

VERJETNOSTNI RAČUN IN STATISTIKA

Aleksandar Jurišić, FRI

26. november 2018

Štetje	Možnosti	Vzrok in posledica	Oblika podatkov
Sredine, momenti	Zvonasta krivulja	Večdim. 2D, 3D	Zveza
Neizbež. (CLI)	Opis. stat. histogram	Vzorec, cenilka	Ocena parametra
Preveri domnevo	Napoved	Model	Uporaba

Kombinatorika	Verjetnost	Pogojna verjetnost	Slučajne spremen.
Sredine, momenti	Bin(n, p) $N(\mu, \sigma)$	Večrazs. porazd.	$Cov(X, Y)$
CLI	Opisna statistika	Vzorci, cenilke	Intervali zaupanja
Prever. domnev	Regresija	Pregled porazd.	Uporaba

Seznam poglavij

1. Uvod	1
I. VERJETNOST	7
2. Poskusi, dogodki in definicija verjetnosti	9
3. Pogojna verjetnost	25
4. Bernoullijevo zaporedje neodvisnih poskusov	37
5. Slučajne spremenljivke in porazdelitve	45
6. Slučajni vektorji	69
7. Funkcije slučajnih spremenljivk in vektorjev	83
8. Momenti in kovarianca	91
9. Karakteristične funkcije in limitni izreki	99
II. STATISTIKA	105
10. Opisna statistika	111
11. Vzorčenje	123
12. Intervali zaupanja	147
13. Preverjanje statističnih domnev	155
14. Bivariatna analiza in regresija	179
15. Časovne vrste in trendi	193
III. KAM NAPREJ	201
16. Nekaj primerov uporabe verjetnosti	205
17. Uporaba statistike	215
A. Matematične osnove (ponovitev)	221
B. Vadnica	241
C. Zgodovina in biografije iz verjetnosti in statistike	263
D. Program R	281
Stvarno kazalo	??

Kazalo

1	Uvod	1
I	VERJETNOST	5
2	Poskusi, dogodki in definicija verjetnosti	9
2.1	Poskusi in dogodki	10
2.2	Računanje z dogodki	11
2.3	Definicija verjetnosti	13
2.4	Osnovne lastnosti verjetnosti	22
2.5	Aksiomi Kolmogorova*	24
3	Pogojna verjetnost	25
3.1	Intriga (po Kvarkadabri)	25
3.2	Definicija pogojne verjetnosti	27
3.3	Dvostopenjski poskusi in izrek o popolni verjetnosti	31
4	Bernoullijevo zaporedje neodvisnih poskusov	37
4.1	Računanje verjetnosti $P_n(k)$	39
4.2	Usojena krivulja v Pascalovem trikotniku	40
5	Slučajne spremenljivke in porazdelitve	45
5.1	Slučajne spremenljivke	45
5.2	Diskretne slučajne spremenljivke	46
5.2.1	Enakomerna diskretna porazdelitev – $U(n)$	47
5.2.2	Binomska porazdelitev – $B(n, p)$	48
5.2.3	Poissonova porazdelitev – $P(\lambda)$	50
5.2.4	Negativna binomska oziroma Pascalova porazdelitev – $P(m, p)$	52
5.2.5	Hipergeometrijska porazdelitev $H(n; M, N)$	54
5.3	Ponovitev: integrali	55
5.4	Zvezne slučajne spremenljivke	56
5.4.1	Enakomerna zvezna porazdelitev – $U(a, b)$	57
5.4.2	Normalna ali Gaussova porazdelitev – $N(\mu, \sigma)$	57
5.4.3	EkspONENTNA porazdelitev – $\text{Exp}(\lambda)$	61
5.4.4	Gama porazdelitev – $\Gamma(k, \lambda)$	62
5.4.5	Hi-kvadrat porazdelitev – $\chi^2(n)$	63
5.4.6	Cauchyeva porazdelitev	65

6 Slučajni vektorji	69
6.1 Diskretne večrazsežne porazdelitve – polinomska	73
6.2 Ponovitev: dvojni integral	73
6.3 Zvezne večrazsežne porazdelitve	76
6.4 Neodvisnost slučajnih spremenljivk	80
7 Funkcije slučajnih spremenljivk in vektorjev	83
7.1 Funkcije slučajnih spremenljivk	83
7.2 Funkcije in neodvisnost	86
7.3 Funkcije slučajnih vektorjev	88
7.4 Pogojne porazdelitve	89
8 Momenti in kovarianca	91
8.1 Pričakovana vrednost	91
8.2 Disperzija	93
8.3 Standardizirane spremenljivke	94
8.4 Kovarianca	95
8.5 Kovariančna matrika	96
8.6 Višji momenti	97
9 Karakteristične funkcije in limitni izreki	99
9.1 Karakteristična funkcija	99
9.2 Limitni izreki	100
9.3 Centralni limitni izrek (CLI)	102
II STATISTIKA	103
10 Opisna statistika	111
10.1 Vrste slučajnih spremenljivk oziroma podatkov	111
10.2 Grafična predstavitev kvantitativnih podatkov	113
10.3 Mere za lokacijo in razpršenost	117
10.4 Standardizacija	122
11 Vzorčenje	123
11.1 Osnovni izrek statistike	124
11.2 Vzorčne ocene	125
11.3 Porazdelitve vzorčnih povprečij	126
11.4 Vzorčna statistika	129
11.4.1 (A) Vzorčno povprečje	130
11.4.2 (B) Vzorčna disperzija	131
11.5 Nove porazdelitve	133
11.5.1 Studentova t -porazdelitev	134
11.5.2 Fisherjeva porazdelitev	135
11.6 Cenilke	136
11.6.1 Osnovni pojmi	136
11.6.2 Rao-Cramérjeva ocena*	138
11.6.3 Učinkovitost cenilk*	139

11.6.4 Metoda momentov	141
11.6.5 Metoda največjega verjetja	142
11.7 Vzorčna statistika (nadaljevanje)	143
11.7.1 (C) Vzorčna povprečja	143
11.7.2 (D) Vzorčni deleži	144
11.7.3 (E) Razlika vzorčnih povprečij	145
11.7.4 (F) Razlika vzorčnih deležev	146
12 Intervali zaupanja	147
12.1 Pomen stopnje tveganja	148
12.2 Intervalsko ocenjevanje parametrov	148
12.2.1 Pričakovana vrednost μ z znanim odklonom σ	150
12.2.2 Velik vzorec za pričakovano vrednost μ	150
12.2.3 Majhen vzorec za pričakovano vrednost μ	150
12.2.4 Razlika pričakovanih vrednosti $\mu_1 - \mu_2$ z znanima odklonoma σ_1 in σ_2	151
12.2.5 Velika vzorca za razliko pričakovanih vrednosti $\mu_1 - \mu_2$	151
12.2.6 Majhna vzorca za razliko pričakovanih vrednosti $\mu_1 - \mu_2$ z neznanima $\sigma_1 = \sigma_2$	151
12.2.7 Majhna vzorca za razliko pričakovanih vrednosti $\mu_1 - \mu_2$	151
12.2.8 Velik vzorec za razliko $\mu_d = \mu_1 - \mu_2$ ujemajočih se parov	152
12.2.9 Majhen vzorec za razliko $\mu_d = \mu_1 - \mu_2$ ujemajočih se parov	152
12.2.10 Delež populacije π z znanim odklonom σ	153
12.2.11 Velik vzorec za delež populacije π	153
12.2.12 Razlika deležev $\pi_1 - \pi_2$ z znanim odklonom $\sigma_{\pi_1 - \pi_2}$	153
12.2.13 Velik vzorec za razliko deležev $\pi_1 - \pi_2$	153
12.2.14 Majhen vzorec za varianco σ^2	154
12.2.15 Majhen vzorec za kvocient varianc σ_1^2/σ_2^2	154
12.3 Izbira velikosti vzorca	154
13 Preverjanje statističnih domnev	155
13.1 Ilustrativni primeri	156
13.2 Alternativna domneva in definicije	158
13.3 Predznačni test	162
13.4 Wilcoxonov predznačni-rang test	163
13.5 Formalen postopek za preverjanje domnev	165
13.6 Domneve za pričakovano vrednost $H_0 : \mu = \mu_0$	166
13.6.1 Znan odklon σ	166
13.6.2 Neznan odklon σ in velik vzorec	167
13.6.3 Neznan odklon σ , normalna populacija in majhen vzorec	167
13.7 Domneve za razliko povprečij $H_0 : \mu_1 - \mu_2 = D_0$	169
13.7.1 Znana odklona σ_1 in σ_2	169
13.7.2 Neznana odklona σ_1 in/ali σ_2 , $n_1 \geq 30$ in/ali $n_2 \geq 30$	170
13.7.3 Neznana σ_1 in/ali σ_2 , norm. pop., $\sigma_1 = \sigma_2$, $n_1 < 30$ ali $n_2 < 30$	170
13.7.4 Neznana σ_1 in/ali σ_2 , norm. pop., $\sigma_1 \neq \sigma_2$, $n_1 < 30$ ali $n_2 < 30$	170
13.8 Domneve za povprečje razlik $H_0 : \mu_d = D_0$	171
13.8.1 Velik vzorec	171
13.8.2 Normalna populacija razlik in majhen vzorec	171

13.9 Domneve za delež $H_0 : \pi = \pi_0$	172
13.9.1 Velik vzorec	172
13.10 Razlika deležev dveh populacij $H_0 : \pi_1 - \pi_2 = D_0$	172
13.10.1 Velik vzorec in $D_0 = 0$	173
13.10.2 Velik vzorec in $D_0 \neq 0$	174
13.11 Analiza variance	174
13.11.1 Domneva o varianci $\sigma^2 = \sigma_0^2$	175
13.11.2 Domneva o kvocientu varianc $\sigma_1^2/\sigma_2^2 = 1$	176
13.12 Domneva o porazdelitvi spremenljivke	176
13.12.1 Domneva o enakomerni porazdelitvi	176
13.12.2 Domneva o normalni porazdelitvi	178
14 Bivariatna analiza in regresija	179
14.1 Povezanost dveh imenskih (nominalnih) spremenljivk	179
14.2 Koeficienti asociacije	181
14.3 Povezanost dveh ordinalnih spremenljivk	182
14.4 Povezanost dveh številskeih spremenljivk	184
14.5 Parcialna korelacija	185
14.6 Regresijska analiza	186
14.7 Linearni model	188
14.7.1 Statistično sklepanje o regresijskem koeficientu	190
14.7.2 Pojasnjena varianca (ang. ANOVA)	191
15 Časovne vrste in trendi	193
15.1 Primerljivost členov v časovni vrsti	194
15.2 Grafični prikaz časovne vrste	194
15.3 Indeksi	195
15.4 Sestavine dinamike v časovnih vrstah	196
III KAM NAPREJ	199
16 Nekaj primerov uporabe verjetnosti	205
16.1 Zamenjalna šifra	205
16.2 Kakšno naključje!!! Mar res?	208
16.3 Ramseyjeva teorija	209
16.4 Teorija kodiranja	213
17 Uporaba statistike	215
17.1 Načrtovanje eksperimentov	215
A MATEMATIČNE OSNOVE (ponovitev)	221
A.1 Računala nove dobe	221
A.2 Funkcije/preslikave	224
A.3 Permutacije	224
A.4 Kombinacije	230
A.5 Vrsta za e	233
A.6 Stirlingov obrazec	234

A.7	Normalna krivulja v prostoru	235
A.8	Sredine nenegativnih števil a_1, \dots, a_n	237
A.9	Cauchyjeva neenakost	238
B	Vadnica	241
B.1	Vaje za uvod	241
B.2	Poskusi, dogodki in definicija verjetnosti	246
B.3	Pogojna verjetnost	248
B.4	Bernoullijevo zaporedje neodvisnih poskusov	250
B.5	Slučajne spremenljivke in porazdelitve	250
B.6	Slučajni vektorji	250
B.7	Funkcije slučajnih spremenljivke in vektorjev	250
B.8	Momenti in kovarianca	250
B.9	Karakteristične funkcije in limitni izreki	251
B.10	Opisna statistika	251
B.11	Vzorčenje	253
B.12	Cenilke	253
B.13	Intervali zaupanja	253
B.14	Preverjanje statističnih domnev	253
B.15	Bivariatna analiza in regresija	261
B.16	Časovne vrste in trendi	261
C	Zgodovina in biografije	263
C.1	Pierre de Fermat (1601-1665)	264
C.2	Blaise Pascal (1623-1662)	265
C.3	Matematični Bernoulli	266
C.4	Abraham de Moivre (1667-1754)	267
C.5	Thomas Bayes (1702-1761)	267
C.6	Leonhard Euler (1707-1783)	268
C.7	Joseph Louis Lagrange (1736-1813)	270
C.8	Pierre-Simon Laplace (1749–1827)	270
C.9	Johann Carl Friedrich Gauss (1777–1855)	272
C.10	Simeon Poisson (1781–1840)	274
C.11	Augustin Louis Cauchy (1789-1857)	274
C.12	Arthur Cayley (1821-1895)	275
C.13	Emile Borel (1871–1956)	275
C.14	William Sealy Gosset (1876-1937)	276
C.15	Sir Ronald A. Fisher (1890-1962)	277
C.16	Egon Sharpe Pearson (1895-1980)	277
C.17	Frank Plumpton Ramsey (1903-1930)	277
C.18	Andrei Kolmogorov (1903-1987)	278
C.19	Paul Erdős (1913-1996)	278
D	PROGRAM R (Martin Raič)	281
D.1	Izvajanje programa	281
D.2	Aritmetika	282
D.3	Najosnovnejše o spremenljivkah	282

D.4	Uporabnikove funkcije	282
D.5	Numerično računanje	283
D.6	Podatkovne strukture	283
D.6.1	Vektorji	283
D.6.2	Matrike	284
D.6.3	Tabele	286
D.6.4	Vektorji, matrike in tabele z označenimi indeksi	286
D.6.5	Zapisi	287
D.6.6	Kontingenčne tabele in vektorji s predpisanimi vrednostmi	287
D.6.7	Preglednice	287
D.7	Osnove programiranja	288
D.7.1	Izvajanje programov, shranjenih v datotekah	288
D.7.2	Najosnovnejši programski ukazi	289
D.7.3	Krmilni stavki	289
D.7.4	Nekaj več o funkcijah	290
D.8	Knjižnice z dodatnimi funkcijami	290
D.9	Vhod in izhod	290
D.9.1	Pisanje	290
D.9.2	Delo z datotekami	291
D.9.3	Branje	292
D.9.4	Izvoz grafike	292
D.10	Verjetnostne porazdelitve	293
D.11	Simulacije	294
D.12	Statistično sklepanje	294
D.12.1	Verjetnost uspeha poskusa/delež v populaciji	294
D.12.2	Primerjava verjetnosti dveh poskusov/deležev v dveh populacijah	295
D.12.3	Primerjava verjetnosti več poskusov/deležev več populacijah	295
D.12.4	Populacijsko povprečje – T -test	296
D.12.5	Test mediane	296
D.12.6	Primerjava porazdelitev dveh spremenljivk	297
D.12.7	Koreliranost	298

Predgovor

*Smer v katero izobrazba zapelje čoveka,
določi njegovo bodočnost in življenje.*

PLATON

Zapiski trenutno še vedno predstavljajo nekakšno razširjeno verzijo prosojnic za predmeta

- *Osnove verjetnosti in statistike* ter
- *Verjetnost in statistika*.

