisen
 - enako porazdeljen (kot celotna populacija),
- vzorčenje iz zvezne porazdelitve,
- verjetnostna porazdelitev ima mediano.

Postavitev statistične domneve

- ničelna domneva $H_0 : \tau = 6,0$ (mediana populacije τ_0)
- alternativna domneva $H_a : \tau < 6,0$

Izbira testne statistike (TS)

- S_+ = število vzorcev, ki so **večji** od mediane τ_0 iz domneve.
- S_- = število vzorcev, ki so **manjši** od mediane τ_0 iz domneve.

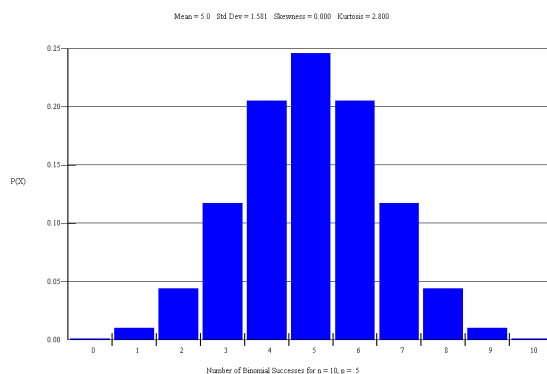
Porazdelitev testne statistike

- vsak poskus je bodisi uspeh ali neuspeh,
- fiksen vzorec, velikosti n ,
- naključni vzorci
 - neodvisni poskusi,
 - konstantna verjetnost uspeha.

Torej gre za

- binomsko porazdelitev: $S_+ \approx B(n, p)$,
- s parameteri $n = 10$ in $p = 0,5$ in
- pričakovano vrednostjo (matematičnim upanjem): $E(X) = np = 5$.

Porazdelitev testne statistike



Določimo zavrnitveni kriterij

x	$P(X = x)$	$F(x)$
0	0,000977	0,00098
1	0,009766	0,01074
2	0,043945	0,05469
3	0,117188	0,17188
4	0,205078	0,37695
5	0,246094	0,62305
6	0,205078	0,82813
7	0,117188	0,94531
8	0,043945	0,98926
9	0,009766	0,99902
10	0,000977	1,00000

- Stopnja značilnosti testa $\alpha = 0,01074$,

- Kritična vrednost:
- $-S_+ = 1$,
- Območje zavrnitve: $S_+ \leq 1$,

Izberemo naključni vzorec

Predpostavimo, da je dejanska mediana (τ) pH iz določene regije 6,0. Da bi preverili to trditev, smo izbrali 10 vzorcev zemlje iz te regije in jih podvrgli kemični analizi in na ta način določili pH vrednost za vsak vzorec. [Ali empirični podatki podpirajo trditev, da je dejanska mediana manjša ali enaka 6,0?](#)

5,93; 6,08; 5,86; 5,91; 6,12; 5,90; 5,95; 5,89; 5,98; 5,96.

Izračunaj vrednost iz testne statistike

pH	predznak
5,93	–
6,08	+
5,86	–
5,91	–
6,12	+
5,90	–
5,95	–
5,89	–
5,98	–
5,96	–

$S_+ = 2$, P -vrednost = $P(S_+ \geq 2 \mid \tau = 6,0) = 0,05469$.

$S_- = 8$, P -vrednost = $P(S_- \geq 8 \mid \tau = 6,0) = 0,05469$.

Odločitev in zaključek

- P -vrednost $> \alpha = 0,01074$.
- Ni osnove za zavrnitev domneve H_0 .
- Zavrni ničelno domnevo.
 - Zaključimo, da empirični podatki sugerirajo, da velja alternativna trditev.
- Ni osnove za zavrnitev ničelne domneve (angl. fail to reject, kratica FTR).
 - Zaključimo, da nimamo dovolj osnov, da bi dokazali, da velja alternativna trditev.
- Premalo podatkov, da bi pokazali, da je dejanska mediana pH manjša od 6,0.
- Privzemimo, da je pH enaka 6,0 v tej konkretni regiji. ◇

15.5 Preverjanje predznaka

- test: $H_0 : \tau = \tau_0$ (mediana populacije)
- predpostavke
 - naključno vzorčenje
 - vzorčenje iz zvezne porazdelitve

Preverjanje domneve z enim vzorcem

Testi za mero centralne tendence???

SASA diagram....

15.6 Wilcoxonov predznačen-rang test

- test
 - $H_0 : \tau = \tau_0$ (mediana populacije)
 - $H_0 : \mu = \mu_0$ (povprečje populacije)
- Predpostavke
 - naključni vzorec iz zvezne porazdelitve.
 - porazdelitev populacije ima simetrično obliko.
 - verjetnostna porazdelitev ima povprečje (mediano).

Testna statistika

- S_+ = vsota rangov, ki ustrezajo pozitivnim številom.
- S_- = vsota rangov, ki ustrezajo negativnim številom.

Primer:

- Naj bo $H_0 : \tau = 500$ in $H_a : \tau > 500$
- Postopek:
 - izračunaj odstopanje od τ_0
 - razvrsti odstopanja glede na velikost absolutne vrednosti (tj., brez upoštevanja predznaka).
 - seštej range, ki ustrezajo bodisi pozitivnemu ali negativnemu predznaku.

meritve	odstopanje	abs. vrednost	rang	+	-
499,2	-0,8	0,8	1		1
498,5	-1,5	1,5	2		2
502,6	2,6	2,6	3	3	
497,3	-2,7	2,7	4		4
496,9	-3,1	3,1	5		5
				S₊ = 3	S₋ = 12

Porazdelitev testne statistike

- 2^n enako verjetnih zaporedij,
- največji rang = $n(n + 1)/2$.

S+	1	2	S-	p	F
0	-	-	3	0,25	0,25
1	+	-	2	0,25	0,5
2	-	+	1	0,25	0,75
3	+	+	0	0,25	1

S+	1	2	3	S-	p	F
0	-	-	-	6	0,125	0,125
1	+	-	-	5	0,125	0,25
2	-	+	-	4	0,125	0,375
3	-	-	+	3	0,125	0,5
3	+	+	-	3	0,125	0,625
4	+	-	+	2	0,125	0,75
5	-	+	+	1	0,125	0,875
6	+	+	+	0	0,125	1

S+	1	2	3	4	S-	p	F
0	-	-	-	-	10	0,0625	0,0625
1	+	-	-	-	9	0,0625	0,125
2	-	+	-	-	8	0,0625	0,1875
3	+	+	-	-	7	0,0625	0,25
3	-	-	+	-	7	0,0625	0,3125
4	+	-	+	-	6	0,0625	0,375
4	-	-	-	+	6	0,0625	0,4375
5	+	-	-	+	5	0,0625	0,5
5	-	+	+	-	5	0,0625	0,5625
6	-	+	-	+	4	0,0625	0,625
6	+	+	+	-	4	0,0625	0,6875
7	+	+	-	+	3	0,0625	0,75
7	-	-	+	+	3	0,0625	0,8125
8	+	-	+	+	2	0,0625	0,875
9	-	+	+	+	1	0,0625	0,9375
10	+	+	+	+	0	0,0625	1

S+	1	2	3	4	5	S-	p	F
0	-	-	-	-	-	15	0,03125	0,03125
1	+	-	-	-	-	14	0,03125	0,0625
2	-	+	-	-	-	13	0,03125	0,09375
3	-	-	+	-	-	12	0,03125	0,125
3	+	+	-	-	-	12	0,03125	0,15625
4	-	-	-	+	-	11	0,03125	0,1875
4	+	-	+	-	-	11	0,03125	0,21875
5	-	-	-	-	+	10	0,03125	0,25
5	+	-	-	+	-	10	0,03125	0,28125
5	-	+	+	-	-	10	0,03125	0,3125
6	+	-	-	-	+	9	0,03125	0,34375
6	-	+	-	+	-	9	0,03125	0,375
6	+	+	+	-	-	9	0,03125	0,40625
7	-	+	-	-	+	8	0,03125	0,40625
7	-	-	+	+	-	8	0,03125	0,4375
7	+	+	-	+	-	8	0,03125	0,46875
8	-	-	+	-	+	7	0,03125	0,5

P-vrednost

- Sprejemljivost domneve H_0 na osnovi vzorca

- možnost za opazovanje vzorca (ali bolj ekstremno podatkov), če je domneva H_0 pravilna.
- Najmanjši α pri katerem zavrnemo domnevo H_0
 - Če je P -vrednost $> \alpha$, potem FTR H_0 .
 - Če je P -vrednost $< \alpha$, potem zavrne H_0 .
 - Če je P -vrednost $= (2)P(Z > 1,278) = (2)(0,1003) = 0,2006$.

Odločitev in zaključek

- odločitev
 - P -vrednost $= P(S_+?3)$ ali $P(S_-?12)$
 - P -vrednost $= 0,15625$
 - P -vrednost $> \alpha = 0,1$
 - FTR H_0
- zaključek
 - privzemimo $\tau = 500$
 - ni osnov, da bi pokazali $\tau > 500$

◇

15.7 Naloga

Pascal je visoko-nivojski programski jezik, ki smo ga nekoč pogosto uporabljali na miniračunalnikih in microprocesorjih. Narejen je bil eksperiment, da bi ugotovili delež Pascalovih spremenljivk, ki so tabelarične spremenljivke (v kontrast skalarim spremenljivkam, ki so manj učinkovite, glede na čas izvajanja). 20 spremenljivk je bilo naključno izbranih iz množice Pascalovih programov, pri tem pa je bilo zabeleženo število tabelaričnih spremenljivk Y . Predpostavimo, da želimo testirati domnevo, da je Pascal bolj učinkovit jezik kot Agol, pri katerem je 20% spremenljivk tabelaričnih. To pomeni, da bomo testirali $H_0 : p = 0,20$, proti $H_a : p > 0,20$, kjer je p verjetnost, da imamo tabelarično spremenljivko na vsakem poskusu. Predpostavimo, da je 20 poskusov neodvisnih.

(a) Določi α za območje zavrnitve $y > 8$. Izračunati želimo verjetnost, da se bo zgodila napaka 1. vrste, torej da bomo zavrnilo pravilno domnevo. Predpostavimo, da je domneva H_0 pravilna, tj. $Y : B(20; 0,2)$. Če se bo zgodilo, da bo Y pri izbranem vzorcu večji ali enak 8, bom domnevo zavrnilo, čeprav je pravilna. Torej velja:

$$\begin{aligned}\alpha &= P(Y \geq 8) = 1 - P(Y \leq 7) = 1 - \sum_{i=0}^7 P(Y = i) \\ &= 1 - \sum_{i=0}^7 \binom{20}{i} 0,2^i 0,2^{20-i} = 1 - 0,9679 = 0,0321 = 3,21\%.\end{aligned}$$



(b) Določi α za območje zavrnitve $y \geq 5$. Do rezultata pridemo na enak način kot v prejšnji točki:

$$\begin{aligned}\alpha &= P(Y \geq 5) = 1 - P(Y \leq 4) = 1 - \sum_{i=0}^4 P(Y = i) \\ &= 1 - \sum_{i=0}^4 \binom{20}{i} 0,2^i 0,2^{20-i} = 1 - 0,6296 = 0,3704 = 37,04\%.\end{aligned}$$

(c) Določi β za območje zavrnitve $Y \geq 8$, če je $p = 0,5$. Izračunati želimo verjetnost, da se bo zgodila napaka 2. vrste, torej da bomo sprejeli napačno domnevo. Ker vemo, da je $p = 0,5$, velja $Y \sim B(20; 0,5)$. Napačno domnevo bomo sprejeli, če bo y pri izbranem vzorcu manjši od 8.

$$\beta = P(y \leq 7) = \sum_{i=0}^7 \binom{20}{i} 0,5^i 0,5^i = 0,1316 = 13,16\%.$$

(d) Določi β za območje zavrnitve $y \geq 5$, če je $p = 0,5$. Do rezultata pridemo na enak način kot v prejšnji točki:

$$\beta = P(y \leq 4) = \sum_{i=0}^4 \binom{20}{i} 0,5^i 0,5^i = 0,0059 = 0,59\%.$$

(e) Katero območje zavrnitve $y \geq 8$ ali $y \geq 5$ je bolj zaželeno, če želimo minimizirati verjetnost napake 1. stopnje oziroma če želimo minimizirati verjetnost napake 2. stopnje.

Napako 1. stopnje minimiziramo z izbiro območja $y \geq 8$, napako 2. stopnje pa z izbiro območja $y \geq 5$.

(f) Določi območje zavrnitve $y \geq a$ tako, da je α približno 0,01. Na osnovi točke (e) zaključimo, da se z večanjem števila a manjša verjetnost α in s poskušanjem (ki ga pričnemo na osnovi izkušnje iz točke (a) pri 9) pridemo do $a = 9$.

(g) Za območje zavrnitve določeno v točki (f) določi moč testa, če je v resnici $p = 0,4$. Moč testa je $1 - \beta$. Verjetnost β izračunamo enako kot v točkah (c) in (d). Velja $Y \sim B(20; 0,4)$ in

$$\beta = P(y \leq 8) = \sum_{i=0}^8 \binom{20}{i} 0,4^i 0,6^i = 0,5956 = 59,56\%.$$

Moč testa znaša 0,4044.

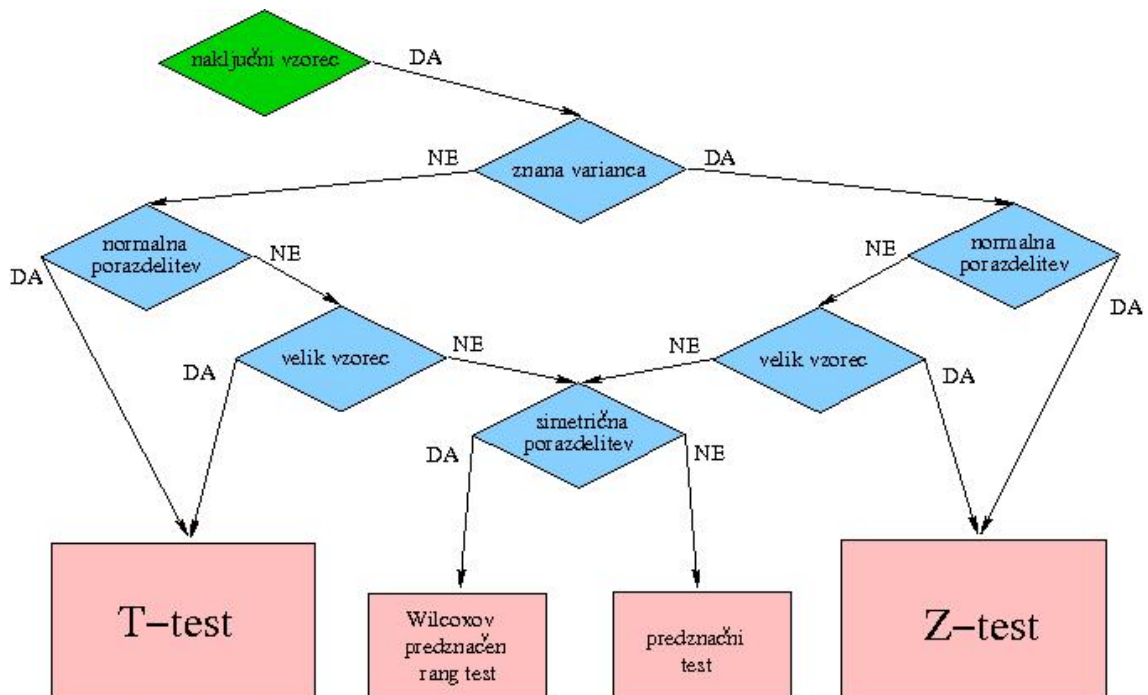
(h) Za območje zavrnitve določeno v točki (f) določi moč testa, če je v resnici $p = 0,7$. Tokrat velja $Y \sim B(20; 0,4)$ in

$$\beta = P(y \leq 8) = \sum_{i=0}^8 \binom{20}{i} 0,7^i 0,3^i = 0,0051 = 0,51\%.$$

Moč testa znaša 0,995.

15.8 Formalen postopek za preverjanje domnev

1. Postavi domnevi o parametrih (ničelno H_0 in alternativno H_1).
2. Za parameter poiščemo kar se da dobro cenilko (npr. nepristransko) in njeno porazdelitev ali porazdelitev ustrezne statistike (izraz, v katerem nastopa cenilka).
3. Določi odločitveno pravilo. Izberemo stopnjo značilnosti (α). Na osnovi stopnje značilnosti in porazdelitve statistike določimo kritično območje;
4. Zberi/manipuliraj podatke ter na vzorčnih podatkih izračunaj (eksperimentalno) vrednost testne statistike.
5. Primerjaj in naredi zaključek.
 - če eksperimentalna vrednost pade v kritično območje, ničelno domnevo zavrni in sprejmi osnovno domnevo ob stopnji značilnosti α .
 - če eksperimentalna vrednost ne pade v kritično območje, pa pravimo da vzorčni podatki kažejo na statistično neznačilne razlike med parametrom in vzorčno oceno.



15.8.1 $\mu = \mu_0$ z znanim σ

$$\text{T.S.} = \frac{\bar{y} - \mu_0}{\frac{\sigma}{\sqrt{n}}} \text{ sledi } z\text{-porazdelitev.}$$



Primer: Proizvajalec omake za špagete da v vsako posodo 28 unče omake za špagete. Količina omake, ki je v vsaki posodi, je porazdeljena normalno s standardnim odklonom 0,005 unče. Podjetje ustavi proizvodni trak in popravi napravo za polnenje, če so posode bodisi premalo napolnjene (to razjezi kupce), ali preveč napolnjene (kar seveda pomeni manjši profit).

Ali naj na osnovi vzorca iz 15ih posod ustavimo proizvodno linijo?

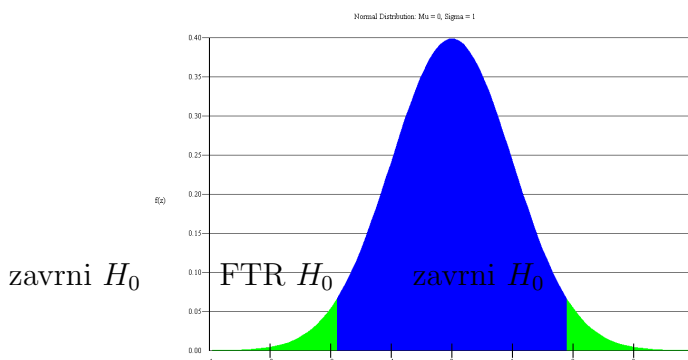
Uporabi stopnjo značilnosti 0,05. Postavimo domnevo

- ničelna domneva $H_0 : \mu = 28$
- alternativna domneva $H_a : \mu \neq 28$

Izberemo testno statistiko: Z-Test

- test: $H_0 : \mu = \mu_0$ (povprečje populacije)
- predpostavke
 - naključno vzorčenje
 - poznamo varianco populacije
 - izbiramo vzorce iz normalne porazdelitve in/ali imamo vzorec pri katerem je n velik.

Določimo zavrnitveni kriterij



Rezultati testiranja

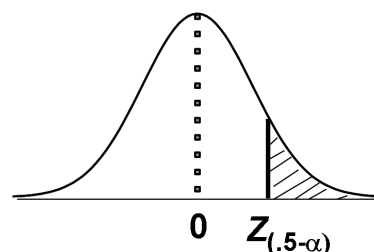
- vzami naključni vzorec – zanj dobimo vzorčno povprečje: 28,0165
- izračunaj vrednost testne statistike: $Z = (28,0165 - 28)/0,0129 = 1,278$.
- naredi odločitev: FTR H_0 .
- zaključek: privzemi $\mu = 28$

P-vrednost

- Sprejemljivost domneve H_0 na osnovi vzorca (možnost za opazovanje vzorca ali bolj ekstremno podatkov, če je domneva H_0 pravilna):
 - P -vrednost = $(2)P(Z > 1,278) = (2)(0,1003) = 0,2006$.
- Najmanjši α pri katerem zavrni domnevo H_0
 - P -vrednost $> \alpha$, zato FTR H_0 .



Za $H_a : \mu > \mu_0$ je **odločitveno pravilo**: zavrni H_0 , če je **T.S.** $\geq z_{(0,5-\alpha)}$



Za $H_a : \mu < \mu_0$ **odločitveno pravilo**: zavrni H_0 , če je **T.S.** $\leq z_{(0,5-\alpha)}$

Za $H_a : \mu \neq \mu_0$ **odločitveno pravilo**: zavrni H_0 če je **T.S.** $\leq -z_{(0,5-\alpha)}$ ali če je **T.S.** $\geq z_{(0,5-\alpha)}$. ◇

15.8.2 $\mu = \mu_0$ z neznanim σ , $n \geq 30$

T.S. = $\frac{\bar{y} - \mu_0}{\frac{s}{\sqrt{n}}}$ sledi t -porazdelitev z $n - 1$ prostostnimi stopnjami.

(Velja omeniti še, da se pri tako velikem n z - in t -porazdelitev tako ne razlikujeta kaj dosti.)

15.8.3 $\mu = \mu_0$, neznan σ , populacija normalna in $n < 30$

T.S. = $\frac{\bar{y} - \mu_0}{\frac{s}{\sqrt{n}}}$ sledi t -porazdelitev z $n - 1$ prostostnimi stopnjami.

Primer: Za slučajni vzorec: 16-ih odraslih Slovencev smo izračunali povprečno število in variance priznanih let šolanja: $\bar{X} = 9$ in $s^2 = 9$. Predpostavljamo, da se spremenljivka na populaciji porazdeljuje normalno. **Ali lahko sprejmemo domnevo, da imajo odrasli**

Sloenci v povprečju več kot osemletko pri 5% stopnji značilnosti? Postavimo najprej ničelno in osnovno domnevo:

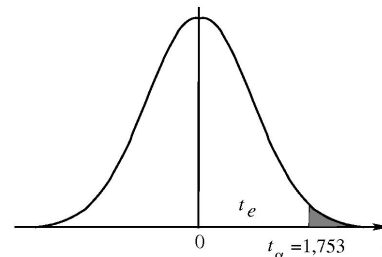
$$H_0 : \mu = 8 \quad \text{in} \quad H_1 : \mu > 8.$$

Ustrezna statistike je

$$t = \frac{\bar{X} - \mu_H}{s} \sqrt{n},$$

ki se porazdeljuje po t -porazdelitvi s 15 prostostnimi stopnjami. Ker gre za enostranski test, je glede na osnovno domnevo kritično območje na desni strani porazdelitve in kritična vrednost $t_{0,05}(15) = 1,753$. Izračunajmo eksperimentalno vrednost statistike:

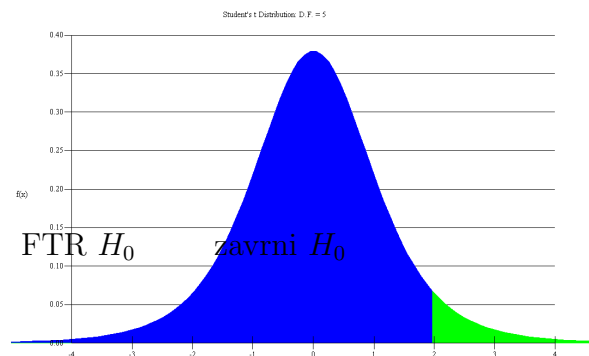
$$t_e = \frac{9 - 8}{3} \sqrt{16} = 1,3$$



Eksperimentalna vrednost ne pade v kritično območje. Zato ničelne domneve ne moremo zavrniti in sprejeti osnovne domneve, da imajo odrasli Slovenci več kot osemletko. \diamond

Primer: Ravnatelj bežigradske gimnazije trdi, da imajo najboljši PT program v Sloveniji s povprečjem APFT 240. Predpostavi, da je porazdelitev rezultatov testov približno normalna. Uporabi $\alpha = 0,05$ za določitev **ali je povprečje APFT rezultatov šestih naključno izbranih dijakov iz bežigradske gimnazije statistično večje od 240**. Postavimo domnevi: $H_0 : \mu = 240$ in $H_a : \mu > 240$ in izberemo testno statistiko T -test.

- test: $H_0 : \mu = \mu_0$ (povprečje populacije)
- predpostavke
 - naključno vzorčenje
 - ne poznamo varianco populacije
 - izbiramo vzorce iz normalne porazdelitve in/ali imamo vzorec pri katerem je n velik.

Določimo zavrnitveni kriterij**Rezultati testov**

- naredi naključni vzorec:
 - vzorčno povprečje: 255,4
 - vzorčni standardni odklon: 40,07
- izračunaj vrednost testne statistike: $T = (255,4 - 240)/16,36 = 0,9413$.
- sprejmi odločitev: FTR H_0
- zaključek: Bežigrajska gimnazija ne more pokazati, da imajo višje povprečje APFT rezultatov, kot slovensko povprečje.

***P*-vrednost**

- Sprejemljivost domneve H_0 na osnovi vzorca
 - možnost za opazovanje vzorca (ali bolj ekstremno podatkov), če je domneva H_0 pravilna
 - P -vrednost = $P(T > 0,9413) = 0,1949$.
- Najmanjši α pri katerem zavrնemo domnevo H_0
 - P -vrednost $> \alpha$, zato FTR H_0 .

Vstavimo podatke v Minitab (Ex9-23.MTV)

C1:

2610

2750

2420

2510

2540

2490

2680



T-test povprečja

Test of $\mu = 2500.0$ vs $\mu > 2500.0$

	N	MEAN	STDEV	SE MEAN
C1	7	2571.4	115.1	43.5

	T	p-VALUE
	1.64	0.076



Razlaga P -vrednosti

1. Izberi največjo vrednost za α , ki smo jo pripravljene tolerirati.
2. Če je P -vrednost testa manjša kot maksimalna vrednost parametra α , potem zavrne ničelno domnevo. ◇

15.9 Razlika povprečij $\mu_1 - \mu_2 = D_0$



15.9.1 Znana σ_1 in σ_2

Vzorci jemljemo neodvisno, zato

$$\text{T.S.} = \frac{(\bar{y}_1 - \bar{y}_2) - D_0}{\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}} \text{ sledi } z\text{-porazdelitev.}$$

Primer: Preveriti želimo domnevo, da so dekleta na izpitu boljša od fantov. To domnevo preverimo tako, da izberemo slučajni vzorec 36 deklet in slučajni vzorec 36 fantov, za katere imamo izpitne rezultate, na katerih izračunamo naslednje statistične karakteristike:

$$\bar{X}_F = 7,0, \quad s_F = 1$$

$$\bar{X}_D = 7,2, \quad s_D = 1$$

Domnevo preverimo pri 5% stopnji značilnosti. Postavimo ničelno in osnovno domnevo:

$$H_0 : \mu_D = \mu_F \quad \text{ozioroma} \quad \mu_D - \mu_F = 0,$$

$$H_1 : \mu_D > \mu_F \quad \text{ozioroma} \quad \mu_D - \mu_F > 0.$$

Za popularijsko razliko aritmetičnih sredin na vzorcih računamo vzorčno razliko aritmetičnih sredin, ki se za dovolj velike vzorce porazdeljuje normalno

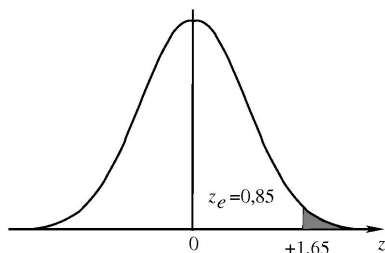
$$\bar{X}_D - \bar{X}_F : N\left(\mu_D - \mu_F, \sqrt{\frac{s_D^2}{n_D} + \frac{s_F^2}{n_F}}\right).$$

ozioroma statistika

$$z = \frac{\bar{X}_D - \bar{X}_F - (\mu_D - \mu_F)_H}{\sqrt{\frac{s_D^2}{n_D} + \frac{s_F^2}{n_F}}}$$

standardizirano normalno $N(0, 1)$. Osnovna domneva kaže enostranski test: možnost napake 1. vrste je le na desni strani normalne porazdelitve, kjer zavračamo ničelno domnevo. Zato je kritično območje določeno z vrednostmi večjimi od 1,65. Eksperimentalna vrednost statistike je

$$z_e = \frac{7,2 - 7 - 0}{\sqrt{\frac{1}{36} + \frac{1}{36}}} = 0,852.$$



Eksperimentalna vrednost ne pade v kritično območje. Ničelne domneve ne moremo zavrniti. Povprečna uspešnost deklet in fantov ni statistično značilno različna. \diamond

15.9.2 Neznana σ_1 in/ali σ_2 , $n_1 \geq 30$ in/ali $n_2 \geq 30$

Ker vzorce jemljemo neodvisno, velja

$$\text{T.S.} = \frac{(\bar{y}_1 - \bar{y}_2) - D_0}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}} \text{ sledi } z\text{-porazdelitev.}$$

15.9.3 Neznana σ_1 in/ali σ_2 , pop. norm., $\sigma_1 = \sigma_2$, $n_1 < 30$ ali $n_2 < 30$

Ker vzorce jemljemo neodvisno, velja:

$$\text{T.S.} = \frac{(\bar{y}_1 - \bar{y}_2) - D_0}{\sqrt{\frac{(n_1 - 1)s_1^2 + (n_2 - 1)s_2^2}{n_1 + n_2 - 2} \left(\frac{1}{n_1} + \frac{1}{n_2} \right)}}$$

sledi t -porazdelitev z $n_1 + n_2 - 2$ prostostnimi stopnjami.

Privzeli smo:

1. Populaciji iz katerih jemljemo vzorce imata obe približno **normalno** relativno porazdelitev frekvenc.
2. Varianci obeh populacij sta **enaki**.
3. Naključni vzorci so izbrani **neodvisno** iz obeh populacij.

15.9.4 Neznana σ_1 in/ali σ_2 , pop. norm., $\sigma_1 \neq \sigma_2$, $n_1 < 30$ ali $n_2 < 30$

Ker jemljemo vzorce neodvisno, velja

$$\text{T.S.} = \frac{(\bar{y}_1 - \bar{y}_2) - D_0}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}} \quad \text{sledi } t\text{-porazdelitev z } \nu \text{ prostostnimi stopnjami,}$$

kjer je

$$\nu = \frac{\left(\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}\right)^2}{\frac{\left(\frac{s_1^2}{n_1}\right)^2}{n_1 - 1} + \frac{\left(\frac{s_2^2}{n_2}\right)^2}{n_2 - 1}}$$

Če ν ni naravno število, zaokroži ν navzdol do najbližjega naravnega števila za uporabo t -tabele.

15.9.5 Povprečje $\mu_d = D_0$ in $n \geq 30$

Ker vzorce jemljemo neodvisno velja

$$\text{T.S.} = \frac{\bar{d} - D_0}{\frac{s_d}{\sqrt{n}}} \quad \text{sledi } z\text{-porazdelitev.}$$

15.9.6 Povprečje $\mu_d = D_0$, populacija razlik normalna, in $n \leq 30$

Ker vzorce ne jemljemo neodvisno, velja

$$\text{T.S.} = \frac{\bar{d} - D_0}{\frac{s_d}{\sqrt{n}}} \quad \text{sledi } t\text{-porazdelitev z } n - 1 \text{ prostostnimi stopnjami.}$$

naloga	clovek. urnik	avtomatizirana metoda	razlika
1	185,4	180,4	5,0
2	146,3	248,5	-102,2
3	174,4	185,5	-11,1
4	184,9	216,4	-31,5
5	240,0	269,3	-29,3
6	253,8	249,6	-4,2
7	238,8	282,0	-43,2
8	263,5	315,9	-52,4

Vstavimo podatke v Minitab (Ex9-40.MTV)

C1: 185,4 146,3 174,4 184,9 240,0 253,8 238,8 263,5

C2: 180,4 248,5 185,5 216,4 269,3 249,6 282,0 315,9

Test za parjenje in interval zaupanja



Parjen T za C1-C2

	N	povpr.	StDev	SE povpr.
C1	8	210,9	43,2	15,3
C2	8	243,4	47,1	16,7
Razlika	8	032,6	35,0	12,4

95% interval zaupanja za razliko povprečja: $(-61,9; -3,3)$

T -test za razliko povpr. = 0 (proti $\neq 0$):

T -vrednost=-2,63

P -vrednost=0,034.