Študentke in študentje ste od samega začetka predavanj spodbujani, da postavljate vprašanja in aktivno sodelujete pri predmetu.

*Še vedno mi služi šest vdnih služabnikov.
Od njih sem se naučil vsega, kar danes vem.
Njihova imena so: Kaj, Zakaj, Kdaj, Kako, Kje in Kdo.*

RLJDYARD KIPLING, *The Elephant Child*

Edino neumno vprašanje je tisto, ki ga ne zastavite.
PAUL MACCREADY, *izumitelj*

Priznati moram, da redki mojo spodbudo vzamete resno. Morda zaradi želje po anonimnosti (predvsem v odnosu do sošolk in sošolcev, saj pri občasnih spisih – Kaj lahko pričakujem od predmeta – pogosto izvem marsikaj in sem nad vašo odkritostjo celo presenečen). Ali pa ste nesodelovanja navajeni že od prej.¹ Vendar sem si zadal nalogu to spremeniti.²³ Če bo po sreči, si bomo pri tem pomagali s tehnologijo (aplikacijo eQuiz, ki jo že šesto leto razvijajo študenti na FRI, ob moji skromni pomoči) in boste vsako uro vsaj enkrat odgovorili na kakšno moje vprašanje.⁴

Naš cilj ni slediti vsakega posameznika, saj imamo vsako leto blizu 500 študentov in bito znalo biti prenaporno (pa tudi smiselno ni). Morda boste pomislili, da so ocene vendarle individualne, a so, če dobro premislimo, predvsem (nekoliko zapoznela) informacija vam samim. V resnici bi radi razumeli, kaj se dogaja z večino, pa naj si bo to malo čez 50%, (2/3)-sko ali celo 95%. Hkrati upamo, da bodo rezultati zanimali tudi vas same, posebej vaše mesto glede na večino, še bolje, glede na izbrani krog vaših prijateljev. Tako se hočeš nočeš srečamo z osnovno potrebo statistike - razumeti podatke (kadar so na voljo - in danes bi moralno biti tako), da gremo lahko pravočasno v akcijo. Z malo sreče, boste kmalu uporabljali novo verzijo aplikacije in boste prek nje dobili navodila, kako izboljšati svoje znanje.

¹ *Svet, v katerem bodo živeli naši otroci, se spreminja štirikrat hitreje kot naše šole.* Dr. WILLARD DAGGETT, direktor Mednarodnega centra za vodenje in izobraževanje, v nagovoru šole v Coloradu l. 1992.

² *Tradicionalni izobraževalni sistem je preživet.* RICHARD L. MEASELLE v Morton Egol, *Transforming Education: Breakthrough Quality at Lower Cost*, Arthur Andersen

³ *Spremeniti je treba način učenja v svetu.* Izjava poslanstva, Izjava je rezultat petdnevnega "umika" članov fakultete Lansdowne, ki je bila maja 1996 v Soto Grande v Španiji.

⁴ *Današnja generacija šolarjev je prva, ki jo res obdajajo digitalni mediji.* DON TAPSCOTT, *Growing Up Digital*, McGraw Hill, New York

Resnična moč računalnikov se bo pokazala, ko bodo postali učno orodje v rokah učencev. PAT NOLAN, v avtorjem intervjuju, Univerza Massey, Palmerston North, Nova Zelandija

Ste se že kdaj vprašali, pod kakšnimi pogoji športnik, glasbenik ali raziskovalec doseže najboljše rezultate. Za začetek potrebuje talent in veselje do svoje panoge. Vendar to še zdaleč ni dovolj. Obstajati mora močna želja za preseganje samega sebe in pripravljenost na trdo delo, vztrajnost ter še kaj. Ste že slišali, da je za vsako mojstrstvo potreben nek minimum vloženega dela? Za posamezno področje to lahko znaša tudi do 10.000 ur in več. Pomembna je tudi sredina, okolje, ki mora biti naravnano spodbujevalno. Sledi vloga trenerja/mentorja/učitelja, ki osebo nenehno spodbuja in usmerja. Postavljati mu mora vmesne izzive, na poti do končnega cilja.

Nekoč so bila najpomembnejša tekmovanja enkrat letno, ali celo na štiri leta (olimpijada), danes pa je tekmovanj vse več. Enako je tudi v šolstvu. Namesto končnega preverjanja ob diplomi, so bili najprej organizirani letni izpiti, sedaj pa že semestralni. Sprotni študij smo poudarili z domačimi nalogami in še posebej s kolokviji – najprej s štirimi letno, sedaj pa imamo lahko že 5 kolokvijev v enem semestru. Najbolj idealno bi bilo, če bi dobivali izzive dnevno. Enako je s primerjavo. Tudi ta bi lahko bila nenehna - tako kot recimo merjenje časa na vsakem treningu, izmerjena višina pri vsakem skoku itd. Seveda se razume, da boksar ne more imeti borbe vsak dan.

Cilj naše spletne aplikacije eQuiz je, da bi študentom omogočila doseganje čim boljših rezultatov in doseganje čim višjih ciljev. Pri semestralnih predmetih VIS (UNI-FRI) in OVS (VŠS-FRI) potekata pilotska projekta.⁵ V ta namen želimo izbrati za vsako posamezno skupino približno 1000 preverjenih nalog. Res pa je tudi, da bomo skušali naloge nenehno izboljševati. Reševalci (študentje) jih bodo med reševanjem ocenjevali (zahtevnost, duhovitost/izvirnost, poučnost). Tudi aplikacija sama bo z ratingom ugotavljala tako znanje študentov, kot nivo/zahtevnost nalog). Hkrati bi radi povečali količino povratne informacije, ki jo ob reševanju dobi študent. Ravno pomanjkanje teh informacij predstavlja pri veliki količini študentov veliko oviro, katero jim bomo s skupnimi močmi, pomagali prestopiti.

Različne vrste razuma - inteligence

Osebna in profesionalna uporaba: 2. Logično-matematična inteligence.

Ponavadi jo najdemo pri matematikih, znanstvenikih, inženirjih, kriminalistih, pravnikih in računovodjih.⁶

⁵Se en pilot se izvaja na MF, pri predmetu Interna medicina.

⁶ Znana osebnost: Marian Diamond, profesorica nevroanatomije na univerzi Berkley v Kaliforniji.

Najverjetnejše lastnosti:

- Uživa v abstraktnem razmišljanju
- Ljubi natančnost
- Rad šteje
- Rad ima organiziranost
- Uporablja logične strukture
- Rad uporablja računalnike
- Uživa v reševanju problemov
- Uživa v poskušanju in logičnem sklepanju
- Prizadeva si za urejeno zapisovanje podrobnosti

*Tisti, ki ima rad prakso brez teorije,
je kot mornar, ki zapljuje z ladjo
brez krmila in kompasa
ter nikoli ne ve, kje lahko nasede.*
LEONARDO DA VINCI (1452-1519)

*Matematika ni samo realna,
pač pa je edina realnost. Očitno je,
da se celotno vesolje sestoji iz materije.
Le-ta je sestavljena iz delcev. Sestavljena je
iz elektronov, neutronov in protonov.
Torej celotno vesolje je sestavljeno iz delcev.
In iz česa so potem sestavljeni delci?
Niso sestavljeni iz ničesar.
Vse kar lahko nekdo pove o realnosti elektrona,
je da navede njegove matematične lastnosti.
V tem smislu se je materija popolnoma razkrojila
in vse kar je ostalo je matematična struktura.*
MARTIN GARDNER

*Ker možgani ne morejo biti pozorni na vse, ...
si nezanimivih, čustveno praznih lekcij
enostavno ne zapomnite.*
LAUNA ELLISON^a

*Učimo se: 10% z branjem
20% s poslušanjem
30% z opazovanjem
50% s poslušanjem in opazovanjem
70% z govorjenjem
90% z govorjenjem in delom*
VERNON A. MAGNESEN^b

^aWhat Does The Brain Have To Do With Learning?
članek v Holistic Education Review, jeseni 1991

^bcitat v knjigi Bobbi DePorter, Mark Reardon, Sarah Singer, Nourie: Quantum Teaching, Allyn and Bacon, Needham Heights

Pomoč pri učenju:

- Spodbujajte k reševanju problemov
- Igrajte se matematične igrice
- Analizirajte in interpretirajte podatke
- Uporablajte sklepanje
- Spodbujajte samostojnost
- Spodbujajte lastno raziskovanje
- Uporablajte napovedovanje prihodnosti
- Uporablajte organizacijske metode in matematična orodja za poučevanje drugih področij
- Vsaka stvar mora imeti svoje mesto
- Dovolite postopno reševanje problema
- Uporablajte deduktivno razmišljanje
- Intenzivno uporablajte računalnike

*Otroci se najbolje učijo takrat,
kadar jim samo pomagamo, da se sami
dokopljejo do iskanih zakonitosti.*

PETER KLINE, *The everyday genius*

*Učenje je čudovita igra,
ki je povrhu še zelo zabavna.
Vsi otroci se rodijo s tem prepričanjem
in ga ohranijo toliko časa,
dokler jih ne prepričamo,
da je učenje naporno in neprijetno delo.*

*So otroci, ki se jih tega ne da naučiti.
Zanje je učenje do konca življenja zabavno in
njim je učenje edina res pozornosti vredna igra.
Za take ljudi imamo prav posebno besedo.*

Rečemo jim geniji.
GLENN DOMAN *Teach Your Baby Math*

*Razum ni čaša, ki čaka, da bi jo napolnili,
je ogenj, ki čaka na iskro.*
PLUTARH, grški biograf,
ki je misel zapisal pred skoraj 3.000 leti.

*Predlog sistema ocenjevanja 21. stoletja
50% samoocenjevanja
30% ocenjevanja vrstnikov
20% učiteljevo (šefovo) ocenjevanje*
JEANNETTE VOS povzeto po njenih seminarjih

Šola bi morala biti najboljša zabava v mestu.
PETER KLINE, *The Everyday Genius,*
Great Ocean Publishers Inc. Arlington.

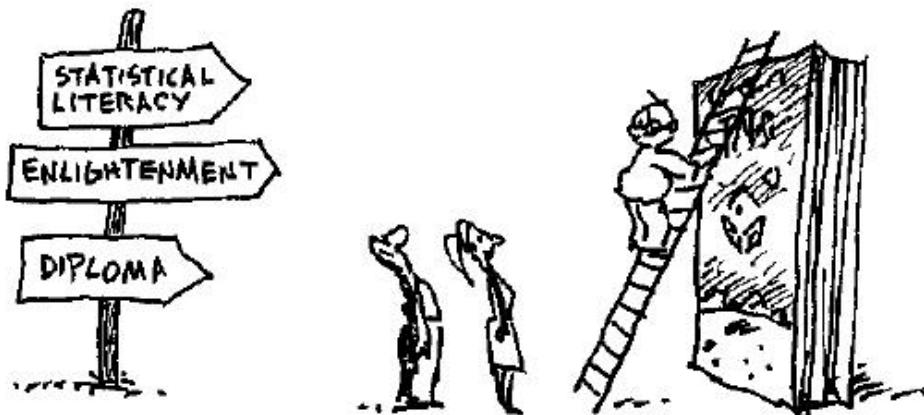
Za to, da ne blestite, ni opravičila.
TOM PETERS, *The Circle of Innovation,*
Alfred A. Knoff, New York, ZDA

Poglavlje 1

Uvod



Z veseljem bi predstavil učbenik, a ga nimamo (vsaj ne v slovenščini). Prav, pa predstavimo skripto in hkrati tudi naš predmet. Le-ta govori o dogodkih in podatkih.¹



Iz učnega načrta: *“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.”*²

¹ Vsekakor ni samo prepreka na poti do diplome, pač pa se uporablja pri večini drugih predmetov na FRI. Če ne verjamete, me pridite negirat – moja vrata so vam vedno odprta (R3.06), posebej v času govorilnih ur, ki so ta semester ob torkih 17:15-18:00.

²VAGABUND s foruma: huh, zveni preveč zapleteno in suhoparno, a saj to je pisano bolj za prfokse.

V nekem članku o računalniških virusih lahko preberemo debato o napadu 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 okužnih 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 predvsem numeričnih dejstev (ni pa nujno), ki jih imenujemo *podatki*. Vsakodnevno smo s podatki takorekoč bombar dirani. Večina ljudi povezuje "statistiko" z biti podatkov, ki izhajajo v dnevni časopisu, novicah, reportažah: povprečna temperatura na današnji dan, procenti pri košarkaških prostih metih, procent tujih vlaganj na našem trgu, in anketa popularnosti politikov. 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 navzoči pri delu mnogih, zato je izobraževanje na področju statistike izredno pomembno pri številnih poklicih. Ekonomisti, finančni svetovalci, vodstveni kader v politiki in gospodarstvu, vsi 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 nočeš izkorišča statistične metode.

Nič lažje 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.

Zakaj verjetnost in statistika?

Statistika 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**, ki jemlje vzorce podatkov in na osnovi njih naredi zaključke (inferenčnost) o populaciji (ekstrapolacija). Slednjo bi po domače morda lahko poimenovali kar *napovedna statistika*, saj ne moremo nikoli biti popolnoma prepričani o zakjučkih.

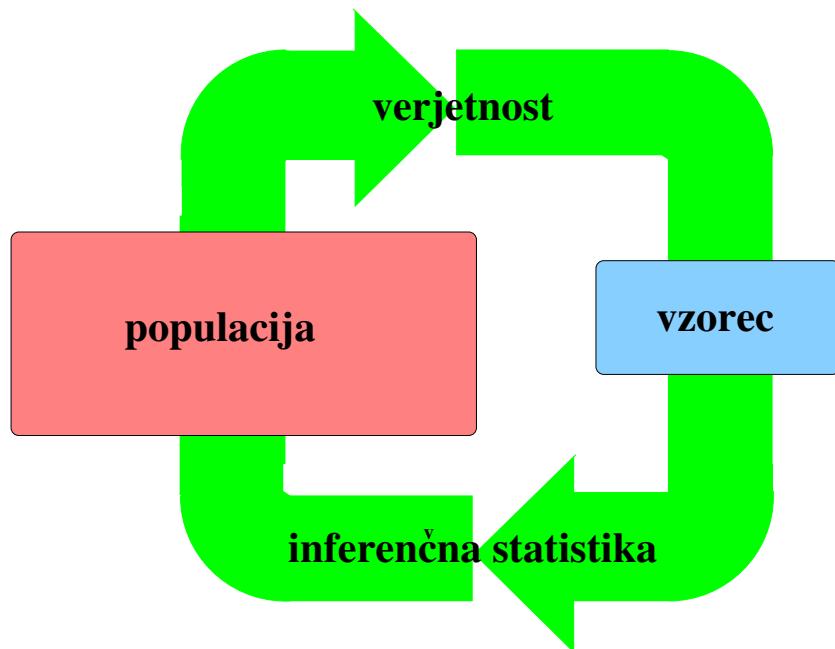


Populacija je podatkovna množica, ki ji je namenjena naša pozornost (vsi objekti, ki jih opazujemo).