15.9.7 Delež $p = p_0$ z dovolj velikim vzorcem

$$\text{T.S.} = \frac{\hat{p} - p_0}{\sqrt{\frac{p_0 q_0}{n}}} \quad \text{sledi } z\text{-porazdelitev.}$$

Kot splošno pravilo bomo zahtevali, da velja

$$n\hat{p} \geq 4 \quad \text{in} \quad n\hat{q} \geq 4.$$

Primer: Postavimo domnevo o vrednosti parametra, npr. π – delež enot z določeno lastnostjo na populaciji. Denimo, da je domneva $H : \pi_H = 0,36$

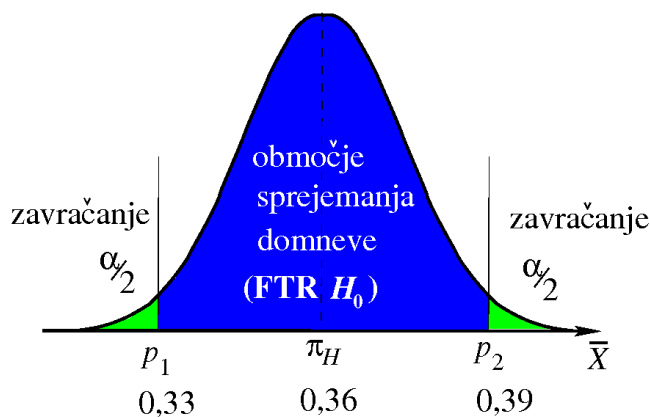
Tvorimo slučajne vzorce npr. velikosti $n = 900$ in na vsakem vzorcu določimo vzorčni delež p (delež enot z določeno lastnostjo na vzorcu). Ob predpostavki, da je domneva pravilna, vemo, da se vzorčni deleži porazdeljujejo približno normalno

$$N\left(\pi_H, \sqrt{\frac{\pi_H(1 - \pi_H)}{n}}\right)$$

Vzemimo en slučajni vzorec z vzorčnim deležem p . Ta se lahko bolj ali manj razlikuje od π_H . Če se zelo razlikuje, lahko podvomimo o resničnosti domneve π_H . Zato okoli π_H naredimo območje sprejemanja domneve in izven tega območja območje zavračanja domneve. Denimo, da je območje zavračanja določeno s 5% vzorcev, ki imajo ekstremne vrednosti deležev (2,5% levo in 2,5% desno). Deleža, ki ločita območje sprejemanja od območja zavračanja lahko izračunamo takole:

$$p_{1,2} = \pi_H \pm z_{\alpha/2} \sqrt{\frac{\pi_H(1 - \pi_H)}{n}},$$

$$\begin{aligned} p_{1,2} &= 0,36 \pm 1,96 \sqrt{\frac{0,36 \times 0,64}{900}} \\ &= 0,36 \pm 0,03. \end{aligned}$$



Kot smo že omenili, je sprejemanje ali zavračanje domnev po opisanem postopku lahko napačno v dveh smislih:

Napaka 1. vrste (α): Če vzorčna vrednost deleža pade v območje zavračanja, domnevo π_H zavrnamo. Pri tem pa vemo, da ob resnični domnevi π_H obstajajo vzorci, ki imajo vrednosti v območju zavračanja. Število α je verjetnost, da vzorčna vrednost pade v območje zavračanja ob predpostavki, da je domneva resnična. Zato je α verjetnost, da zavrnamo pravilno domnevo – **napaka 1. vrste**. Ta napaka je merljiva in jo lahko poljubno manjšamo.

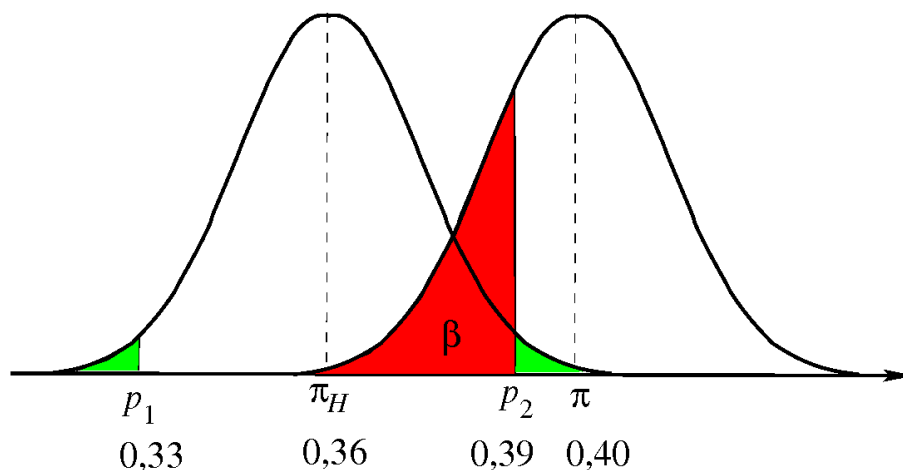
Napaka 2. vrste (β): Vzorcna vrednost lahko pade v območje sprejemanja, čeprav je domnevna vrednost parametra napačna. V primeru, ki ga obravnavamo, naj bo prava vrednost deleža na populaciji $\pi = 0,40$. Tedaj je porazdelitev vzorčnih deležev

$$N\left(\pi, \sqrt{\frac{\pi(1 - \pi)}{n}}\right) = N(0,40; 0,0163)$$

Ker je območje sprejemanja, domneve v intervalu $0,33 < p < 0,39$, lahko izračunamo verjetnost, da bomo sprejeli napačno domnevo takole:

$$\beta = P(0,33 < p < 0,39) = 0,27$$

Napako 2. vrste lahko izračunamo le, če imamo znano resnično vrednost parametra π . Ker ga ponavadi ne poznamo, tudi ne poznamo napake 2. vrste. Zato ne moremo sprejemati domnev. \diamond



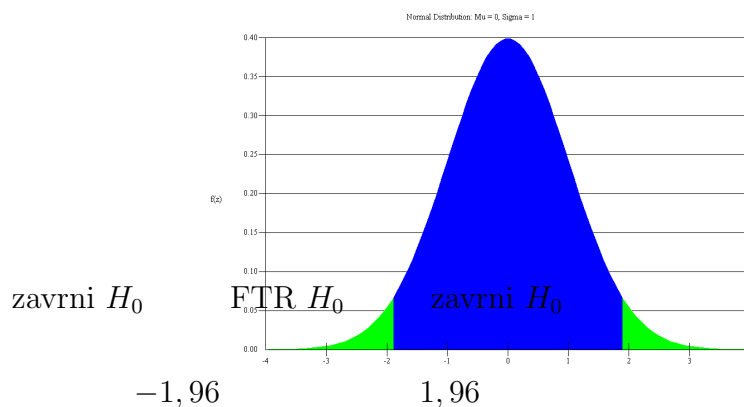
Primer: Državni zapisi indicirajo, da je od vseh vozil, ki gredo skozi testiranje izpušnih plinov v preteklem letu, 70% uspešno opravilo testiranje v prvem poskusu. Naključni vzorec 200ih avtomobilov testiranih v določeni pokrajni v tekočem letu je pokazalo, da jih je 156 šlo čez prvi test. Ali to nakazuje, da je dejanski delež populacije za to pokrajno v tekočem letu različno od preteklega državnega deleža? Pri testiranju domneve uporabi $\alpha = 0,05$. \diamond

15.10 Preverjanje domneve za delež

- Ničelna domneva $H_0 : p = 0,7$
- Alternativna domneva $H_a : p \neq 0,7$
- Test: $H_0 : p = p_0$ (delež populacije)
- Predpostavke
 - naključni vzorec
 - izbiranje vzorca iz binomske porazdelitve
- T.S.: $Z = \frac{\hat{p} - p_0}{\sqrt{p_0 q_0 / n}}$, upoštevamo pa tudi $n\hat{p} \geq 4$ in $n\hat{q} \geq 4$.

Primer:

Določimo zavrnitveni kriterij



Rezultati testiranja

- Naredi naključni vzorec: dobimo delež vzorca: $156/200 = 0,78$
- Izračunaj vrednost testne statistike: $Z = (0,78 - 0,7)/0,0324 = 2,4688$.
- Naredi odločitev: zavrni domnevo H_0
- Zaključek: pokrajna ima drugačen kriterij.

P-vrednost

- Sprejemljivost domneve H_0 na osnovi vzorca

- možnost za opazovanje vzorca (ali bolj ekstremno podatkov), če je domneva H_0 pravilna
- P -vrednost = $(2) * P(Z > 2,469) = (2) * (0,0068) = 0,0136$
- Najmanjši α pri katerem zavrnilo domnevo H_0
 - P -vrednost $< \alpha$, zato zavrni domnevo H_0 ◇

15.11 Razlika deležev dveh populaciji

Velik vzorec za testiranje domneve o $p_1 - p_2$



Kot splošno pravilo bomo zahtevali, da velja

$$n_1 \hat{p}_1 \geq 4, \quad n_1 \hat{q}_1 \geq 4,$$

$$n_2 \hat{p}_2 \geq 4 \quad \text{in} \quad n_2 \hat{q}_2 \geq 4.$$

15.11.1 Velik vzorec za testiranje domneve o $p_1 - p_2$, kadar je $D_0 = 0$.

$$\text{T.S.} = \frac{\hat{p}_1 - \hat{p}_2}{\sqrt{\hat{p}\hat{q}\left(\frac{1}{n_1} + \frac{1}{n_2}\right)}} \quad \text{sledi } z\text{-porazdelitev,}$$

kjer je $\hat{p} = \frac{y_1 + y_2}{n_1 + n_2}$.

Primer: Želimo preveriti, ali je predsedniški kandidat različno priljubljen med mestnimi in vaškimi prebivalci. Zato smo slučajni vzorec mestnih prebivalcev povprašali, ali bi glasovali za predsedniškega kandidata. Od 300 vprašanih (n_1) jih je 90 glasovalo za kandidata (k_1). Od 200 slučajno izbranih vaških prebivalcev (n_2) pa je za kandidata glasovalo 50 prebivalcev (k_2). Domnevo, da je kandidat različno priljubljen v teh dveh območjih preverimo pri 10% stopnji značilnosti.

$$H_0 : \pi_1 = \pi_2 \quad \text{ozirama} \quad \pi_1 - \pi_2 = 0,$$

$$H_1 : \pi_1 \neq \pi_2 \quad \text{oziroma} \quad \pi_1 - \pi_2 \neq 0.$$

Vemo, da se razlika vzorčnih deležev porazdeljuje približno normalno:

$$p_1 - p_2 : N\left(\pi_1 - \pi_2, \sqrt{\frac{\pi_1(1 - \pi_1)}{n_1} + \frac{\pi_2(1 - \pi_2)}{n_2}}\right).$$

Seveda π_1 in π_2 nista znana. Ob predpostavki, da je ničelna domneva pravilna, je matematično upanje razlike vzorčnih deležev hipotetična vrednost razlike deležev, ki je v našem primeru enaka 0. Problem pa je, kako oceniti standardni odklon. Ker velja domneva $\pi_1 = \pi_2 = \pi$, je disperzija razlike vzorčnih deležev

$$\frac{\pi_1(1 - \pi_1)}{n_1} + \frac{\pi_2(1 - \pi_2)}{n_2} = \frac{\pi(1 - \pi)}{n_1} + \frac{\pi(1 - \pi)}{n_2} = \pi(1 - \pi)\left(\frac{1}{n_1} + \frac{1}{n_2}\right).$$

Populacijski delež π ocenimo z obteženim povprečjem vzorčnih deležev p_1 in p_2

$$p = \frac{n_1 p_1 + n_2 p_2}{n_1 + n_2} = \frac{k_1 + k_2}{n_1 + n_2}.$$

Vrnimo se na primer. Vzorčna deleža sta:

$$p_1 = \frac{90}{300} = 0,30, \quad p_2 = \frac{50}{200} = 0,25.$$

Ocena populacijskega deleža je

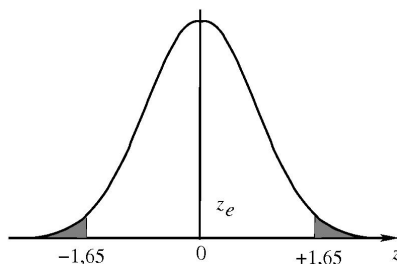
$$P = \frac{50 + 90}{200 + 300} = 0,28.$$

Kot smo že omenili, se statistika

$$z = \frac{p_1 - p_2 - (\pi_1 - \pi_2)_H}{\sqrt{p(1 - p)\left(\frac{1}{n_1} + \frac{1}{n_2}\right)}}$$

porazdeljuje približno standardizirano normalno $N(0, 1)$. Ker gre za dvostranski test, sta kritični vrednosti $\pm z_{\alpha/2} = \pm 1,65$. Eksperimentalna vrednost statistike pa je

$$z_e = \frac{0,30 - 0,25 - 0}{\sqrt{0,28(1 - 0,28)\left(\frac{1}{300} + \frac{1}{200}\right)}} = 1,22.$$



Eksperimentalna vrednost ne pade v kritično območje. Zato ničelne domneve ne moremo zavrniti. Priljubljenost predsedniškega kandidata ni statistično značilno različna med mestnimi in vaškimi prebivalci. \diamond

15.11.2 Velik vzorec za testiranje domneve o $p_1 - p_2$, kadar je $D_0 \neq 0$.

$$\text{T.S.} = \frac{(\hat{p}_1 - \hat{p}_2) - D_0}{\sqrt{\frac{\hat{p}_1 \hat{q}_1}{n_1} + \frac{\hat{p}_2 \hat{q}_2}{n_2}}} \quad \text{sledi } z\text{-porazdelitev.}$$

Primer: Neka tovarna cigaret proizvaja dve znamki cigaret. Ugotovljeno je, da ima 56 od 200 kadilcev raje znamko A in da ima 29 od 150 kadilcev raje znamko B . Preveri domnevo pri 0,06 stopnji zaupanja, da bo prodaja znamke A boljša od prodaje znamke B za 10% proti alternativni domnevi, da bo razlika manj kot 10%. \diamond



15.12 Analiza variance

Če opravljamo isti poskus v nespremenjenih pogojih, kljub temu v rezultatu poskusa opazamo spremembe (variacije) ali odstopanja. Ker vzrokov ne poznamo in jih ne moremo kontrolirati, spremembe pripisujemo *slučajnim vplivom* in jih imenujemo *slučajna odstopanja*. Če pa enega ali več pogojev v poskusu spreminjamo, seveda dobimo dodatna odstopanja od povprečja. Analiza tega, ali so odstopanja zaradi sprememb različnih faktorjev ali pa zgolj slučajna, in kateri faktorji vplivajo na variacijo, se imenuje *analiza variance*.

Zgleda:

- (a) Namesto dveh zdravil proti nespečnosti kot v Studentovem primeru lahko preskušamo učinkovitost več različnih zdravil A, B, C, D, \dots in s preskušanjem domneve $H_0 : \mu_1 = \mu_2 = \dots = \mu_r$ raziskujemo, ali katero od zdravil sploh vpliva na rezultat. Torej je to posplošitev testa za $H_0 : \mu_1 = \mu_2$

- (b) Raziskujemo hektarski donos pšenice. Nanj vplivajo različni faktorji: različne sorte pšenice, različni načini gnojenja, obdelave zemlje itd., nadalje klima, čas sejanja itd

Analiza variance je nastala prav v zvezi z raziskovanjem v kmetijstvu. Glede na število faktorjev, ki jih spreminjamo, ločimo t.i. *enojno klasifikacijo* ali *enofaktorski eksperiment*, *dvojno klasifikacijo* ali *dvofaktorski eksperiment*, itd. V izrazu

$$Q_v^2 = \sum_{i=1}^r \sum_{k=1}^{n_i} (X_{ik} - \bar{X}_i)^2 = \sum_{i=1}^r (n_i - 1) S_i^2$$

je

$$S_i^2 = \frac{1}{n_i - 1} \sum_{k=1}^{n_i} (X_{ik} - \bar{X}_i)^2$$

nepristranska cenilka za disperzijo v i -ti skupini; neodvisna od S_j^2 , za $i \neq j$. Zato ima

$$\frac{Q_v^2}{\sigma^2} = \sum_{i=1}^r (n_i - 1) \frac{S_i^2}{\sigma^2}$$

porazdelitev $\chi^2(n - r)$, saj je ravno $\sum_{i=1}^r (n_i - 1) = n - r$ prostostnih stopenj. Ker je $E \frac{Q_v^2}{\sigma^2} = n - r$, je tudi $S_v^2 = \frac{1}{n - r} Q_v^2$ nepristranska cenilka za σ^2 . Izračunajmo še Q_m^2 pri predpostavki o veljavnosti osnovne domneve H_0 . Dobimo

$$Q_m^2 = \sum_{i=1}^r n_i (\bar{X}_i - \mu)^2 - n (\bar{X} - \mu)^2.$$

Torej je

$$\frac{Q_m^2}{\sigma^2} = \sum_{i=1}^r n_i \left(\frac{\bar{X}_i - \mu}{\sigma / \sqrt{n_i}} \right)^2 - n \left(\frac{\bar{X} - \mu}{\sigma / \sqrt{n}} \right)^2$$

od tu sprevidimo, da je statistika Q_m^2 / σ^2 porazdeljena po $\chi^2(r - 1)$. Poleg tega je $S_m^2 = Q_m^2 / (r - 1)$ nepristranska cenilka za σ^2 , neodvisna od S_v^2 . Ker sta obe cenilki za varianco σ^2 , pri domnevi H_0 , njuno razmerje $F = S_m^2 / S_v^2$ ne more biti zelo veliko. Iz

$$F = \frac{S_m^2}{S_v^2} = \frac{Q_m^2 / (r - 1)}{Q_v^2 / (n - r)} = \frac{Q_m^2 / (r - 1)}{Q_v^2 / (n - r)}$$

vidimo, da gre za Fisherjevo (Snedecorjevo) porazdelitev $F(r - 1, n - r)$. Podatke zapišemo v *tabelo analize variance*

VV	VK	PS	PK	F
faktor	Q_m^2	$r - 1$	S_m^2	F
slučaj	Q_v^2	$n - r$	S_v^2	
	Q^2	$n - 1$		

Analiza variance v R-ju

Zgled: Petnajst enako velikih njiv je bilo posejanih z isto vrsto pšenice, vendar gnojeno na tri različne načine – z vsakim po pet njiv.

```
> ena <- c(47,47,40,32,40)
> dva <- c(76,68,71,46,54)
> tri <- c(49,40,34,36,44)
> d <- stack(list(e=ena,d=dva,t=tri))
> names(d)
[1] "values" "ind"
> oneway.test(values ~ ind, data=d, var.equal=TRUE)

      One-way analysis of means

data:  values and ind
F = 10.5092, num df = 2, denom df = 12, p-value = 0.002304
> av <- aov(values ~ ind, data=d)
> summary(av)
      Df  Sum Sq Mean Sq F value    Pr(>F)
ind      2 1628.93   814.47  10.509 0.002304 **
Residuals 12  930.00    77.50
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

Domnevo H_0 zavrnemo.

15.12.1 Preverjanje domneve $\sigma^2 = \sigma_0^2$

$$\text{T.S.} = \frac{(n-1)s^2}{\sigma_0^2} \text{ sledi } \chi^2\text{-porazd.}$$



Če je

- $H_a : \sigma^2 > \sigma_0^2$, potem je **odločitveno pravilo**:
zavrni ničelno domnevo, če je test statistike večji ali enak $\chi_{(\alpha, n-1)}^2$.
- $H_a : \sigma^2 < \sigma_0^2$ potem je **odločitveno pravilo**:
zavrni ničelno domnevo, če je test statistike manjši ali enak $\chi_{(\alpha, n-1)}^2$.
- $H_a : \sigma^2 \neq \sigma_0^2$, potem je **odločitveno pravilo**:
zavrni ničelno domnevo, če je test statistike manjši ali enak $\chi_{(\alpha, n-1)}^2$
ali če je test statistike večji ali enak $\chi_{(\alpha, n-1)}^2$.

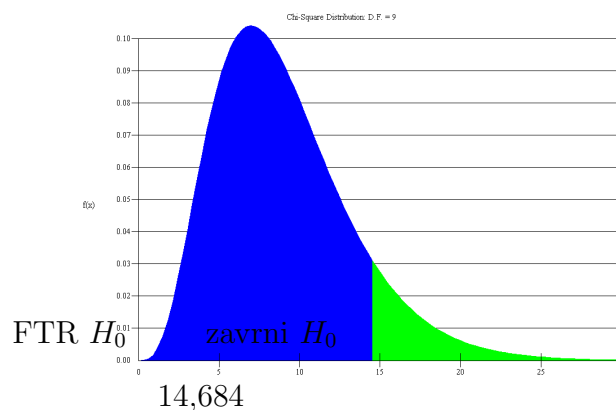
Primer: Količina pijače, ki jo naprava za mrzle napitke zavrže je normalno porazdeljena s povprečjem 12 unčev in standardnim odklonom 0,1 unče. Vsakič, ko servisirajo napravo, si izberejo 10 vzorcev in izmerijo zavrženo tekočino. Če je razpršenost zavržene količine prevelika, potem mora naprava na servis. **Ali naj jo odpeljejo na servis?** Uporabi $\alpha = 0,1$.

◇

15.12.2 Preverjanje domneve za varianco

- Ničelna domneva $H_0 : \sigma^2 = 0,01$,
- Alternativna domneva $H_a : \sigma^2 > 0,01$,
- Predpostavke
 - naključni vzorec
 - vzorčenje iz normalne porazdelitve.
- Testna statistika $\chi^2_{\nu=n-1} = S^2(n-1)/\sigma_0^2$.

Določimo zavrnitveni kriterij



Rezultati testiranja

- naredi naključni vzorec izračunamo naslednjo varianco vzorca: 0,02041,
- izračunaj vrednost testne statistike $\chi^2 = (0,02041)(9)/(0,01) = 18,369$,
- naredi odločitev: zavrni H_0 ,
- zaključek popravi napravo.

P-vrednost

- Sprejemljivost domneve H_0 na osnovi vzorca
 - možnost za opazovanje vzorca (ali bolj ekstremno podatkov), če je domneva H_0 pravilna
 - *P*-vrednost = $P(\chi^2 > 18,369) = 0,0311$
- Najmanjši α pri katerem zavrnemo domnevo H_0
 - *P*-vrednost $< \alpha$, zato zavrnemo domnevo H_0 .

15.12.3 Preverjanje domneve $\sigma_1^2/\sigma_2^2 = 1$

Če velja

$$H_a : \sigma_1^2/\sigma_2^2 > 1,$$

potem je **test statistike** enak s_1^2/s_2^2 , **odločitveno pravilo** pa je: zavrni ničelno domnevo, če velja

$$\text{T.S.} \geq F_{\alpha, n_1-1, n_2-1}.$$

Če velja $H_a : \sigma_1^2/\sigma_2^2 < 1$, potem je **test statistike** enak

$$\frac{\text{varianca večjega vzorca}}{\text{varianca manjšega vzorca}},$$

odločitveno pravilo pa je: zavrni ničelno domnevo, če velja $s_1^2 > s_2^2$ in

$$\text{T.S.} \geq F_{\alpha, n_1-1, n_2-1}$$

oziroma zavrni ničelno domnevo, če velja $s_1^2 < s_2^2$ in

$$\text{T.S.} \geq F_{\alpha, n_2-1, n_1-1}.$$

15.13 Preverjanje domnev o porazdelitvi spremenljivke

Do sedaj smo ocenjevali in preverjali domnevo o parametrih populacije kot μ , σ in π . Sedaj pa bomo preverjali, če se spremenljivka porazdeljuje po določeni porazdelitvi. Test je zasnovan na dejstvu, kako dobro se prilegajo empirične (eksperimentalne) frekvence vrednosti spremenljivke hipotetičnim (teoretičnim) frekvencam, ki so določene s predpostavljeno porazdelitvijo.

15.13.1 Preverjanje domneve o enakomerni porazdelitvi

Za primer vzemimo met kocke in za spremenljivko število pik pri metu kocke. Preizkusimo domnevo, da je kocka poštena, kar je enakovredno domnevi, da je porazdelitev spremenljivke enakomerna. Tedaj sta ničelna in osnovna domneva

H_0 : spremenljivka se porazdeljuje enakomerno,

H_1 : spremenljivka se ne porazdeljuje enakomerno.

Denimo, da smo 120-krat vrgli kocko ($n = 120$) in štejemo kolikokrat smo vrgli posamezno število pik. To so empirične ali opazovane frekvence, ki jih označimo s f_i . Teoretično, če je kocka poštena, pričakujemo, da bomo dobili vsako vrednost z verjetnostjo $1/6$ oziroma 20 krat. To so teoretične ali pričakovane frekvence, ki jih označimo s f'_i . Podatke zapišimo v naslednji tabeli

x_i	1	2	3	4	5	6
p_i	1/6	1/6	1/6	1/6	1/6	1/6
f'_i	20	20	20	20	20	20
f_i	20	22	17	18	19	24

S primerjavo empiričnih frekvenc z ustreznimi teoretičnimi frekvencami se moramo odločiti, če so razlike posledica le vzorčnih učinkov in je kocka poštena ali pa so razlike prevelike, kar kaže, da je kocka nepoštena. Statistika, ki meri prilagojenost empiričnih frekvenc teoretičnim je

$$\chi^2 = \sum_{i=1}^k \frac{(f_i - f'_i)^2}{f'_i},$$

ki se porazdeljuje po χ^2 porazdelitvi z $m = k - 1$ prostostnimi stopnjami, ki so enake številu vrednosti spremenljivke ali celic (k) minus število količin dobljenih iz podatkov, ki so uporabljene za izračun teoretičnih frekvenc.

V našem primeru smo uporabili le eno količino in sicer skupno število metov kocke ($n = 120$). Torej število prostostnih stopenj je $m = k - 1 = 6 - 1 = 5$. Ničelna in osnovna

domneva sta tedaj

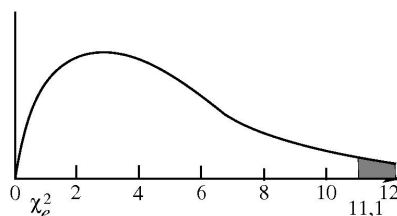
$$H_0 : \chi^2 = 0 \quad \text{in} \quad H_1 : \chi^2 > 0.$$

Doomnevo preverimo pri stopnji značilnosti $\alpha = 5\%$. Ker gre za enostranski test, je kritična vrednost enaka

$$\chi_{1-\alpha}^2(k-1) = \chi_{0,95}^2(5) = 11,1.$$

Eksperimentalna vrednost statistike pa je

$$\begin{aligned} \chi_e^2 &= \frac{(20-20)^2}{20} + \frac{(22-20)^2}{20} + \frac{(17-20)^2}{20} + \frac{(18-20)^2}{20} + \frac{(19-20)^2}{20} + \frac{(24-20)^2}{20} \\ &= \frac{4+9+4+1+16}{20} = \frac{34}{20} = 1,7. \end{aligned}$$



Ker eksperimentalna vrednost statistike ne pade v kritično območje, ničelne domneve ne moremo zavrnila. Empirične in teoretične frekvence niso statistično značilno različne med seboj.

15.13.2 Preverjanje domneve o normalni porazdelitvi

Omenjeni test najpogosteje uporabljamo za preverjanje ali se spremenljivka porazdeljuje normalno. V tem primeru je izračun teoretičnih frekvenc potrebno vložiti malo več truda.

Primer: Preizkusimo domnevo, da se spremenljivka telesna višina porazdeljuje normalno $N(177, 10)$. Domnevo preverimo pri 5% stopnji značilnosti. Podatki za 100 slučajno izbranih oseb so urejeni v frekvenčni porazdelitvi takole:

	f_i
nad 150-160	2
nad 160-170	20
nad 170-180	40
nad 180-190	30
nad 190-200	8
	100

Ničelna in osnovna domneva sta tedaj

$$H_0 : \chi^2 = 0 \quad \text{in} \quad H_1 : \chi^2 \neq 0.$$

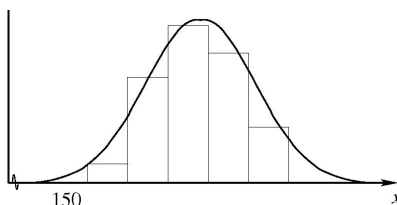
Za test uporabimo statistiko

$$\chi^2 = \sum_{i=1}^k \frac{(f_i - f'_i)^2}{f'_i},$$

ki se porazdeljuje po χ^2 porazdelitvi z $m = 5 - 1$ prostostnimi stopnjami. Kritična vrednost je

$$\chi_{0,95}^2(4) = 9,49.$$

V naslednjem koraku je potrebno izračunati teoretične frekvence. Najprej je potrebno za vsak razred izračunati verjetnost p_i , da spremenljivka zavzame vrednosti določenega intervala, če se porazdeljuje normalno. To lahko prikažemo na sliki:



Tako je na primer verjetnost, da je višina med 150 in 160 cm:

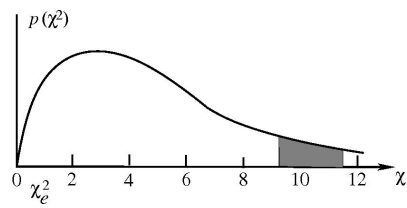
$$\begin{aligned} P(150 < X < 160) &= P\left(\frac{150 - 177}{10} < Z < \frac{160 - 177}{10}\right) \\ &= P(-2,7 < Z < -1,7) = H(2,7) - H(1,7) = 0,4965 - 0,4554 \\ &= 0,0411. \end{aligned}$$

Podobno lahko izračunamo ostale verjetnosti. Teoretične frekvence so $f'_i = n \times p_i$. Izračunane verjetnosti p_i in teoretične frekvence f'_i so

	f'_i	p_i	f'_i
nad 150-160	2	0,0411	4,11
nad 160-170	20	0,1974	19,74
nad 170-180	40	0,3759	37,59
nad 180-190	30	0,2853	28,53
nad 190-200	8	0,0861	8,61
	100		98,58

Eksperimentalna vrednost statistike je tedaj

$$\chi_e^2 = \frac{(2 - 4,11)^2}{4,11} + \frac{(20 - 19,74)^2}{19,74} + \frac{(40 - 37,59)^2}{37,59} + \frac{(30 - 28,53)^2}{28,53} + \frac{(8 - 8,61)^2}{8,61} \approx 1$$



Ker ekperimentalna vrednost ne pade v kritično območje, ne moremo zavrni ničelne domneve, da je spremenljivka normalno porazdeljena. \diamond

Obstajajo tudi drugi testi za preverjanje porazdelitve spremenljivke, npr. Kolmogorov-Smirnov test.

Pri objavi anketiranih rezultatov je potrebno navesti:

1. naročnika in izvajalca,
2. populacijo in vzorčni okvir,
3. opis vzorca,
4. velikost vzorca in velikost realiziranega vzorca (stopnja odgovorov)
5. čas, kraj in način anketiranja,
6. anketno vprašanje,
7. vzorčno napako.

Preizkus Kolmogorov-Smirnova v R-ju

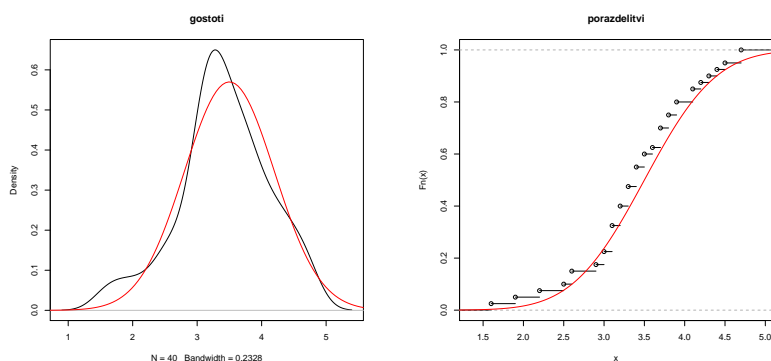
```
> t <- c(1.6,1.9,2.2,2.5,2.6,2.6,2.9,3.0,3.0,3.1,3.1,3.1,
+ 3.1,3.2,3.2,3.2,3.3,3.3,3.3,3.3,3.4,3.4,3.4,3.5,3.5,3.6,3.7,
+ 3.7,3.7,3.8,3.8,3.9,3.9,4.1,4.1,4.2,4.3,4.4,4.5,4.7,4.7)
> z <- sample(t)
> z
[1] 4.1 2.2 4.7 3.8 4.7 3.5 3.3 3.6 4.4 2.9 4.3 4.2 3.9 3.1
[15] 3.1 2.5 4.5 3.0 2.6 3.4 1.9 3.5 3.2 3.1 3.7 2.6 3.1 3.2
[29] 3.4 3.7 3.4 3.3 4.1 1.6 3.9 3.3 3.0 3.7 3.2 3.8
> "lot(density(z),main="gostoti")
> "urve(dnorm(x,mean=3.5,sd=0.7),add=TRUE,col="red")
> "lot(ecdf(z),main="porazdelitvi")
> curve(pnorm(x,mean=3.5,sd=0.7),add=TRUE,col="red")
> ks.test(z,"pnorm",mean=3.5,sd=0.7)
```

One-sample Kolmogorov-Smirnov test