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



Pravijo, da je slika vredna 1 000 besed, zato si še nazorno predočimo, 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 osnovi vzorcev.



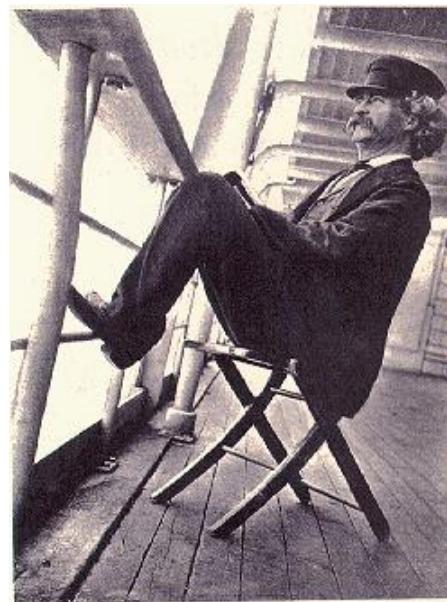
Za konec pa še tole. Matematika je dobra osnova, odličen temelj. Če pri konkretnem računalniškem ustvarjanju (pri tem mislim na razvoj programske opreme in reševanje logističnih problemov, ne pa prodajo in popravilo računalniške opreme) nimamo matematične izobrazbe/znanja, potem je skoraj tako kot če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

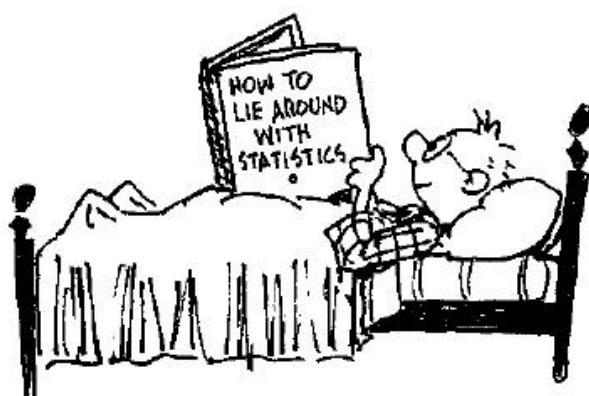
*“Obstajajo tri vrste laži:
laži,
preklete laži in
statistika.”*

iz avtobiografije Marka Twaina



Be good & you will be lonesome.
Mark Twain

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



Vaš cilj je lahko na primer, da se ne pustite drugim vleči za nos.

Del I

VERJETNOST

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



Slika 1.1: Ruleta je igra na srečo, ki izvira iz Francije in se je prvič pojavila v 17. ali zgodnjem 18. stoletju. Najbolj znan center, kjer igrajo to igro, je Monte Carlo, ki leži v kneževini Monako. Pri igri sunemo kroglico na vrteče se kolo, ki ima po obodu 37 ali 38 enakih predalov, v katere se kroglica lahko ujame (slika 1). Predali so izmenjajo rdeče in črne barve, označeni s številkami od 1 do 36 in ne nujno v pravem vrstnem redu. En od predalov je poseben, označen je z zeleno barvo in nosi številko 0 (predvsem v Ameriki dodajo še en predal z oznako 00). Pri igri stavimo določen znesek na določeno število, množico števil ali pa barvo predalčka, v katerem se bo kroglica ustavila. Igro dobimo, kadar se kroglica ustavi na enem od polj, ki jih naša stava pokrije. Koliko bomo dobili nazaj, je odvisno od zneska, ki ga stavimo, in načina stave. Več o tej igri si lahko preberete v članku Aleša Mohoriča, Ali lahko ukanemo ruleto?, *Presek 33*, 8–9.

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 stotinov 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.

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.

Poglavlje 2

Poskusi, dogodki in definicija verjetnosti



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

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

Na področju Francije je leta 1662 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 z matematikom Pierre Fermatom in iz tega so nastali začetki verjetnostnega računa.



Cardano

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

Leta 1713 je švicarski matematik Jacob 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 iskustvene (empirične) narave in njegovi osnovni pojmi so povzeti iz izkušnje. Osnovni pojmi verjetnosti 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. \mathcal{X} , \mathcal{Y} , \mathcal{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

Dogodek je **slučajen**, če so posamezni izidi negotovi, vendar pa je na dolgi rok vzorec velikega števila posameznih izidov napovedljiv.

Za statistika *slučajen* ne pomeni *neurejen*. Za slučajnostjo je neke vrste red, ki se pokaže šele na dolgi rok, po velikem številu ponovitev. Veliko pojavov, naravnih in tistih, ki so delo človeka, je slučajnih. Življenska doba zavarovancev in barva las otrok sta primera naravne slučajnosti. Res, kvantna mehanika zagotavlja, da je na subatomskem nivoju v naravo vgrajena slučajnost. Teorija verjetnosti, matematični opis slučajnosti, je bistvenega pomena za prenekatero sodobno znanost.

Igre naključij so primeri slučajnosti, ki jo namenoma povzroči človek. Kocke v igralnicah

so skrbno izdelane in izvrtane luknje, ki služijo označevanju pik, so zapolnjene z materialom, ki ima enako gostoto kot ostali del kocke. S tem je zagotovljeno, da ima stran s šestimi pikami enako težo kot nasprotna stran, na kateri je le ena pika. Tako je za vsako stran enako verjetno, da bo končala zgoraj. Vse verjetnosti in izplačila pri igrah s kockami temeljijo na tej skrbno načrtovani slučajnosti.

Izpostavimo nekatere posebne dogodke:

(a) **gotov** dogodek – oznaka $\textcolor{red}{G}$: 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 – oznaka $\textcolor{red}{N}$: nikoli se ne zgodi.

Primer: dogodek, da vržemo 7 pik pri metu igralne kocke. \diamond

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

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

2.2 Računanje z dogodki

Dogodek A je **poddogodek** ali **način** dogodka B , kar zapišemo $\textcolor{red}{A} \subseteq B$, 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. \diamond

Če je dogodek A način dogodka B in sočasno dogodek B način dogodka A , sta dogodka **enaka**: $(A \subseteq B) \wedge (B \subseteq A) \iff A = B$. **Vsota** dogodkov A in B je dogodek, označimo jo z $\textcolor{red}{A} \cup \textcolor{red}{B}$ ali $\textcolor{red}{A} + \textcolor{red}{B}$, ki se zgodi, če se zgodi *vsaj* eden izmed dogodkov A in B .

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

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 <http://lkrv.fri.uni-lj.si>, skrajšana verzija pa je v dodatku A.1 na strani 221).

Trditev 2.1. Za poljubna dogodka A in B velja:

- (i) $A \cup B = B \cup A$ (tj. za vsoto velja pravilo o *zamenjavi*),
- (ii) $A \cup A = A$ (tj. vsi dogodki so *idempotenti* za vsoto),
- (iii) $A \cup N = A$ (tj. nemogoč dogodek N je za vsoto *nevtralen element*, tj. enota),
- (iv) $A \cup G = G$ (tj. gotov dogodek G za vsoto absorbira, je univerzalen),
- (v) $B \subseteq A \iff A \cup B = A$

(z vsoto znamo na drug način povedati kdaj je dogodek A način dogodka B),

- (vi) $A \cup (B \cup C) = (A \cup B) \cup C$ (tj. za vsoto velja pravilo o *združevanju*). \square

Produkt dogodkov A in B je dogodek, označimo ga z $A \cap B$ ali AB , ki 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. \diamond

Trditev 2.2. Za poljubna dogodka A in B velja:

- (i) $A \cap B = B \cap A$ (tj. za produkt velja pravilo o *zamenjavi*),
- (ii) $A \cap A = A$ (vsi dogodki so *idempotenti* za produkt),
- (iii) $A \cap N = N$ (tj. nemogoč dogodek N absorbira, je univerzalen za produkt),
- (iv) $A \cap G = A$ (tj. gotov dogodek G je za produkt *neutralen element*, tj. enota),
- (v) $B \subseteq A \iff A \cap B = B$ (tudi s produktom znamo na drug način povedati kdaj je dogodek A način dogodka B),
- (vi) $A \cap (B \cap C) = (A \cap B) \cap C$ (tj. za produkt velja pravilo o *zdrževanju*),
- (vii) $A \cap (B \cup C) = (A \cap B) \cup (A \cap C) \quad \text{in} \quad A \cup (B \cap C) = (A \cup B) \cap (A \cup C)$
(tj. za produkt in vsoto veljata pravili o *porazdelitvi*). \square

Dogodka A in B sta **nezdružljiva**, če se ne moreta zgoditi hkrati. To se zgodi natanko tedaj, ko je njun produkt nemogoč dogodek.

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

Dogodku A **nasproten** dogodek \bar{A} , je tisti, ki se zgodi natanko takrat, ko se dogodek A ne zgodi, in ga imenujemo tudi **negacija** 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:

- (i) $A \cap \bar{A} = N$ (dogodek in njegov nasprotni dogodek sta nezdružljiva),
- (ii) $A \cup \bar{A} = G$, (dogodek in njegov nasprotni dogodek se dopolnita do gotovega dogodka),
- (iii) $\bar{N} = G$, $\bar{G} = N$,
- (iv) $\bar{\bar{A}} = A$,
- (iv) $\overline{A \cup B} = \bar{A} \cap \bar{B} \quad \text{in} \quad \overline{A \cap B} = \bar{A} \cup \bar{B}$ (de Morganovi pravili). \square

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 pik, E_2 , da padeta 2 pik, \dots , 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 , tj.

- $A_i \neq N$, za $i = 1, 2, \dots, n$ (noben med njimi ni nemogoč),
- $A_i A_j = N$ za $i \neq j$ (so paroma nezdružljivi) in
- $A_1 + A_2 + \dots + A_n = G$ (njihova vsota je gotov dogodek).

Primer: Popoln sistem dogodkov pri metu kocke sestavlja 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

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

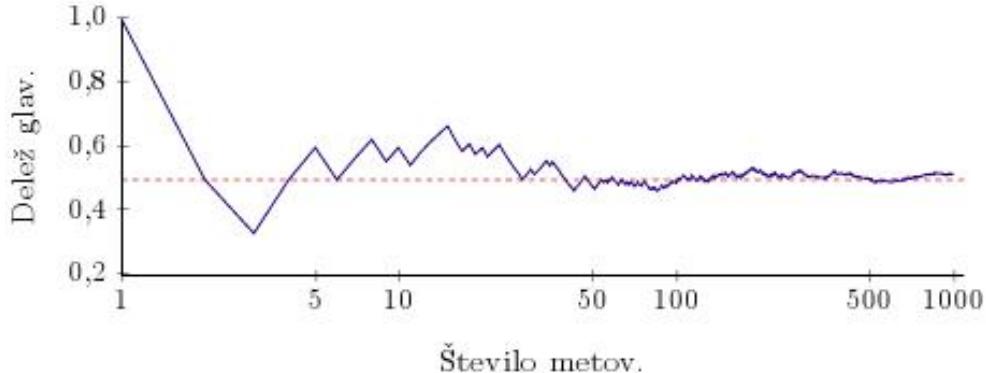
pa je **relativna frekvenca** (pogostost) dogodka A v opravljenih poskusih. Statistični zakon, ki ga kaže izkušnja, je:

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

To temeljno zakonitost so empirično preverjali na več načinov.

Primer: (Metanje kovanca) Ko vržemo kovanec, sta le dva možna izida: slika (glava, grb,...) ali cifra. Na sliki 2.1 je prikazan rezultat 1 000 metov kovanca. Za vsako število metov med 1 in 1 000 smo narisali delež tistih metov, katerih rezultat je bila glava. Prvi met je bila glava,

zato je začetni delež glav enak 1. Pri drugem metu je padla cifra, zato se je delež glav po dveh metih zmanjšal na 0·5. Tretji met je bila spet cifra, sledili pa sta ji dve glavi, zato je bil delež glav po petih metih enak 3/5 ali 0·6.



Slika 2.1: Delež glav v odvisnosti od števila metov kovanca. Delež se sčasoma ustali pri verjetnosti za glavo.

Delež metov, pri katerih pade glava, se na začetku precej spreminja, vendar pa se ustali, ko število metov narašča. Nazadnje pride ta delež blizu 0·5 in tam obstane. Pogovorno rečemo, da se glava pojavi z *verjetnostjo* 0·5. Ta verjetnost je prikazana na grafu s črtkano črto. \diamond

Najbolj znani poskusi s kovanci, kjer so določali relativno frekvenco cifre ($f(A)$), pa so bili še bolj vstajni:

- Buffon (cca. 1750) je v 4040 metih dobil $f(A) = 0\cdot5069$,
- Pearson (cca. 1900) je v 12000 metih dobil $f(A) = 0\cdot5016$,
- Pearson (cca. 1900) je v 24000 metih dobil $f(A) = 0\cdot5005$.

Ti in 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 frekvanca dogodka A v velikem številu ponovitev tega poskusa.

Ker je relativna frekvanca 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 $A \subseteq B \Rightarrow P(A) \leq P(B)$,
3. če sta dogodka A in B nezdružljiva, potem je $P(A + B) = P(A) + P(B)$. \square

Klasični pristop k verjetnosti

Igralcem je že dolgo časa znano, da se meti kovanca, kart ali kock sčasoma ustalijo v točno določene vzorce. Verjetnostna matematika ima začetke v Franciji 17. stoletja, ko so hazarderji začeli prihajati k matematikom po nasvete (več v dodatku C.2, na strani 266). Ideja

verjetnosti temelji na dejstvu, da lahko povprečni rezultat velikega števila slučajnih izidov poznamo z veliko gotovostjo. Vendar pa je definicija verjetnosti z izrazom ‐na dolgi rok‐ nejasna. Kdo ve, kaj ‐dolgi rok‐ je? Namesto tega matematično opišemo *kako se verjetnosti obnašajo*, pri čemer temeljimo na našem razumevanju deležev, ki se pojavljajo na dolgi rok. Da bi nadaljevali, si najprej zamislimo zelo preprost slučajni pojav, en sam met kovanca. Ko vržemo kovanec, ne vemo vnaprej, kakšen bo izid. Kaj pa vemo? Pripravljeni smo priznati, da bo izid bodisi glava bodisi cifra. Verjamemo, da se vsak od teh rezultatov pojavi z verjetnostjo $1/2$. Opis meta kovanca sestoji iz dveh delov:

- seznama vseh možnih izidov in
- verjetnosti za vsakega od teh izidov.

Tak opis je osnova vseh verjetnostnih modelov. Tule je slovarček besed, ki jih pri tem uporabljamo:

Vzorčni prostor S slučajnega pojava je množica vseh možnih izidov.

V tem kontekstu je *dogodek* katerikoli izid ali množica izidov slučajnega pojava. Dogodek je torej podmnožica verjetnostnega prostora.

Verjetnostni model je matematični opis slučajnega pojava, sestavljen iz dveh delov: verjetnostnega prostora S in predpisa, ki dogodkom priredi verjetnosti.

Vzorčni prostor S je lahko zelo preprost ali pa zelo zapleten. Ko vržemo kovanec enkrat, sta le dva možna izida, glava in cifra. Vzorčni prostor je torej $S = \{G, C\}$. Če izbiramo slučajni vzorec 1500 polnoletnih Američanov kot v Gallupovih raziskavah, pa vzorčni prostor vsebuje vse možne izbire 1500 izmed več kot 200 milijonov odraslih prebivalcev. Ta S je hudo velik. Vsak element vzorčnega prostora S je možni vzorec za Gallupovo raziskavo, od koder ime *vzorčni prostor*.

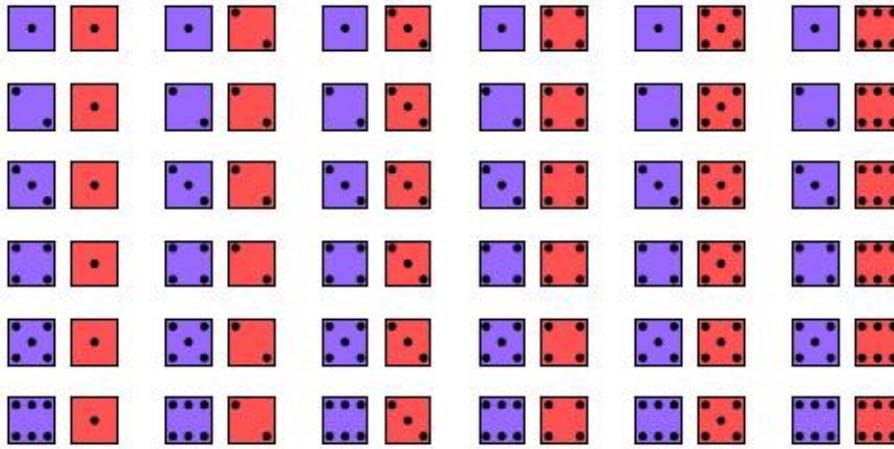
Enako verjetni izidi in *klasična definicija verjetnosti*:

Pri enostavnem slučajnem vzorčenju imajo vsi vzorci enake možnosti, da so izbrani. Kadar je slučajnost produkt človekovega načrtovanja, se velikokrat zgodi, da so izidi iz vzorčnega prostora enako verjetni.

Primer: (Metanje kocke) Metanje dveh kock hkrati je običajni način za izgubljanje denarja v igralnicah. Ko vržemo dve kocki, je možnih 36 izidov, če med kockama razlikujemo. Ti možni izidi so prikazani na sliki 2.2. Tvorijo vzorčni prostor S .

‐Pade 5‐ je nek dogodek, označimo ga z A . Vsebuje 4 od 36 možnih izidov:

$$A = \{ \begin{array}{c} \text{•} \\ \text{—} \end{array} \begin{array}{c} \text{—} \\ \text{•} \end{array}, \begin{array}{c} \text{—} \\ \text{•} \end{array} \begin{array}{c} \text{•} \\ \text{—} \end{array}, \begin{array}{c} \text{—} \\ \text{•} \end{array} \begin{array}{c} \text{—} \\ \text{•} \end{array}, \begin{array}{c} \text{—} \\ \text{•} \end{array} \begin{array}{c} \text{—} \\ \text{—} \end{array} \},$$



Slika 2.2: Možni izidi pri metu dveh kock.

Opozorimo, da ločimo izida $1 + 4$ in $4 + 1$, kakor tudi $2 + 3$ in $3 + 2$. V primeru sode vsote (npr. 6) pa imamo en sam izid, ko kocki pokažeta enako (tj. $3 + 3$). Če so kocke dobro izdelane, izkušnje kažejo, da se vsak od 36 izidov s slike 2.2 pojavi enako pogosto. Razumen verjetnostni model torej pripisuje vsakemu od izidov verjetnost $1/36$. \diamond

Trditev 2.4 nam v primeru enako verjetnih izidov pove, kolikšne so te verjetnosti.

Vzemimo $n \in \mathbb{N}$ in predpostavimo, da so dogodki iz popolnega sistema dogodkov $\{E_1, E_2, \dots, E_n\}$ enako verjetni: $P(E_1) = P(E_2) = \dots = P(E_n) = p$. Tedaj je $P(E_i) = 1/n$, $1 \leq i \leq n$. Če je nek dogodek A sestavljen iz m dogodkov tega popolnega sistema dogodkov, potem je njegova verjetnost $P(A) = m/n$.

Z drugimi besedami:

Če je pri slučajnem pojavu možnih n izidov, ki so vsi enako verjetni, potem je verjetnost vsakega od izidov enaka $1/n$. Za poljubni dogodek A , ki je sestavljen iz m izidov, pa je njegova verjetnost enaka

$$P(A) = \frac{\text{število izidov v } A}{\text{število izidov v } S} = \frac{m}{n},$$

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

Primer: (Naključna števila) Z generatorjem naključnih števil generiramo števila iz množice $\{0, 1, \dots, 9\}$ za neko zaporedje. Če je generator dober, potem se vsak element iz zaporedja z enako verjetnostjo katerikoli od desetih možnih. Torej je verjetnost vsakega od desetih izidov enaka $1/10$, verjetnostni model pa:

Izid	0	1	2	3	4	5	6	7	8	9
Verjetnost	0·1	0·1	0·1	0·1	0·1	0·1	0·1	0·1	0·1	0·1

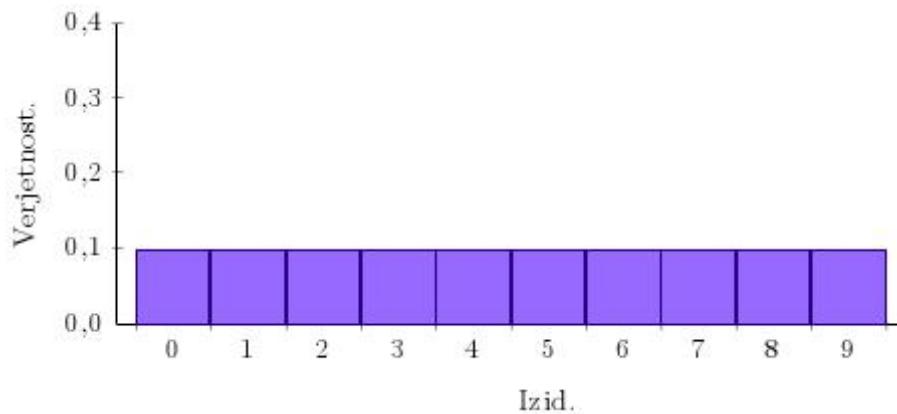
Na sliki 2.3 je ta model prikazan z verjetnostnim histogramom. Verjetnost poljubnega dogodka poiščemo tako, da preštejemo, koliko izidov vsebuje. Naj bo na primer

$$A = \{\text{lih izid}\} = \{1, 3, 5, 7, 9\} \text{ in } B = \{\text{izid je manjši ali enak 3}\} = \{0, 1, 2, 3\}.$$

Vidimo, da je $P(A) = 0,5$ in $P(B) = 0,4$. Dogodek A ali B vsebuje 7 izidov,

$$A + B = \{0, 1, 2, 3, 5, 7, 9\},$$

zato ima verjetnost 0·7. Ta verjetnost *ni* vsota verjetnosti $P(A)$ in $P(B)$, ker A in B *nista* disjunktna dogodka. Izida 1 in 3 pripadata tako A kot B . \diamond



Slika 2.3: Verjetnostni histogram, ki prikazuje verjetnosti pri naključnem generirajuju števke med 0 in 9.

Primer: (Verjetnost pri metanju dveh kock) V igrah kot je *craps*¹ je pomembna samo *vsota* pik, ki jih vržemo. Spremenimo torej izide, ki nas zanimajo, takole: vržemo dve kocki in preštejemo število pik. V tem primeru je možnih le 11 izidov, od vsote 2, ki jo dobimo, če vržemo dve enici, do vsote 12, ki jo dobimo z dvema šesticama. Vzorčni prostor je

$$S = \{2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12\}.$$

Če primerjamo ta S z vzorčnim prostorom na sliki 2.2, vidimo, da se S lahko spremeni, če spremenimo podrobnosti, ki jih pri nekem pojavu opazujemo. Izidi v tem novem vzorčnem prostoru *niso* več enako verjetni, ker lahko 7 dobimo na šest načinov, 2 ali 12 pa samo na enega. Za igralniške kocke je smiselno vsakemu od omenjenih 36 izidov predpisati isto verjetnost. Ker mora biti verjetnost vseh 36 dogodkov skupaj enaka 1, mora imeti vsak od izidov verjetnost 1/36. **Kolikšna je verjetnost, da je vsota pik na obeh kockah enaka 5?** Ker je ta dogodek sestavljen iz štirih možnih izidov, ki smo jih zapisali v prejšnjem primeru, nam **pravilo vsote**, Trditev 2.4(3), pove, da je

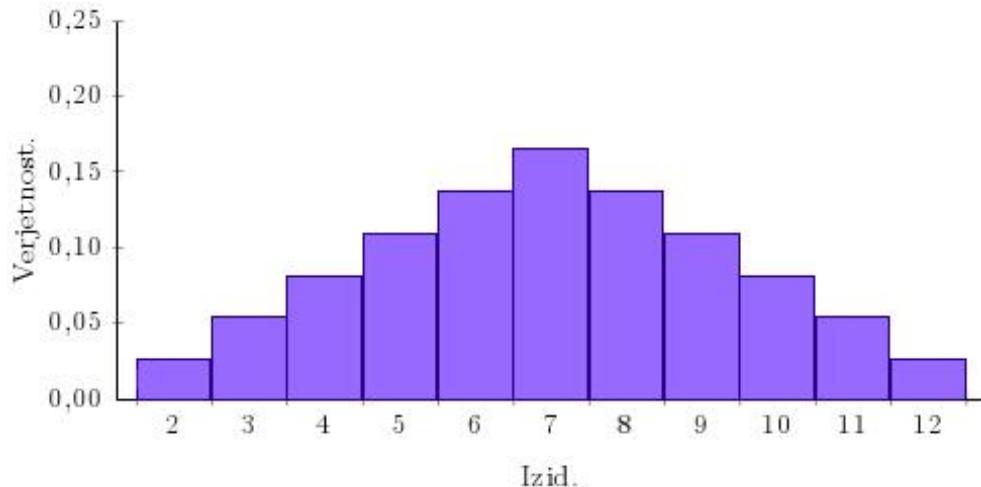
¹Ameriška igra z dvema kockama.

$$\begin{aligned}
 P(\text{pade 5}) &= P(\begin{array}{|c|c|}\hline \bullet & \blacksquare \\ \hline \end{array}) + P(\begin{array}{|c|c|}\hline \bullet & \bullet \\ \hline \end{array}) + P(\begin{array}{|c|c|}\hline \bullet & \bullet \\ \hline \end{array}) + P(\begin{array}{|c|c|}\hline \bullet & \bullet \\ \hline \end{array}) = \\
 &= \frac{1}{36} + \frac{1}{36} + \frac{1}{36} + \frac{1}{36} = \frac{4}{36} = 0.111.
 \end{aligned}$$

Kaj pa verjetnost, da dobimo 7? Na sliki 2.2 najdemo šest izidov, pri katerih je vsota pik enaka 7. Verjetnost za 7 je torej $6/36$, kar je približno 0.167. Na ta način nadaljuj z računanjem, da dobiš celoten verjetnostni model (vzorčni prostor in predpise verjetnosti) za met dveh kock, pri katerem opazujemo vsoto pik. Tole je rezultat:

Izid	2	3	4	5	6	7	8	9	10	11	12
Verjetnost	$\frac{1}{36}$	$\frac{2}{36}$	$\frac{3}{36}$	$\frac{4}{36}$	$\frac{5}{36}$	$\frac{6}{36}$	$\frac{5}{36}$	$\frac{4}{36}$	$\frac{3}{36}$	$\frac{2}{36}$	$\frac{1}{36}$

Na sliki 2.4 je *verjetnostni histogram* tega verjetnostnega modela. Višina vsakega stolpca prikazuje verjetnost dogodka, ki ga ta stolpec predstavlja. Ker so višine ravno verjetnosti, se seštejejo v 1. Slika 2.4 je idealizirana podoba rezultatov velikega števila metov kock. Kot idealizacija je popolnoma simetrična. \diamond



Slika 2.4: Verjetnostni histogram za met dveh kock, ki prikazuje, kakšne so verjetnosti za posamezno vsoto pik.

Primer: V 17. stoletju je bila popularna igra, pri kateri se je stavilo na vsoto pik pri metu treh kock. Stavili so večinoma na 9 ali 10.

Kockarji so napisali vse možnosti na naslednji način:

vsota 9	vsota 10
1 2 6	1 3 6
1 3 5	1 4 5
1 4 4	2 2 6
2 2 5	2 3 5
2 3 4	2 4 4
3 3 3	3 3 4

Sklepali so, da sta stavi enakovredni. [So imeli prav?](#) Problem je kockarjem rešil sam Galileo Galilei (1564-1642). Na list papirja jim je napisal vse možne trojice padlih pik na kockah. [Koliko jih je?](#) \diamond

Z zgornjimi primeri smo spoznali enega od načinov, kako dogodkom priredimo verjetnosti: predpišemo verjetnosti vsakemu od izidov, nato pa jih seštejemo, da dobimo verjetnost dogodka. Da bi tako prirejanje zadoščalo pravilom verjetnosti, se morajo verjetnosti posameznih izidov sešteeti v 1.

Verjetnostni model za (končen) vzorčni prostor podamo tako, da predpišemo verjetnost vsakemu posameznemu izidu. Te verjetnosti morajo biti števila med 0 in 1 in njihova vsota mora biti enaka 1. Verjetnost poljubnega dogodka je vsota verjetnosti izidov, ki ga sestavljajo.