```
data: z
D = 0.1068, p-value = 0.7516
alternative hypothesis: two.sided
```

```
Warning message:
cannot compute correct p-values with ties in:
ks.test(z, "pnorm", mean = 3.5, sd = 0.7)
```

Preizkus Kolmogorov-Smirnova v R-ju



V R-ju so pri preizkusih izpisane vrednosti p -value = Π
(preizkusna statistika ima pri veljavnosti osnovne domneve vrednost vsaj tako ekstremno, kot je zračunana).

[0, 0.001] – izjemno značilno (***);
(0.001, 0.01] – zelo značilno (**);
(0.01, 0.05] – statistično značilno (*);
(0.05, 0.1] – morda značilno;
(0.1, 1] – neznačilno.

Osnovno domnevo zavrnamo, če je p -value pod izbrano stopnjo značilnosti.

Preizkus Kolmogorov-Smirnova v R-ju

```
> p <- pnorm(t, mean=3.5, sd=0.7)
> s <- (1:40)/40; d <- abs(s-p)
> options(digits=3)
> cbind(t,p,s,d)
      t      p      s      d
[1,] 1.6 0.00332 0.025 0.02168 [21,] 3.4 0.44320 0.525 0.08180
[2,] 1.9 0.01114 0.050 0.03886 [22,] 3.4 0.44320 0.550 0.10680
[3,] 2.2 0.03165 0.075 0.04335 [23,] 3.5 0.50000 0.575 0.07500
[4,] 2.5 0.07656 0.100 0.02344 [24,] 3.5 0.50000 0.600 0.10000
[5,] 2.6 0.09927 0.125 0.02573 [25,] 3.6 0.55680 0.625 0.06820
[6,] 2.6 0.09927 0.150 0.05073 [26,] 3.7 0.61245 0.650 0.03755
[7,] 2.9 0.19568 0.175 0.02068 [27,] 3.7 0.61245 0.675 0.06255
[8,] 3.0 0.23753 0.200 0.03753 [28,] 3.7 0.61245 0.700 0.08755
[9,] 3.0 0.23753 0.225 0.01253 [29,] 3.8 0.66588 0.725 0.05912
[10,] 3.1 0.28385 0.250 0.03385 [30,] 3.8 0.66588 0.750 0.08412
[11,] 3.1 0.28385 0.275 0.00885 [31,] 3.9 0.71615 0.775 0.05885
[12,] 3.1 0.28385 0.300 0.01615 [32,] 3.9 0.71615 0.800 0.08385
[13,] 3.1 0.28385 0.325 0.04115 [33,] 4.1 0.80432 0.825 0.02068
[14,] 3.2 0.33412 0.350 0.01588 [34,] 4.1 0.80432 0.850 0.04568
[15,] 3.2 0.33412 0.375 0.04088 [35,] 4.2 0.84134 0.875 0.03366
[16,] 3.2 0.33412 0.400 0.06588 [36,] 4.3 0.87345 0.900 0.02655
[17,] 3.3 0.38755 0.425 0.03745 [37,] 4.4 0.90073 0.925 0.02427
[18,] 3.3 0.38755 0.450 0.06245 [38,] 4.5 0.92344 0.950 0.02656
[19,] 3.3 0.38755 0.475 0.08745 [39,] 4.7 0.95676 0.975 0.01824
[20,] 3.4 0.44320 0.500 0.05680 [40,] 4.7 0.95676 1.000 0.04324
> options(digits=7)
> max(d)
[1] 0.1067985
```

Preizkus χ^2 v R-ju

```
> a <- rbind(c(80,5,15), c(40,20,20), c(20,30,20))
> rownames(a) <- c("za", "proti", "neodlocen")
> colnames(a) <- c("do 25", "25-50", "nad 50")
> a
      do 25 25-50 nad 50
za      80     5    15
proti   40    20    20
neodlocen 20    30    20
> chisq.test(a)
```

Pearson's Chi-squared test

data: a

X-squared = 51.4378, df = 4, p-value = 1.808e-10

Poleg `ks.test` in `chisq.test` obstaja v R-ju še več drugih preizkusov: `prop.test`, `binom.test`, `t.test`, `wilcox.test`, `var.test`, `shapiro.test`, `cor.test`, `fisher.test`, `kruskal.test`.

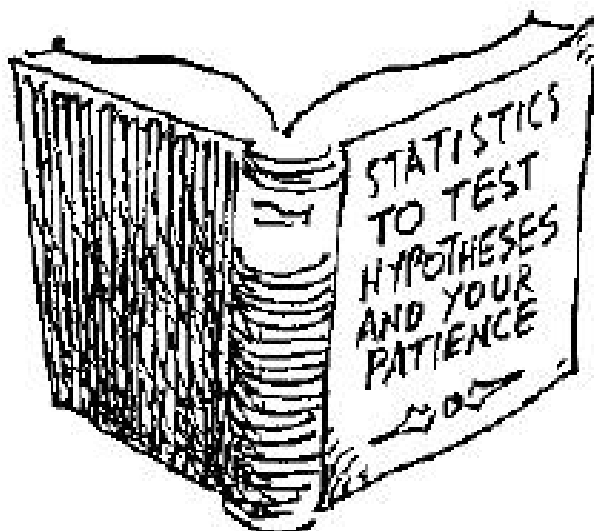
Opise posameznega preizkusa dobimo z zahtevo `help(preizkus)`.

Spearmanov preizkus povezanosti v R-ju

```
> slo <- c(5,3,1,2,4)
> mat <- c(5,2,3,1,4)
> slo
[1] 5 3 1 2 4
> mat
[1] 5 2 3 1 4
> cor.test(slo,mat,method="spearman")
```

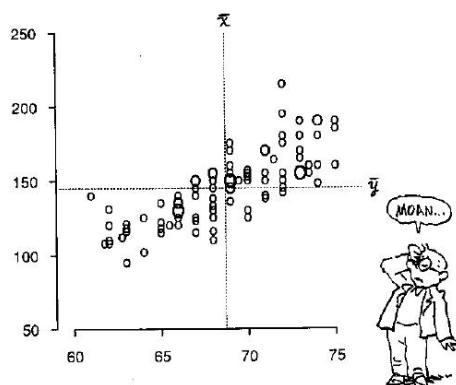
Spearman's rank correlation rho

```
data: slo and mat
S = 6, p-value = 0.2333
alternative hypothesis: true rho is not equal to 0
sample estimates:
rho
0.7
```



Poglavje 16

Bivariatna analiza in regresija



Bivariatna analiza

$X \longleftrightarrow Y$ povezanost

$X \longrightarrow Y$ odvisnost

Mere povezanosti ločimo glede na tip spremenljivk:

1. NOMINALNI tip para spremenljivk (ena od spremenljivk je nominalna): χ^2 , kontingenčni koeficienti, koeficienti asociacije;
2. ORDINALNI tip para spremenljivk (ena spremenljivka je ordinalna druga ordinalna ali boljša) koeficient korelacije rangov;
3. ŠTEVILSKI tip para spremenljivk (obe spremenljivki sta številski): koeficient korelacije.

16.1 Preverjanje domneve o povezanosti dveh nominalnih spremenljivk

Vzemimo primer:

- ENOTA: dodiplomski študent neke fakultete v letu 1993/94;
- VZOREC: slučajni vzorec 200 študentov;
- 1. SPREMENLJIVKA: spol;
- 2. SPREMENLJIVKA: stanovanje v času študija.

Zanima nas ali študentke drugače stanujejo kot študentje oziroma: ali sta spol in stanovanje v času študija povezana. V ta namen podatke študentov po obeh spremenljivkah uredimo v dvorazsežno frekvenčno porazdelitev. To tabelo imenujemo kontingenčna tabela. Denimo, da so podatki za vzorec urejeni v naslednji kontingenčni tabeli:

	starši	št. dom	zasebno	skupaj
moški	16	40	24	80
ženske	48	36	36	120
skupaj	64	76	60	200

Ker nas zanima ali študentke drugače stanujejo v času študija kot študentje, moramo porazdelitev stanovanja študentk primerjati s porazdelitvijo študentov. Ker je število študentk različno od števila študentov, moramo zaradi primerjave izračunati relativne frekvence:

	starši	št. dom	zasebno	skupaj
moški	20	50	30	100
ženske	40	30	30	100
skupaj	32	38	30	100

Če med spoloma ne bi bilo razlik, bi bili obe porazdelitvi (za moške in ženske) enaki porazdelitvi pod "skupaj". Naš primer kaže, da se odstotki razlikujejo: npr. le 20% študentov in kar 40% študentk živi med študijem pri starših. Odstotki v študentskih domovih pa so ravno obratni. Zasebno pa stanuje enak odstotek deklet in fantov. Že pregled relativnih frekvenc (po vrsticah) kaže, da sta spremenljivki povezani med seboj. Relativne frekvence lahko računamo tudi po stolpcih:

	starši	št. dom	zasebno	skupaj
moški	25	56,6	40	40
ženske	75	43,4	60	60
skupaj	100	100	100	100

Relativno frekvenco lahko prikažemo s stolpci ali krogi. Kontingenčna tabela kaže podatke za slučajni vzorec. Zato nas zanima, ali so razlike v porazdelitvi tipa stanovanja v času študija po spolu statistično značilne in ne le učinek vzorca.

H_0 : spremenljivki nista povezani

H_1 : spremenljivki sta povezani

Za preverjanje domneve o povezanosti med dvema nominalnima spremenljivkama na osnovi vzorčnih podatkov, podanih v dvo-razsežni frekvenčni porazdelitvi, lahko uporabimo χ^2 test. Ta test sloni na primerjavi empiričnih (dejanskih) frekvenc s teoretičnimi frekvencami, ki so v tem primeru frekvence, ki bi bile v kontingenčni tabeli, če spremenljivki ne bi bili povezani med seboj. To pomeni, da bi bili porazdelitvi stanovanja v času študija deklet in fantov enaki. Če spremenljivki nista povezani med seboj, so verjetnosti hkratne zgoditve posameznih vrednosti prve in druge spremenljivke enake produktu verjetnosti posameznih vrednosti. Npr., če označimo moške z M in stanovanje pri starših s S , je:

$$P(M) = \frac{80}{200} = 0,40, \quad P(S) = \frac{64}{200} = 0,32,$$

$$P(M \cap S) = P(M) \cdot P(S) = \frac{80}{200} \cdot \frac{64}{200} = 0,128.$$

Teoretična frekvenca je verjetnost $P(M \cap S)$ pomnožena s številom enot v vzorcu:

$$f'(M \cap S) = n \cdot P(M \cap S) = 200 \cdot \frac{80}{200} \cdot \frac{64}{200} = 25,6.$$

Podobno izračunamo teoretične frekvence tudi za druge celice kontingenčne tabele. Če teoretične frekvence zaokrožimo na cela števila, je tabela izračunanih teoretičnih frekvenc f'_i naslednja:

	starši	št. dom zasebno	skupaj
moški	26	30	80
ženske	38	46	120
skupaj	64	76	200

Spomnimo se tabel empiričnih (dejanskih) frekvenc f_i : χ^2 statistika, ki primerja dejanske in teoretične frekvence je

$$\chi^2 = \sum_{i=1}^k \frac{(f_i - f'_i)^2}{f'_i},$$

kjer je k število celic v kontingenčni tabeli. Statistika χ^2 se porazdeljuje po χ^2 porazdelitvi s $(s-1)(v-1)$ prostostnimi stopnjami, kjer je s število vrstic v kontingenčni tabeli in v število stolpcev. Ničelna in osnovna domneva sta v primeru tega testa

$H_0: \chi^2 = 0$ (spremenljivki nista povezani)

$H_1: \chi^2 > 0$ (spremenljivki sta povezani)

Iz tabele za porazdelitev χ^2 lahko razberemo kritične vrednost te statistike pri 5% stopnji značilnosti:

$$\chi_{1-\alpha}^2[(s-1)(v-1)] = \chi_{0,95}^2(2) = 5,99.$$

Eksperimentalna vrednost statistike χ^2 pa je:

$$\chi_e^2 = \frac{(16-26)^2}{26} + \frac{(40-30)^2}{30} + \frac{(24-24)^2}{24} + \frac{(48-38)^2}{38} + \frac{(36-46)^2}{46} + \frac{(36-36)^2}{36} = 12.$$

Ker je eksperimentalna vrednost večja od kritične vrednosti, pomeni, da pade v kritično območje. To pomeni, da ničelno domnevo zavrnamo. Pri 5% stopnji značilnosti lahko sprejmemo osnovno domnevo, da sta spremenljivki statistično značilno povezani med seboj. Statistika χ^2 je lahko le pozitivna. Zavzame lahko vrednosti v intervalu $[0, \chi_{\max}^2]$, kjer je $\chi_{\max}^2 = n(k-1)$, če je $k = \min(v, s)$. χ^2 statistika v splošnem ni primerljiva. Zato je definiranih več **kontingenčnih koeficientov**, ki so bolj ali manj primerni. Omenimo naslednje:

1. **Pearsonov koeficient:**

$$\Phi = \frac{\chi^2}{n},$$

ki ima zgornjo mejo $\Phi_{\max}^2 = k - 1$.

2. **Cramerjev koeficient:**

$$\alpha = \sqrt{\frac{\Phi^2}{k-1}} = \sqrt{\frac{\chi^2}{n(k-1)}},$$

ki je definiran na intervalu $[0, 1]$.

3. **Kontingenčni koeficient:**

$$C = \sqrt{\frac{\chi^2}{\chi^2 + n}},$$

ki je definiran na intervalu $[0, C_{\max}]$, kjer je $C_{\max} = \sqrt{k/(k-1)}$.

16.2 Koeficienti asociacije

Denimo, da imamo dve nominalni spremenljivki, ki imata le po dve vrednosti (sta dihoto-
mni). Povezanost med njima lahko računamo poleg kontingenčnih koeficientov s **koefici-**

enti asociacije na osnovi frekvenc iz kontingenčne tabele 2×2 :

$Y \setminus X$	x_1	x_2	
y_1	a	b	$a + b$
y_2	c	d	$c + d$
	$a + c$	$b + d$	N

kjer je $N = a + b + c + d$. Na osnovi štirih frekvenc v tabeli je definiranih več koeficientov asociacije:

- **Yulov koeficient asociacije:**

$$Q = \frac{ad - bc}{ad + bc} \in [-1, 1].$$

- **Sokal Michenerjev koeficient:**

$$S = \frac{a + d}{a + b + c + d} = \frac{a + d}{N} \in [0, 1].$$

- **Pearsonov koeficient:**

$$\phi = \frac{ad - bc}{\sqrt{(a + b)(c + d)(a + c)(b + d)}} \in [-1, 1].$$

Velja

$$\chi^2 = N \cdot \phi^2.$$

- **Jaccardov koeficient:**

$$J = \frac{a}{a + b + c} \in [0, 1],$$

- in še več drugih.

Primer: Vzemimo primer, ki kaže povezanost mod kaznivimi dejanji in alkoholizmom.

Tabela kaže podatke za $N = 10.750$ ljudi

alk. \ kaz. d.	DA	NE	skupaj
DA	50	500	550
NE	200	10.000	10.200
skupaj	250	10.500	10.750

Izračunajmo koeficiente asociacije:

$$Q = \frac{50 \times 10000 - 200 \times 500}{50 \times 10000 + 200 \times 500} = 0,67,$$

$$S = \frac{10050}{10750} = 0,93, \quad \text{in} \quad J = \frac{50}{50 + 500 + 200} = 0,066.$$

Izračunani koeficienti so precej različni. Yulov in Sokal Michenerjev koeficient kažeta na zelo močno povezanost med kaznjivimi dejanji in alkoholizmom, medtem kot Jaccardov koeficient kaže, da med spremenljivkama ni povezanosti. **Pri prvih dveh koeficientih povezanost povzroča dejstvo, da večina alkoholiziranih oseb ni naredila kaznivih dejanj in niso alkoholiki (frekvenca d).** Ker Jaccardov koeficient upošteva le DA DA ujemanje, je lažji za interpretacijo. V našem primeru pomeni, da oseba, ki je naredila kaznivo dejanje, sploh ni nujno alkoholik.

16.3 Preverjanje domneve o povezanosti dveh ordinalnih spremenljivk

V tem primeru gre za študij povezanosti med dvema spremenljivkama, ki sta vsaj ordinalnega značaja.

Primer: Vzemimo slučajni vzorec šestih poklicev in ocenimo, koliko so odgovorni (O) in koliko fizično naporni (N). V tem primeru smo poklice uredili od najmanj odgovornega do najbolj odgovornega in podobno od najmanj fizično napornega do najbolj napornega. Poklicem smo torej priredili range po odgovornosti (R_0) in po napornosti (R_N) od 1 do 6. Podatki so podani v tabeli:

poklic	R_0	R_N
<i>A</i>	1	6
<i>D</i>	2	4
<i>C</i>	3	5
<i>D</i>	4	2
<i>E</i>	5	3
<i>F</i>	6	1

Povezanost med spremenljivkama lahko merimo s koeficientom korelacije rangov r_s (Spearman), ki je definiran takole:

$$r_s = 1 - \frac{6 \cdot \sum_{i=1}^n d_i^2}{n \cdot (n^2 - 1)}.$$

kjer je d_i razlika med **rangoma** v i -ti enoti. Koeficient korelacije rangov lahko zavzame vrednosti v intervalu $[-1, 1]$. Če se z večanjem rangov po prvi spremenljivki večajo rangi tudi po drugi spremenljivki, gre za pozitivno povezanost. Tedaj je koeficient pozitiven in blizu 1. Če pa se z večanjem rangov po prvi spremenljivki rangi po drugi spremenljivki

manjšajo, gre za negativno povezanost. Koeficient je tedaj negativen in blizu -1 . V našem preprostem primeru gre negativno povezanost. Če ne gre za pozitivno in ne za negativno povezanost, rečemo, da spremenljivki nista povezani. Izračunajmo koeficient korelacije rangov za primer šestih poklicev:

poklic	R_0	R_N	d_i	d_i^2
<i>A</i>	1	6	-5	25
<i>B</i>	2	4	-2	4
<i>C</i>	3	5	-2	4
<i>D</i>	4	2	2	4
<i>E</i>	5	3	2	4
<i>F</i>	6	1	5	25
vsota			0	66

$$r_s = 1 - \frac{6 \cdot 66}{6 \cdot 35} = 1 - 1,88 = -0,88.$$

Res je koeficient blizu, kar kaže na močno negativno povezanost teh 6-ih poklicev. Omenili smo, da obravnavamo 6 slučajno izbranih poklicev. Zanima nas, ali lahko na osnovi tega vzorca posplošimo na vse poklice, da sta odgovornost in fizična napornost poklicev (negativno) povezana med seboj. Upoštevajmo 5% stopnjo značilnosti. Postavimo torej ničelno in osnovno domnevo:

$$H_0: \rho_s = 0 \text{ (spremenljivki nista povezani)}$$

$$H_1: \rho_s \neq 0 \text{ (spremenljivki sta povezani)}$$

kjer populacijski koeficient označimo s ρ_s . Pokaže se, da se statistika

$$t = \frac{r_s \cdot \sqrt{n-2}}{\sqrt{1-r_s^2}}$$

porazdeljuje približno po porazdelitvi z $m = (n-2)$ prostostnimi stopnjami. Ker gre za dvostranski test, sta kritični vrednosti enaki

$$\pm t_{\alpha/2} = \pm t_{0,025}(4) = \pm 2,776.$$

Eksperimentalna vrednost statistike je za naš primer

$$t_e = \frac{-0,88 \times 2}{\sqrt{1 - (-0,88)^2}} = \frac{-1,76}{0,475} = -3,71.$$

Eksperimentalna vrednost pade v kritično območje. Pri 5% stopnji značilnosti lahko rečemo, da sta odgovornost in fizična napornost (negativno) povezani med seboj. Če je ena od obeh spremenljivk številska, moramo vrednosti pred izračunom d_i rangirati. Če so kakšne vrednosti enake, zanje izračunamo povprečne pripadajoče range.

16.4 Preverjanje domneve o povezanosti dveh številskih spremenljivk

Vzemimo primer dveh številskih spremenljivk:

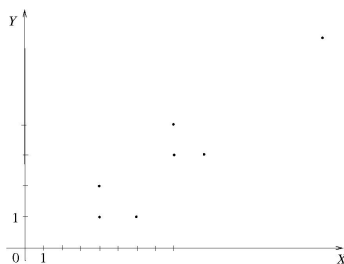
X - izobrazba (število priznanih let šole)

Y - število ur branja dnevnih časopisov na teden

Podatki za 8 slučajno izbranih oseb so:

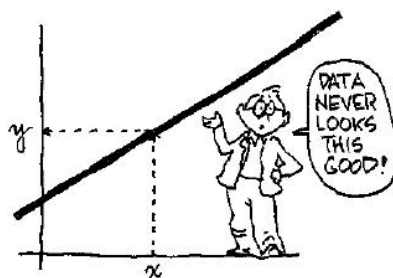
X	10	8	16	8	6	4	8	4
Y	3	4	7	3	1	2	3	1

Grafično lahko ponazorimo povezanost med dvema številskima spremenljivkama z razsevnim grafikonom. To je, da v koordinatni sistem, kjer sta koordinati obe spremenljivki, vrišemo enote s pari vrednosti. V našem primeru je izgleda razsevni grafikon takole:

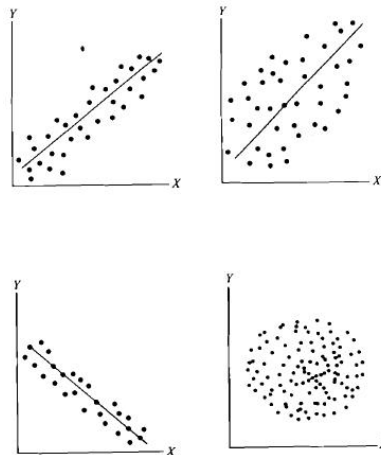


Tipi povezanosti:

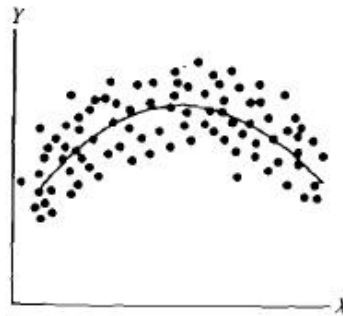
- **funkcijska** povezanost: vse točke ležijo na krivulji:
- **korelacijska** (stohastična) povezanost: točke so od krivulje bolj ali manj odklanjajo (manjša ali večja povezanost).



Tipični primeri linearne povezanosti spremenljivk:



Primer nelinearne povezanosti spremenljivk:



Kovarianca

$$\text{Cov}(X, Y) = \frac{1}{N} \sum_{i=1}^N (x_i - \mu_X) \cdot (y_i - \mu_Y)$$

meri linearno povezanost med spremenljivkama.

$\text{Cov}(X, Y) > 0$ pomeni pozitivno linearno povezanost,

$\text{Cov}(X, Y) = 0$ pomeni da ni linearne povezanosti,

$\text{Cov}(X, Y) < 0$ pomeni negativno linearno povezanost.

(Pearsonov) koeficient korelacije je

$$\rho_{XY} = \frac{\text{Cov}(X, Y)}{\sigma_X \sigma_Y} = \frac{\sum_{i=1}^N (x_i - \mu_X) \cdot (y_i - \mu_Y)}{\sqrt{\sum_{i=1}^N (x_i - \mu_X)^2 \cdot \sum_{i=1}^N (y_i - \mu_Y)^2}}$$

Koeficient korelacije lahko zavzame vrednosti v intervalu $[-1, 1]$. Če se z večanjem vrednosti prve spremenljivke večajo tudi vrednosti druge spremenljivke, gre za **pozitivno povezanost**. Tedaj je koeficient povezanosti blizu 1. Če pa se z večanjem vrednosti prve spremenljivke vrednosti druge spremenljivke manjšajo, gre za **negativno poveza-**

nost. Koeficient je tedaj negativen in blizu -1 . Če ne gre za pozitivno in ne za negativno povezanost, rečemo da spremenljivki nista povezani in koeficient je blizu 0.

Statistično sklepanje o korelacijski povezanosti:

Postavimo torej ničelno in osnovno domnevo:

$H_0: \rho = 0$ (spremenljivki nista linearno povezani)

$H_1: \rho \neq 0$ (spremenljivki sta linearno povezani)

Pokaže se, da se statistika

$$t = \frac{r \cdot \sqrt{n-2}}{\sqrt{1-r^2}}$$

porazdeljuje po t porazdelitvi z $m = (n-2)$ prostostnimi stopnjami. Z r označujemo koeficient korelacije na vzorcu in z ρ koeficient korelacije na populaciji.

Primer: Preverimo domnevo, da sta izobrazba (število priznanih let šole) in število ur branja dnevnih časopisov na teden povezana med seboj pri 5% stopnji značilnosti. Najprej izračunajmo vzorčni koeficient korelacije:

x_i	y_i	$x_i - \mu_x$	$y_i - \mu_y$	$(x_i - \mu_x)^2$	$(y_i - \mu_y)^2$	$\frac{(x_i - \mu_x) \cdot (y_i - \mu_y)}{(y_i - \mu_y)}$
10	3	2	0	4	0	0
8	4	0	1	0	1	0
16	7	8	4	64	16	32
8	3	0	0	0	0	0
6	1	-2	-2	4	4	4
4	2	-4	-1	16	1	4
8	3	0	0	0	0	0
4	1	-4	-2	16	4	8
64	24	0	0	104	26	48

$$r = \frac{\sum_{i=1}^n (x_i - \mu_x)(y_i - \mu_y)}{\sqrt{\sum_{i=1}^n (x_i - \mu_x)^2 \sum_{i=1}^n (y_i - \mu_y)^2}} = \frac{48}{\sqrt{104 \cdot 26}} = 0,92.$$

Ker gre za dvostranski test, je kritično območje določeno s kritičnima vrednostima

$$\pm t_{\alpha/2}(n-2) = \pm t_{0,025}(6) = \pm 2,447.$$

Eksperimentalna vrednost statistike pa je:

$$t_e = \frac{0,92\sqrt{8-2}}{\sqrt{1-0,92^2}} = 2,66.$$

Eksperimentalna vrednost pade v kritično območje. **Zaključek:** ob 5% stopnji značilnosti lahko rečemo, da je izobrazba linearno povezana z branjem dnevnih časopisov.

16.5 Parcialna korelacija

Včasih je potrebno meriti zvezo med dvema spremenljivkama in odstraniti vpliv vseh ostalih spremenljivk. To zvezo dobimo s pomočjo koeficienta parcialne korelacije. Pri tem seveda predpostavljamo, da so vse spremenljivke med seboj linearno povezane. Če hočemo iz zveze med spremenljivkama X in Y odstraniti vpliv tretje spremenljivke Z , je **koeficient parcialne korelacije**:

$$r_{XY,Z} = \frac{r_{XY} - r_{XZ}r_{YZ}}{\sqrt{1 - r_{XZ}^2} \sqrt{1 - r_{YZ}^2}}.$$

Tudi ta koeficient, ki zavzema vrednosti v intervalu $[-1, 1]$, interpretiramo podobno kot običajni koeficient korelacije. S pomočjo tega obrazca lahko razmišljamo naprej, kako bi izločili vpliv naslednjih spremenljivk.

Primer: V neki ameriški raziskavi, v kateri so proučevali vzroke za kriminal v mestih, so upoštevali naslednje spremenljivke:

- X : % nebelih prebivalcev,
- Y : % kaznivih dejanj,
- Z : % revnih prebivalcev,
- U : velikost mesta.

Izračunali so naslednje koeficiente korelacije:

	X	Z	U	Y
X	1	0,51	0,41	0,36
Z		1	0,29	0,60
U			1	0,49
Y				1

Zveza med nebelim prebivalstvom in kriminalom je

$$r_{XY} = 0,36.$$

Zveza je kar močna in lahko bi mislili, da nebeli prebivalci povzročajo več kaznivih dejanj. Vidimo pa še, da je zveza med revščino in kriminalom tudi precejšna

$$r_{YZ} = 0,60.$$

Lahko bi predpostavili, da revščina vpliva na zvezo med nebelci in kriminalom, saj je tudi zveza med revnimi in nebelimi precejšna $r_{XZ} = 0,51$. Zato poskusim odstraniti vpliv revščine iz zveze: “nebelo prebivalstvo : kazniva dejanja”:

$$r_{XY,Z} = \frac{0,36 - 0,51 \cdot 0,60}{\sqrt{1 - 0,51^2} \sqrt{1 - 0,60^2}}.$$

Vidimo, da se je linearna zveza zelo zmanjšala. Če pa odstranimo še vpliv velikosti mesta, dobimo parcialno korelacijo $-0,02$ oziroma zveze praktično ni več. \diamond

16.6 Regresijska analiza

Regresijska funkcija $Y' = f(X)$ kaže, kakšen bi bil vpliv spremenljivke X na Y , če razen vpliva spremenljivke X ne bi bilo drugih vplivov na spremenljivko Y . Ker pa so ponavadi še drugi vplivi na proučevano spremenljivko Y , se točke, ki predstavljajo enote v razsevnem grafikonu, odklanjajo od idealne regresijske krivulje

$$Y = Y' + E = f(X) + E$$

kjer X imenujemo neodvisna spremenljivka, Y odvisna spremenljivka in E člen napake (ali motnja, disturbanca). Če je regresijska funkcija linearna:

$$Y' = f(X) = a + bX,$$

je regresijska odvisnost

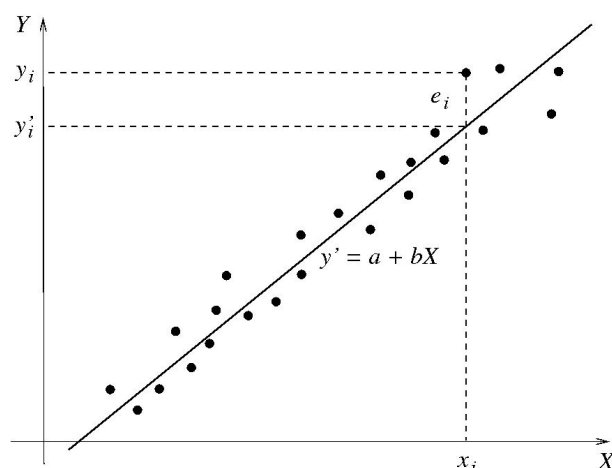
$$Y = Y' + E = a + bX + E$$

oziroma za i to enoto

$$y_i = y'_i + e_i = a + bx_i + e_i.$$



Regresijsko odvisnost si lahko zelo nazorno predstavimo v razsevnem grafikonu:



Regresijsko funkcijo lahko v splošnem zapišemo

$$Y' = f(X, a, b, \dots),$$

kjer so a, b, \dots parametri funkcije. Ponavadi se moramo na osnovi pregleda razsevnega grafikona odločiti za tip regresijske funkcije in nato oceniti parametre funkcije, tako da se regresijska krivulja kar se da dobro prilega točkam v razsevnom grafikonu.

Pri dvorazsežno normalno porazdeljenem slučajnem vektorju $(X, Y) : N(\mu_x, \mu_y, \sigma_x, \sigma_y, \rho)$ je, kot vemo

$$E(Y|x) = \mu_y + \rho \frac{\sigma_y}{\sigma_x} (x - \mu_x).$$

Pogojna porazdelitev Y glede na X je tudi normalna:

$$N\left(\mu_y + \rho \frac{\sigma_y}{\sigma_x} (x - \mu_x), \sigma_y \sqrt{1 - \rho^2}\right).$$

Regresija je linearna in regresijska krivulja premica, ki gre skozi točko (σ_x, σ_y) . Med Y in X ni linearne zveze, sta le 'v povprečju' linearno odvisni. Če označimo z $\beta = \rho \frac{\sigma_y}{\sigma_x}$ *regresijski koeficient*, $\alpha = \mu_y - \beta \mu_x$ in $\sigma^2 = \sigma_y \sqrt{1 - \rho^2}$, lahko zapišemo zvezo v obliki

$$y = \alpha + \beta x.$$

Preizkušanje regresijskih koeficientov

Po metodi momentov dobimo cenilki za α in β :

$$B = R \frac{C_y}{C_x} = \frac{C_{xy}}{C_x^2} \quad \text{in} \quad A = \bar{Y} - B\bar{X},$$

kjer so $C_x^2 = \sum_{i=1}^n (X_i - \bar{X})^2$, $C_y^2 = \sum_{i=1}^n (Y_i - \bar{Y})^2$ in $C_{xy} = \sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y})$.

Kako sta cenilki B in A porazdeljeni?

$$B = \frac{C_{xy}}{C_x^2} = \sum_{i=1}^n \frac{X_i - \bar{X}}{C_x^2} (Y_i - \bar{Y}) = \sum_{i=1}^n \frac{X_i - \bar{X}}{C_x^2} Y_i.$$

Ker proučujemo pogojno porazdelitev Y glede na X (torej so vrednost X poznane), obravnavamo spremenljivke X_1, \dots, X_n kot konstante. Ker je B linearna funkcija spremenljivk Y_1, \dots, Y_n , ki so normalno porazdeljene $Y_i : N(\alpha + \beta X_i, \sigma)$, je tudi B normalno porazdeljena. Določimo parametra te porazdelitve:

$$\mathbb{E}B = \sum_{i=1}^n \frac{X_i - \bar{X}}{C_x^2} \mathbb{E}Y_i = \sum_{i=1}^n \frac{X_i - \bar{X}}{C_x^2} (\alpha + \beta X_i) = \sum_{i=1}^n \frac{X_i - \bar{X}}{C_x^2} \beta (X_i - \bar{X}) = \beta.$$

Pri tem upoštevamo, da je $\sum_{i=1}^n (X_i - \bar{X}) = 0$ in da sta α ter \bar{X} konstanti.

$$\mathbb{D}B = \sum_{i=1}^n \frac{(X_i - \bar{X})^2}{C_x^4} \mathbb{D}Y_i = \frac{\sigma^2}{C_x^2}.$$

Torej je $B : N\left(\beta, \frac{\sigma}{C_x}\right)$, oziroma $\frac{B - \beta}{\sigma} C_x : N(0, 1)$.

Podobno dobimo

$$\mathbb{E}A = \alpha \quad \text{in} \quad \mathbb{D}A = \sigma^2 \left(\frac{1}{n} + \frac{\bar{X}^2}{C_x^2} \right).$$

Težje se je dokopati do cenilke za parameter σ^2 . Označimo $Q^2 = \sum_{i=1}^n (Y_i - A - BX_i)^2$.

Po nekaj računanja se izkaže, da velja $\mathbb{E} \frac{Q^2}{\sigma^2} = n - 2$.

Torej je $S^2 = \frac{Q^2}{n - 2} = \frac{\sigma^2}{n - 2} \frac{Q^2}{\sigma^2}$ nepristranska cenilka za σ^2 .

S^2 je neodvisna od A in B . Testni statistiki za A in B sta tedaj

$$T_A = \frac{A - \mathbb{E}A}{\sqrt{\mathbb{D}A}} = \frac{A - \alpha}{S} \sqrt{\frac{nC_x^2}{C_x^2 + n\bar{X}^2}} = \frac{A - \alpha}{S} C_x \sqrt{\frac{n}{\sum_{i=1}^n X_i^2}},$$

$$T_B = \frac{B - \mathbb{E}B}{\sqrt{\mathbb{D}B}} = \frac{B - \beta}{S} C_x,$$

ki sta obe porazdeljeni po Studentu $S(n - 2)$. Statistika za σ^2 pa je spremenljivka $\frac{Q^2}{\sigma^2} = (n - 2) \frac{S^2}{\sigma^2}$, ki je porazdeljena po $\chi^2(n - 2)$. Pokazati je mogoče tudi, da velja

$$Q^2 = C_y^2 - B^2 C_x^2 = C_y^2 (1 - R^2).$$

To nam omogoča S v statistikah zapisati z C_y in R . Te statistike uporabimo tudi za določitev intervalov zaupanja za parametre α , β in σ^2 .

16.7 Linearni model

Pri proučevanju pojavov pogosto teorija postavi določeno funkcijsko zvezo med obravnavanimi spremenljivkami. Oglejmo si primer *linernega modela*, ko je med spremenljivkama x in y linearna zveza

$$y = \alpha + \beta x$$

Za dejanske meritve se pogosto izkaže, da zaradi različnih vplivov, ki jih ne poznamo, razlika $u = y - \alpha - \beta x$ v splošnem ni enaka 0, čeprav je model točen. Zato je ustreznejši *verjetnostni linearni model*

$$Y = \alpha + \beta X + U,$$

kjer so X , Y in U slučajne spremenljivke in $EU = 0$ – model je vsaj v povprečju linearen.

Slučajni vzorec (meritve) $(X_1, Y_1), \dots, (X_n, Y_n)$ je realizacija slučajnega vektorja. Vpeljimo spremenljivke

$$U_i = Y_i - \alpha - \beta X_i$$

in predpostavimo, da so spremenljivke U_i med seboj neodvisne in enako porazdeljene z matematičnim upanjem 0 in disperzijo σ^2 . Torej je:

$$EU_i = 0, \quad DU_i = \sigma^2 \quad \text{in} \quad E(U_i U_j) = 0, \quad \text{za } i \neq j.$$

Običajno privzamemo še, da lahko vrednosti X_i točno določamo – X_i ima vedno isto vrednost. Poleg tega naj bosta vsaj dve vrednosti X različni. Težava je, da (koeficientov) premice $y = \alpha + \beta x$ ne poznamo. Recimo, da je približek zanjo premica $y = a + bx$. Določimo jo po *načelu najmanjših kvadratov* z minimizacijo funkcije

$$f(a, b) = \sum_{i=1}^n (y_i - (bx_i + a))^2.$$

Naloga zadošča pogojem izreka. Iz pogoja $\nabla P = 0$ dobimo enačbi

$$\begin{aligned} \frac{\partial f}{\partial a} &= \sum_{i=1}^n 2(y_i - (bx_i + a)) = 0, \\ \frac{\partial f}{\partial b} &= \sum_{i=1}^n 2(y_i - (bx_i + a))x_i = 0, \end{aligned}$$

z rešitvijo

$$b = \frac{n \sum xy - \sum x \sum y}{n \sum x^2 - (\sum x)^2}, \quad a = \frac{1}{n} \left(\sum y - b \sum x \right).$$

oziroma, če vpeljemo oznako $\bar{z} = \frac{1}{n} \sum z$:

$$b = \frac{\overline{xy} - \bar{x}\bar{y}}{\overline{x^2} - \bar{x}^2}, \quad a = \bar{y} - b\bar{x}.$$

Poglejmo še Hessovo matriko

$$H = \begin{bmatrix} \frac{\partial^2 f}{\partial a^2} & \frac{\partial^2 f}{\partial a \partial b} \\ \frac{\partial^2 f}{\partial b \partial a} & \frac{\partial^2 f}{\partial b^2} \end{bmatrix} = 2 \begin{bmatrix} n & \sum x \\ \sum x & \sum x^2 \end{bmatrix}.$$

Ker je $\Delta_1 = 2 \sum x^2 > 0$ in

$$\Delta_2 = 4 \left(n \sum x^2 - \left(\sum x \right)^2 \right) = 2 \sum \sum (x_i - x_j)^2 > 0,$$

je matrika H pozitivno definitna in zato funkcija P strogo konveksna. Torej je *regresijska premica* enolično določena. Seveda sta parametra a in b odvisna od slučajnega vzorca – torej slučajni spremenljivki. Iz dobljenih zvez za a in b dobimo že znani cenilki za koeficients α in β

$$B = \frac{C_{xy}}{C_x^2} \quad \text{in} \quad A = \bar{Y} - B\bar{X}.$$

Iz prej omenjenih predpostavk lahko (brez poznavanja porazdelitve Y in U) pokažemo

$$EA = \alpha \quad \text{in} \quad DA = \sigma^2 \left(\frac{1}{n} + \frac{\bar{X}^2}{C_x^2} \right), \quad EB = \beta \quad \text{in} \quad DB = \frac{\sigma^2}{C_x^2}, \quad K(A, B) = -\sigma^2 \frac{\bar{X}}{C_x^2}.$$

Cenilki za A in B sta najboljši linearni nepristranski cenilki za α in β .

```
> x <- c(3520, 3730, 4110, 4410, 4620, 4900, 5290, 5770, 6410, 6920, 7430)
> y <- c(166, 153, 177, 201, 216, 208, 227, 238, 268, 268, 274)
> l <- 1947:1957
> plot(y ~ x); abline(lm(y ~ x), col="red")
> m <- lm(y ~ x)
> summary(m)
```

Call:

```
lm(formula = y ~ x)
```

Residuals:

```
      Min       1Q   Median       3Q      Max
-19.2149  -5.4003   0.3364   6.8453  16.0204
```

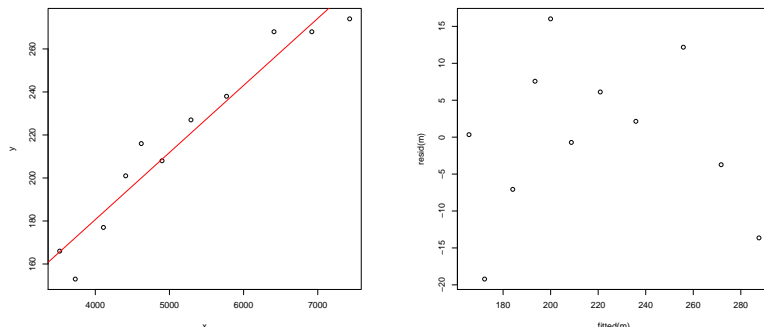
Coefficients:

```
              Estimate Std. Error t value Pr(>|t|)
(Intercept) 55.852675   14.491253   3.854  0.00388 **
x            0.031196    0.002715  11.492 1.11e-06 ***
```

```
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

```
Residual standard error: 11.18 on 9 degrees of freedom
Multiple R-Squared:  0.9362,    Adjusted R-squared:  0.9291
F-statistic: 132.1 on 1 and 9 DF,  p-value: 1.112e-06
```

```
> plot(fitted(m), resid(m))
```



```
> coef(m)
(Intercept)      x
55.8526752    0.0311963
```

To metodo ocenjevanja parametrov regresijske funkcije imenujemo **metoda najmanjših kvadratov**.

Če izračunana parametra vstavimo v regresijsko funkcijo, dobimo:

$$Y = \mu_Y + \frac{\text{Cov}(X, Y)}{\sigma_X^2}(X - \mu_X).$$

To funkcijo imenujemo tudi **prva** regresijska funkcija. Podobno bi lahko ocenili linearno regresijsko funkcijo

$$X = a^* + b^*Y.$$

Če z metodo najmanjših kvadratov podobno ocenimo parametra a^* in b^* , dobimo:

$$X = \mu_X + \frac{\text{Cov}(X, Y)}{\sigma_Y^2}(Y - \mu_Y).$$

To funkcijo imenujemo **druga** regresijska funkcija,

Primer: Vzemimo primer 8 oseb, ki smo ga obravnavali v poglavju o povezanosti dveh številskih spremenljivk. Spremenljivki sta bili:

X - izobrazba (število priznanih let šole),

Y - št. ur branja dnevnih časopisov na teden.

Spomnimo se podatkov za teh 8 slučajno izbranih oseb:

X	10	8	16	8	6	4	8	4
Y	3	4	7	3	1	2	3	1

Zanje izračunajmo obe regresijski premici in ju vpišimo v razsevni grafikon. Ko smo računali koeficient korelacije smo že izračunali aritmetični sredini

$$\mu_X = \frac{64}{8} = 8, \quad \mu_Y = \frac{24}{8} = 3,$$

vsoti kvadratov odklonov od aritmetične sredine za obe spremenljivki

$$\sum_{i=1}^n (x_i - \mu_X)^2 = 104, \quad \sum_{i=1}^n (y_i - \mu_Y)^2 = 26$$

in vsoto produktov odklonov od obeh aritmetičnih sredin

$$\sum_{i=1}^n (x_i - \mu_X)(y_i - \mu_Y) = 48.$$

Potem sta regresijski premici

$$Y = \mu_Y + \frac{\sum_{i=1}^N (x_i - \mu_X)(y_i - \mu_Y)}{\sum_{i=1}^N (x_i - \mu_X)^2} (X - \mu_X),$$

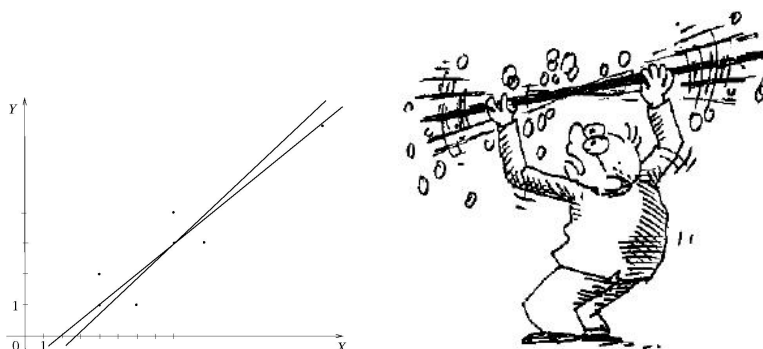
$$X = \mu_X + \frac{\sum_{i=1}^N (x_i - \mu_X)(y_i - \mu_Y)}{\sum_{i=1}^N (y_i - \mu_Y)^2} (Y - \mu_Y),$$

oziroma

$$Y = 3 + \frac{48}{104}(X - 8) = -0,68 + 0,46X,$$

$$X = 8 + \frac{48}{26}(Y - 3) = -2,46 + 1,85Y.$$

Obe regresijski premici lahko vrisemo v razsevni grafikon in preverimo, če se res najboljše prilegata točkam v grafikonu: ◇



Regresijski premici se sečeta v točki, določeni z aritmetičnima sredinama spremenljivk X in Y . Dokažite, da se premici vedno sečeta v tej točki.

16.7.1 Statistično sklepanje o regresijskem koeficientu

Vpeljmo naslednje oznake:

$Y = \alpha + \beta X$ regresijska premica na populaciji,

$Y = a + bX$ regresijska premica na vzorcu.

Denimo, da želimo preveriti domnevo o regresijskem koeficientu β . Postavimo ničelno in osnovno domnevo takole:

$H_0: \beta = \beta_0$,

$H_1: \beta \neq \beta_0$.

Nepristranska cenilka za regresijski koeficient β je $b = \text{Cov}(X, Y)/s_X^2$, ki ima matematično upanje in standardno napako:

$$E b = \beta; \quad \text{SE}(b) = \frac{s_Y \sqrt{1 - r^2}}{s_X \sqrt{n - 2}}.$$

Testna statistika za zgornjo ničelno domnevo je:

$$t = \frac{s_Y \sqrt{n - 2}}{s_X \sqrt{1 - r^2}} (b - \beta_0),$$

ki se porazdeljuje po t -porazdelitvi z $m = (n - 2)$ prostostnimi stopnjami.

Primer: Vzemimo primer, ki smo ga že obravnavali. Spremenljivki sta

X - izobrazba (število priznanih let šole),

Y - št. ur branja dnevnih časopisov na teden.

Podatke za slučajno izbrane enote ($n = 8$) najdemo na prejšnjih prosojnicah.

Preverimo domnevo, da je regresijski koeficient različen od 0 pri $\alpha = 5\%$.

Postavimo najprej ničelno in osnovno domnevo:

$H_0: \beta = 0$,

$H_1: \beta \neq 0$.

Gre za dvostranski test. Zato je ob 5% stopnji značilnosti kritično območje določeno s kritičnima vrednostima:

$$\pm t_{\alpha/2}(n - 2) = \pm t_{0,025}(6) = \pm 2,447.$$

Eksperimentalna vrednost statistike pa je:

$$t_e = \sqrt{\frac{104 \cdot (8 - 2)}{26 \cdot (1 - 0,92^2)}} \cdot (0,46 - 0) = 5,8.$$

Regresijski koeficient je statistično značilno različen od 0. ◇

16.7.2 Pojasnjena varianca (ang. ANOVA)

Vrednost odvisne spremenljivke Y_i lahko razstavimo na tri komponente:

$$y_i = \mu_Y + (y'_i - \mu_Y) + (y_i - y'_i),$$

kjer so pomeni posameznih komponent

μ_Y : rezultat splošnih vplivov,

$(y'_i - \mu_Y)$: rezultat vpliva spremenljivke X (regresija),

$(y_i - y'_i)$: rezultat vpliva drugih dejavnikov (napake/motnje).

Če zgornjo enakost najprej na obeh straneh kvadriramo, nato seštejemo po vseh enotah in končno delimo s številom enot (N), dobimo:

$$\frac{1}{N} \sum_{i=1}^N (y_i - \mu_Y)^2 = \frac{1}{N} \sum_{i=1}^N (y'_i - \mu_Y)^2 + \frac{1}{N} \sum_{i=1}^N (y_i - y'_i)^2.$$

To lahko zapišemo takole:

$$\sigma_Y^2 = \sigma_{Y'}^2 + \sigma_e^2,$$

kjer posamezni členi pomenijo:

σ_Y^2 : celotna varianca spremenljivke Y ,

$\sigma_{Y'}^2$: pojasnjena varianca spremenljivke Y ,

σ_e^2 : nepojasnjena varianca spremenljivke Y .

Delež pojasnjene variance spremenljivke Y s spremenljivko X je

$$R = \frac{\sigma_{Y'}^2}{\sigma_Y^2}.$$

Imenujemo ga **determinacijski koeficient** in je definiran na intervalu $[0, 1]$. Pokazati se da, da je v primeru linearne regresijske odvisnosti determinacijski koeficient enak

$$R = \rho^2,$$

kjer je ρ koeficient korelacije. Kvadratni koren iz nepojasnjene variance σ_e imenujemo **standardna napaka regresijske ocene**, ki meri razpršenost točk okoli regresijske krivulje. Standardna napaka ocene meri kakovost ocenjevanja vrednosti odvisne spremenljivke z regresijsko funkcijo. V primeru linearne regresijske odvisnosti je standardna napaka enaka:

$$\sigma_e = \sigma_Y \sqrt{1 - \rho^2}.$$

Primer: Vzemimo spremenljivki

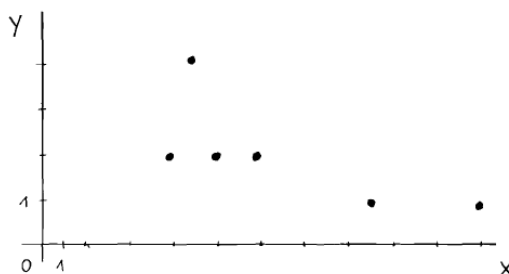
X - število ur gledanja televizije na teden

Y - število obiskov kino predstav na mesec

Podatki za 6 oseb so:

X	10	15	6	7	20	8
Y	2	1	2	4	1	2

Z linearno regresijsko funkcijo ocenimo, kolikokrat bo šla oseba v kino na mesec, če gleda 18 ur na teden televizijo. **Kolikšna je standardna napaka? Kolikšen delež variance obiska kinopredstav lahko pojasnimo z gledanjem televizije?** Najprej si podatke predstavimo v razsevnem grafikonu:



Za odgovore potrebujemo naslednje izračune:

x_i	y_i	$x_i - \mu_x$	$y_i - \mu_y$	$(x_i - \mu_x)^2$	$(y_i - \mu_y)^2$	$\frac{(x_i - \mu_x) \cdot (y_i - \mu_y)}{(y_i - \mu_y)}$
10	2	-1	0	1	0	0
15	1	4	-1	16	1	-4
6	2	-5	0	25	0	0
7	4	4	2	16	4	-8
20	1	9	-1	81	1	-9
8	2	-3	0	9	0	0
66	12	0	0	148	6	21

$$Y' = 2 - \frac{21}{148} (X - 11) = 3,54 - 0,14X$$

$$y'(18) = 3,54 - 0,14 \cdot 18 = 1,02$$

$$\rho = \frac{21}{146 \cdot 6} = -0,70$$

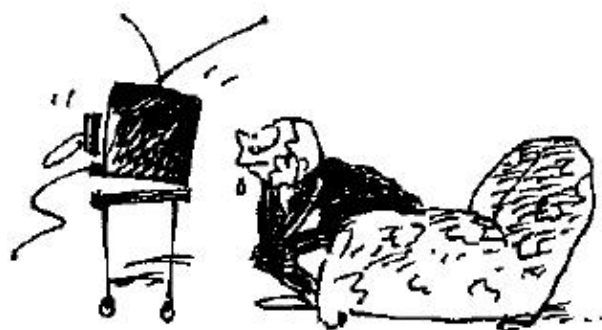
$$\sigma_e^2 = \frac{6}{6} \sqrt{1 - (-0,70)^2} = \sqrt{0,51} = 0,71$$

$$R = (0,70)^2 = 0,49$$

Če oseba gleda 18 ur na teden televizijo, lahko pričakujemo, da bo 1-krat na mesec šla v kino, pri čemer je standardna napaka 0,7. 49% variance obiska kino predstav lahko pojasnimo z gledanjem televizije. \diamond

Poglavje 17

Časovne vrste



Družbeno-ekonomski pojavi so časovno spremenljivi. Spremembe so rezultat delovanja najrazličnejših dejavnikov, ki tako ali drugače vplivajo na pojave. Sliko dinamike pojavov dobimo s časovnimi vrstami. *Časovna vrsta* je niz istovrstnih podatkov, ki se nanašajo na zaporedne časovne razmike ali trenutke. Osnovni namen analize časovnih vrst je

- opazovati časovni razvoj pojavov,
- iskati njihove zakonitosti in
- predvidevati nadaljni razvoj.

Seveda to predvidevanje ne more biti popolnoma zanesljivo, ker je skoraj nemogoče vnaprej napovedati in upoštevati vse faktorje, ki vplivajo na proučevani pojav. Napoved bi veljala strogo le v primeru, če bi bile izpolnjene predpostavke, pod katerimi je napoved izdelana. Časovne vrste prikazujejo individualne vrednosti neke spremenljivke v času. Čas lahko interpretiramo kot trenutek ali razdobje; skladno s tem so časovne vrste

- trenutne, npr. število zaposlenih v določenem trenutku:
- intervalne, npr. družbeni proizvod v letu 1993.

Časovne vrste analiziramo tako, da opazujemo spreminjanje vrednosti členov v časovih vrstah in iščemo zakonitosti tega spreminjanja. Naloga enostavne analize časovnih vrst je primerjava med členi v isti časovni vrsti. Z metodami, ki so specifične za analizo časovnih vrst, analiziramo zakonitosti dinamike ene same vrste, s korelacijsko analizo pa zakonitosti odvisnosti v dinamiki več pojavov, ki so med seboj v zvezi.

Primer: Vzemimo število nezaposlenih v Sloveniji v letih od 1981 do 1990. V metodoloških pojasnilih v Statističnem letopisu Republike Slovenije 1991, so nezaposlni (spremenljivka X) opredeljeni takole:

Brezposelna oseba je oseba, ki je sposobna in voljna delati ter je pripravljena sprejeti zaposlitev, ki ustreza njeni strokovni izobrazbi oz. z delom pridobljeni delovni zmožnosti, vendar brez svoje krivde nima dela in možnosti, da si z delom zagotavlja sredstva za preživetje in se zaradi zaposlitve prijavi pri območni enoti Zavoda za zaposlovanje (do leta 1989 skupnosti za zaposlovanje).

leto	X_k
1981	12.315
1982	13.700
1983	15.781
1984	15.300
1985	11.657
1986	14.102
1987	15.184
1988	21.311
1989	28.218
1990	44.227

◇

17.1 Primerljivost členov v časovni vrsti

Kljub temu, da so členi v isti časovni vrsti istovrstne količine, dostikrat niso med seboj neposredno primerljivi.

Osnovni pogoj za primerljivost členov v isti časovni vrsti je pravilna in nedvoumna opredelitev pojava, ki ga časovna vrsta prikazuje.

Ta opredelitev mora biti vso dobo opazovanja enaka in se ne sme spreminjati.

Ker so spremembe pojava, ki ga časovna vrsta prikazuje bistveno odvisne od časa, je zelo koristno, če so **časovni razmiki med posameznimi členi enaki**. Na velikost

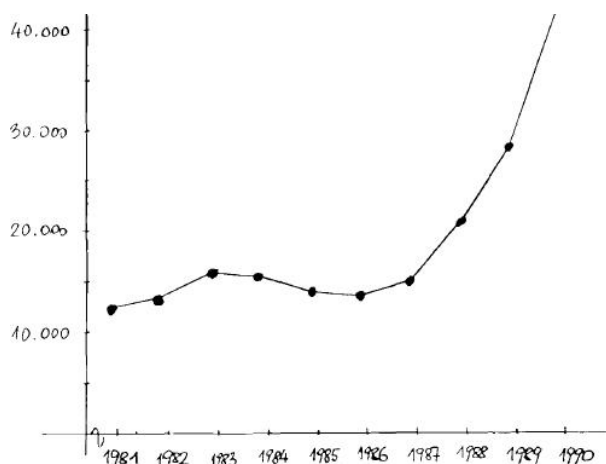
pojavov dostikrat vplivajo tudi **administrativni** ukrepi, ki z vsebino proučevanja nimajo neposredne zveze.

En izmed običajnih vzrokov so upravnoteritorialne spremembe, s katerimi se spremeni geografska opredelitev pojava, ki onemogoča primerljivost podatkov v časovni vrsti. V tem primeru je potrebno podatke časovne vrste za nazaj preračunati za novo območje.

17.2 Grafični prikaz časovne vrste

Kompleksen vpogled v dinamiko pojavov dobimo z grafičnim prikazom časovnih vrst v koordinatnem sistemu, kjer nanašamo na abscisno os čas in na ordinatno vrednosti dane spremenljivke. V isti koordinatni sistem smemo vnašati in primerjati le istovrstne časovne vrste.

Primer: Grafično prikažimo število brezposelnih v Sloveniji v letih od 1981 do 1990.



◇

17.3 Indeksi

Denimo, da je časovna vrsta dana z vrednostmi neke spremenljivke v časovnih točkah takole:

$$X_1, X_2, \dots, X_n$$

o indeksih govorimo, kadar z relativnimi števili primerjamo istovrstne podatke. Glede na to, kako določimo osnovo, s katero primerjamo člene v časovni vrsti, ločimo dve vrsti indeksov:

- **Indeksi s stalno osnovo.** Člene časovnih vrst primerjamo z nekim stalnim členom v časovni vrsti, ki ga imenujemo osnova X_0

$$I_{k/0} = \frac{X_k}{X_0} \cdot 100.$$

- **Verižni indeksi.** Za dano časovno vrsto računamo vrsto verižnih indeksov tako, da za vsak člen vzamemo za osnovo predhodni člen

$$I_k = \frac{X_k}{X_{k-1}} \cdot 100.$$

člene časovne vrste lahko primerjamo tudi z absolutno in relativno razliko med členi:

- **Absolutna razlika**

$$D_k = X_k - X_{k-1}.$$

- **Stopnja rasti** (relativna razlika med členi)

$$T_k = \frac{X_k - X_{k-1}}{X_{k-1}} \cdot 100 = I_k - 100.$$

Interpretacija indeksov

indeks	pojav		
	raste	stagnira	pada
s stalno osnovo	$I_{k+1/0} > I_{k/0}$	$I_{k+1/0} = I_{k/0}$	$I_{k+1/0} < I_{k/0}$
verižni indeks	$I_k > 100$	$I_k = 100$	$I_k < 100$
stopnja rasti	$T_k > 0$	$T_k = 0$	$T_k < 0$

Primer: Izračunajmo omenjene indekse za primer brezposelnih v Sloveniji:

leto	X_k	$I_{k/0}$	I_k	T_k
1981	12.315	100	—	—
1982	13.700	111	111	11
1983	15.781	128	115	15
1984	15.300	124	97	-3
1985	11.657	119	96	-4
1986	14.102	115	97	-3
1987	15.184	124	107	7
1988	21.311	173	141	41
1989	28.218	229	132	32
1990	44.227	359	157	57

Rezultati kažejo, da je bila brezposenost v letu 1990 kar 3,5 krat večja kot v letu 1981 (glej indeks s stalno osnovo). Iz leta 1989 na leto 1990 je bil prirast nezposlenih 57% (glej stopnjo rasti).