Primer: (Rangiranje srednješolcev) Naključno izbiramo študente in študentke prvih letnikov in jih vprašamo, kje so se v srednji šoli nahajali glede na učni uspeh. Tule so verjetnosti, ki jih dobimo iz deležev v velikem vzorcu študentov:

Razred	Zgornjih 20%	Drugih 20%	Tretjih 20%	Četrthih 20%	Spodnjih 20%
Verjetnost	0·41	0·23	0·29	0·06	0·01

Prepričaj se, da se te verjetnosti seštejejo v 1. Zdaj lahko izračunamo še naslednji dve verjetnosti:

$$\begin{aligned} P(\text{študent je v zgornjih } 40\%) &= P(\text{zgornjih } 20\%) + P(\text{drugih } 20\%) \\ &= 0\cdot41 + 0\cdot23 = 0\cdot64, \end{aligned}$$

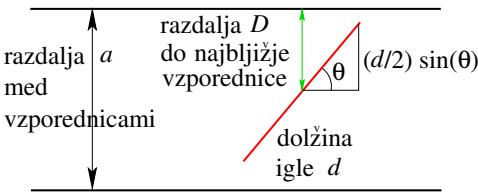
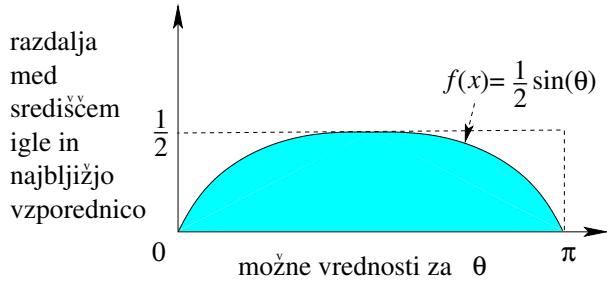
$$P(\text{študent ni v zgornjih } 40\%) = 0\cdot23 + 0\cdot29 + 0\cdot06 + 0\cdot01 = 0\cdot59.$$

[Ali veš, zakaj je verjetnost, da študent ni v zgornjih 20%, enaka 1 minus verjetnost, da izbrani študent je v zgornjih 20%?](#) \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. Verjetnost torej lahko računamo s pomočjo geometrije. Lahko pa se zgodi tudi obratno, da nam verjetnost pomaga določiti nekaj geometrijskega, npr. število π , ki ima neskončen neperiodičen zapis $3\cdot14159265\dots$ S tem številom je povezana ena od znamenitih nalog klasičnega verjetnostnega računa iz leta 1777, ki jo je Buffon zastavil ter rešil z enostavno Monte Carlo metodo in jo predstavimo v naslednjem primeru.

Primer: Na večji papir narišemo na razdalji a vrsto vzporednih premic. Nato vzamemo iglo dolžine d (zaradi enostavnosti naj bo $d \leq a$) in jo na slepo večkrat vržemo na papir. **Zanima nas verjetnost dogodka A , da igla seče eno od vzporednih premic.** Igla seče najbljižjo vzporednico, če je $D \leq (d/2) \sin(\theta)$.

Slika 2.5: (b) Splošna a in d .(a) Vzporednice na razdalji a .(c) Za naše potrebe bi lahko vzeli kar $d = 1$ in $a = 1$.

Ploščino pod krivuljo funkcije $(d/2) \sin(\theta)$ izračunamo z določenim integralom te funkcije od 0 do π , in je enaka d . Ploščina pravokotnika $\pi \times (a/2)$ pa je $\pi a/2$. Kvocient teh dveh ploščin pa nam da $P(A) = 2d/(\pi a)$, saj je na sliki (c) vsaka točka enako verjetna. \diamond

Pri zgornjem primeru je najbolj zanimivo, da v rezultatu nastopa število π , ki ga srečamo na primer pri obsegu ali pa ploščini kroga. Ravno ta prisotnost števila π je vzpodbudila vrsto ljudi k naslednjemu poskusu. Poskus so ponovili zelo velikokrat (n -krat) in prešeli, kolikokrat je sekala eno od črt (m -krat). Relativno frekvenco $f(A) = m/n$ so vzeli za približek zgornje verjetnosti in od tod ocenili približek za število π : $\pi \doteq 2dn/(am)$.

eksperimentator	d	leto	število poskusov	ocena za π
Wolf	0.8	1850	5000	3.1596
Smith	0.6	1855	3204	3.1553
deMorgan	1.0	1860	600	3.1370
Fox	0.75	1894	1120	3.1419

n	m	ocena za π
98	57	3.43860
234	150	3.14094
239	152	3.14474
357	220	3.24545
355	226	3.14159

Tabela 2: Rezultati (a) iz raznih knjig (za $a = 1$), (b) našega programa (na 5 decimalnih mest), $a = d = 1$.

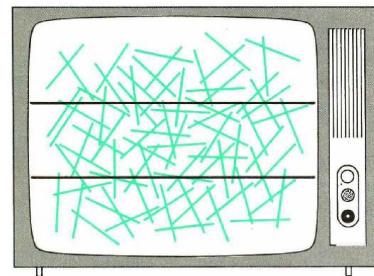
Za računanje čim večjega števila decimalnih mest števila π je znanih vrsta drugih uspešnejših metod. Zgornji poskus z metanjem igle ne prinese velike natančnosti, prestavili smo ga bolj kot zanimivost. Pravzaprav je tu še en pomemben pogled. Metanje igle lahko zelo nazorno posnemamo z računalnikom. Zaradi enostavnosti vzemimo, da je $d = a$. E. Kramar je v Preseku 12/4 (1984/85), glej <http://www.presek.si/12/731-Kramar.pdf>, predstavil zgled programa v basicu za hišni računalnik ZX Spectrum. Vzel je $a = 50$ enot in na zaslonu

narisal dve vzporednici na tej razdalji. Slučajnost pri metu je generiral s funkcijo RND, ki pomeni na slepo izbrano realno število med 0 in 1. Lego igle je določil s slučajnim kotom med 0° in 180° ($\text{fi} = \text{PI} * \text{RND}$) in slučajno razdaljo ene konice igle od zgornje vzporednice ($b = \text{INT}(a * \text{RND})$). Dodal je še slučajni pomik v levo ali desno (stavek 85), da igla pada nekako po celem ekranu, vendar to ne vpliva na rezultat.

```

5 REM metanje igle
6 REM in računanje približka za število pi
7 REM
10 INPUT "število metov"; n
15 PRINT "število metov ";n
20 LET a=50: LET m=0
30 PLOT 0,70: DRAW 255,0
40 PLOT 0,120: DRAW 255,0
50 FOR i =1 TO n
60 LET b=INT(a*RND): LET fi=PI*RND
70 LET x1=a*COS(fi): LET y1=a*0.5*SIN(fi)
80 LET z1=70+b-y1: LET z2=70+b+y1
85 LET c=INT(150*RND)+50
90 PLOT c,z1: DRAW x1,2*y1
100 IF z1 < 70 AND z2 > 120 THEN GO TO 140
110 LET m=m+1
120 PRINT AT 3,3; m; " " ;i
130 NEXT i
140 PRINT "približek za pi = " ; 2*n/m
150 STOP

```



Tule je še nekaj Kramarjevih komentarjev: "Bralec, ki ima dostop do računalnika, bo gotovo tudi sam poskusil posnemati metanje igle (ali igel) na računalniku, pri tem bo spremenjal program po svoji želji. Če vzamemo večje število poskusov (n), lahko izključimo risanje, saj nas zanima samo rezultat na koncu. Velja opozoriti, da prevelikega števila n pravzaprav nima smisla vzeti, saj rezultati ne bodo bistveno boljši, zaradi vsakokratnega računanja nekaterih funkcij pa bo računalnik kar precej časa "mlel". Prav tako nima smisla iti na lov za čim boljšim približkom števila π . Ker poznamo veliko njegovih decimalnih mest, bi lahko poskus privedli tako, da bi ga prekinili v trenutku za nas najbolj sprejemljivega približka. To pa ni poskus v zgornjem smislu, ko je število metov vnaprej izbrano. Drug način grobe potvorbe rezultatov pa bi bil, da bi vzeli ulomek, ki je dober približek za število π , kot je na primer star kitajski približek $355/113$, in bi nekomu sporočili, da nam je pri 355-kratnem metu igle ta sekala vzporednico 226-krat. V resnici je prav malo verjetno (da se ugotoviti celo verjetnost tega dogodka, in sicer je manjša od 0·034), da se nam bo ravno to zgodilo. Omeniti velja še dejstvo, da smo zgoraj privzeli, da je generator slučajnih števil (RND) na računalniku dobro narejen, da dobro posnema slučajno izbiro realnega števila med 0 in 1, kar pa je zopet le približek."

Za konec pa omenimo možnost eksperimentiranja s programi na internetu, npr.:

<http://www.angelfire.com/wa/hurben/buff.html> in
<http://demonstrations.wolfram.com/BuffonsNeedleProblem/>.

2.4 Osnovne lastnosti verjetnosti

Trditev 2.5. Za poljubna dogodka A in B velja:

$$P(A + B) = P(A) + P(B) - P(AB).$$

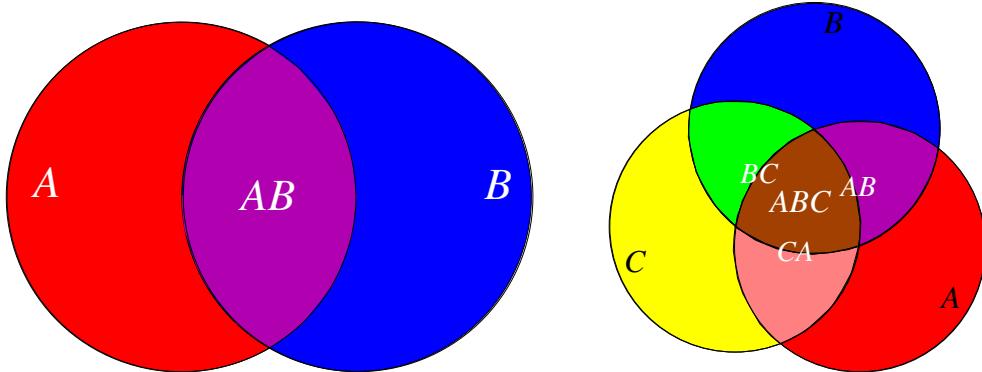
Dokaz. Pri statističnem pristopu je dovolj, da se prepričamo o veljavnosti zveze

$$k_{A+B} + k_{AB} = k_A + k_B$$

(glej sliko 2.6(b)), ki jo nato na obeh stranih delimo z n in pošljemo $n \rightarrow \infty$. \square



(a) John Venn, angleški matematik 1834-1923.



Slika 2.6: Ne to ni reklama za Mastercard, pač pa Vennova diagrama za (b) dve množici in (c) za tri množice. Vennov diagram za štiri množice pa raje rišemo v prostoru. [Zakaj?](#)

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 $P(S + P) = 4/5$, kolikšna je verjetnost, da naredi oba izpita?

$$P(S \cap P) = P(S) + P(P) - P(S + P) = \frac{2}{3} + \frac{5}{9} - \frac{4}{5} \doteq 0.42. \quad \diamond$$

Posledica 2.6. $P(\overline{A}) = 1 - P(A)$. \square

Primer: Iz kupa 32 kart slučajno povlečemo 3 karte. **Kolikšna je verjetnost, da je med tremi kartami vsaj ena as (dogodek A)?** Pomagamo si z nasprotnim dogodkom \overline{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

$$P(\overline{A}) = \frac{\binom{28}{3}}{\binom{32}{3}} \doteq 0.66; P(A) = 1 - P(\overline{A}) \doteq 1 - 0.66 = 0.34. \quad \diamond$$

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:

$$P(A + B + C) = P(A) + P(B) + P(C) - (P(AB) + P(AC) + P(BC)) + P(ABC). \quad \square$$

Kako lahko to pravilo posplošimo še na več dogodkov?

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

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

pri čemer predznaki alternirajo.

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

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

Velja tudi za števno neskončne množice dogodkov. \square

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

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

$$\begin{array}{ccc} (A, B) & \xrightarrow{+} & A + B \\ \downarrow (P(\cdot), P(\cdot)) & & \downarrow P(\cdot) \\ (P(A), P(B)) & \xrightarrow{+} & P(A) + P(B) = P(A + B) \end{array}$$

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

S simulacijo lahko pridemo do števila e .

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 $\bar{A} = A_1 + \dots + A_n$, dogodki v slednji vsoti 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, tj. $S_1 = n P(A_1)$. Poskus je v našem primeru razdeljevanje n dopisov po ovojnicih, 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. $S_1 = n(n-1)!/n! = 1$. 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

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

in končno, če upoštevamo še Trditev A.5, dobimo $\lim_{n \rightarrow \infty} P(A) = 1/e$, tako da imamo že pri razmeroma majhnih n $P(A) \doteq 1/e \doteq 0.369$. \diamond

2.5 Aksiomi Kolmogorova*

Poglejmo še, kako na verjetnost gledajo matematiki. 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** (tudi σ -algebra), če velja:

$$A \in \mathcal{D} \Rightarrow \bar{A} \in \mathcal{D} \quad \text{in} \quad A, B \in \mathcal{D} \Rightarrow A + B \in \mathcal{D}.$$

Pri neskončnih množicah dogodkov moramo drugo zahtevo posplošiti:

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

Naj bo \mathcal{D} algebra dogodkov v G . **Verjetnost na G** je preslikava

$P : \mathcal{D} \rightarrow \mathbb{R}$ z lastnostmi:

1. $P(A) \geq 0$ za vsak $A \in \mathcal{D}$.
2. $P(G) = 1$.
3. Če so dogodki $A_i, i \in I$ (I je množica indeksov), paroma nezdružljivi, je

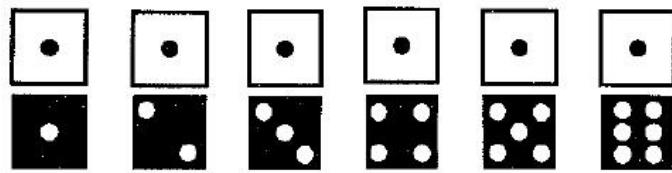
$$P\left(\sum_{i \in I} A_i\right) = \sum_{i \in I} P(A_i).$$

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



Poglavlje 3

Pogojna verjetnost



3.1 Intriga (po Kvarkadabri)

Ne podcenujmo 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 priloga 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 enim se skriva avto, za drugima dvema pa koza. Recimo, da izberete vrata številka 3, voditelj igre, ki ve, kaj se nahaja za posameznimi vrti, 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 splača zamenjati izbor vrat?”¹

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 splača zamenjati, saj se tako verjetnost, da bomo zadeli avto, poveča za dvakrat. Tole je njen odgovor: **Seveda se splača zamenjati vrata.** Prva vrata imajo le **1/3** verjetnosti za zmago, medtem ko imajo druga verjetnost **2/3**.

¹ 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).

Morda si je najlažje vse skupaj predstavljate takole. Predpostavimo, da je na voljo milijon vrat in vi izberete prva. Nato voditelj, ki ve, kaj se nahaja za posameznimi vrti, odpre vsa vrata razen prvih vrat in vrat številka 777 777. V tem primeru bi zelo hitro zamenjali svoj izbor, kajne? Se najinteligenčnejša ženska na planetu moti?

Sledila je ploha kritik (več kot 10 000 pisem jeznih bralcev, med katerimi je bilo ogromno učiteljev matematike). Skoraj 1 000 pisem je bilo podpisanih z imeni (dr. nazivi, napisana na papirju z glavo katere od ameriških univerz - www.marilynvossavant.com). Marylin bralce zavaja, **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 obožil, da je ona sama koza.

Polemika je pristala celo na naslovni 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 splača spremeniti izbiro vrat. Kdo ima prav?

Razprava o Monty Hall problemu spada na področje, ki mu matematiki pravijo **pogojna 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 vrti je avtomobil.