17.4 Sestavine dinamike v časovnih vrstah

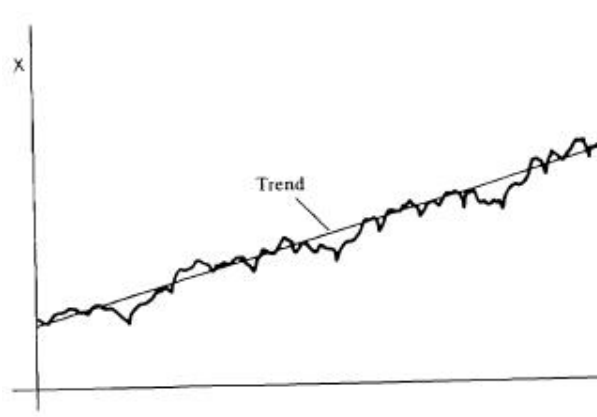
Posamezne vrednosti časovnih vrst so rezultat številnih dejavnikov, ki na pojav vplivajo. Iz časovne vrste je moč razbrati skupen učinek dejavnikov, ki imajo širok vpliv na pojav, ki ga proučujemo. Na časovni vrsti opazujemo naslednje vrste sprememb:

1. Dolgoročno gibanje ali trend - X_T podaja dolgoročno smer razvoja. Običajno ga je mogoče izraziti s preprostimi rahlo ukrivljenimi krivuljami.
2. Ciklična gibanja - X_C , so oscilarijo okoli trenda. Poriode so ponavdi daljše od enega leta in so lahko različno dolge.
3. Sezonske oscilacije - X_S so posledice vzrokov, ki se pojavljajo na stalno razdobje. Poriode so krajše od enega leta, ponavadi sezonskega značaja.
4. Naključne spremembe - X_E so spremembe, ki jih ne moremo razložiti s sistematičnimi gibanji (1, 2 in 3).

Časovna vrsta ne vsebuje nujno vseh sestavin. Zvezo med sestavinami je mogoče prikazati z nekaj osnovnim modeli. Npr.:

$$X = X_T + X_C + X_S + X_E \text{ ali } X = X_T \cdot X_C \cdot X_S \cdot X_E; \text{ ali } X = X_T \cdot X_C \cdot X_S + X_E.$$

Primer časovne vrste z vsemi štirimi sestavinami:



Ali je v časovni vrsti trend?

Obstaja statistični test, s katerim preverjamo ali trend obstaja v časovni vrsti. Med časom in spremenljivko izračunamo koeficient korelacije rangov

$$r_s = 1 - \frac{6 \cdot \sum_{i=1}^n d_i^2}{n \cdot (n^2 - 1)},$$

kjer je d_i , razlika med rangoma i tega časa in pripadajoče vrednosti spremenljivke. Ničelna in osnovna domneva sta:

$H_0: \rho_e = 0$ trend ne obstaja

$H_1: \rho_e \neq 0$ trend obstaja

Ustrezna statistika je

$$t = \frac{r_s \sqrt{n-2}}{\sqrt{1-r_s^2}},$$

ki se porazdeljuje približno po t porazdelitvi z $(n-2)$ prostostnimi stopnjami.

Metode določanja trenda

- Prostorčno
- Metoda drsečih sredin
- Metoda najmanjših kvadratov
- Druge analitične metode

Drseče sredine

Metoda drsečih sredin lahko pomaga pri določitvi ustreznega tipa krivulje trenda. V tem primeru namesto člena časovne vrste zapišemo povprečje določenega števila sosednjih članov. Če se odločimo za povprečje treh členov, govorimo o tričlenski vrsti drsečih sredin. Tedaj namesto članov v osnovni časovni vrsti X_k : tvorimo tričlenske drseče sredine X :

$$X'_k = \frac{X_{k-1} + X_k + X_{k+1}}{3}.$$

V tem primeru prvega in zadnjega člena časovne vrste moramo izračunati.

- Včasih se uporablja obtežena aritmetična sredina, včasih celo geometrijska za izračun drsečih sredin.

- Če so v časovni vrsti le naključni vplivi, dobimo po uporabi drsečih sredin ciklična gibanja (učinek Slutskega).
- Če so v časovni vrsti stalne periode, lahko drseče sredine zabrišejo oscilacije v celoti.
- V splošnem so drseče sredine lahko dober približek pravemu trendu.

Primer: Kot primer drsečih sredin vzemimo zopet brezposelne v Sloveniji. Izračunajmo tričlensko drsečo sredino:

T	X_k	tričl. drs. sred.
1981	12.315	–
1982	13.700	13.032
1983	15.781	14.030
1984	15.240	15.249
1985	15.300	14.710
1986	14.657	14.678
1987	14.102	15.184
1988	21.341	21.581
1989	28.218	31.262
1990	44.227	–

◇

Analitično določanje trenda

Trend lahko obravnavamo kot posebni primer regresijske funkcije, kjer je neodvisna spremenljivka čas (T). Če je trend

$$X_T = f(T),$$

lahko parametre trenda določimo z metoda najmanjših kvadratov

$$\sum_{i=1}^n (X_i - X_{iT})^2 = \min.$$

V primeru linearnega trenda

$$X_T = a + bT,$$

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

dobimo naslednjo oceno trenda

$$X_T = \bar{X} + \frac{\sum_{i=1}^n (X_i - \bar{X})(T_i - \bar{T})}{\sum_{i=1}^n (T_i - \bar{T})^2} (T - \bar{T}).$$

Ponavadi je čas T transformiran tako, da je $t = 0$. Tedaj je ocena trenda

$$X_T = \bar{X} + \frac{\sum_{i=1}^n (X_i - \bar{X}) \cdot t_i}{\sum_{i=1}^n t_i^2} t.$$

Standardna napaka ocene, ki meri razpršenost točk okoli trenda, je

$$\sigma_e = \sqrt{\frac{1}{n} \sum_{i=1}^n (X_i - X_{iT})^2}.$$

Primer: Kot primer ocenimo število doktoratov znanosti v Sloveniji v razdobju od leta 1986 do 1990. Z linearnim trendom ocenimo koliko doktorjev znanosti je v letu 1991. Izračunajmo tudi standardno napako ocene. Izračunajmo najprej trend:

T	Y_i	t_i	$Y_i - \bar{Y}$	$(Y_i - \bar{Y})t_i$	t_i^2
1986	89	-2	-19,8	39,6	4
1987	100	-1	-8,8	8,8	1
1988	118	0	9,2	0	0
1989	116	1	7,2	7,2	1
1990	121	2	12,2	24,4	4
	544	0		80	10

$$\bar{Y} = \frac{544}{4} = 108,8,$$

$$Y_T = 108,8 + \frac{80}{10} t = 108,8 + 8t,$$

$$Y_T(1991) = 108,8 + 8 \cdot 3 = 132,8.$$

Ocena za leto 1991 je približno 133 doktorjev znanosti. Zdaj pa izračunajmo standardno napako ocene. Za vsako leto je potrebno najprej izračunati

T	Y_i	Y_{iT}	$Y_i - Y_{iT}$	$(Y_i - Y_{iT})^2$
1986	89	92,8	-3,8	14,44
1987	100	100,8	-0,8	0,64
1988	118	108,8	9,2	84,64
1989	116	116,8	-0,8	0,64
1990	121	124,8	-3,8	14,44
	544	544	0	114,8

$$\sigma_e = \sqrt{\frac{114,8}{5}} = 4,8.$$

◇

Poglavje 18

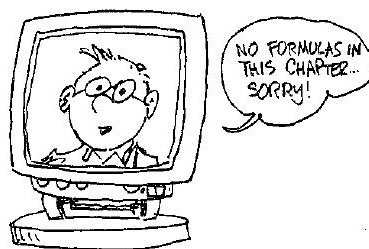
Uporaba

(Testiranje PRNG, Teorija informacij in entropija)

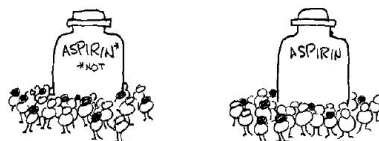
18.1 Načrtovanje eksperimentov



Načrtovanje eksperimentov se pogosto neposredno prevede v uspeh oziroma neuspeh. V primeru parjenja lahko statistik spremeni svojo vlogo iz pasivne v aktivno. Predstavimo samo osnovne ideje, podrobno numerično analizo pa prepustimo statistični programski opremleni.



Elementi načrta so eksperimentalne enote ter terapije, ki jih želimo uporabiti na enotah.



- medicina: bolniki (enote) in zdravila (terapije),
- optimizacija porabe: taxi-ji (enote) in različne vrste goriva (terapije),
- agronomija: območja na polju in različne vrste kulture, gnojiva, špricanja,...

Danes uporabljamo ideje načrtovanja eksperimentov na številnih področjih:

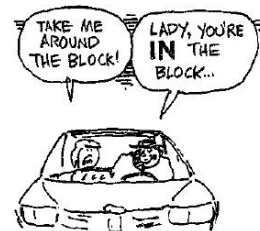
- optimizacija industrijskih procesov,
- medicina,
- sociologija.



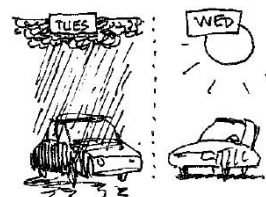
Na primeru bomo predstavili tri osnovne principe načrtovanja eksperimentov:

1. **Ponavljanje**: enake terapije pridružimo različnim enotam, saj ni mogoče oceniti naravno spremenljivost (ang. natural variability) in napake pri merjenju.
2. **Lokalna kontrola** pomeni vsako metodo, ki zmanjša naravno spremenljivost.

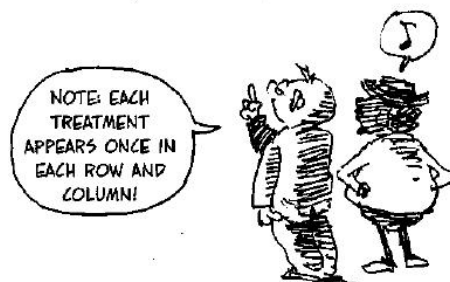
En od načinov grupira podobne enote eksperimentov v **bloke**. V primeru taxijev uporabimo obe vrsti goriva na vsakem avtomobilu in rečemo, da je avto blok.



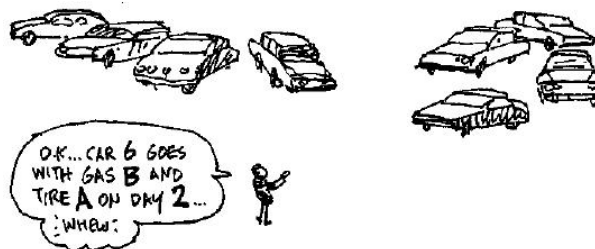
3. **Naključna izbira** je bistven korak povsod v statistiki! Terapije za enote izbiramo naključno. Za vsak taksi izberemo vrsto goriva za torek oziroma sredo z metom kovanca. Če tega ne bi storili, bi lahko razlika med torkom in sredo vplivala na rezultate.



		DAY			
		1	2	3	4
CAB	1	a	b	c	d
	2	b	c	d	a
	3	c	d	a	b
	4	d	a	b	c



Latinski kvadrati



Latinski kvadrat reda v je $v \times v$ -razsežna matrika, v kateri vsi simboli iz množice

$$\{1, \dots, v\}$$

nastopajo v vsaki vrstici in vsakem stolpcu.

Trije paroma ortogonalni latinski kvadrati reda 4, tj. vsak par znak-črka ali črka-barva ali barva-znak se pojavi natanko enkrat.

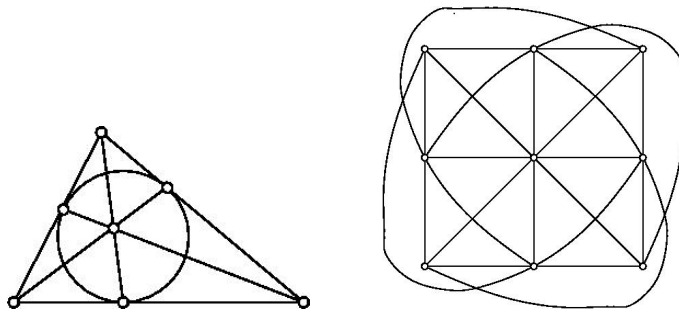
Projektivni prostor $PG(d, q)$ (razsežnosti d nad q) dobimo iz vektorskega prostora $[GF(q)]^{d+1}$, tako da naredimo kvocient po 1-razsežnih podprostorih.

Projektivna ravnina $PG(2, q)$ je incidenčna struktura z 1- in 2-dim. podprostori prostora $[GF(q)]^3$ kot **točkami** in **premicami**, kjer je “ \subset ” incidenčna relacija. To je $2-(q^2 + q + 1, q + 1, 1)$ -design, tj.,

- $v = q^2 + q + 1$ je število točk (in število premic b),
- vsaka premica ima $k = q + 1$ točk (in skozi vsako točko gre $r = q + 1$ premic),
- vsak par točk leži na $\lambda = 1$ primicah (in vsaki premici se sekata v natanko eno točki).

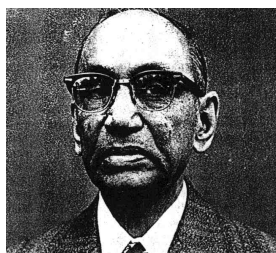
Primeri:

1. Projektivno ravnino $PG(2, 2)$ imenujemo **Fano ravnina** (7 točk in 7 premic).



2. $PG(2, 3)$ lahko skonstruiramo iz 3×3 mreže oziroma afine ravnine $AG(2, 3)$.

Bose in Shrikhande



Prva sta konec tridesetih let prejšnjega stoletja vpeljala **asociativne sheme Bose** in **Nair** a potrebe statistike.

Toda **Delsarte** je pokazal, da nam lahko služijo kot povezava med številnimi področji matematike, naprimer teorijo kodiranja in teorijo načrtov.

Del III
KAM NAPREJ

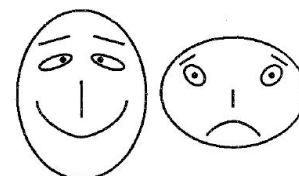
Osnovni principi in orodja, ki smo jih spoznali pri VIS, lahko posplošimo in razširimo do te mere, da se dajo z njimi rešiti tudi bolj kompleksni problemi.



Spoznali smo kako predstaviti **eno** spremenljivko (dot-plot, histogrami,...) in **dve** spremenljivki (razsevni diagram).

Kako pa predstavimo več kot dve spremenljivki na ravnem listu papirja?

Med številnimi možnostmi moramo omeniti idejo **Hermana Chernoffa**, ki je uporabil človeški obraz, pri čemer je vsako lastnost povezal z eno spremenljivko. Oglejmo si Chernoffov obraz: X =naklon obrvi, Y =velikost oči, Z =dolžina nosu, T =dolžina ust, U =višino obraza, itd.



Multivariantna analiza

Širok izbor multivariantnih modelov nam omogoča analizo in ponazoritev n -razsežnih podatkov.

Združevalna/grozdna tehnika (ang. cluster technique):

Iskanje delitve populacije na homogene podskupine, npr. z analizo vzorcev senatorskih glasovanj v ZDA zaključimo, da *jug* in *zahod* tvorita dva različna grozda.



Diskriminacijska analiza

je obraten proces. Npr. odbor/komisija za sprejem novih študentov bi rad našel podatke, ki bi že vnaprej opozorili ali bodo prijavljeni kandidati nekega dne uspešno zaključili program (in finančno pomagali šoli - npr. z dobrodelnimi prispevki) ali pa ne bodo uspešni (gre delati dobro po svetu in šola nikoli več ne sliši zanj(o)).



Analiza faktorjev

išče poenostavljeno razlago večrazsežnih podatkov z manjšo skupino spremenljivk. Npr. Psihiater lahko postavi 100 vprašanj, skrivoma pa pričakuje, da so odgovori odvisni samo od nekaterih faktorjev:

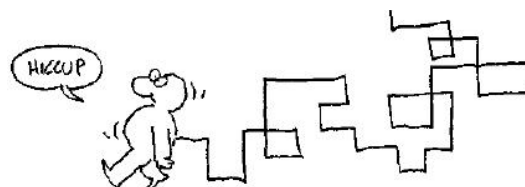
ekstravertiranost, avtoritativnost, alutarizem, itd. Rezultate testa lahko potem povzamemo le z nekaterimi sestavljenimi rezultati v ustreznih dimenzijah.



Naključni sprehodi

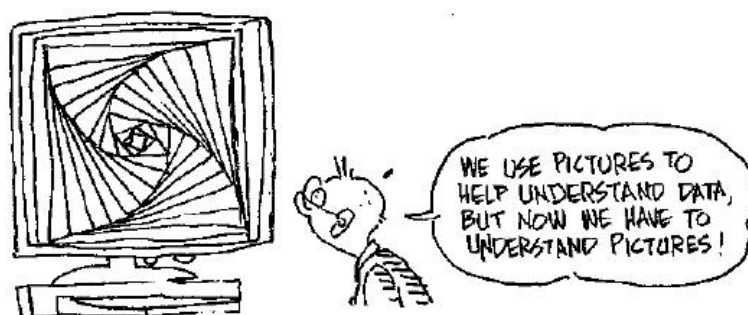
pričnejo z metom kovanca, recimo, da se pomaknemo korak nazaj, če pade grb, in korak naprej, če pade cifra. (z dvema kovancema se lahko gibljemo v 2-razsežnem prostoru - tj. ravnini). Če postopek ponavljamo, pridemo do *stohastičnega procesa*, ki ga imenujemo naključni sprehod (ang. random walk).

Modeli na osnovi naključnih sprehodov se uporabljajo za nakup/prodajo delnic in portfolio management.



Vizualizacija in analiza slik

Sliko lahko sestavlja 1000×1000 pikselov, ki so predstavljeni z eno izmed 16,7 milijonov barv. Statistična analiza slik želi najti nek pomen iz "informacije" kot je ta.



Ponovno vzorčenje

Pogosto ne moremo izračunati standardne napake in limite zaupanja. Takrat uporabimo tehniko ponovnega vzorčenja, ki tretira vzorec, kot bi bila celotna populacija. Za takšne tehnike uporabljamo pod imeni: randomization Jackknife, in Bootstrapping.



Kvaliteta podatkov

Navidezno majhne napake pri vzorčenju, merjenju, zapisovanju podatkov, lahko povzročijo katastrofalne učinke na vsako analizo. R. A. Fisher, genetik in ustanovitelj moderne statistike ni samo načrtoval in analiziral eksperimentalno rejo, pač pa je tudi čistil kletke in pazil na živali. Zavedal se je namreč, da bi izguba živali vplivala na rezultat.



Moderni statistiki, z njihovimi računalniki in podatkovnimi bazami ter vladnimi projekti (beri denarjem) si pogosto ne umažejo rok.

Inovacija

Najboljše rešitve niso vedno v knjigah (no vsaj najti jih ni kar tako). Npr. mestni odpad je najel strokovnjake, da ocenijo kaj sestavljajo odpadki, le-ti pa so se znašli pred zanimivimi problemi, ki se jih ni dalo najti v standardnih učbenikih.



Komunikacija

Še tako uspešna in bistroumna analiza je zelo malo vredna, če je ne znamo jasno predstaviti, vključujoč stopnjo statistične značilnosti? v zaključku.



Npr. V medijih danes veliko bolj natančno poročajo o velikosti napake pri svojih anketah.

Timsko delo

V današnji kompleksni družbi. Reševanje številnih problemov zahteva *timsko delo*. Inženirji, statistiki in delavci sodelujejo, da bi izboljšali kvaliteto produktov. Biostatistiki, zdravniki, in AIDS-aktivisti združeno sestavljajo klinične poiskuse, ki bolj učinkovito ocenijo terapije.



Literatura

- [1] A. Ferligoj: *Osnove statistike na prosojnicah*. Samozaložba, Ljubljana 1995.
- [2] L. Gonick in W. Smith, *The Cartoon guide to Statistics*, 1993.
- [3] M. Hladnik: *Verjetnost in statistika*. Založba FE in FRI, Ljubljana 2002.
- [4] W. Mendenhall in T. Sincich, *Statistics for engineering and the sciences*, 4th edition, Prentice Hall, 1995.
- [5] D. S. Moore (Purdue University), *Statistika: znanost o podatkih* (5. izdaja prevedena v slovenščino leta 2007).

Obstaja obilna literatura na spletu in v knjižnicah.

Gradiva bodo dosegljiva preko internetne učilnice (moodle).

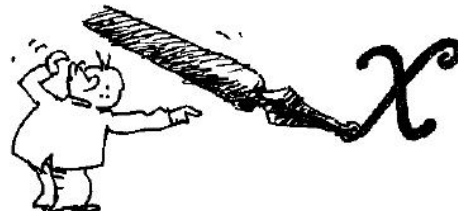
Pri delu z dejanskimi podatki se bomo v glavnem naslonili na prosti statistični program R.

Program je prosto dostopen na:

<http://www.r-project.org/>

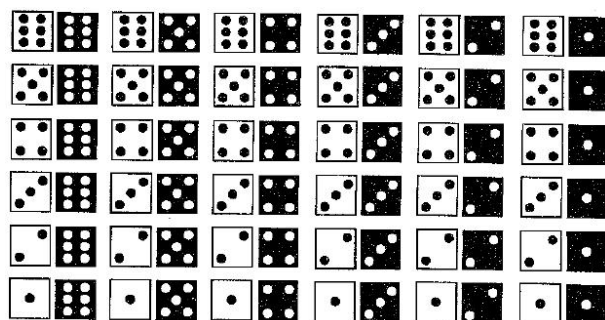
Proti koncu semestra pa morda tudi Minitab.

OUR HUMBLE OPINION IS THAT LEARNING A LITTLE MORE ABOUT THE SUBJECT MIGHT NOT BE SUCH A BAD IDEA... AND THAT'S WHY WE WROTE THIS BOOK!



Dodatek A

MATEMATIČNE OSNOVE (ponovitev)



A.1 Računala nove dobe

Ste že kdaj razmišljali o računanju (aritmetiki), ki ga uporabljamo v vsakdanjem življenju? Večina ljudi jo zamenjuje kar za celotno matematiko. Na kakšen način računajo računalniki ter ostale digitalne naprave (digit je angl. beseda za število), ki nas obkrožajo v času informacijske dobe? Nekateri se sicer skušajo prilagajati našemu načinu računanja, vse več pa je takih, ki so jim časovna in prostorska učinkovitost ter preciznost ključnega pomena. Take naprave računajo na malce drugačen način. V tem razdelku se bomo poskusili s pomočjo osnovnošolskega računanja približati računalom, ki jih preko številnih naprav, kot so osebni računalniki, diskmani in pametne kartice, uporabljamo v vsakdanji praksi.

Poleg seštevanja in množenja pa se v prvih razredih osnovne šole naučimo tudi odštevati in deliti. Seveda začnemo najprej odštevati manjša števila od večjih. Če želimo izračunati $a - b$ in je $a \geq b$, se lahko vprašamo

b plus koliko je a ?

Šele nekoliko kasneje se naučimo, da moramo v primeru, ko želimo odšteti večje število od manjšega, števili najprej zamenjati, na koncu pa dobljeni razliki spremeniti predznak. Zaradi tega smo povečali množico naravnih števil $\mathbb{N} = \{1, 2, \dots\}$ do množice celih števil $\mathbb{Z} = \{\dots, -2, -1, 0, 1, 2, \dots\}$.

Deljenje ni tako preprosto. Če želimo a deliti z b , se lahko prav tako kot prej vprašamo “ b krat koliko je a ?” Vendar se pogosto zgodi, da število a sploh ni deljivo s številom b . Množico števil lahko sicer povečamo do množice ulomkov \mathbb{Q} , kjer se da deliti s poljubnim od nič različnim številom, a potem nastopijo druge težave. Najdemo lahko različne ulomke, ki so si poljubno blizu, tudi tako blizu, da jih računalnik ne more več ločiti. Ker pa si želimo, da bi se računalniki čim manj motili, se vprašajmo po množicah, v katerih bi lahko brez problemov tudi delili, po možnosti na enak način kot znamo odšteti. Pravzaprav se je potrebno vprašati, na katera pravila se želimo pri računanju opreti. Naštejmo jih nekaj.

1. Običajno je prvo pravilo *zaprtost*, rezultat, ki ga dobimo po opravljeni operaciji med dvema številoma, je tudi v množici, iz katere smo izbrali števili. Množica naravnih števil je zaprta za seštevanje in množenje, saj v tabelah 1a in 1b nastopajo samo naravna števila. Ni pa množica naravnih števil zaprta za odštevanje. To lastnost ima na primer množica celih števil.

2. V množici celih števil igra pomembno vlogo število 0; pa ne samo zato, ker loči pozitivna števila od negativnih, pač pa tudi zato, ker se nobeno število s prištevanjem števila 0, ne spremeni. Tudi pri množenju najdemo nekaj podobnega. Če pomnožimo katerokoli od nič različno število z 1, dobimo zopet isto število. Takemu številu pravimo *neutralni element* ali pa tudi *enota* za ustrezno operacijo.

3. V množici celih števil sta poljubni števili $-a$ in a povezani z enoto za seštevanje na naslednji način: $a + (-a) = 0$. Pravimo, da je $-a$ *nasprotni* element števila a . Celo število b je *obratni element* celega števila a , če je $ab = 1$. Od tod sledi $a = b = 1$, tj. v množici celih števil imata le števili 1 in -1 obratni element.

4. Če si izberemo poljubna števila a , b in c , potem velja $a + (b + c) = (a + b) + c$ in $a(bc) = (ab)c$. O drugi enakosti se lahko prepričamo z računanjem prostornine kvadra s stranicami a , b in c . Tem lastnostim pravimo *zakon o združevanju* za seštevanje oziroma za množenje (ali tudi *asociativnost*). Le-ta nam pove, da je vseeno, ali začnemo računati z leve ali z desne. To seveda ne drži za odštevanje ali deljenje.

Če v neki množici G z binarno (dvočleno) operacijo \circ , tj. operacijo, ki vsakemu urejenemu paru elementov iz G priredi natanko določen element, veljajo naslednja pravila:

(G1) za vsaka $a, b \in G$ je $a \circ b \in G$,

(G2) obstaja tak element $e \in G$, da za vsak $g \in G$ velja $e \circ g = g \circ e = g$,

(G3) za vsak element $g \in G$ obstaja tak $f \in G$, da je $g \circ f = f \circ g = e$,

(G4) za vse $a, b, c \in G$ velja $(a \circ b) \circ c = a \circ (b \circ c)$,

potem pravimo, da je par (G, \circ) **grupa**. Elementu e pravimo **enota** grupe, elementu f pa **inverz** elementa g . Množica celih števil je grupa za seštevanje, ni pa grupa za množenje, saj ni izpolnjeno pravilo (G3) (le 1 in -1 imata inverzni element za množenje).

Morda bo kdo pomislil, da je prišla definicija grupe iz glave enega samega matematika, pa temu sploh ni tako. Matematiki so potrebovali več kot 100 let trdega dela, da so končno (eksplicitno) zapisali zgornja pravila (*aksiome*). *Joseph Louis Lagrange* (1736-1813) je leta 1771 postavil prvi pomembnejši izrek. *Augustin Louis Cauchy* (1789-1857) je študiral grupe permutacij, medtem, ko je *Niels Henrik Abel* (1802-1829) s teorijo grup pokazal, da enačba 5. stopnje ni rešljiva z radikali (tj. rešitve ne znamo zapisati s formulami kot v primeru enačb nižjih stopenj). Po njem pravimo grupam, v katerih velja pravilo zamenjave, tudi *Abelove grupe* (ali komutativne grupe). Pravi pionir abstraktnega pristopa pa je bil *Evariste Galois* (1811-1832), ki je leta 1823 prvi uporabil besedo "grupa". Proces poudarka na strukturi se je nadaljeval vse do leta 1854, ko je *Arthur Cayley* (1821-1895) pokazal, da je grupo moč definirati ne glede na konkretno naravo njenih elementov.

Galois je vpeljal tudi naslednji pojem. Če za neko množico \mathcal{O} z binarnima operacijama, ki ju bomo označili s $+$ in $*$ (četudi ne predstavljata nujno običajnega seštevanja in množenja), velja

(O1) par $(\mathcal{O}, +)$ je grupa z enoto 0,

(O2) par $(\mathcal{O} \setminus \{0\}, *)$ je grupa z enoto 1,

(O3) za vse $a, b, c \in \mathcal{O}$ je $a * (b + c) = a * b + b * c$ in $(b + c) * a = b * a + c * a$,

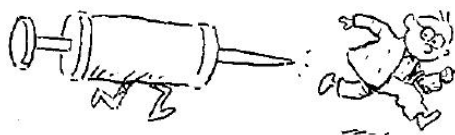
potem imenujemo trojico $(\mathcal{O}, +, *)$ **obseg**. Množica ulomkov z običajnim seštevanjem in množenjem je primer obsega. O lastnosti (O3), ki jo imenujemo *zakon o razčlenjevanju* oziroma *distributivnost*, se lahko prepričamo z računanjem površine pravokotnika s stranicama a in $b + c$.

Primer: Za cilj postavimo iskanje obsega s končno mnogo elementi, v katerem bo računanje v nekem smislu še udobnejše kot v obsegih, ki jih srečamo v osnovni ali srednji šoli (racionalna števila \mathbb{Q} , realna števila \mathbb{R} ali celo kompleksna števila \mathbb{C}).

Gotovo ste hitro ugotovili, da mora imeti grupa zaradi aksioma (G2) vsaj en element, enoto e namreč, obseg pa vsaj dva, enoto za operacijo “+” in enoto za operacijo “*”. Potem se ni več težko prepričati, da je en element v primeru grupe že dovolj, saj nam $e \circ e = e$ zadovolji vse aksiome (G1)-(G4). V primeru obsega z dvema elementoma je enako z multiplikativno grupo: $1 * 1 = 1$ zadovolji aksiom (O2). Ne pozabite, da operaciji “+” in “*” ne predstavljata (nujno) običajnega seštevanja in množenja. V tem sestavku bomo spoznali kar nekaj takih obsegov. V vsakem obsegu je produkt poljubnega elementa a z aditivno enoto 0 enak 0 , saj je $0 * a = (0 + 0) * a = 0 * a + 0 * a$ (upoštevali smo (G2) in (O3)) in po krajšanju z $0 * a$ res dobimo $0 * a = 0$. Opozoriti je treba, da *pravilo krajšanja* velja v poljubni grupi, saj v resnici na obeh straneh “dodamo” inverzni element in nato upoštevamo zakon o združevanju (G4) ter (G2) (do tega ste se gotovo dokopali že sami pri reševanju 3. naloge iz prejšnjega razdelka). Seveda velja enako tudi, kadar vrstni red zamenjamo: $a * 0 = 0$. Torej tudi v primeru najmanjšega obsega velja $0 * 0 = 0$ in $0 * 1 = 0 = 1 * 0$, kjer je 1 multiplikativna enota. Kako pa je z grupo, ki ima dva elementa, npr. enoto e in a ? Poleg $e \circ e = e$ in $e \circ a = a = a \circ e$ mora zaradi aksioma (G3), pravila krajšanja in $e \neq a$ veljati še $a \circ a = e$ in že so izpolnjeni vsi aksiomi (G1)-(G4). Torej velja za obseg z dvema elementoma in pravkar odkrito aditivno grupo tudi aksiom (O1). Zlahka preverimo še (O3) in že smo ugnali tudi najmanjši obseg. \diamond

A.2 Funkcije/preslikave

Funkcija f iz množice A v množico B je predpis, ki vsakemu elementu iz množice A priredi natanko določen element iz množice B , oznaka $f : A \rightarrow B$.



Funkcija $f : A \rightarrow B$ je:

- **injektivna** (angl. one to one) če za $\forall x, y \in A \quad x \neq y \Rightarrow f(x) \neq f(y)$,
- **surjektivna** (angl. on to), če za $\forall b \in B \quad \exists a \in A$, tako da je $f(a) = b$.

Injektivni in surjektivni funkciji pravimo **bijekcija**. Množicama med katerima obstaja bijekcija pravimo **bijektivni** množici. Bijektivni množici imata enako število elementov (npr. končno, števno neskončno, itd).

Trditev A.1. Če sta množici A in B končni ter je $f : A \rightarrow B$ funkcija, iz injektivnosti funkcije f sledi surjektivnost, in obratno, iz surjektivnosti funkcije f sledi injektivnost. \square

A.3 Permutacije

Permutacija elementov $1, \dots, n$ je bijekcija, ki slika iz množice $\{1, \dots, n\}$ v množico $\{1, \dots, n\}$. Npr. permutacija kart je običajno premešanje kart (spremeni se vrstni red, karte pa ostanejo iste). Število permutacij n elementov, tj. razvrstitev n -tih različnih elementov, je enako $n! := 1 \cdot 2 \cdot \dots \cdot n$ (oziroma definirano rekurzivno $n! = (n-1)!n$ in $0! = 1$). Permutacijo lahko opišemo z zapisom:

$$\pi = \begin{pmatrix} 1 & 2 & \dots & n \\ a_1 & a_2 & \dots & a_n \end{pmatrix},$$

kjer je $\{1, 2, \dots, n\} = \{a_1, a_2, \dots, a_n\}$. To pomeni $\pi(1) = a_1, \pi(2) = a_2, \dots, \pi(n) = a_n$.

Primer: $n = 11$,

$$\pi_1 = \begin{pmatrix} 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 & 10 & 11 \\ 3 & 4 & 5 & 10 & 2 & 1 & 7 & 9 & 11 & 6 & 8 \end{pmatrix} \quad \diamond$$

Naj bo A neka množica. Permutacije množice A med seboj množimo po naslednjem pravilu: $\pi = \pi_1 \circ \pi_2$ je permutacija množice A , ki preslika $a \in A$ v $\pi_2(\pi_1(a))$.

Primer:

$$\begin{aligned} \pi_1 &= \begin{pmatrix} 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 & 10 & 11 \\ 3 & 4 & 5 & 10 & 2 & 1 & 7 & 9 & 11 & 6 & 8 \end{pmatrix} \\ \pi_2 &= \begin{pmatrix} 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 & 10 & 11 \\ 8 & 2 & 1 & 3 & 10 & 9 & 4 & 5 & 7 & 6 & 11 \end{pmatrix} \end{aligned}$$

Potem je

$$\pi = \pi_1 \circ \pi_2 = \begin{pmatrix} 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 & 10 & 11 \\ 1 & 3 & 10 & 6 & 2 & 8 & 4 & 7 & 11 & 9 & 5 \end{pmatrix} \quad \diamond$$

Cikel je permutacija, za katero je

$$\pi(a_1) = a_2, \pi(a_2) = a_3, \dots, \pi(a_r) = a_1,$$

ostale elementi pa so fiksni (tj. $\pi(a) = a$). Na kratko jo zapišemo z $(a_1 a_2 \dots a_r)$.

Trditev A.2. *Vsako permutacijo lahko zapišemo kot produkt disjunktnih ciklov.* □

Primer:

$$\begin{aligned} \pi_1 &= \begin{pmatrix} 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 & 10 & 11 \\ 3 & 4 & 5 & 10 & 2 & 1 & 7 & 9 & 11 & 6 & 8 \end{pmatrix} \\ \pi_2 &= \begin{pmatrix} 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 & 10 & 11 \\ 8 & 2 & 1 & 3 & 10 & 9 & 4 & 5 & 7 & 6 & 11 \end{pmatrix} \\ \pi &= \begin{pmatrix} 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 & 10 & 11 \\ 1 & 3 & 10 & 6 & 2 & 8 & 4 & 7 & 11 & 9 & 5 \end{pmatrix} \end{aligned}$$

Potem je

$$\pi_1 = (1\ 3\ 5\ 2\ 4\ 10\ 6)(8\ 9\ 11), \quad \pi_2 = (1\ 8\ 5\ 10\ 6\ 9\ 7\ 4\ 3), \quad \pi = (2\ 3\ 10\ 9\ 11\ 5\ 2)(4\ 6\ 8\ 7) \quad \diamond$$

Transpozicija je cikel dolžine 2. Vsak cikel pa je produkt transpozicij:

$$(a_1 a_2 a_3 \dots a_r) = (a_1 a_2) \circ (a_2 a_3) \circ \dots \circ (a_{r-1} a_r),$$

torej je tudi vsaka permutacija produkt transpozicij. Seveda ta produkt ni nujno enolično določen, vseeno pa velja:

Trditev A.3. *Nobena permutacija se ne da zapisati kot produkt sodega števila in kot produkt lihega števila permutacij.*

Dokaz. Naj bodo x_1, x_2, \dots, x_n različna realna števila. Poglejmo si produkt:

$$P = \prod_{i < j} (x_i - x_j).$$

Izberimo indeksa a in b , $a < b$, in pogledimo v katerih razlikah se pojavita:

$x_1 - x_a, \dots, x_{a-1} - x_a,$		$x_a - x_{a+1}, \dots, x_a - x_{b-1},$	$x_a - x_b,$	$x_a - x_{b+1}, \dots, x_a - x_n,$
$x_1 - x_b, \dots, x_{a-1} - x_b,$	$x_a - x_b,$	$x_{a+1} - x_b, \dots, x_{b-1} - x_b,$		$x_b - x_{b+1}, \dots, x_b - x_n.$



Razliko $x_a - x_b$ smo navedli dvakrat, a se v produktu P pojavi samo enkrat. Če na množici indeksov opravimo transpozicijo $(a b)$, razlika $x_a - x_b$ preide v razliko $x_b - x_a$, torej zamenja predznak, razlike iz prvega in zadnjega stolpca se med seboj zamenjajo, razlike iz srednjega stolpca pa tudi zamenjajo predznake (vendar je le-teh sodo mnogo in zato ne vplivajo na produkt P). Sedaj pa napravimo na množici indeksov permutacijo π . V tem primeru je produkt

$$P_\pi = \prod_{i < j} (x_{\pi(i)} - x_{\pi(j)}).$$

enak $\pm P$. Če uporabimo sodo število transpozicij, potem je $P_\pi = P$, sicer pa $P_\pi = -P$. \square

Glede na sodo oziroma liho število transpozicij imenujemo permutacijo **soda** oziroma **liha** permutacija.

Permutacije s ponavljanjem

Permutacije s ponavljanjem so nekakšne permutacije, pri katerih pa ne ločimo elementov v skupinah s k_1, k_2, \dots, k_r elementi, torej imamo $n = k_1 + k_2 + \dots + k_r$ elementov - zato delimo število vseh permutacij n elementov s številom njihovih vrstnih redov, tj. permutacij:

$$P_n^{k_1, k_2, \dots, k_r} = \frac{n!}{k_1! k_2! \dots k_r!}.$$

Primer: 8 vojakov je potrebno poslati na stražo v štiri postojanke. Recimo, da želimo vsak dan izbrati drugačno razporeditev. **Na koliko načinov lahko to storimo?**

Odgovor: $P_8^{2,2,2,2} = 8 \cdot 7 \cdot 6 \cdot 5 \cdot 4 \cdot 3 \cdot 2 / 2^4 = 7 \cdot 6 \cdot 5 \cdot 4 \cdot 3 = 42 \cdot 60 = 2520$. Torej je načinov vsekakor preveč, da bi vojaki odšli na stražo na vse možne načine, četudi bi služili vojaški rok celo leto in šli na stražo prav vsak dan po šestkrat. \diamond

Če je $r = 1$, je število permutacij s ponavljanjem enako 1, če je $r = n$, pa gre za čisto navadne permutacije. Če je $r = 2$, ločimo elemente v dve skupini. Če je $k = k_1$, je $n - k = k_2$ in pišemo

$$\binom{n}{k} := P_n^{k, n-k}.$$

Ta primer bomo obravnavali posebej v naslednjem razdelku.

Primer: Na koliko načinov lahko med sošolke Aleksandro, Evo in Rebeko razdelimo pet knjig, če dobi vsaka vsaj eno knjigo?

Naj bo m število sošolk in n število knjig, iskano število pa označimo s $S(n, m)$. Potem je očitno $S(n, 1) = 1$. Poglejmo si sedaj primer $m = 2$, tj. primer, ko knjige dobita le dve sošolki. Za vsako knjigo se odločimo, kateri sošolkii jo damo in hitro vidimo, da imamo 2^n različnih možnosti, vendar pa dve možnosti nista pravi (tisti pri katerih bi vse knjige dali prvi oziroma drugi sošolki), tj. iskano število je enako

$$S(n, 2) = 2^n - 2.$$

Morda bo sedaj kaj lažje ugnati primer $m = 3$. Za $n = 5$ bi morda lahko izpisali vsa petmestna števila v trojiškem sistemu (le teh je natanko $3^5 = 3^2 \cdot 3^2 \cdot 3 = 9 \cdot 9 \cdot 3 = 81 \cdot 3 = 243$), nato pa označili z * tista števila, ki imajo vse števke enake (teh je ravno 3), z x , y in z pa zaporedoma še tista preostala števila, ki ne vsebujejo nobene dvojke, enico oziroma ničlo, vendar iz prejšnjega primera že vemo, da je število oznak x (oziroma y oziroma z) je ravno $S(n, 2) = 2^5 - 2$. Torej je iskano število enako:

$$S(n, 3) = 3^n - \binom{3}{2}(2^n - 2) - \binom{3}{1} = 3^n - 3(2^n - 2) - 3.$$

Za $n = 5$ pa dobimo $S(5, 3) = 3(3^4 - 2^5 + 2 - 1) = 3(81 - 32 + 2 - 1) = 3 \cdot 50 = 150$.

Preverimo dobljeno formulo za $S(n, 3)$ še s formulo za permutacije s ponavljanjem:

n	$S(n, 3)$		
3	$3^3 - 3(2^3 - 2) - 3 = 6$	$= 3!$	$= P_3^{111}$
4	$3^4 - 3(2^4 - 2) - 3 = 36$	$= 3 \cdot 4! / 2!$	$= 3 \cdot P_4^{112}$
5	$3^5 - 3(2^5 - 2) - 3 = 150$	$= 3 \cdot \left(\frac{5!}{3!} + \frac{5!}{2!2!} \right)$	$= 3 \cdot (P_5^{113} + P_5^{122})$
6	$3^6 - 3(2^6 - 2) - 3 = 540$	$= \frac{3 \cdot 6!}{4!} + \frac{6 \cdot 6!}{2!3!} + \frac{6!}{2!2!2!}$	$= 3 \cdot P_6^{114} + 6 \cdot P_6^{123} + 1 \cdot P_6^{222}$
7	$3^7 - 3(2^7 - 2) - 3 = 1806$	$= \frac{3 \cdot 7!}{5!} + \frac{6 \cdot 7!}{2!4!} + \frac{3 \cdot 7!}{3!3!} + \frac{3 \cdot 7!}{2!2!3!}$	$= 3P_7^{115} + 6P_7^{124} + 3P_7^{133} + 3P_7^{223}$

V primeru $n = 3$ je bil izračun otročje lahek (pa tudi za zelo majhno število gre), že v naslednjem primeru pa smo morali upoštevati tri možnosti, tj. katera od deklet dobi dve knjigi (tudi to število v resnici šteje permutacije s ponavljanjem: $P_3^{21} = 3$). Primer $n = 5$ je zelo podoben prejšnjemu primeru: če dobi eno dekle tri knjige, ostali dve morata dobiti po eno, če pa damo eni dve knjigi, bo dobila še ena dve, za tretjo pa ostane ena sama knjiga. Za $n = 6$ in $n = 7$ pa omenimo še $P_3^{111} = 6$ in $P_3^3 = 1$.

0000*	0100 <i>x</i>	0200 <i>y</i>	1000 <i>x</i>	1100 <i>x</i>	1200	2000 <i>y</i>	2100	2200 <i>y</i>
00001 <i>x</i>	01001 <i>x</i>	02001	10001 <i>x</i>	11001 <i>x</i>	12001	20001	21001	22001
00002 <i>y</i>	01002	02002 <i>y</i>	10002	11002	12002	20002 <i>y</i>	21002	22002 <i>y</i>
00010 <i>x</i>	01010 <i>x</i>	02010	10010 <i>x</i>	11010 <i>x</i>	12010	20010	21010	22010
00011 <i>x</i>	01011 <i>x</i>	02011	10011 <i>x</i>	11011 <i>x</i>	12011	20011	21011	22011
00012 <i>y</i>	01012	02012	10012	11012	12012	20012	21012	22012
00020	01020	02020 <i>y</i>	10020	11020	12020	20020 <i>y</i>	21020	22020 <i>y</i>
00021	01021	02021	10021	11021	12021	20021	21021	22021
00022 <i>y</i>	01022	02022 <i>y</i>	10022	11022	12022	20022 <i>y</i>	21022	22022 <i>y</i>
00100 <i>x</i>	01100 <i>x</i>	02100	10100 <i>x</i>	11100 <i>x</i>	12100	20100	21100	22100
00101 <i>x</i>	01101 <i>x</i>	02101	10101 <i>x</i>	11101 <i>x</i>	12101	20101	21101	22101
00102	01102	02102	10102	11102	12102	20102	21102	22102
00110 <i>x</i>	01110 <i>x</i>	02110	10110 <i>x</i>	11110 <i>x</i>	12110	20110	21110	22110
00111 <i>x</i>	01111 <i>x</i>	02111	10111 <i>x</i>	11111*	12111 <i>z</i>	20111	21111 <i>z</i>	22111 <i>z</i>
00112	01112	02112	10112	11112 <i>z</i>	12112 <i>z</i>	20112	21112 <i>z</i>	22112 <i>z</i>
00120	01120	02120	10120	11120	12120	20120	21120	22120
00121	01121	02121	10121	11121 <i>z</i>	12121 <i>z</i>	20121	21121 <i>z</i>	22121 <i>z</i>
00122	01122	02122	10122	11122 <i>z</i>	12122 <i>z</i>	20122	21122 <i>z</i>	22122 <i>z</i>
00200 <i>y</i>	01200	02200 <i>y</i>	10200	11200	12200	20200 <i>y</i>	21200	22200 <i>y</i>
00201	01201	02201	10201	11201	12201	20201	21201	22201
00202 <i>y</i>	01202	02202 <i>y</i>	10202	11202	12202	20202 <i>y</i>	21202	22202 <i>y</i>
00210	01210	02210	10210	11210	12210	20210	21210	22210
00211	01211	02211	10211	11211 <i>z</i>	12211 <i>z</i>	20211	21211 <i>z</i>	22211 <i>z</i>
00212	01212	02212	10212	11212 <i>z</i>	12212 <i>z</i>	20212	21212 <i>z</i>	22212 <i>z</i>
00220 <i>y</i>	01220	02220 <i>y</i>	10220	11220	12220	20220 <i>y</i>	21220	22220 <i>y</i>
00221	01221	02221	10221	11221 <i>z</i>	12221 <i>z</i>	20221	21221 <i>z</i>	22221 <i>z</i>
00222 <i>y</i>	01222	02222 <i>y</i>	10222	11222 <i>z</i>	12222 <i>z</i>	20222 <i>y</i>	21222 <i>z</i>	22222*
7 <i>x</i> 7 <i>y</i>	8 <i>x</i>	8 <i>y</i>	8 <i>x</i>	7 <i>x</i> 7 <i>z</i>	8 <i>z</i>	8 <i>y</i>	8 <i>z</i>	7 <i>y</i> 7 <i>z</i>
27-1-14	27-8	27-8	27-8	27-1-14	27-8	27-8	27-8	27-1-14
12	19	19	19	12	19	19	19	12
	50			50			50	

Tabela. Nejeverni Tomaži si lahko res izpišejo vseh 3^5 možnosti delitve knjig in prečrtajo (označijo) napačne, pa bodo zopet prišli do $S(5, 3) = 150$. Naj pa bodo pozorni, da bi bilo dovolj izpisati prve tri stolpce, saj drugi trije in zadnji trije izgledajo precej podobno (v resnici jih dobimo iz prvega s permutacijo oznak: (012) oziroma (021)). To pa pomeni, da bi lahko z enako truda izračunali tudi $S(6, 3)$. V resnici lahko v ta namen uporabimo kar zgornjo tabelo - le ničlo si moramo prestavljati na začetku vsake peterice. To pa pomeni, da šesterice označenih z ni potrebno več odševati, dve šesterici, ki sta označeni z zvezdico pa bi morali označiti z oziroma y . Torej je $S(6, 3) = (150 + 30)3 = 540$. V splošnem pa dobimo na ta način rekurzivno zvezo $S(n + 1, 3) = 3(S(n, 3) + S(n, 2))$.

Če nadaljujemo izračun števila $S(n, 3)$ s uporabo formule za permutacije s ponavljanjem, pa stvar postane že skoraj rutinirano dolgočasna:

$$n = 8: \quad 3^8 - 3(2^8 - 2) - 3 = 5796 = 3P_8^{116} + 6P_8^{125} + 6P_8^{134} + 3P_8^{224} + 3P_8^{233},$$

$$n = 9: \quad 3^9 - 3(2^9 - 2) - 3 = 18150 = 3P_9^{117} + 6P_9^{126} + 6P_9^{135} + 3P_9^{144} + 3P_9^{225} + 6P_9^{234} + P_9^{333}.$$

Vse več je sumandov na desni strani, kar pomeni, da postaja za večje m prvi način bolj praktičen/učinkovit. Da pa se ne bi preveč dolgočasili, je morda čas, da rešimo še kakšen

primer, npr. ko je $m = 4$: Zapišimo zvezo iz katerih smo izračunali $S(n, 2)$ in $S(n, 3)$ za splošen m :

$$m^n = S(n, m) \binom{m}{m} + S(n, m-1) \binom{m}{m-1} + \cdots + S(n, 2) \binom{m}{2} + S(n, 1) \binom{m}{1}.$$

Le-ta nam da rekurzivno formulo za $S(n, m)$. V primeru $m = 4$ dobimo

$$S(n, 4) = 4^n - 4 \cdot S(n, 3) - 6 \cdot S(n, 2) - 4 \cdot S(n, 1)$$

oziroma, če upoštevamo še formule za $S(n, 3)$, $S(n, 2)$ in $S(n, 1)$:

$$S(n, 4) = 4^n - 4(3^n - 3(2^n - 2) - 3) - 6(2^n - 2) - 4 = 4^n - 4 \cdot 3^n + 6 \cdot 2^n - 4.$$

Bralcu prepuščamo, da pravkar dobljeno formulo testira (npr. bodisi s permutacijami s ponavljanjem ali pa kar štetjem izračuna $S(5, 4)$ (glej spodnjo tabelo) ter $S(6, 4)$ in $S(7, 4)$). \diamond

0000	0200	1000	1200	2000	2200	3000	3200
0001	0201	1001	1201	2001	2201	3001	3201
0002	0202	1002	1202	2002	2202	3002	3202
0003	0203	1003	1203	2003	2203	3003	3203
0010	0210	1010	1210	2010	2210	3010	3210
0011	0211	1011	1211	2011	2211	3011	3211
0012	0212	1012	1212	2012	2212	3012	3212
0013	0213	1013	1213	2013	2213	3013	3213
0020	0220	1020	1220	2020	2220	3020	3220
0021	0221	1021	1221	2021	2221	3021	3221
0022	0222	1022	1222	2022	2222	3022	3222
0023	0223	1023	1223	2023	2223	3023	3223
0030	0230	1030	1230	2030	2230	3030	3230
0031	0231	1031	1231	2031	2231	3031	3231
0032	0232	1032	1232	2032	2232	3032	3232
0033	0233	1033	1233	2033	2233	3033	3233
0100	0300	1100	1300	2100	2300	3100	3300
0101	0301	1101	1301	2101	2301	3101	3301
0102	0302	1102	1302	2102	2302	3102	3302
0103	0303	1103	1303	2103	2303	3103	3303
0110	0310	1110	1310	2110	2310	3110	3310
0111	0311	1111	1311	2111	2311	3111	3311
0112	0312	1112	1312	2112	2312	3112	3312
0113	0313	1113	1313	2113	2313	3113	3313
0120	0320	1120	1320	2120	2320	3120	3320
0121	0321	1121	1321	2121	2321	3121	3321
0122	0322	1122	1322	2122	2322	3122	3322
0123	0323	1123	1323	2123	2323	3123	3323
0130	0330	1130	1330	2130	2330	3130	3330
0131	0331	1131	1331	2131	2331	3131	3331
0132	0332	1132	1332	2132	2332	3132	3332
0133	0333	1133	1333	2133	2333	3133	3333

Tabela. Vsa štirimestna števila s števki 0, 1, 2, 3.

Slika 1: To so v obliki trikotnika zapisani binomski simboli, vsaka vrstica pa ustreza enemu binoskemu obrazcu.

Primer: Na polico bi radi postavili štiri mat., 6 fizikalnih in 2 kemijski knjigi. Na koliko načinov lahko to storimo:

- (a) če naj knjige iste stroke stojijo skupaj,
- (b) če kemijski knjigi ne smeta stati skupaj,
- (c) če morajo matematične knjige stati na začetku.

(Namig za (b): če želimo ponagajati prijateljici, dajmo tisti dve kemijski knjigi namenoma skupaj.) ◇

Primer: Na koliko načinov lahko sestavimo iz 7ih soglasnikov in 5ih samoglasnikov besedo, ki ima 4 soglasnike in 3 samoglasnike? (Namig: nalogo lahko rešimo tako s kmpbinacijami in permutacijami, kot tudi z variacijami in permutacijami s pomavljanjem.) ◇

A.5 Vrsta za e

Število e , ki predstavlja osnovo za naravni logaritem, pogosto definiramo s formulo

$$e = \lim_{n \rightarrow \infty} \left(1 + \frac{1}{n}\right)^n. \quad (\text{A.1})$$

Leta 1683 je Jakob Bernoulli poskušal izračunati limito $(1+1/n)^n$, ko gre n proti neskončno. Uporabil je binomski obrazec:

$$\left(1 + \frac{1}{n}\right)^n = 1 + \binom{n}{1} \frac{1}{n} + \binom{n}{2} \frac{1}{n^2} + \binom{n}{3} \frac{1}{n^3} + \cdots + \binom{n}{n} \frac{1}{n^n}.$$

The k -ti sumand na desni strani zgornje relacije je enak

$$\binom{n}{k} \frac{1}{n^k} = \frac{1}{k!} \cdot \frac{n(n-1)(n-2)\cdots(n-k+1)}{n^k}.$$

Za $n \rightarrow \infty$, gre slednji ulomek na desni proti 1, tj.

$$\lim_{n \rightarrow \infty} \binom{n}{k} \frac{1}{n^k} = \frac{1}{k!},$$

kar pomeni, da lahko e zapišemo kot vrsto.

$$e = 1 + \frac{1}{1!} + \frac{1}{2!} + \frac{1}{3!} + \cdots. \quad (\text{A.2})$$

O tem se prepričamo zato, ker je vsak člen v binomski razširitvi naraščujoča funkcija od n , sledi iz izreka o monotoni konvergenci za vrste, da je vsota te neskončne vrste enaka e . Bernoulli se je na ta način prepričal, da število e leži med 2 in 3. To sta torej v nekem smislu prva približka za e , vendar pa Bernoulli nikoli ni povezal tega števila z logaritmom¹. Za kaj takega je bilo potrebno izkristalizirati pojem funkcije in dognati, da sta eksponentna in logaritemska funkcija inverzni. Euler je v resnici prvi dokazal zgornjo zvezo (A.2), hkrati pa izračunal prvih 18 decimalk števila e : $e \approx 2,718281828459045235$

Še splošnejšo vrsto pa dobimo z uporabo Taylorjeve vrste:

Trditev A.5.

$$e^x = 1 + \frac{x}{1!} + \frac{x^2}{2!} + \frac{x^3}{3!} + \cdots.$$

¹Glej <http://www.gap-system.org/history/HistTopics/e.html#s19>

A.6 Stirlingov obrazec

Stirlingovo aproksimacijo (ozriroma formulo ali obrazec) uporabljamo za učinkovito računanje/ocenjevanje velikih faktorjelov in je poimenovana po škotskem matematiku Jamesu Stirlingu (1692–1770):

$$n! \approx \sqrt{2\pi n} \left(\frac{n}{e}\right)^n. \quad (\text{A.3})$$

Zavedati se moramo, da je naivno računanje zgornjega faktorjela $1 \cdot 2 \cdot 3 \dots (n-1) \cdot n$ eksponentne časovne zahtevnosti v odvisnosti od dolžina zapisa števila n (le-ta je seveda enaka naravnemu številu k , ki je zelo blizu $\log n$, pri čemer je logaritemska osnova odvisna od številske osnove, v kateri zapišemo število n). V čem smo torej na boljšem, ko računamo izraz na desni strani (A.3)? Računanje potence lahko izvedemo tako, da najprej izračunamo naslednje potence $(n/e)^0, (n/e)^1, (n/e)^2, \dots, (n/e)^k$, nato pa zmnožimo med seboj tiste katerih eksponenti ustrezajo mestom enic v binarni predstavitvi števila n , kar pomeni da smo opravili največ $2k$ množenj.

Primer: Namesto, da bi izračunali $P = a^{21}$ z dvajsetimi množenji ($P = a$, dvajsetkrat ponavljaj $P := P * a$), raje izračunamo potence a, a^2, a^4, a^8, a^{16} , nato pa zaradi $21 = 2^4 + 2^2 + 2^1$ še $P = a^{16} \cdot a^4 \cdot a^1$, kar znese samo 4 kvadriranja in 2 množenji. \diamond

Formulo (A.3) je prvi odkril Abraham de Moivre v naslednji obliki

$$n! \sim [\text{constant}] \cdot n^{n+1/2} e^{-n}.$$

pri čemer je konstanto izrazil s hiperboličnim logaritmom. James Stirlingov prispevek pa je bil, da je konstanta v resnici enaka $\sqrt{2\pi}$. Formulo tipično uporabljamo v aplikacijah v obliki

$$\ln n! \approx n \ln n - n.$$

V zgornji verziji manjka faktor $\frac{1}{2} \ln(2\pi n)$, ki ga lahko za velike n zanemarimo v primerjavi z drugimi faktorji. Zapišimo $\ln(n!) = \ln 1 + \ln 2 + \dots + \ln n$. pri čemer lahko na desno stran zgornje relacije gledamo kot na približek za integral

$$\int_1^n \ln(x) dx = n \ln n - n + 1.$$

Od tu naprej pa si lahko pomagamo z Euler–Maclaurinovo formulo in uporabo Bernoullijevih števil. Glej npr. http://en.wikipedia.org/wiki/Stirling's_approximation.

A.7 Normalna krivulja v prostoru

Želimo pokazati naslednjo identiteto.²

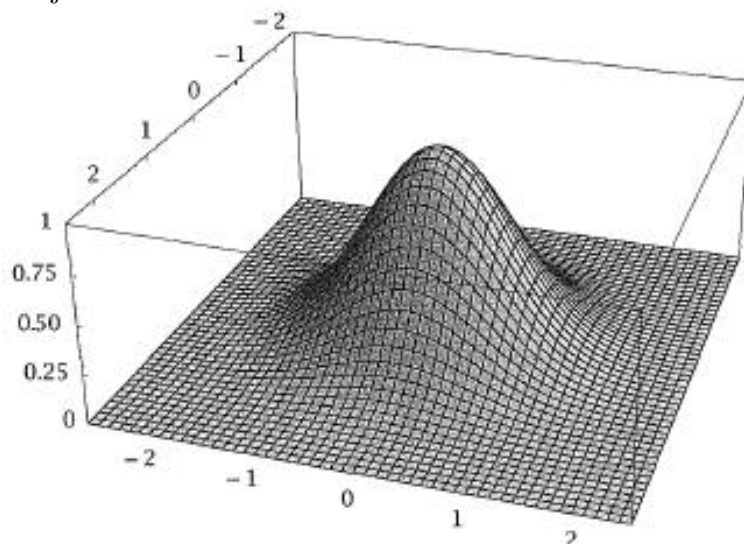
Izrek A.6.

$$\int_{-\infty}^{\infty} e^{-t^2} dt = \sqrt{\pi}. \quad (\text{A.4})$$

Dokaz. Označimo vrednost integrala na levi strani (A.4) z I . Funkcija

$$g(s, t) = e^{-(s^2+t^2)} = e^{-s^2} e^{-t^2}$$

je narisana na spodnji sliki.



Slika: Normalna gora.

Sedaj pa prerežimo zgornjo ploskev z ravnino $s = 1$. Prerez seveda izgleda kot normalna krivulja, ploščina pod dobljeno krivuljo pa je enaka ploščini pod normalno krivuljo, ki je pomnožena z e^{-1} :

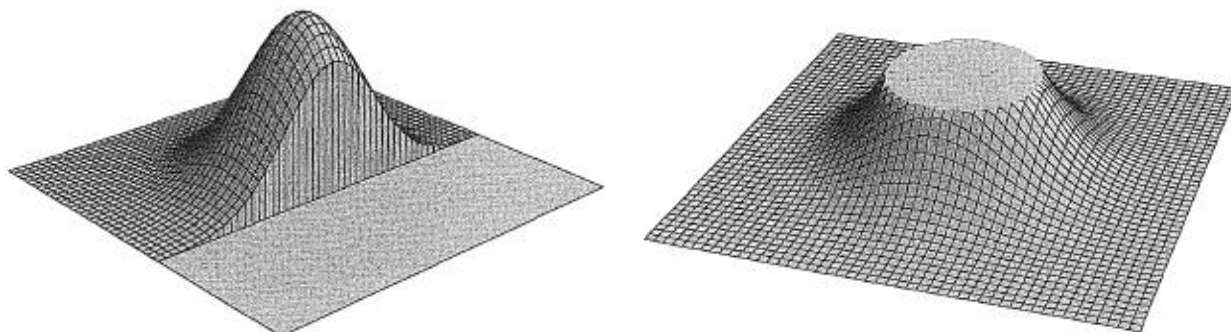
$$\int_{-\infty}^{\infty} e^{-1} e^{-t^2} dt = e^{-1} I.$$

Podobno je za katerokoli drugo vrednost števila s ploščina ustrezne krivulje enaka $e^{-s^2} I$. Sedaj lahko izračunamo prostornino normalne gore z naslednjim integralom

$$V = \int_{-\infty}^{\infty} e^{-s^2} I ds = I^2.$$

²Stara zgodba pravi, da je Lord Kelvin nekoč dejal, da je matematik nekdo, za katerega je ta identiteta očitna.

Preostane nam le še, da dokažemo, da je $V = \pi$. Tokrat preseka jmo normalno ploskev z ravnino $z = h$.



Slika: Vertikalni in horizontalni prerez normalne ploskve.