Ko v drugem delu odpre vrata, za katerimi se pokaže koza, vnaprej ve, da za temi vrti 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 1 000 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 enim ali za drugimi vrti, ko so dvoja 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 vrti 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 tekmovalec 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 tekmovalec 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 tekmovalec 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 \mathcal{X} , ki je realizacija kompleksa (seznama) 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 \mathcal{X}' , verjetnost dogodka A v tem poskusu pa označimo s $P_B(A)$ oziroma $P(A|B)$ (ali celo $P(A|B)$), če želimo poudariti, da ne gre za nobeno deljenje).² Le-ta se z verjetnostjo $P(A)$ ujema ali pa tudi ne. Pravimo, da je poskus \mathcal{X}' poskus \mathcal{X} s pogojem B in verjetnost $P(A|B)$ **pogojna verjetnost** dogodka A glede na dogodek B . Je torej tudi verjetnost, le obravnavani kompleks pogojev, ki mora biti izpolnjen, se je spremenil.

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

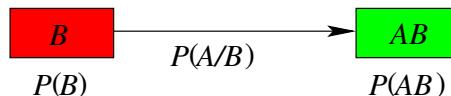
$$P(A|B) = \frac{P(AB)}{P(B)}.$$

Dokaz. Denimo, da smo n -krat ponovili poskus \mathcal{X} in da se je ob tem k_B -krat zgodil dogodek B . To pomeni, da smo v n ponovitvah poskusa \mathcal{X} napravili k_B -krat poskus \mathcal{X}' . Dogodek A se je zgodil ob poskusu \mathcal{X}' le, če se je zgodil tudi B , tj. AB . Denimo, da se je dogodek AB zgodil ob ponovitvi poskusa k_{AB} -krat. Potem je relativna frekvenca dogodka A v opravljenih ponovitvah poskusa \mathcal{X}' :

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

□

Pogojna verjetnost P_B ima prav take lastnosti kot 'brezpogojna'.³



Slika 3.3: Zvezo iz Trditve 3.1 lahko predstavimo grafično: verjetnost dogodka AB (desna škatla) je enaka produktu verjetnosti dogodka B (leva škatla) in pogojni verjetnosti dogodka A glede na dogodek B . Slednjo verjetnost zato postavimo ob povezavo med temi škatlami.

Primeri: (1) 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**:

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

² Opozorilo: A/B ne mešati z razliko množic $A \setminus B$.

³ Trojica (B, \mathcal{D}_B, P_B) , $\mathcal{D}_B = \{AB \mid A \in \mathcal{D}\}$ je zopet verjetnostni prostor.

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

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

ali neposredno iz kontingenčne tabele $P(M/Z) = 460/700$.

(2) 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 natanko 3-krat belo kroglo?** Dogodek A je, da izvlečem belo kroglo. Potem je

$$p = P(A) = \frac{8}{10} = 0.8, \quad q = 1-p = 1-0.8 = 0.2 \quad \text{in} \quad P_5(3) = \binom{5}{3} 0.8^3 \cdot (1-0.8)^{5-3} = 0.205,$$

slednja je verjetnost, da v petih poskusih izberemo 3-krat belo kroglo. Mati ima dva otroka.

(3) **Kolikšna je verjetnost, da je drugi otrok moškega spola, če vemo, da je vsaj en otrok moškega spola?** Naj bo A dogodek, da je drugi otrok moškega spola, B pa dogodek, da je vsaj en otrok moškega spola. Radi bi izračunali $P(A/B)$. V ta namen si oglejmo popoln sistem naslednjih elementarnih dogodkov, ki jih sestavlja pari [spol,spol]: $\Omega_1 = \{[M, M], [M, \check{Z}], [\check{Z}, M], [\check{Z}, \check{Z}]\}$. Le-ti so zaradi simetrije enako verjetni, tj. zgodijo se z verjetnostjo $\frac{1}{4}$ (pri tem pa zanemarimo enojajčne dvojčke). Dogodek B je sestavljen iz treh elementarnih dogodkov, dogodek AB pa le iz dveh, tj. $P(A/B) = 2/3$. Lahko pa najprej izračunamo $P(B) = \frac{3}{4}$, $P(AB) = \frac{1}{2}$ in po formuli za pogojno verjetnost pa dobimo

$$P(A/B) = \frac{P(AB)}{P(B)} = \frac{1}{2} \cdot \frac{4}{3} = \frac{2}{3}.$$

(4) **Kolikšna je verjetnost, da je drugi otrok moškega spola, če vemo, da je vsaj en otrok moškega spola in rojen na torek?** Definiramo še dogodek: C - vsaj en otrok moškega spola in rojen na torek. Množica Ω_2 je sestavljena iz parov $[(\text{spol,dan}), (\text{spol,dan})]$, kar pomeni, da imamo $(2 \times 7)^2 = 14^2$ enakovrednih elementarnih dogodkov. Dogodek C predstavlja

14 parov oblike $[(\text{spol,dan}), (\text{M,torek})]$ in 14 parov oblike $[(\text{M,torek}), (\text{spol,dan})]$,

vendar moramo upoštevati, da smo par $[(\text{M,torek}), (\text{M,torek})]$ šteli dvakrat, torej velja $P(C) = (14 + 14 - 1)/14^2 = 27/14^2$. Dogodek AC predstavlja

7 parov oblike $[(\text{M,torek}), (\text{M,dan})]$ in 14 parov oblike $[(\text{spol,dan}), (\text{M,torek})]$,

vendar smo tudi v tem primeru par $[(\text{M,torek}), (\text{M,torek})]$ šteli dvakrat, zato je $P(AC) = (7 + 14 - 1)/14^2 = 20/14^2$. Sedaj pa po formuli za pogojno verjetnost dobimo

$$P(A/C) = \frac{P(AC)}{P(C)} = \frac{20}{14^2} \cdot \frac{14^2}{27} = \frac{20}{27}.$$

Tudi tu bi lahko prišli do enakega rezultata neposredno, tj.

$$P(A/C) = k_{AC}/k_C = (7 + 14 - 1)/(14 + 14 - 1) = 20/27. \quad \diamond$$

Iz formule za pogojno verjetnost sledita naslednji zvezi:

$$P(AB) = P(B)P(A/B), \quad (3.1)$$

$$P(AB) = P(A)P(B/A). \quad (3.2)$$

(Drugo oz. ‘dualno’ enakost smo dobili iz prve tako, da smo v prvi zamenjali vlogi dogodkov A in B , seveda pa zanjo potrebujemo še pogoj $P(A) \neq 0$.) Torej velja:

$$P(A)P(B/A) = P(B)P(A/B). \quad (3.3)$$

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

$$P(AB) = P(A) \cdot P(B).$$

Opozorilo: za par nezdružljivih dogodkov A in B pa velja $P(A/B) = 0$.

Primer: Za neodvisna dogodka A in B , za katera velja $P(A) = 0.3$ in $P(B) = 0.6$, poišči $P(A + B)$. Iz $P(A + B) = P(A) + P(B) - P(AB)$ in za neodvisne dogodke še $P(AB) = P(A)P(B)$ dobimo $P(A + B) = 0.3 + 0.6 - 0.18 = 0.72$. \diamond

Trditev 3.2. *Naj bodo A , B in C poljubni dogodki. Potem velja*

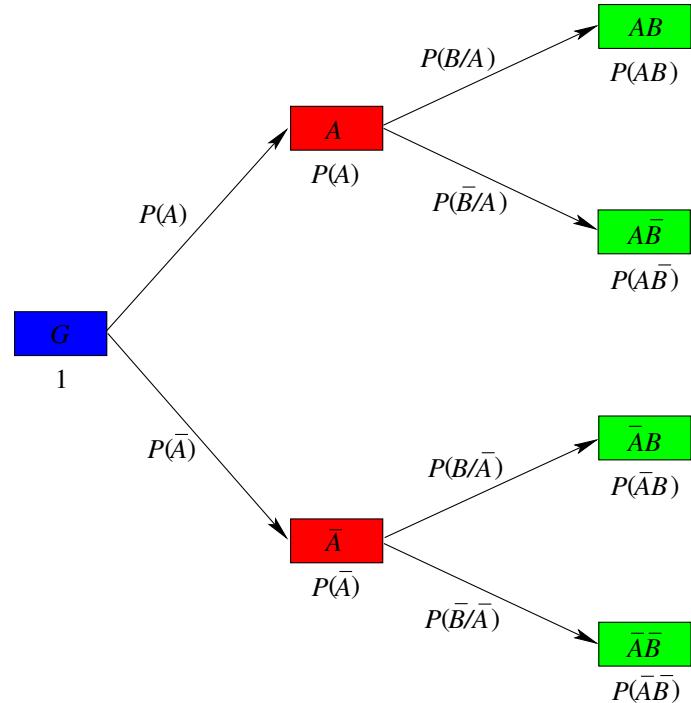
(i) *Dogodka A in B sta neodvisna natanko tedaj, ko je $P(A/B) = P(A)$ oziroma natanko tedaj, ko je $P(A/B) = P(A/\bar{B})$.*

(ii) $P(ABC) = P(A) \cdot P(B/A) \cdot P(C/(AB))$.

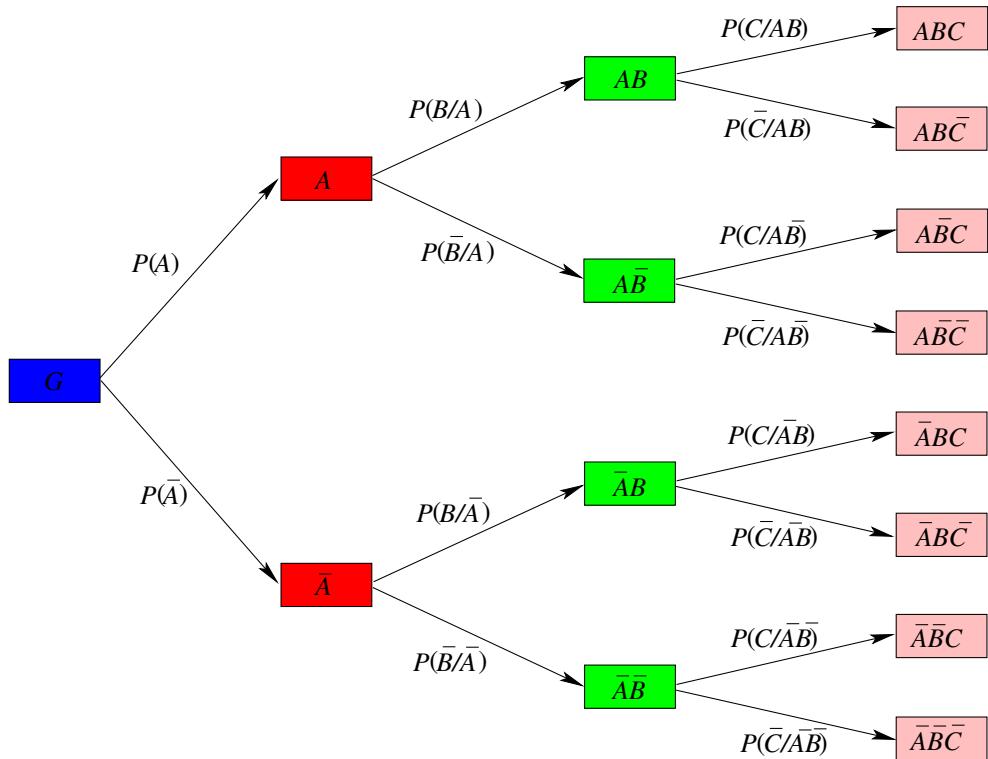
(iii) $P(B/A) + P(\bar{B}/A) = 1$.

Dokaz. (i) in (ii) sledita iz Trditve 3.1 (pri (ii) najprej združimo dogodka A in B).

(iii) Dogodki, ki ustrezajo tem verjetnostim, sestavljajo namreč popoln sistem dogodkov (zavadati pa se je potrebno za kateri verjetnostni prostor gre). Seveda pa bi se o tej identiteti lahko prepričali tudi neposredno s formulo za pogojno verjetnost in upoštevanjem, da sta AB in $A\bar{B}$ nezdružljiva dogodka. \square



Slika 3.4: Verjetnost vsakega izmed dogodkov AB , $A\bar{B}$, $\bar{A}B$ in $\bar{A}\bar{B}$ je enaka produktu verjetnosti na puščicah od začetka (koren na levi) pa do samega dogodka, kar nam zagotavlja identiteta (3.1). Primerjaj te dogodke s polji kontingenčne tabele 2×2 . Le-ti sestavljajo popoln sistem dogodkov, vsota prvega in tretjega pa je ravno dogodek B . Kakor velja $P(A) + P(\bar{A}) = 1$, velja tudi $P(B/A) + P(\bar{B}/A) = 1$, tj. vsota verjetnosti na izhodnih puščicah iz A , je enaka 1. Glej Trditev 3.2.



Slika 3.5: Binarno drevo, ki nas pripelje do vseh osmih produktov med tremi dogodki A , B ter C in njihovimi nasprotnimi dogodki \bar{A} , \bar{B} in \bar{C} , pri čemer mora nastopati v produktu natanko en izmed nasprotnih dogodkov X in \bar{X} za vsak $X \in \{A, B, C\}$.

Dogodki A_i , $i \in I$ so **neodvisni**, če je

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

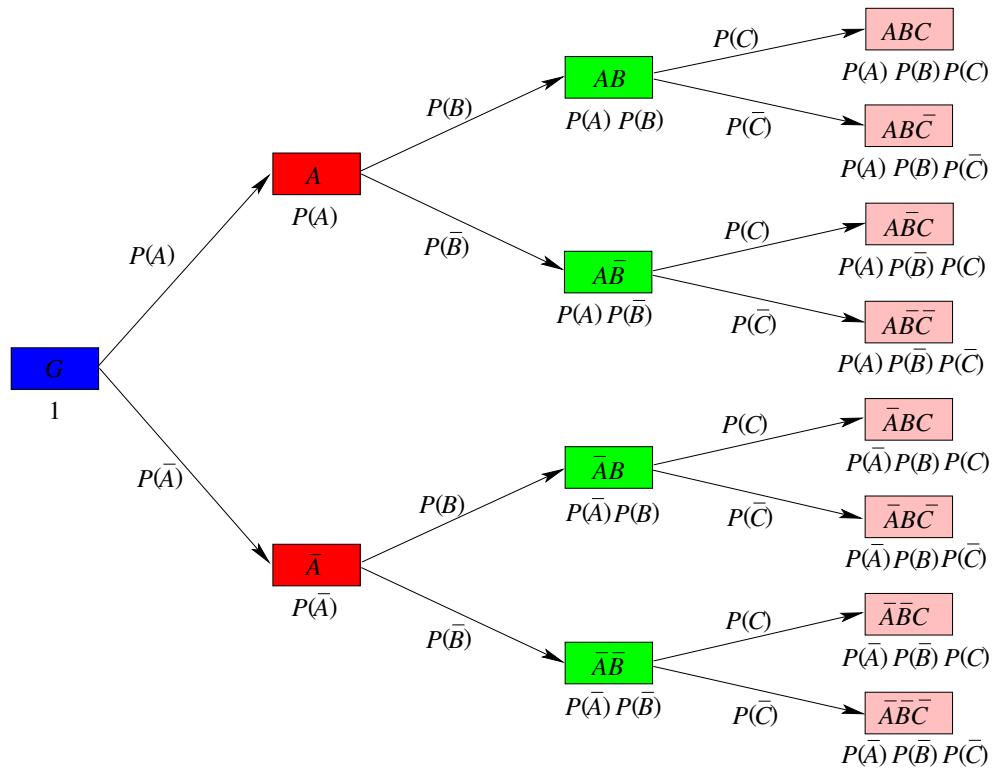
Pomembno je poudariti razliko med *nezdružljivostjo* in *neodvisnostjo*, zato za šalo zopet “recitirajmo”:

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

$$\begin{array}{ccc}
 (A, B) & \xrightarrow{*} & A B \\
 \left(P(.), P(.) \right) \downarrow & & \downarrow P(.) \\
 (P(A), P(B)) & \xrightarrow{*} & P(A)P(B) = P(AB)
 \end{array}$$

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

Za neodvisne dogodke A_i , $i \in I$ velja $P(A_j) = P(A_j / \prod_{i=1}^{j-1} A_i)$ za vsak $j \in I$.

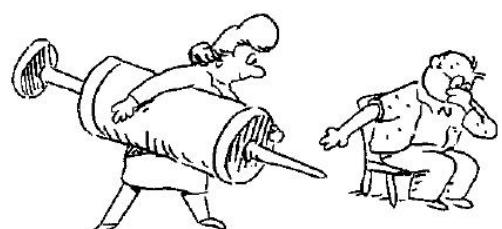
Slika 3.6: Binarno drevo za tri neodvisne dogodke A , B in C .

3.3 Dvostopenjski poskusi in izrek o popolni verjetnosti

Primer: Redko nalezljivo bolezen dobi ena oseba na 1 000.

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 1 000 se naze),
 $P(B/A) = 0.99$ (test pravilno označi okuženega),
 $P(B/\bar{A}) = 0.02$ (test napačno označi zdravega).

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

1. način (algebraični pristop): Iz (3.3) dobimo: $P(A/B) = P(A)P(B/A)/P(B)$, od koder je razvidno, da potrebujemo 'le še' $P(B)$. Spomnimo se popolnega sistema dogodkov, ki ga v našem primeru lahko predstavljata nasprotna dogodka A in \bar{A} . Potem velja:

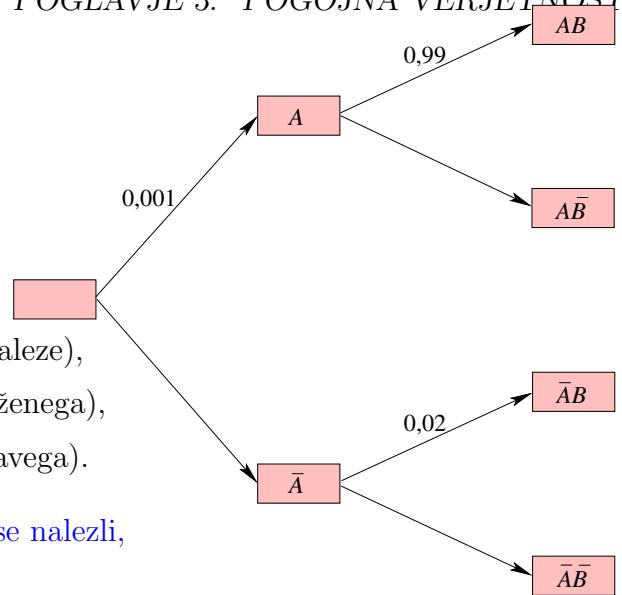
$$P(B) = P(BG) = P(B(A + \bar{A})) = P(BA) + P(B\bar{A})$$

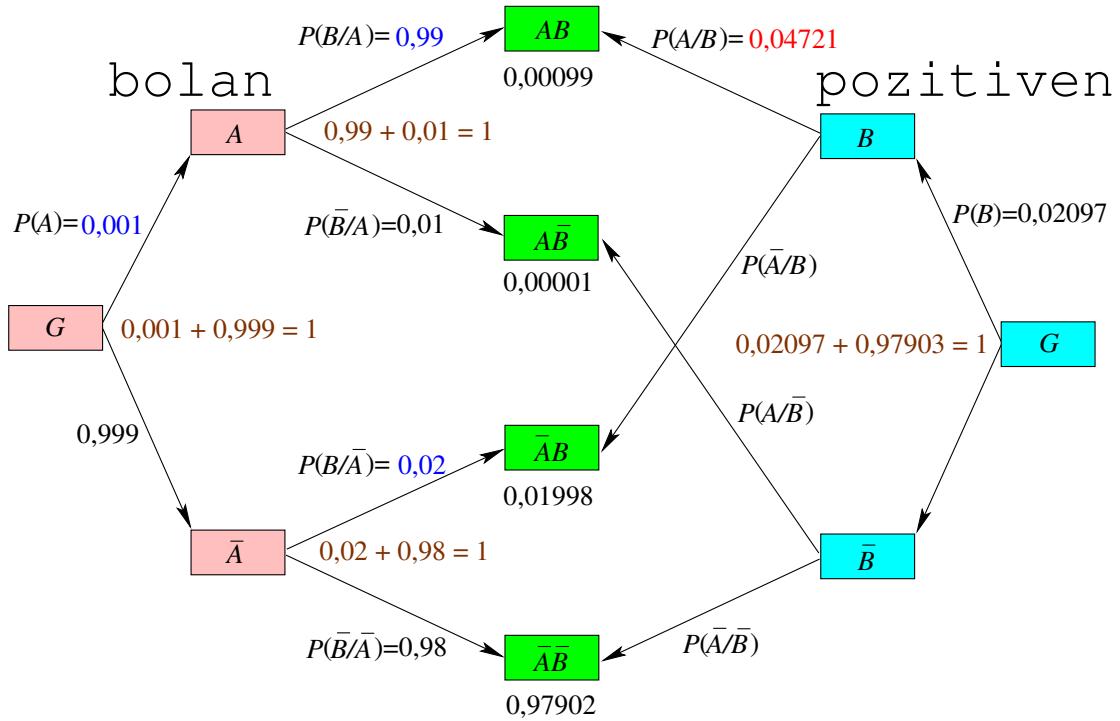
(zadnji enačaj valja zato, ker sta dogodka BA in $B\bar{A}$ nezdružljiva). Ali bi znali izračunati $P(BA)$? Pomagamo si lahko s formulo za pogojno verjetnost (ne morda prvo, ki nam pride na misel, pač pa tisto drugo, tj. 'dualno', glej (3.2)):

$$P(BA) = P(A)P(B/A) = 0.001 \times 0.99 = 0.00099.$$

Spomnimo se tudi, da je $P(\bar{A}) = 1 - P(A) = 0.999$ in že na enak način kot v prejšnjem primeru, ko smo računali $P(BA)$, dobimo: $P(B\bar{A}) = P(\bar{A})P(B/\bar{A}) = 0.999 \times 0.02 = 0.01998$, kar nam da $P(B) = 0.00099 + 0.01998 = 0.02097$ in $P(A/B) = 0.00099/0.02097 = 0.04721$.

2. način (grafični pristop): Tokrat si pomagamo z binarnim drevesom na prejšnji strani. Za izračun $P(\bar{A})$ uporabimo Posledico 2.6. Enako sklepamo tudi na naslednjem nivoju (z leve proti desni), le da se tokrat sklicemo na Trditev 3.2, in že imamo verjetnosti na vseh puščicah na naslednjem nivoju. Verjetnost vsakega izmed dogodkov v zelenih škatlah je enaka produktu verjetnosti na puščicah od začetka (koren na lev) pa do samega dogodka, kar nam zagotavlja identiteta (3.2). Glej sliko 3.9.



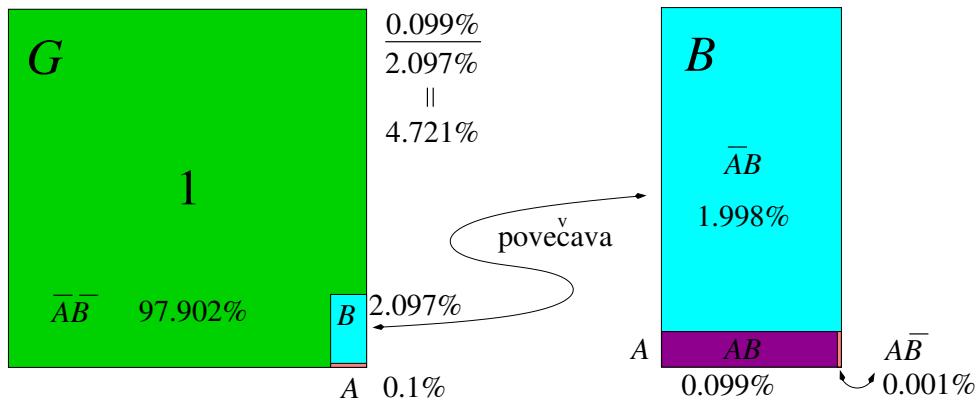


Slika 3.9: Popoln sistem dogodkov AB , $A\bar{B}$, $\bar{A}B$ in $\bar{A}\bar{B}$ je označen z zeleno, njihove verjetnosti pa so $P(AB) = 0 \cdot 001 \times 0 \cdot 99 = 0 \cdot 00099$, $P(A\bar{B}) = 0 \cdot 001 \times 0 \cdot 01 = 0 \cdot 00001$, $P(\bar{A}B) = 0 \cdot 999 \times 0 \cdot 02 = 0 \cdot 01998$ in $P(\bar{A}\bar{B}) = 0 \cdot 999 \times 0 \cdot 98 = 0 \cdot 97902$. in so zapisane pod njihove škatle. Za vsak slučaj lahko preverimo, da je vsota pravkar izračunanih verjetnosti res enaka 1: $0 \cdot 00099 + 0 \cdot 00001 + 0 \cdot 01998 + 0 \cdot 97902 = 1$.

Verjetnost dogodka B je torej enaka $0 \cdot 00099 + 0 \cdot 01998 = 0 \cdot 02097$. Končno si pogledamo še pot, ki nas pripelje iz desnega korena do dogodka AB in upoštevamo še identiteto (3.1): $0 \cdot 00099 / 0 \cdot 02097 = 0 \cdot 04721$.

Verjetnost $P(\bar{A}/B)$, tj. da smo zdravi, kljub temu, da je bil test pozitiven, je $1 - 0 \cdot 04721 = 0 \cdot 95279$. Naredimo še en preiskus: $0 \cdot 02097 \times (1 - 0 \cdot 04721) = 0 \cdot 01998$, a premislite sami, kaj smo v resnici preverili. Verjetnost $P(\bar{B})$, tj. da bo test negativen, je $0 \cdot 97903$, verjetnost $P(\bar{A}/\bar{B})$, tj. da smo zdravi, če je bil test negativen, pa približno $1 - 0 \cdot 00001 = 0 \cdot 99999$, kar je naravnost odlično.

Premislimo še malo o problematični pogojni verjetnosti $P(A/B)$. Ali je res možno, da je verjetnost $P(A/B)$, tj. da smo se nalezli, če je test pozitiven, tako majhna? Kako si lahko to razlagamo?



Slika 3.10: Čeprav imata oba dogodka A in B zelo majhno vejetnost, je pogojna verjetnost dogodka A glede na dogodek B lahko tudi majhna.

Če je testiranje drago, potem lahko najprej izvedemo cenejši test in glede na to, kaj pokaže, gremo na dražjega. \diamond

Naj bo H_i , $i \in I$ (kjer je I množica indeksov, običajno kar $1, 2, \dots, n$), popoln sistem dogodkov, tj. **razbitje** gotovega dogodka: $\sum_{i \in I} H_i = G$, na paroma nezdružljive dogodke: $H_i H_j = N$, $i \neq j$. Gotov dogodek smo torej kot hlebec narezali z domnevami 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 $i \in I$:

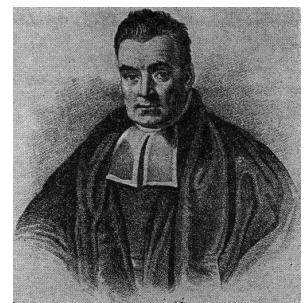
$$A = A(H_1 + H_2 + \dots + H_n) = A H_1 + \dots + A H_n.$$

Torej si bomo ogledali verjetnost dogodka A na posameznih ‘kosih’, skoraj tako kot pri marmornem kolaču. Ker so tudi dogodki $A H_i$ paroma nezdružljivi, velja:

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

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

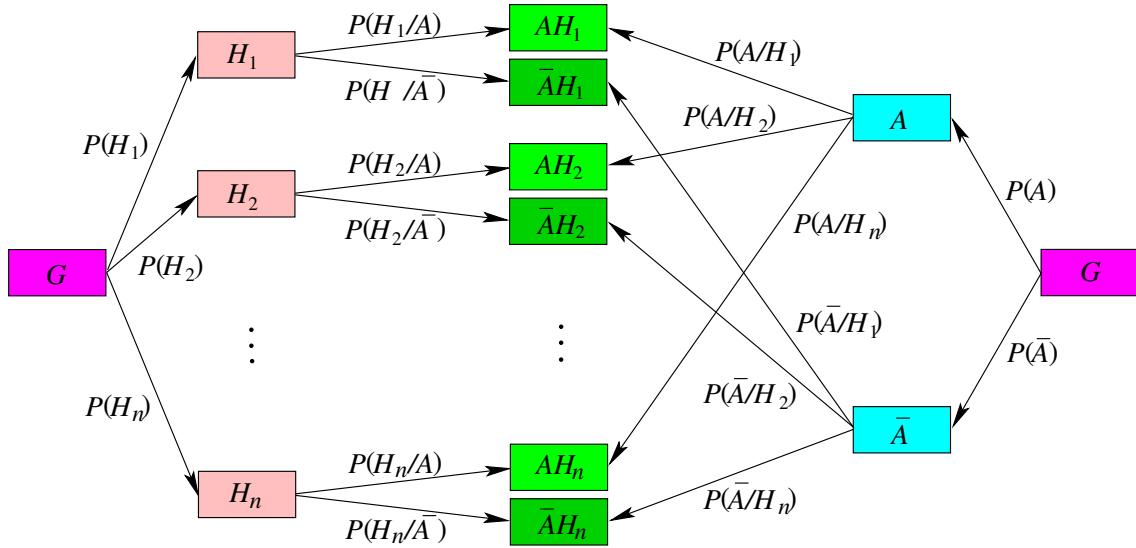
Zgornji trditvi pravimo tudi *izrek o popolni verjetnosti*, formuli pa *formula za popolno verjetnost* (tudi obrazec za razbitje). Na to lahko pogledamo tudi kot na večstopenjski poskus: v prvem koraku se zgodi natanko eden od dogodkov H_i , ki ga imenujemo domneva (hipoteza) (domneve sestavlja popoln sistem dogodkov). Šele izidi na prejšnjih stopnjah določajo, kako bo potekal poskus na naslednji stopnji. Omejimo se na poskus z dvema stopnjama oz. t.i. **dvo-stopenjski poskus**. 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.