Potem za točke preseka velja

$$e^{-s^2-t^2} \geq h \quad \text{oziroma} \quad s^2 + t^2 \leq -\ln h.$$

To pa je ravno krog s središčem $(0,0)$ in polmerom $r = \sqrt{-\ln h}$. Ploščina tega kroga je $\pi r^2 = \pi(-\ln h)$, prostornino V pa dobimo z integriranjem od $h = 0$ do $h = 1$:

$$V = \int_0^1 \pi(-\ln h) dh = -\pi \int_0^1 \ln h dh.$$

Postavimo $u = \ln x$, $dv = dx$, $du = dx/x$ in $v = \int dx = x$. Z integriranjem po delih dobimo

$$\int_0^1 \ln x dx = x \ln x \Big|_0^1 - \int_0^1 \frac{x dx}{x} = -1$$

in od tod $V = \pi$. □

A.9 Cauchyjeva neenakost

Skalarni produkt vektorjev \vec{u} in \vec{v} iz \mathbb{R}^n je po definiciji enak

$$\vec{u} \cdot \vec{v} = |\vec{u}| \operatorname{proj}_{\vec{u}} \vec{v},$$

kjer je $\operatorname{proj}_{\vec{u}} \vec{v}$ pravokotna projekcija vektorja \vec{v} na vektor \vec{u} . Od tod sledi $\vec{u} \cdot \vec{v} = |\vec{u}| |\vec{v}| \cos \varphi$ (glej sliko). Naj bo $\vec{u} = (u_1, \dots, u_n)$, $\vec{v} = (v_1, \dots, v_n)$, potem lahko skalarni produkt izračunamo z naslednjo vsoto

$$\vec{u} \cdot \vec{v} = \sum_{i=1}^n u_i v_i = u_1 v_1 + \dots + u_n v_n.$$

Potem je dolžina vektorja \vec{v} enaka $|\vec{v}| = \sqrt{\vec{v} \cdot \vec{v}}$. Neposredno iz $|\cos \varphi| \leq 1$ sledi za realne vektorje naslednja neenakost.



Trditev A.7. (Cauchyjeva neenakost) Za poljubna vektorja \vec{u} in \vec{v} velja

$$|\vec{u} \cdot \vec{v}| \leq |\vec{u}| |\vec{v}|.$$

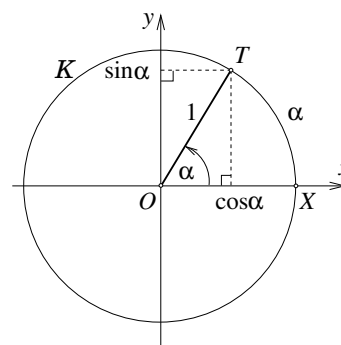
Enakost velja natanko tedaj, ko je kota med vektorjema \vec{u} in \vec{v} enak $k\pi$, za $k \in \mathbb{N}$, to je natanko tedaj, ko sta vektorja \vec{u} in \vec{v} kolinearna.

Pravimo ji tudi Cauchy-Schwarzova neenakost ali neenakost Bunjakovskega, glej http://en.wikipedia.org/wiki/Cauchy-Schwarz_inequality. To je ena izmed najpomembnejših neenakosti. Omogoča nam, da preverimo kolinearnost dveh vektorjev tako, da izračunamo vrednosti na levi in desni strani zgornje neenakosti in preverimo, če sta enaki.

Za $n = 2$ Cauchyjeva neenakost sledi tudi iz identitete

$$(a^2 + b^2)(c^2 + d^2) = (ad - bc)^2 + (ac + bd)^2, \quad a, b, c, d \in \mathbb{R},$$

ki je že naša stara znanka in pove, da za kompleksni števili $z = a + ib$ in $w = c + id$ produkt absolutnih vrednosti dveh kompleksnih števil enak absolutni vrednosti ustreznega produkta, tj. $|z| \cdot |w| = |zw|$.



Slika A.4.1: Definiciji sinusa in kosinusa. V ravnini narišemo enotsko krožnico \mathcal{K} s središčem O v izhodišču koordinatnega sistema. Iz točke $X = (1, 0)$ se v nasprotni smeri od urinega kazalca poda na pot po krožnici \mathcal{K} točka T . Ko ima za seboj "prehojen" lok dolžine α (takrat je kot $\angle XOT$ enak α radianov), ima točka T koordinati $(\cos \alpha, \sin \alpha)$. Funkcija sinus je pozitivna v prvem in drugem kvadrantu, funkcija kosinus pa v prvem in četrtem. Obe funkciji sta periodični s periodo 2π (tj. 360°).

Za $n = 3$ lahko dokažemo Cauchyjevo neenakost s pomočjo vektorskega produkta in Lagrangeove identitete. **Vektorski produkt** vektorjev $\vec{u} = (u_1, u_2, u_3)$ in $\vec{v} = (v_1, v_2, v_3)$ iz \mathbb{R}^3 je vektor v \mathbb{R}^3 podan s formulo: $\vec{u} \times \vec{v} = (u_2v_3 - u_3v_2, -u_1v_3 + u_3v_1, u_1v_2 - u_2v_1)$. Dolžina vektorja $\vec{u} \times \vec{v}$ je enaka $|\vec{u} \times \vec{v}| = |\vec{u}||\vec{v}|\sin\varphi$, geometrično pa to pomeni, da je dolžina vektorskega produkta enaka ploščini paralelograma, ki ga razpenjata vektorja \vec{u} in \vec{v} . **Lagrangeova identiteta** $|\vec{u}|^2|\vec{v}|^2 = (\vec{u}\vec{v})^2 + |\vec{u} \times \vec{v}|^2$, ni nič drugega kot na drugačen način zapisana relacija $\sin^2\varphi + \cos^2\varphi = 1$ (oziroma Pitagorjev izrek).

Za splošen $n \in \mathbb{N}$ lahko Cauchyjevo neenakost zapišemo tudi v naslednji obliki:

$$(a_1^2 + \cdots + a_n^2)(b_1^2 + \cdots + b_n^2) \geq (a_1b_1 + \cdots + a_nb_n)^2,$$

kjer so $a_1, \dots, a_n, b_1, \dots, b_n$ poljubna realna števila. Lagrangeova identiteta za splošen n pa izgleda takole:

$$\left(\sum_{i=1}^n a_i^2\right)\left(\sum_{i=1}^n b_i^2\right) = \left(\sum_{i=1}^n a_ib_i\right)^2 + \sum_{i<j} (a_ib_j - a_jb_i)^2.$$

Zadnja vsota ima $(n-1) + (n-2) + \cdots + 2 + 1 = (n-1)n/2$ členov.

Raba v verjetnosti

Za slučajni spremenljivki X in Y je matematično upanje njunega produkta skalarni produkt, tj.

$$\langle X, Y \rangle := E(XY)$$

zadovoljuje tri aksiome iz naslednje škatle. (V tem primeru velja $\langle X, X \rangle = 0$ natanko tedaj, ko je $P(X=0) = 1$.) Potem iz Cauchyjeve neenakosti sledi

$$|E(XY)|^2 \leq E(X^2)E(Y^2).$$

Naj bo $\mu = E(X)$ in $\nu = E(Y)$. Potem po Cauchyjevi neenakosti velja

$$\begin{aligned} |\text{Cov}(X, Y)|^2 &= |E((X - \mu)(Y - \nu))|^2 = |\langle X - \mu, Y - \nu \rangle|^2 \\ &\leq \langle X - \mu, X - \mu \rangle \langle Y - \nu, Y - \nu \rangle = E((X - \mu)^2)E((Y - \nu)^2) = D(X)D(Y), \end{aligned}$$

kjer je D disperzija, Cov pa kovarianca.

Formalno je **vektorski prostor s skalarnim produktom** (rečemo tudi *unitarni vektorski prostor*) vektorski prostor V nad poljubnim obsegom \mathbb{F} s skalarnim produktom, tj. s preslikavo

$$\langle \cdot, \cdot \rangle : V \times V \rightarrow \mathbb{F}$$

ki zadovoljuje naslednje tri aksiome za poljubne vektorje $x, y, z \in V$ in skalarje $a \in \mathbb{F}$:

- **Konjugirana simetrija:** $\langle x, y \rangle = \overline{\langle y, x \rangle}$.
- **Linearnost na prvi koordinati:** $\langle ax, y \rangle = a\langle x, y \rangle$. $\langle x + y, z \rangle = \langle x, z \rangle + \langle y, z \rangle$.
- **Positivna-definitnost:** $\langle x, x \rangle \geq 0$, kjer velja enakost, če in samo če je $x = 0$.

(zgoraj smo opustili vektorske oznake). Glej

http://en.wikipedia.org/wiki/Inner_product_space.

Predstavimo še dokaz Cauchyjeve neenakosti za vektorski prostor s skalarnim produktom.

Dokaz. Naj bosta u in v poljubna vektorja vektorskega prostora V nad obsegom \mathbb{F} . Neenakost je očitna za $v = 0$, zato predpostavimo, da je $\langle v, v \rangle \neq 0$. Naj bo $\delta \in \mathbb{F}$. Potem velja

$$0 \leq |u - \delta v|^2 = \langle u - \delta v, u - \delta v \rangle = \langle u, u \rangle - \bar{\delta}\langle u, v \rangle - \delta\langle u, v \rangle + |\delta|^2\langle v, v \rangle.$$

Sedaj pa izberimo $\delta = \langle u, v \rangle \cdot \langle v, v \rangle^{-1}$, in dobimo

$$0 \leq \langle u, u \rangle - |\langle u, v \rangle|^2 \cdot \langle v, v \rangle^{-1} \quad \text{ozirama} \quad |\langle u, v \rangle|^2 \leq \langle u, u \rangle \cdot \langle v, v \rangle,$$

in končno po korenjenju neenakost, ki smo jo želeli pokazati. □

Trikotniško neenakost za vektorske prostore s skalarnim produktom pogosto pokažemo kot posledico Cauchyjeve neenakosti na naslednji način: za vektorja x in y velja

$$|x + y|^2 = \langle x + y, x + y \rangle = |x|^2 + \langle x, y \rangle + \langle y, x \rangle + |y|^2 \leq |x|^2 + 2|x||y| + |y|^2 = (|x| + |y|)^2.$$

Po korenjenju dobimo trikotniško neenakost.

Leonhard Euler (1707–1783)

Swiss mathematician who was tutored by Johann Bernoulli. He worked at the Petersburg Academy and Berlin Academy of Science. He had a phenomenal memory, and once did a calculation in his head to settle an argument between students whose computations differed in the fiftieth decimal place. Euler lost sight in his right eye in 1735, and in his left eye in 1766. Nevertheless, aided by his phenomenal memory (and having practiced writing on a large slate when his sight was failing him), he continued to publish his results by dictating them. Euler was the most prolific mathematical writer of all times finding time (even with his 13 children) to publish over 800 papers in his lifetime. He won the Paris Academy Prize 12 times. When asked for an explanation why his memoirs flowed so easily in such huge quantities, Euler is reported to have replied that his pencil seemed to surpass him in intelligence. François Arago said of him "He calculated just as men breathe, as eagles sustain themselves in the air" (Beckmann 1971, p. 143; Boyer 1968, p. 482).

Euler systematized mathematics by introducing the symbols e , i , and $f(x)$ for f a function of x . He also made major contributions in optics, mechanics, electricity, and magnetism. He made significant contributions to the study of differential equations. His *Introductio in analysin infinitorum* (1748) provided the foundations of analysis. He showed that any complex number to a complex power can be written as a complex number, and investigated the beta and gamma functions. He computed the Riemann zeta function to for even numbers.

He also did important work in number theory, proving that that the divergence of the harmonic series implied an infinite number of Primes, factoring the fifth Fermat number (thus disproving Fermat's conjecture), proving Fermat's lesser theorem, and showing that e was irrational. In 1772, he introduced a synodic coordinates (rotating) coordinate system to the study of the three-body problem (especially the Moon). Had Euler pursued the matter, he would have discovered the constant of motion later found in a different form by Jacobi and known as the Jacobi integral.

Euler also found the solution to the two fixed center of force problem for a third body. Finally, he proved the binomial theorem was valid for any rational exponent. In a testament to Euler's proficiency in all branches of mathematics, the great French mathematician and celestial mechanic Laplace told his students, "Lisez Euler, lisez Euler, c'est notre maître à tous" ("Read Euler, read Euler, he is our master in everything" (Beckmann 1971, p. 153).

<http://scienceworld.wolfram.com/biography/Euler.html>

Blaise Pascal (1623-1662)

French mathematician, philosopher, and religious figure. He studied the region above the mercury in a barometer, maintaining that it was a vacuum. In his investigations of the barometer, he found that the height to which the mercury rose was the same regardless of shape. Based on his double vacuum experiment, he formulated Pascal's principle, which states that the pressure is constant throughout a static fluid. He performed an experiment in which he convinced his brother-in-law to climb the Puy-de-Dôme Mountain in France. He found the height of the mercury dropped with altitude, indicating pressure decreases with altitude.

Pascal also designed and built mechanical adding machines, and incorporated a company in 1649 to produce and market them. Unfortunately, the machines were rather highly priced and had reliability problems. Only seven of Pascal's devices survive today.

Pascal suffered from serious health problems, and spent most of his final years writing on religious philosophy. <http://scienceworld.wolfram.com/biography/Pascal.html>

Augustin Louis Cauchy(1789 – 1857)

<http://www-history.mcs.st-and.ac.uk/Mathematicians/Cauchy.html>

Cauchy pioneered the study of analysis, both real and complex, and the theory of permutation groups. He also researched in convergence and divergence of infinite series, differential equations, determinants, probability and mathematical physics.

<http://scienceworld.wolfram.com/biography/Cauchy.html>

French mathematician who wrote 789 papers, a quantity exceeded only by Euler and Cayley, which brought precision and rigor to mathematics. He invented the name for the determinant and systematized its study and gave nearly modern definitions of limit, continuity, and convergence. Cauchy founded complex analysis by discovering the Cauchy-Riemann equations (although these had been previously discovered by d'Alembert).

Cauchy also presented a mathematical treatment of optics, hypothesized that ether had the mechanical properties of an elasticity medium, and published classical papers on wave propagation in liquids and elastic media. After generalizing Navier's equations for isotropic media, he formulated one for anisotropic media. Cauchy published his first elasticity theory in 1830 and his second in 1836. Both were rather ad hoc and were riddled with problems, and Cauchy proposed a third theory in 1839. Cauchy also studied the reflection from metals and dispersion relationships.

Cauchy extended the polyhedral formula in a paper which was criticized by Malus. His theory of substitutions led to the theory of finite groups. He proved that the order of any subgroup is a divisor of the order of the group. He also proved Fermat's three triangle theorem. He refereed a long paper by Le Verrier on the asteroid Pallas and invented techniques which allowed him to redo Le Verrier's calculations at record speed. He was a man of strong convictions, and a devout Catholic. He refused to take an oath of loyalty, but also refused to leave the French Academy of Science.

Dodatek B

PROGRAM R (Martin Raič)

Predstavili bomo osnovne ukaze v programu R, ki so povezani z našim predmetom.

Informacije

Dokumentacija na Linuxu: `/usr/share/doc/r-doc-html/manual` .

Pomoč v R-ovem pozivniku:

- `help(točno določena funkcija)`
- `help.search("nekaj približnega")`

B.1 Izvajanje programa

Iz pozivnika našega operacijskega sistema program zaženemo z ukazom `R`. Če ni določeno drugače, se znajdemo v R-ovem pozivniku. Program lahko zaženemo z obilo opcijami. Njihov seznam izpiše ukaz `R -h` ali `R --help` (in nato takoj konča). R ne zažene pozivnika in le izpiše rezultat, nakar konča, če mu na vhodu predpišemo ukazni niz. To pa je lahko:

- Standardni vhod, če R zaženemo kot cevovod ali pa z dostavkom `< datoteka`.
V tem primeru mu moramo predpisati še opcijo `--save`, `--no-save` ali `--vanilla`.
- Vsebina datoteke z imenom, ki sledi opciji `-f` oziroma `--file`.
- Ukazni niz, ki sledi opciji `-e`.

Pri izvajanju R-a na ta način pride prav opcija `--slave`, ki izključi ves nenujni izhod.

Zadnjo vrnjeno vrednost dobimo z ukazom `.Last.value`.

Izhod iz programa dosežemo z ukazom `q()`.

B.2 Aritmetika

Elementarne binarne operacije: `+`, `-`, `*`, in `**`

Elementarne funkcije: `sqrt`, `exp`, `log`, `sin`, `cos`, `tan`, `asin`, `acos`, `atanr`

Konstanta: `pi`

Izpis na določeno število (n) signifikantnih decimalk: `print(x, digits=n)`

(Več o izpisovanju kasneje).

Zaokrožitvene funkcije: `trunc`, `round`, `floor`, `ceiling`

Fakulteta: `factorial`

Funkcija gama: `gamma`

Binomski simbol: `choose(n, k)` vrne n nad k , tj. $\binom{n}{k}$.

Naključna števila: `runif(1)` vrne psevdonaključno število med 0 in 1 z enakomerno porazdelitvijo. Več o naključnih številih kasneje.

B.3 Najosnovnejše o spremenljivkah

Prireditvev: `x <- nekaj` ali `nekaj -> x`

Izbris: `rm(x)` ali tudi `rm(x, y)`.

Ukaz `ls` vrne seznam vseh simbolov, ki so trenutno definirani, razen tistih, katerih imena se začenjajo s piko. Če želimo vključiti še te, vnesemo `ls(all=TRUE)`.

Ukaz `rm(list=ls(all=TRUE))` izbrši vse trenutno definirane simbole.

Osnovni podatkovni tipi:

- števila: cela (npr. `integer(-42)`), realna (npr. `1.23`) in kompleksna (npr. `2 + 1i`);
- nizi znakov (npr. `"Zemlja"`);
- logične vrednosti: `TRUE`, `FALSE`;
- prazna vrednost: `NULL`.

B.4 Uporabnikove funkcije

Anonimna funkcija: `function(parametri) telo`

Funkcija z imenom: `ime <- function(parametri) telo`

Ukaz `plot(funkcija, sp. meja, zg. meja)` nariše graf funkcije.

B.5 Numerično računanje

`uniroot(f, c(a, b))` vrne ničlo zvezne funkcije f na intervalu $[a, b]$. Vrednosti na krajiščih morata imeti nasproten predznak. Vselej vrne le eno ničlo. Pravzaprav vrne celotno poročilo o ničli. Če želimo le ničlo, ukažemo: `uniroot(f, c(a, b))$root`.

`integrate(f, a, b)` numerično integrira funkcijo f od a do b . Spet vrne celo poročilo, če želimo le vrednost, ukažemo: `integrate(f, a, b)$value`.

B.6 Podatkovne strukture

B.6.1 Vektorji

Primer: Konstrukcija vektorja: `c(7, 8, 9)` da isto kot `7:9` ali `seq(7, 9)`. ◇

Pri ukazu `seq` lahko predpišemo tudi korak, recimo `seq(70, 90, 10)`. `runif(n)` vrne naključni vektor dolžine n . Več o tem kasneje. Vektorje lahko tvorimo tudi iz nizov:

```
x <- c("Merkur", "Venera", "Zemlja", "Mars", "Jupiter", "Saturn", "Uran", "Neptun").
```

Ukaz `c` tudi združuje vektorje.

POZOR! Vsi elementi v vektorju morajo biti istega tipa. Če so tipi različni, se nižji tipi pretvorijo v višje.

Primeri:

- `c(1, 2, "Zemlja")` se pretvori v `c("1", "2", "Zemlja")`.
- `c(1, 2, TRUE)` se pretvori v `c(1, 2, 1)`.
- `c(1, 2, TRUE, "Zemlja")` se pretvori v `c("1", "2", "TRUE", "Zemlja")`. ◇

Če ni določeno drugače, ukaz `c` izpusti vrednosti `NULL`. Z ukazom `[...]` dobimo ven posamezne komponente vektorja. Natančneje, če je v vektor, ukaz `v[i]` deluje odvisno od narave objekta i na naslednji način:

- Če je i naravno število, vrne element z indeksom i . Indeksi se štejejo od 1 naprej.
- Če je i negativno celo število, vrne vektor brez elementa z ustreznim indeksom.

- Če je i vektor iz naravnih števil, vrne vektor iz elementov z ustreznimi indeksi (primerno za komponiranje preslikav).
POZOR! Primer tovrstne kode je `v[c(2, 4, 3, 2)]` in ne recimo `v[2, 4, 3, 2]`: slednje pomeni večrazsežni indeks v tabeli – glej kasneje.
- Če je i vektor iz negativnih celih števil, vrne vektor, ki ima ustrezne elemente izpuščene.
- Če je i vektor iz logičnih vrednosti, vrne vektor iz komponent vektorja v , pri katerih je na odgovarjajočem položaju v i vrednost `TRUE`.

Ponovitev elementa ali vektorja: `rep(vrednost ali vektor, kolikokrat)`.

Obrat vektorja: `rev(x)` vrne vektor x od zadaj naprej.

`plot(x)` nariše točke, ki pripadajo koordinatam vektorja x .

To je isto kot `plot(x, type="p")`.

Druge opcije za ukaz `type`:

- `"l"`: nariše lomljenko;
- `"b"`: točke poveže z daljicami;
- `"h"`: nariše histogram iz navpičnih črt;
- `"s"`: nariše stopnice, pri čemer gre posamezna stopnica desno od točke;
- `"S"`: nariše stopnice, pri čemer gre posamezna stopnica levo od točke;

Ukaz `barplot(x)` nariše vektor x v obliki lepega histograma.

B.6.2 Matrike

Matriko:

1	2	3
4	5	6
7	8	9
10	11	12

vnesemo z enim izmed naslednjih ukazov:

- `matrix(c(1, 4, 7, 10, 2, 5, 8, 11, 3, 6, 9, 12), nrow=4, ncol=3)`

- `matrix(1:12, nrow=4, ncol=3, byrow=TRUE)`
- `array(c(1, 4, 7, 10, 2, 5, 8, 11, 3, 6, 9, 12), dim=c(4, 3))`
- `rbind(1:3, 4:6, 7:9, 10:12)`
- `cbind(c(1, 4, 7, 10), c(2, 5, 8, 11), c(3, 6, 9, 12))`

Preprost je tudi vnos diagonalnih matrik, npr.

`diag(3, nrow=5)` ali `diag(c(3, 2, 4, 5, 1))`.

Priklic elementov matrike:

- `A[i, j]` vrne element v i -ti vrstici in j -tem stolpcu matrike A .
- `A[i,]` vrne i -to vrstico matrike A .
- `A[, j]` vrne j -ti stolpec matrike A .
- `A[i]` vrne i -ti element matrike, pri čemer so elementi urejeni po stolpcih.
Tako pri zgornji matriki `A[8]` vrne `11`.

Ukaz `c(A)` ali `as.vector(A)` iz matrike A naredi dolg vektor, pri čemer združuje po stolpcih.

Vse aritmetične operacije delujejo tudi na vektorjih in matrikah – po komponentah. Tako `A*B` zmnoži matriki A in B po komponentah. Ukaz `A + 2` pa vrne matriko A , ki ima vse elemente povečane za 2. Seveda lahko matrike in vektorje posredujemo tudi funkcijam. Recimo `(function(x) x ** 3 + x)(0:3)` vrne isto kot `c(0, 2, 10, 30)`.

Matrične operacije:

- `%*%:` matrično množenje,
- `%o%:` tenzorski produkt – množenje vsake komponente z vsako,
- `t:` transponiranje,
- `det:` determinanta,
- `solve(A, B):` reši sistem $Ax = b$,
- `solve(A):` poišče A^{-1} ,
- `eigen(A):` poišče lastne vrednosti in lastne vektorje (dobro deluje le, če se da matrika diagonalizirati).

Osnovna obdelava podatkov na vektorjih in matrikah:

- `sum(x)`: vsota elementov vektorja ali matrike
- `prod(x)`: produkt elementov vektorja ali matrike
- `mean(x)`: povprečje elementov vektorja ali matrike
- `sd(x)`: popravljeni standardni odklon elementov vektorja ali matrike
- `cumsum(x)`: vektor iz kumulativnih vsot
- `cumprod(x)`: vektor iz kumulativnih produktov
- `min(x)`: minimum elementov
- `max(x)`: maksimum elementov

Posredovanje funkcij po komponentah:

- Ukaz `sapply(vektor ali matrika, funkcija)` posreduje funkcijo posameznim komponentam vektorja ali matrike. Tako npr. ukaz `sapply(c(x1, ... xn), f)` vrne `c(f(x1), ... f(xn))`.
- Ukaz `apply(matrika, 1, funkcija)` posreduje funkcijo posameznim vrsticam.
- Ukaz `apply(matrika, 2, funkcija)` posreduje funkcijo posameznim stolpcem.
- Nasploh ukaz `array` deluje na poljubnorazsežnih tabelah. V drugem argumentu določimo številke razsežnosti, ki ostanejo.
- Ukaz `mapply(f, c(x11, ... x1n), ... c(xm1, ... xmn))` vrne `c(f(x11, ... xm1), ... f(x1n, ... xmn))`.
- Ukaz `outer(c(x1, ... xm), c(y1, ... yn), FUN=f)` vrne matriko z elementi `f(xi, yj)`. Recimo, operacija `%o%` se da definirati z ukazom `outer`, kjer za funkcijo izberemo množenje. Toda pozor: ukaz je implementiran tako, da funkcijo kliče na ustreznih zgeneriranih matrikah. Rešitev:


```
outer(c(x1, ... xm), c(y1, ... yn), FUN=function(v1, v2) mapply(f, v1, v2) ).
```

B.6.3 Tabele

Tabele imajo lahko poljubno mnogo razsežnosti. Dobimo jih lahko iz vektorjev z ukazom: `array(vektor, dim=c(razsežnosti))`. Z ukazom `[...]` lahko podobno kot pri matrikah izluščimo poljubno razsežne komponente tabele.

B.6.4 Vektorji, matrike in tabele z označenimi indeksi

Vektor z označenimi indeksi lahko dobimo z ukazom

```
c(oznaka1 = vrednost1, oznaka2 = vrednost2, ... ).
```

Lahko pa konstruiramo tudi enorazsežno tabelo z ukazom `array`.

Primer: `array(c(20, 50, 30), dimnames=list(c("prvi", "drugi", "tretji")))` ◇

Seveda nastavitve `dimnames` deluje za poljubno razsežno tabelo.

Deluje tudi pri ukazu `matrix`.

Primer: `matrix(1:12, nrow=4, ncol=3, byrow=TRUE, dimnames=list(c("prva", "druga", "tretja", "cetrt"), c("prvi", "drugi", "tretji")))`

Vrsticam in stolpcem lahko imena prirejamo ali spreminjamo tudi naknadno, recimo:

```
dimnames(x) <- list(c("prvi", "drugi", "tretji"))
```

 ◇

POZOR! Vektor ni eno-razsežna tabela, zato mu ne moremo prirejati nastavitve `dimnames`.

Če želimo to, ga moramo najprej z ukazom `array` pretvoriti v enorazsežno tabelo.

Če ima tabela označene indekse, jih lahko uporabimo tudi znotraj oklepajev `[...]`.

Seveda pa lahko še vedno uporabimo tudi številke.

Oznake komponent upoštevata tudi ukaza `plot` in `boxplot`.

B.6.5 Zapisi

Zapisi so objekti, sestavljeni iz več komponent, ki so lahko različnih tipov. Komponentam bomo rekli rubrike. Zapise konstruiramo z ukazom `list`.

Primer: `nas_clovek <- list(ime="Janez", starost=35, zakonec="Marija", starosti_otrok=c(15, 13, 2))`. ◇

Posamezne rubrike dobimo z ukazom `$`, npr. `nas_clovek$ime`.

Lahko uporabimo tudi `nas_clovek[["ime"]]`.

Lahko priklicujemo nadaljnje komponente, npr. `nas_clovek$starosti_otrok[2]`.

Z oglatimi oklepaji izluščimo del zapisa, npr.

```
nas_clovek[ime] ali nas_clovek[c("ime", "starost")].
```

Imena rubrik dobimo in nastavljammo z ukazom `names`.

B.6.6 Kontingenčne tabele in vektorji s predpisanimi vrednostmi

Kontingenčne tabele so 1- ali 2-razsežne tabele z označenimi indeksi, elementi pa so števila. Dobimo jih lahko iz vektorjev z označenimi indeksi z ukazom `as.table`. Ukaz `table(vektor)` pa iz vektorja naredi tabelo zastopanosti njegovih vrednosti. Le-te sortira po abecedi. Če želimo vrednosti ali njihov vrstni red predpisati, uporabimo vektor s predpisanimi vrednostmi, ki ga dobimo z ukazom `factor` in nastavitvijo `levels`: `table(factor(vektor, levels=vrednosti))`.

Pri vektorju s predpisanimi vrednostmi so le-te vedno nizi. Navaden vektor dobimo nazaj z ukazom `as.vector`. POZOR! Ukaz `c` spremeni vrednosti v naravna števila!

Ukaz `table(vektor1, vektor2)` naredi 2-razsežno kontingenčno tabelo. Tabele prav tako rišemo z ukazoma `plot` in `barplot`. Če je t kontingenčna tabela, ukaza `t[...]` in `t[[...]]` delujeta podobno kot pri zapisih.

B.6.7 Preglednice

Preglednice so podobne dvo-razsežnimi tabelam – z razliko, da so lahko stolpci različnih tipov, zato jim bomo tudi tu rekli rubrike. V posamezni rubriki so bodisi sama števila bodisi sami nizi bodisi same logične vrednosti – tako kot pri vektorjih. Preglednice konstruiramo z ukazom `data.frame`.

Primer:

```
nasi_mozje <- data.frame(
  ime=c("Janez", "Franc", "Joze"),
  starost=c(35, 42, 55), zena=c("Marija", "Stefka", "Lojzka"), st_otrok=c(3, 2, 4)
).
```

◇

Dostop do posameznih komponent je podoben kot pri matrikah – s tem, da:

- če priklicujemo vrstice, je rezultat spet preglednica;
- če priklicujemo stolpce, je rezultat vektor, ki ima predpisane vrednosti, če gre za nize;
- če priklicujemo posamezen element v stolpcu z nizi, je rezultat še vedno vektor s predpisanimi vrednostmi.

Posamezne rubrike lahko podobno kot pri zapisih dobimo tudi z ukazoma `$` in `[[...]]`. Imena rubrik tudi tu dobimo in nastavljammo z ukazom `names`.

Primer: `nase_zene <- nasi_mozje[,c("žena", "starost", "ime", "st_otrok")]`
`names(nase_zene) <- c("ime", "starost", "moz", "st_otrok")` ◇

B.7 Osnove programiranja

R je močan programski jezik, ki podpira tudi objektno orientirano programiranje. Pišemo lahko kratke programčke v pozivniku, daljše programe pa lahko shranimo v datoteke.

B.7.1 Izvajanje programov, shranjenih v datotekah

Samostojne programe, shranjene v datotekah, lahko izvajamo z naslednjimi klici:

- `R < datoteka` – v tem primeru mu moramo predpisati še opcijo `--save`, `--no-save` ali `--vanilla`;
- `R -f datoteka` ali `R --file datoteka`;
- `Rscript datoteka`;
- na Linuxu lahko izvedemo tudi samo datoteko, če se le-ta začne z `#!/usr/bin/Rscript`.

Pri prvih dveh načinih R izpiše celoten potek izvajanja programa. To mu lahko preprečimo z opcijo `--slave`.

Morebitne parametre iz ukazne vrstice dobimo z ukazom `commandArgs(TRUE)`, ki vrne vektor iz pripadajočih nizov. Vključevanje datotek iz pozivnika ali med izvajanjem programa

Pomožne programe lahko vključimo v R-ov pozivnik ali drugo programsko kodo z ukazom `source`.

Pregled nad delovnim imenikom (direktorijem):

- Ukaz `getwd()` vrne lokacijo imenika, kjer R bere in piše.
- Ukaz `setwd(imenik)` nastavi lokacijo delovnega.
- Ukaz `list.files()` vrne vsebino delovnega imenika – kot vektor.
Lahko mu predpišemo obilo nastavitev.

B.7.2 Najosnovnejši programski ukazi

Posamezne stavke ločimo s podpičjem ali prehodom v novo vrstico. Slednje lahko storimo tudi znotraj oklepajev ali narekovajev.

Znak `#` označuje komentar: vsebina od tega znaka do konca vrstice se ignorira.

Za splošno izpisovanje uporabimo ukaz `cat`. Če mu podamo več argumentov ali vektor, jih loči s presledkom. To lahko spremenimo z nastavitvijo `sep`. V nasprotju z ukazom `print` ukaz `sep` ne zaključí vrstice. To dosežemo z nizom `"\n"`.

Primer: `cat(1:5, 10, sep=" ", "); cat("\n")` ◇

Zaviti oklepaji `{ ... }` označujejo blok. Blok je ukaz, ki zaporedoma izvaja ukaze znotraj zavutih oklepajev in vrne vrednost zadnjega.

B.7.3 Krmilni stavki

Ukaz `if(pogoj) izraz1 else izraz2` ali `ifelse(pogoj, izraz1, izraz2)` vrne `izraz1`, če je pogoj pravilen, sicer pa `izraz2`. Pri prvi konstrukciji lahko del z `else` izpustimo. V tem primeru, če je pogoj napačen, dobimo vrednost `NULL`.

Zanke:

- Ukaz `for(spremenljivka in vektor)` ukaz zaporedoma izvaja ukaz, pri čemer spremenljivki prireja vrednosti v vektorju.
- Ukaz `while(pogoj)` ukaz izvaja ukaz, dokler je pogoj pravilen.
- Ukaz `repeat` ukaz ponavlja izvajanje ukaza.
- Ukaz `break` prekine izvajanje zanke.
- Ukaz `next` prekine izvajanje tekočega cikla zanke.

Ukazi za zanke vrnejo izhod zadnjega izvedenega ukaza.

B.7.4 Nekaj več o funkcijah

Funkcije lahko sprejemajo izbirne argumente (opcije), ki jim predpišemo privzete vrednosti. To storimo v deklaraciji:

```
function(parametri, opcija1=vrednost1, opcija2=vrednost2 ...)
```

Primer: Če deklariramo: `f <- function(x, pristej=0) { x*x + pristej }`, ukaz `f(2)` vrne 4, ukaz `f(2, pristej=3)` pa vrne 7. ◇

Ukaz `return(vrednost)` prekine izvajanje funkcije in vrne predpisano vrednost.

B.8 Knjižnice z dodatnimi funkcijami

Ukaz `library()` vrne seznam knjižnic, ki so na voljo. Ukaz `library(knjižnica)` naloži ustrezno knjižnico.

B.9 Vhod in izhod

B.9.1 Pisanje

Včasih želimo kaj izpisati še kako drugače, kot to stori R. Poleg tega R ne izpisuje avtomatično znotraj zank, ker pač le-te ne vračajo argumentov. Če želimo to, uporabimo ukaz `print`, ki izpiše vrednost, tako kot bi jo sicer izpisal R (z nekaj dodatnimi nastavitvami, kot je npr. `digits`). Izpisovanje v osnovni obliki dosežemo z ukazom `cat`. Sprejme več argumentov, ki jih loči s presledkom. Ločitveni niz lahko spremenimo z nastavitvijo `sep`.

Primeri:

- `cat("bla", "ble", "bli")`
- `cat("bla", "ble", "bli", sep="*")`
- `cat("bla", "ble", "bli", "\n")`
- `cat("bla", "ble", "bli", sep="\n")`
- `cat("bla", "ble", "bli", sep="*\n")` ◇

POZOR! Če ločitveni niz vsebuje znak za novo vrstico, R slednjega avtomatično doda tudi na konec.

Ukaz `paste` deluje podobno kot `cat`, le da vrednost vrne, namesto da bi jo izpisal. Če mu kot argument podamo več nizov, vrne isto, kot bi izpisal `cat`, le da ne upošteva pravila o novi vrstici.

Ukazu `paste` pa lahko posredujemo tudi več vektorjev. V tem primeru vrne vektor, čigar i -to komponento sestavljajo vse i -te komponente podanih vektorjev, ločene z nizom, podanim s `sep`.

Če ukazu `sep` predpišemo nastavev `collapse=niz`, iz ustvarjenega vektorja naredi enoten niz, pri čemer komponente loči s predpisanim nizom.

Primer: ukaz `paste(c(1, 2, 3), c(4, 5, 6, 7), sep=, collapse="\n")` vrne "1 4\n2 5\n3 6\n1 7". ◇

Ukaz `format` je namenjen izpisovanju števil v predpisanem formatu (denimo na določeno število decimalk), deluje pa tudi za nize. Podobno kot paste tudi ukaz `format` ne izpisuje, temveč vrne niz. Ena od mnogih nastavitev je `justify`, ki deluje za poravnavo nizov (možne vrednosti so `"left"`, `"right"`, `"centre"` in `"none"`). Števila pa so vedno poravnana desno.

Primer: ukaz `format(2/3, trim = TRUE, digits = 3, width = 6, justify = "left")` vrne `"0.667"`, ukaz: `format(format(2/3, digits = 3), width = 6, justify = "left")` pa vrne `"0.667"`. ◇

POZOR! Pri obratni poševnici je hrošč – šteje jo kot dva znaka.

Še nekaj ukazov za pisanje:

- `writeLines` izpiše komponente podanega vektorja kot vrstice.
- `write.table` izpiše preglednico.
- `write.csv` izpiše preglednico v formatu csv.

B.9.2 Delo z datotekami

Pregled nad delovnim imenikom (direktorijem):

- Ukaz `getwd()` vrne lokacijo imenika, kjer R bere in piše.
- Ukaz `setwd(imenik)` nastavi lokacijo delovnega.
- Ukaz `list.files()` vrne vsebino delovnega imenika – kot vektor.
Lahko mu predpišemo obilo nastavitev.

Branje in pisanje navadno izvedemo tako, da ustreznemu ukazu predpišemo opcijo `file`. Če jo nastavimo na `niz`, to pomeni enkratno branje oz. pisanje na datoteko z danim imenom. R pa podpira tudi postopno delo z datotekami. V tem primeru datoteko:

- Najprej odpremo z ukazom `file(datoteka, open = način)`.
Tipične vrednosti načina so `"r"` (branje), `"w"` (pisanje) in `"a"` (dodajanje).
- Funkcija `file` vrne kazalec na datoteko, ki ga podajamo pri opciji `file`.
- Nazadnje datoteko zapremo z ukazom `close(kazalec)`.

Primer: ukaz `cat("blabla", file = "blabla.txt")` naredi isto kot zaporedje ukazov:

```
bla <- file("blabla.txt", "w")
cat("blabla", file = bla)
close(bla).
```

◇

B.9.3 Branje

Osnovni ukazi:

- `readChar(datoteka, n)` prebere n znakov z datoteke in vrne prebrani niz.
- `readLines(datoteka)` prebere datoteko in vrne vektor iz njenih vrstic.
- `scan(datoteka, what = tip)` prebere datoteko in njeno vsebino interpretira kot seznam elementov predpisanega tipa, recimo `"logical"`, `"numeric"`, `"character"` ali `"list"`. Rezultat je vektor.

Ukaz `read.table` prebere preglednico. Prva vrstica so glave, elementi v posamezni vrstici so ločeni s presledki, če imajo nizi presledke, jih damo v narekovaje. Posamezne glave morajo biti brez presledkov. Če ima npr. datoteka `Preglednica.txt` vsebino:

```
Ime Prva Druga
"Novak Janez" 5 35
Jurak 3 37
```

ima ukaz: `read.table("Preglednica.txt", header=TRUE)` isti rezultat kot:

```
data.frame(
  Ime=factor(c("Novak Janez", "Jurak")),
  Prva=factor(c(5, 3))
  Druga=factor(c(35,37))
).
```

Ukaz `read.csv` je namenjen branju preglednic v formatu `csv`.

B.9.4 Izvoz grafike

Datoteko najprej odpremo s klicem ustreznega gonilnika, npr. `postscript(datoteka)` ali `pdf(datoteka)`. Nato kličemo ustrezne ukaze, npr. `plot`. Nazadnje datoteko zapremo z ukazom `dev.off()`.

B.10 Verjetnostne porazdelitve

Porazdelitve v R-u računamo z ukazom: `predpona+porazdelitev(vrednost, parametri)`. Parametri so opcije ustrezne funkcije, tj. oblike `parameter=vrednost` (glej spodaj). Možne predpone so:

- ‘d’ za točkasto verjetnost $P(X = x)$ diskretnih oz. gostoto $p_X(x)$ zveznih porazdelitev;
- ‘p’ za kumulativno porazdelitveno funkcijo $P(X \leq x)$;
- ‘q’ za kvantilno funkcijo;
- ‘r’ za simulacijo.

Primer: `dbinom(3, size=10, prob=0.4)` vrne $P(X = 3)$, kjer je slučajna spremenljivka X porazdeljena binomijsko $\text{Bi}(10, 0.4)$. ◇

Varianta za ‘d’:

- `pporazdelitev(x, parametri, log=TRUE)` vrne $\ln P(X = x)$ oz. $\ln p_X(x)$.

Variante za ‘p’:

- `pporazdelitev(x, parametri, lower.tail=FALSE)` vrne $P(X > x)$.
- `pporazdelitev(x, parametri, log.p=TRUE)` vrne $\ln P(X \leq x)$.
- `pporazdelitev(x, parametri, lower.tail=FALSE, log.p=TRUE)` vrne $\ln P(X > x)$.

Najpogostejše diskretne porazdelitve:

porazdelitev	R-ovo ime	parametri
binomska	<code>binom</code>	<code>size, prob</code>
geometrijska	<code>geom</code>	<code>prob</code>
neg. binomska	<code>nbinom</code>	<code>size, prob</code>
Poissonova	<code>pois</code>	<code>lambda</code>
hipergeometrijska	<code>hyper</code>	<code>m, n, k</code>

Najpogostejše zvezne porazdelitve:

porazdelitev	R-ovo ime	parametri
enakomerna	<code>unif</code>	<code>min, max</code>
normalna	<code>norm</code>	<code>mean, sd</code>
eksponentna	<code>exp</code>	<code>rate</code>
Cauchyjeva	<code>cauchy</code>	<code>location, scale</code>
Studentova	<code>t</code>	<code>df, ncp</code>
Fisherjeva	<code>f</code>	<code>df1, df2, ncp</code>

Določeni parametri imajo privzete vrednosti, npr. `mean=0` in `sd=1` pri normalni porazdelitvi.

B.11 Simulacije

Ukaz `rporazdelitev(n)` naredi vektor iz n realizacij slučajne spremenljivke z dano porazdelitvijo. Recimo funkcija `runif` ustreza funkciji `rnd` ali `random` iz mnogih programskih jezikov.

Vzorčenju je namenjen ukaz `sample`:

- `sample(x)` vrne slučajno permutacijo vektorja x .
- `sample(x, velikost)` vrne vzorec brez ponavljanja ustrezne velikosti iz x .
- `sample(x, velikost, replace = TRUE)` vrne vzorec iz x s ponavljanjem.

Dodatna možna nastavitve: `prob` za vektor verjetnosti, recimo:

```
sample(c("a", "e", "i"), 20, replace=TRUE, prob=c(0.2, 0.5, 0.4))
```

Isto naredi ukaz:

```
sample(c("a", "e", "i"), 20, replace=TRUE, prob=c(20, 50, 40)).
```

B.12 Statistično sklepanje

B.12.1 Verjetnost uspeha poskusa/delež v populaciji

Zanima nas verjetnost, da poskus uspe oziroma delež enot v populaciji z dano lastnostjo. To označimo s p . Izvedemo n poskusov, k jih uspe (oziroma vzamemo vzorec n enot in k jih ima iskano lastnost).

Interval zaupanja za p pri stopnji zaupanja β dobimo z ukazom:


```
binom.test(k, n, conf.level =  $\beta$ ).
```

ali tudi:

```
binom.test(c(k, n - k), conf.level =  $\beta$ ).
```

Hipotezo, da je $p = p_0$, testiramo z ukazom:

```
binom.test(k, n, p = p0, alternative = "two.sided" ali "less" ali "greater").
```

Kot rezultat dobimo p -vrednost. Privzeta vrednost za določilo p je 0.5.

Določilo `alternative` pove, kaj je alternativna hipoteza:

Izbira `"two.sided"` ali `"t"` ali izpustitev določila pomeni, da alternativna hipoteza trdi $p \neq p_0$.

Izbira `"less"` ali `"l"` pomeni, da alternativna hipoteza trdi $p < p_0$. Izbira `"greater"` ali `"g"` pa pomeni, da alternativna hipoteza trdi $p > p_0$.

B.12.2 Primerjava verjetnosti dveh poskusov/deležev v dveh populacijah

Testiramo ničelno hipotezo, da sta verjetnosti dveh različnih poskusov oziroma deleža enot z določeno lastnostjo v dveh različnih populacijah enaka. Označimo ju s p_1 in p_2 . Izvedemo n_1 poskusov prve vrste, od katerih jih uspe k_1 , in n_2 poskusov druge vrste, od katerih jih uspe k_2 . Ekvivalentno, vzamemo vzorec n_1 enot iz prve populacije, izmed katerih jih ima k_1 iskano lastnost, in vzorec n_2 enot iz druge populacije, izmed katerih jih ima k_2 našo lastnost. Test ničelne hipoteze $p_1 = p_2$ izvedemo z ukazom:

```
prop.test(c(k1, k2), c(n1,n2), alternative = "two.sided" ali "less" ali "greater").
```

Določilo `alternative` pove, kaj je alternativna hipoteza:

- Izbira `"two.sided"` ali `"t"` ali izpustitev določila pomeni, da alternativna hipoteza trdi $p_1 \neq p_2$.
- Izbira `"less"` ali `"l"` pomeni, da alternativna hipoteza trdi $p_1 < p_2$.
- Izbira `"greater"` ali `"g"` pa pomeni, da alternativna hipoteza trdi $p_1 > p_2$.

B.12.3 Primerjava verjetnosti več poskusov/deležev več populacijah

Testiramo ničelno hipotezo, da so vse verjetnosti več različnih poskusov oziroma vsi deleži enot z določeno lastnostjo v več različnih populacijah enaki. Izvedemo poskuse oziroma iz vsake populacije vzamemo vzorec. Števila vseh poskusov posamezne vrste oziroma velikosti vzorcev iz posamezne populacije naj tvorijo vektor n , števila vseh uspešnih poskusov posamezne vrste oziroma števila enot iz posamezne populacije z določeno lastnostjo pa vektor k . Test izvedemo z ukazom: `prop.test(k, n)`.

B.12.4 Populacijsko povprečje – T -test

Zanima nas povprečje spremenljivke na populaciji, kjer privzamemo, da ima (vsaj približno) normalno porazdelitev. Označimo ga z μ . Vzamemo vzorec, vrednosti spremenljivke na vzorcu naj tvorijo vektor v . Parameter v pa je lahko tudi kontingenčna tabela (recimo dobljena z ukazom `as.table`).

Interval zaupanja za μ pri stopnji zaupanja β dobimo z ukazom:

```
t.test(v, conf.level =  $\beta$ ).
```

Hipotezo, da je $\mu = \mu_0$, testiramo z ukazom:

```
t.test(v, alternative = "two.sided"ali "less"ali "greater",  $\mu = \mu_0$ ).
```

Privzeta vrednost določila μ je 0.

Določilo `alternative` pove, kaj je alternativna hipoteza:

- Izbira `"two.sided"` ali `"t"` ali izpustitev določila pomeni, da alternativna hipoteza trdi $\mu \neq \mu_0$.
- Izbira `"less"` ali `"l"` pomeni, da alternativna hipoteza trdi $\mu < \mu_0$.
- Izbira `"greater"` ali `"g"` pa pomeni, da alternativna hipoteza trdi $\mu > \mu_0$.

B.12.5 Test mediane

Testiramo hipotezo, da ima dana urejenostna (ordinalna) spremenljivka na populaciji mediano μ_0 . To lahko izvedemo tudi na številskih spremenljivkah namesto T -testa, še posebej, če so njihove porazdelitve daleč od normalne. Vzamemo vzorec, vrednosti naše spremenljivke na vzorcu naj tvorijo vektor v . Obravnavali bomo dva možna testa.

Test z znaki. Pri tem testu je pomembno le, ali je dana vrednost večja, manjša ali enaka μ_0 , zato je njegova uporaba smiselna pri čistih urejenostnih spremenljivkah, ko ni določeno, kateri dve vrednosti imata isti odmik od μ_0 navzgor in navzdol. Test se tako prevede na testiranje deleža in ga izvedemo z ukazom:

```
binom.test(sum(sapply(v, function(x) { x >  $\mu_0$  })), sum(sapply(v, function(x) { x !=  $\mu_0$  })),
alternative = "two.sided" ali "less" ali "greater").
```

Wilcoxonov test z rangi. Ta test izvedemo, če je določeno, kateri dve vrednosti imata isti odmik od μ_0 . Test izvedemo z ukazom:

```
wilcox.test(v, alternative = "two.sided"ali "less"ali "greater").
```

B.12.6 Primerjava porazdelitev dveh spremenljivk

1. Dve (približno) normalni spremenljivk na isti populaciji – T -test

Ta primer se prevede na testiranja povprečja ene same spremenljivke, če obe spremenljivki odštejemo: če je količina tvori vektor v , druga pa vektor w , ukažemo:

```
t.test(v - w, alternative = "two.sided"ali "less"ali "greater").
```

Isto pa dosežemo tudi z ukazom:

```
t.test(v, w, paired = TRUE, alternative = "two.sided"ali "less"ali "greater").
```

2. Dve urejenostni spremenljivki na isti populaciji

Tukaj primerjamo spremenljivki, ki sta bodisi urejenostni ali pa ne moremo privzeti, da sta njuni porazdelitvi blizu normalne. Spet lahko uporabimo bodisi test z znaki bodisi test z rangi. Če sta v in w vektorja podatkov, test z znaki izvedemo z ukazom:

```
binom.test(sum(mapply(quote(>)), x, y)),
sum(mapply(quote("!="), x, y)),
alternative = "two.sided"ali "less"ali "greater").
```

test z rangi pa izvedemo z ukazom:

```
wilcox.test(v, w, paired = TRUE, alternative = "two.sided"ali "less"ali "greater").
```

3. Povprečji dveh spremenljivk na dveh populacijah – T -test

Testiramo hipotezo, da imata spremenljivki, definirani na različnih populacijah, enako povprečje, pri čemer privzamemo, da sta vsaj približno normalno porazdeljeni. Z drugimi

besedami, če povprečji označimo z μ_1 in μ_2 , ničelna hipoteza trdi, da je $\mu_1 = \mu_2$. Iz vsake populacije vzamemo vzorec, vrednosti na prvem naj tvorijo vektor v , vrednosti na drugem pa vektor w . Test izvedemo z ukazom:

```
t.test(v, w, alternative = "two.sided"ali "less"ali "greater").
```

Natančneje, ta ukaz izvede heteroskedastični T -test, ki ne privzame enakosti varianc. Če smemo privzeti, da sta varianci obeh spremenljivk enaki, izvedemo homoskedastični test:

```
t.test(v, w, var.equal = TRUE, alternative = "two.sided"ali "less"ali "greater").
```

Pomen določila `alternative` je tak kot ponavadi:

- Izbira `"two.sided"` ali `"t"` ali izpustitev določila pomeni, da alternativna hipoteza trdi $\mu_1 \neq \mu_2$.
- Izbira `"less"` ali `"l"` pomeni, da alternativna hipoteza trdi $\mu_1 < \mu_2$.
- Izbira `"greater"` ali `"g"` pa pomeni, da alternativna hipoteza trdi $\mu_1 > \mu_2$.

4. Dve urejenostni spremenljivki na dveh populacijah

– Wilcoxon-Mann-Whitneyev test

Spet primerjamo spremenljivki, ki sta bodisi urejenostni ali pa ne moremo privzeti, da sta njuni porazdelitvi blizu normalne. Iz vsake populacije vzamemo vzorec, vrednosti na prvem naj tvorijo vektor v , vrednosti na drugem pa vektor w . Hipotezo, da sta spremenljivki enako porazdeljeni, testiramo z ukazom:

```
wilcox.test(v, w, paired = TRUE, alternative = "two.sided"ali "less"ali "greater").
```

Določilo `alternative` pove, kaj je alternativna hipoteza:

- izbira `"two.sided"` ali `"t"` ali izpustitev določila pomeni, da alternativna hipoteza trdi, da sta porazdelitvi različni.
- Izbira `"less"` ali `"l"` pomeni, da alternativna hipoteza trdi, da je prva spremenljivka stohastično večja od druge.
- Izbira `"greater"` ali `"g"` pa pomeni, da alternativna hipoteza trdi, da je prva spremenljivka stohastično manjša od druge.

Če je eden od vzorcev velik, je računanje zahtevno. Lahko ga poenostavimo z določilom `exact = FALSE` – v tem primeru bo R računal približek.

5. Povprečji spremenljivk na več populacijah – analiza variance (ANOVA) z enojno klasifikacijo

Testiramo hipotezo, da imajo spremenljivke, za katere privzamemo, da so porazdeljene normalno, enaka povprečja. Podatke uvrstimo v dva vektorja, denimo v in s . V vektorju v so vse vrednosti, v vektorju s pa skupine (ta vektor torej pove, iz katere populacije je dani podatek). Elementi vektorja s morajo biti nizi, ne števila! Test izvedemo z ukazom:

```
anova(lm(vr ~ sk, data.frame(vr = v, sk = s))).
```

p -vrednost odčitamo v rubriki $Pr(< F)$ na desni strani R-ovega poročila.

6. Urejenostni spremenljivki na več populacijah – Kruskal-Wallisov test

Ravnamo podobno kot pri analizi variance: podatke uvrstimo v dva vektorja, denimo v in s , pri čemer so v vektorju v vrednosti, v vektorju s pa skupine. Toda pozor: v nasprotju z analizo variance morajo biti tu v vektorju s števila, ne nizi!

Test izvedemo z ukazom: `kruskal.test(v, s)`.

B.12.7 Koreliranost

Zanima nas koreliranost dveh spremenljivk na isti populaciji. Vzamemo vzorec, vrednosti prve spremenljivke naj tvorijo vektor v , vrednosti druga pa vektor w .

Interval zaupanja za korelacijski koeficient pri stopnji zaupanja β poiščemo z ukazom:

```
cor.test(v, w, method = "pearson"ali "kendall"ali "spearman", conf.level = beta).
```

Določilo `method` pove, kateri korelacijski koeficient nas zanima.

- Izbira `"pearson"` ali `"p"` ali opustitev določila pomeni običajni Pearsonov korelacijski koeficient, ki je primeren za številske spremenljivke, porazdeljene (približno) normalno.
- Izbira `"spearman"` ali `"s"` pomeni Spearmanov,
- izbira `"kendall"` ali `"k"` pa Kendallov korelacijski koeficient;

Slednja testa sta primerna za urejenostne spremenljivke. Spearmanov korelacijski koeficient (ρ) je lažji za računanje, o Kendallovem koeficientu (τ) pa dosti statistikov meni, da je verodostojnejši.

Hipotezo, da sta spremenljivki nekorelirani (oziroma neodvisni) pa testiramo z ukazom:

```
cor.test(v, w, method = "pearson"ali "kendall"ali "spearman", alternative = "two.sided" ali  
"less" ali "greater").
```

Določilo `alternative` pove, kaj je alternativna hipoteza:

- Izbira `"two.sided"` ali `"t"` ali izpustitev določila pomeni, da alternativna hipoteza trdi, da sta spremenljivki odvisni (oziroma korelirani).
- Izbira `"less"` ali `"l"` pomeni, da alternativna hipoteza trdi, da sta spremenljivki pozitivno asociirani (oziroma korelirani).
- Izbira `"greater"` ali `"g"` pa pomeni, da alternativna hipoteza trdi, da sta spremenljivki negativno asociirani (oziroma korelirani).

Računanje p -vrednosti za Kendallov korelacijski koeficient je zahtevno in R sam izbere, kdaj bo računal natančno vrednost in kdaj približek. To lahko spremenimo z določilom `exact = TRUE` ali `exact = FALSE`.

Stvarno kazalo

- P*-vrednost, [190](#)
- σ -algebra, [20](#)
- časovna vrsta, [249](#)
- šibki zakon velikih števil, [99](#)

- algebra dogodkov, [20](#)
- analiza variance, [214](#)
- asimetrija, [96](#)

- Bayes, T. (1702-1761), [32](#)
- Bayesov obrazec, [30](#)
- Bernoulli, J. (1654–1705), [43](#)
- Bernoullijev obrazec, [34](#)
- Bernoullijevo zaporedje, [33](#)
- binomski simbol, [281](#)
- bivariantna analiza, [227](#)
- Borel, E., [85](#)
- Borelove množice, [78](#)

- Cauchy, A.L. (1789 – 1857), [292](#)
- cenilka
 - asimptotično nepristranska, [160](#)
 - dosledna, [159](#)
 - nepristranska, [160](#)
 - točkovna, [159](#)
- centil, [134](#)
- Centralni limitni izrek, [101](#), [148](#)
- centralni limitni zakon, [102](#)

- De Moivre, A. (1667–1754), [44](#)

- De Moivrov točkovni obrazec, [41](#)
- disperzija, [90](#)
- dogodek, [12](#)
 - elementaren, [14](#)
 - enak, [13](#)
 - gotov, [12](#)
 - način, [13](#)
 - nasproten, [14](#)
 - nemogoč, [13](#)
 - neodvisna, [28](#)
 - nezdružljiva, [14](#)
 - osnoven, [14](#)
 - produkt, [13](#)
 - sestavljen, [14](#)
 - slučajen, [13](#)
 - vsota, [13](#)
- domneva
 - enostavna, [186](#)
 - formalen postopek za preverjanje, [199](#)
 - neparametrična, [186](#)
 - sestavljena, [186](#)
- domneva (statistična), [186](#)
- enakomerna zvezna, [52](#)
- Erdöseva probabilistična metoda, [111](#)
- Erdős, P. (1913 – 1996), [112](#)
- Euler, L. (1707–1783), [291](#)

- frekvenca razreda, [130](#)
 - relativna, [130](#)

- funkcija, 274
- bijektivna, 275
 - Gama, 56
 - injektivna, 274
 - karakteristična, 97
 - napake, 54
 - pogojna verjetnostna, 83
 - slučajnega vektorja, 80
 - surjektivna, 274
- Gauss, J.C.F. (1777–1855), 60
- gostota verjetnosti, 51
- histogram, 130
- interval zaupanja, 174
- teoretična interpretacija, 176
- Jacobijeva determinanta, 82
- koeficient
- asimetrije, 138
 - Cramerjev, 230
 - Jaccardov, 231
 - korelacijski, 92
 - parcialne korelacije, 237
 - Pearsonov, 230, 231
 - Sokal Michenerjev, 231
 - sploščenosti, 138
 - Yulov, 231
- Kolmogorov, A. (1903-1987), 21
- kombinacije, 281
- kontingenčna tabela, 27
- konvolucija, 81
- koreliranost, 90
- kovariančna matrika, 95
- kovarianca, 91
- kreпки zakon velikih števil, 99
- kumulativa, 129
- kvantil, 96
- kvartil, 134
- kvartilni razmik, 96
- Lagrangeova indentiteta, 289
- Laplace, S.P. (1749–1827), 43
- Laplaceov intervalski obrazec, 53
- Latinski kvadrat, 259
- list, 129
- matematično upanje, 87
- pogojno, 93
 - slučajnega vektorja, 94
- mediana, 96, 133, 134
- moč statističnega testa, 189
- modus, 133
- moment, 95
- centralni, 96
 - vzorčni začetni, 160
 - začetni, 96
- napaka
- vrste, 189
 - vrste, 189
- neenakost
- Čebiševa, 99
 - Cauchyjeva, 288
- ogiva, 130
- osnovni izrek statistike, 144
- parameter, 123
- Pascal, B. (1623-1662), 292
- permutacija, 275
- ciklična (cikel), 276

- liha, 277
- s ponavljanjem, 277
- soda, 277
- transpozicija, 276
- poddogodek, 13
- pogojna gostota, 84
- pogojna porazdelitvena funkcija, 83
- Poisson, S. (1781–1840), 59
- pojasnjena varianca (ANOVA), 246
- popoln sistem dogodkov, 14
- populacija, 4, 122
 - povprečje, 133
- porazdelitev
 - binomska, 47
 - eksponentna, 55
 - Fisherjeva, 157
 - frekvenčna, 128
 - Gama, 56
 - geometrijska, 49
 - hi-kvadrat, 57
 - hipergeometrijska, 50
 - negativna binomska, 49
 - normalna, 52
 - Pascalova, 49
 - Poissonova, 48
 - polinomska, 66
 - standardizirana normalna, 54
 - Studentova, 156
- porazdelitvena funkcija, 45
 - robna, 64
- porazdelitveni zakon, 45
- poskus, 12
- prikaz
 - krožni, 127
 - Q-Q, 134
 - stolpčni, 127
- pristranost, 160
- Ramsey, F. P. (1903–1930), 113
- Ramseyjeva teorija, 108
- ranžirana vrsta, 128
- rang, 128
- razbitje, 29
- razpon, 130, 134
- razred, 128
- regresija, 94, 238
- regresijska premica, 242
 - druga, 243
 - prva, 243
- regresijski koeficient, 239
- relativna frekvenca, 15
- skalarni produkt, 288
- skatla z brki, 134
- skoraj gotova konvergenca, 98
- slučajna odstopanja, 214
- slučajna spremenljivka, 45
 - diskretna, 46
 - zvezna, 46, 51
- slučajni vektor, 63
- slučajni vzorec
 - enostavni, 142
- sploščenost, 96
- spremenljivka, 122
 - številsko, 125
 - absolutna, 126
 - imenska, 126
 - opisna, 125
 - razmernostna, 126

- razmična, 126
- urejenostna, 126
- sredina
 - aritmetična, 287
 - geometrična, 287
 - harmonična, 287
 - kvadratna, 287
 - potenčna, 287
- sredinska mera, 145
 - geometrijska, 145
 - povprečna vrednost, 145
- standardizacija, 91, 139
- standardna deviacija, 91
- standardni odklon, 91
- statistična enota, 122
- statistika, 122, 123
 - analitična, 4
 - inferenčna, 122
 - opisna, 4, 122, 125
- steblo, 129
- Stirlingov obrazec, 284
- stopnja tveganja, 189
- stopnja zaupanja, 189
- stopnja značilnosti, 189

- teorija vzorčenja, 142
- trend, 255

- urejeno zaporedje, 128

- varianca, 90
- vektorski produkt, 289
- verjetnost, 20
 - definicija, 11
 - klasična definicija, 16
 - pogojna, 26
 - statistična definicija, 16
- verjetnostna funkcija, 65
- verjetnostna konvergenca, 98
- verjetnostna tabela, 65
- verjetnostni prostor, 20
- vzorčna disperzija, 145
 - popravljen, 145
- vzorčna statistika, 150
- vzorčne ocene, 145
- vzorčni odklon, 145
- vzorčni razmah, 145
- vzorčno povprečje, 145
- vzorec, 4, 122, 142
 - mediana, 145
 - modus, 145
 - povprečje, 133