REV. T. BAYES
Thomas Bayes, angleški župnik (1702-1761).

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

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



Slika 3.11: Vsak stolpec škatelj predstavlja popoln sistem dogodkov. V stolcu zelenih škatelj svetlejše predstavlja popoln sistem dogodkov za A , preostale pa popoln sistem dogodkov za \bar{A} .

Primer: Trije lovci so hkrati ustrelili na divjega prasiča 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, (c) 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 domnev (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:

$$\begin{aligned} P(\text{nobena krogla ni zadela}) &= 0.192 = P(N'), \\ P(\text{dve krogli sta zadeli}) &= 0.192 + 0.072 + 0.032 = 0.296 = P(D), \\ P(\text{tri krogle so zadele}) &= 0.048 = P(T). \end{aligned}$$

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

$$\begin{aligned} P(H_5/E) &= \frac{P(H_5E)}{P(E)} = \frac{0.048}{0.464} = 0.103 = P(\text{prvi je zadel}/E), \\ P(H_6/E) &= \frac{P(H_6E)}{P(E)} = \frac{0.128}{0.464} = 0.276 = P(\text{drugi je zadel}/E), \\ P(H_7/E) &= \frac{P(H_7E)}{P(E)} = \frac{0.288}{0.464} = 0.621 = P(\text{tretji je zadel}/E). \end{aligned}$$

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 utegnili na hitro pomisliti).

Kako bi si pravično razdelili plen, če bi v divjim prašiču našli dve krogli? ◇

Poglavlje 4

Bernoullijevo zaporedje neodvisnih poskusov



Zanimajo nas oblike histogramov različnih podatkov. Oblike v dveh razsežnostih opisujemo s funkcijami, Eulerjev zapis $y = f(x)$. Najenostavnejša funkcija je konstantna funkcija $y = n$, kjer je n neka konstanta. Le-to predstavimo v koordinatnem sistemu z vodoravno premico, ki gre skozi točko $(0, n)$. Naslednja je linearna funkcija, tj. premica $y = kx + n$, kjer je k smernostni koeficient (tj. tangens naklonskega kota premice), n pa odsek na y osi, ki ga odreže premica. Sledijo polinomi (npr. parabola), a vseeno še nismo srečali zvonaste krivulje, katera naj bi bila najbolj pogosta oblika podatkov. V ta namen se spomnimo metanja kovanca in poskusimo najti obliko histograma.

Primer: Bodimo bolj natančni. Naj bo n neko naravno število. Kovanec vržemo n -krat. Zanimajo nas frekvence, tj. na koliko načinov bo padel grb natanko k -krat ($0 \leq k \leq n$). Torej gre za čisto konkreten primer ne pa teorijo, kot na prvi pogled izgleda. Takšen poskus/simulacijo lahko izvedemo ročno ali z računalnikom.

Če želite bolj ilustrativen primer, si predstavljate prijatelja, ki ima težave z ravnotežjem in mu hoja naravnost dela težave. Namesto, da bi hodil naravnost, na vsakem koraku stopi naključno bodisi pod kotom 45 stopinj v levo ali v desno (če smo bolj konkratni, iz točke $(0, 0)$ stopi bodisi v točko $(1, 1)$ ali pa v točko $(1, -1)$). Po n korakih se nahaja na premici $x = n$, nekje med točkama (n, n) in $(n, -n)$. Verjetnost, da pride v katero izmed pravkar omenjenih točk je izredno majhna, saj lahko pride tja smo tako, da se na vsakem koraku odloči stopiti v desno (ali pa v levo), namreč $1/2^n$. Bolj gremo proti sredini daljice, ki jo določata ti točki, večja je verjetnost, saj je načinov vse več in več. ◇

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 $P(\bar{A}) = 1 - P(A) = 1 - p = q$.



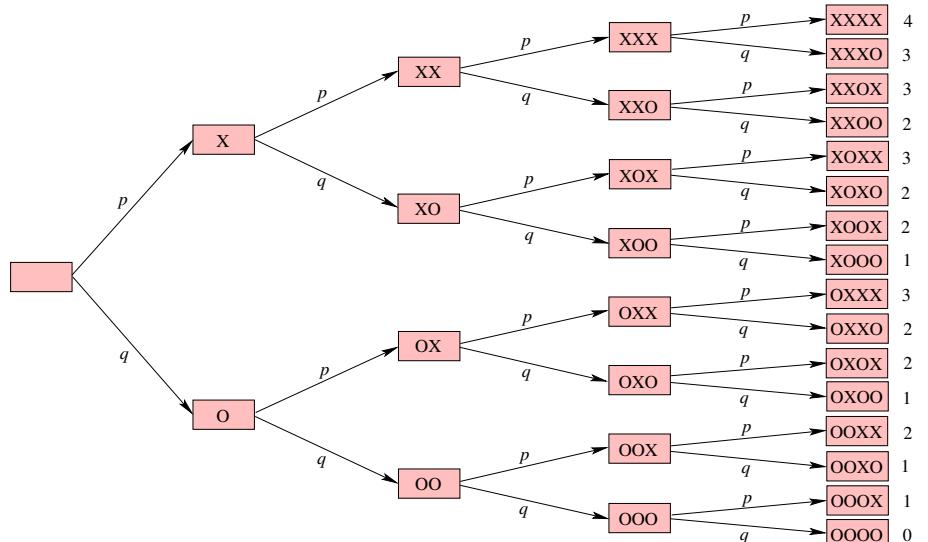
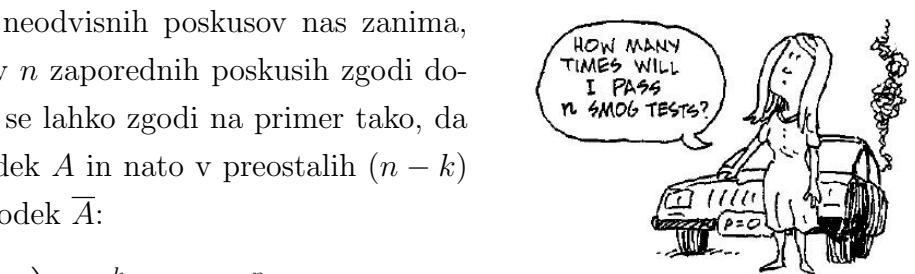
JAKOB BERNoulli um. 1687

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} :

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

Dogodek, 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 toliko kolikor je kombinacij $\binom{n}{k}$. In ker je po binomskem obrazcu $\binom{n}{0} + \binom{n}{1} + \dots + \binom{n}{n} = 2^n$, nismo prav nobenega pozabili. Ker so ti načini nezdružljivi med seboj, dobimo **Bernoullijev obrazec**, tj. verjetnost tega dogodka je enaka



Slika 4.3: Iz binarnega drevesa je očitno, da je število možnih izidov 2^n , če se seveda ustavimo po n zaporednih poskusih. Verjetnost, da pridemo do določenega vozlišča v drevesu je seveda enaka produktu verjetnosti na povezavah od korenja (začetka) pa vse do tega vozlišča.

$$P_n(k) = \binom{n}{k} p^k (1 - p)^{n-k}. \quad (4.1)$$

4.1 Računanje verjetnosti $P_n(k)$

Primer: V povprečju med prehitrimi vozniki, ki jih ujame radar, 32% prekorači hitrost do 10 km/h. Izračunajmo verjetnost, da bo pri nadzoru prometa med 45imi prehitrimi vozniki, natanko 15 voznikov prekoračilo hitrost do 10 km/h:

$$P_{45}(15) = \binom{45}{15} 0.32^{15} 0.68^{45-15} = \frac{45!}{15! \cdot 30!} 0.32^{15} 0.68^{30}.$$

Če nismo pazljivi, se pri računanju pojavi problem, saj je število $45! \doteq 1 \cdot 196 \ 222 \ 209 \cdot 10^{56}$ izredno veliko in lahko našemu računalu zmanjka mest. Glede na to, da je binomski koeficient naravno število, ga lahko izračunamo čisto natančno (pa čeprav z nekaj dela, je pa lažje, če ulomek najprej pokrajšamo)

$$\binom{45}{15} = \frac{45 \cdot 44 \cdot 43 \cdot 42 \cdot 41 \cdot 40 \cdot 39 \cdot 38 \cdot 37 \cdot 36 \cdot 35 \cdot 34 \cdot 33 \cdot 32 \cdot 31}{15 \cdot 14 \cdot 13 \cdot 12 \cdot 11 \cdot 10 \cdot 9 \cdot 8 \cdot 7 \cdot 6 \cdot 5 \cdot 4 \cdot 3 \cdot 2 \cdot 1} = 344 \ 867 \ 425 \ 584.$$

Poskusimo nadaljevati z takšno natančnostjo vse do konca

$$P_{45}(15) = \frac{99 \ 419 \ 052 \ 833 \ 245 \ 159 \ 660 \ 807 \ 119 \ 139 \ 225 \ 127 \ 698 \ 425 \ 137 \ 346 \ 221 \ 902 \ 934 \ 310 \ 912}{807 \ 793 \ 566 \ 946 \ 316 \ 088 \ 741 \ 610 \ 050 \ 849 \ 573 \ 099 \ 185 \ 363 \ 389 \ 551 \ 639 \ 556 \ 884 \ 765 \ 625}.$$

Gotovo niste navdušeni, pa zapišimo ulomek še nekoliko drugače

$$P_{45}(15) = \frac{2^{49} \cdot 3 \cdot 11 \cdot 17^{31} \cdot 19 \cdot 31 \cdot 37 \cdot 41 \cdot 43}{5^{90}} = 2^{139} \cdot 3 \cdot 11 \cdot 17^{31} \cdot 19 \cdot 31 \cdot 37 \cdot 41 \cdot 43 \cdot 10^{-90}.$$

To je res nekoliko lepše, a nas vseeno zanima le ocena in če zaupamo našemu računalu, potem dobimo končno $P_{45}(15) \doteq 0.123$ (do istega rezultata bi prišli tudi z nekaj vodilnimi števkami, npr. 994/8078, a je lahko biti pameten, ko že imamo rezultat).

Sedaj pa izračunajmo verjetnost, da bo pri nadzoru prometa med 45imi prehitrimi vozniki, od (vključno) 12 do 15 voznikov prekoračilo hitrost do 10 km/h. Tudi na to vprašanje lahko hitro odgovorimo:

$$P_{45}(12) + P_{45}(13) + P_{45}(14) + P_{45}(15),$$

vendar pa je tu računanja še več, saj moramo izračunati še tri podobna števila, kot smo jih zgoraj, če želimo priti do konkretnih vrednosti. \diamond

Ena možnost za računanje $P_n(k)$ je uporaba rekurzije:

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

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

Vendar pa je takšno računanje izredno zamudno (eksponentno), tako kot rekurzivno računanje faktorjela, glej razdelek A.6. Kako bi se počutili šele, če bi morali izračunati v zgornjem primeru npr. $P_{45}(5) + \dots + P_{45}(25)$, da ne govorimo o tem, da bi bilo lahko teh sumandov kaj hitro eksponentno mnogo (glede na dolžino zapisa števila n)?

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

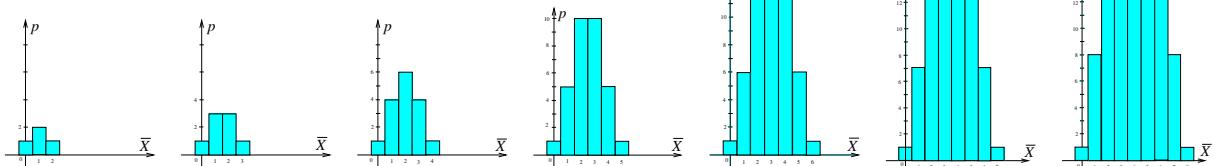
```
> dbinom(50, size=1000, prob=0.05)
[1] 0.05778798
```

Kaj pa, če so bili vsi primeri doslej bolj šolski, in imamo opravka z več kot recimo 6-mestnimi števili? Potem se moramo domisliti kakšnega bolj učinkovitega algoritma in temu je posvečen naslednji razdelek.

4.2 Usojena krivulja v Pascalovemu 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 4.5).

Ali bi znali ugotoviti kakšno pravilo,
ki vlada obliki teh histogramov
(glej sliko 4.5)?



Slika 4.5: Predstavitev 2., ..., 7. in 8. vrstice Pascalovega trikotnika v obliki histogramov. Vsota višin stolpcev v posameznem histogramu je zaporedno 4, 8, 16, 32, 64, 128 in 256, višina najvišjega stolpca pa je 2, 3, 6, ...

Eno je gotovo, višina histograma vrtoglavu 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:

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

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 namesto $\binom{n}{k}$ 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 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 decimalko 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 korakih pride do 0.7979 ali, če smo še bolj potprežljivi do $0.7978646\dots$. 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 z njegovo (študentom dobro znano) oceno 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 zgornje 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 spremenjali, smo končno izpolnili obe zahtevi (konstantna ploščina in konstantna maksimalna višina histograma). 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 rastочim 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 oziroma v levo, če se vozimo po krivulji z leve proti desni. Pri tem si pomagamo z opazovanjem spreminjaanja 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

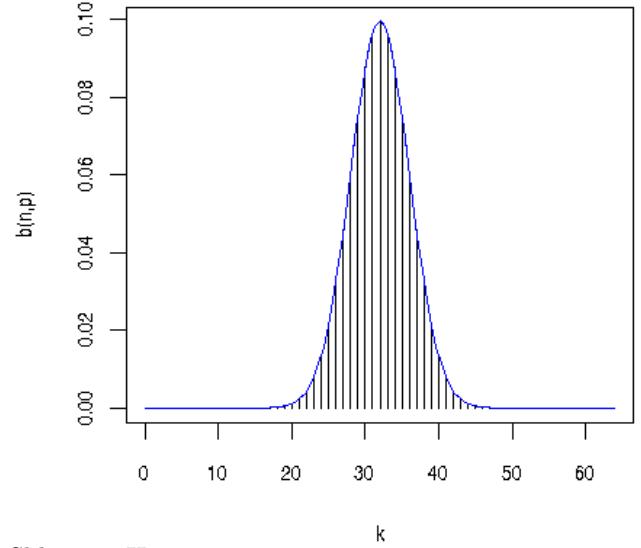
$$d_i = -d_{n-1-i}, \quad 0 \leq i \leq n-1,$$

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

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

oziroma $d_0 < d_1 < \dots < d_{k-1} \leq d_k > d_{k+1} > d_{k+2} > \dots > d_{\lfloor (n-1)/2 \rfloor}$, kjer je $k := \lfloor (n - \sqrt{n+2})/2 \rfloor$. Enačaj velja natanko tedaj, ko je $n+2$ popoln kvadrat in je $h = k$.



Slika 4.7: Histogram za $n = 64$.
Stolpce smo nadomestili kar z intervali.

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

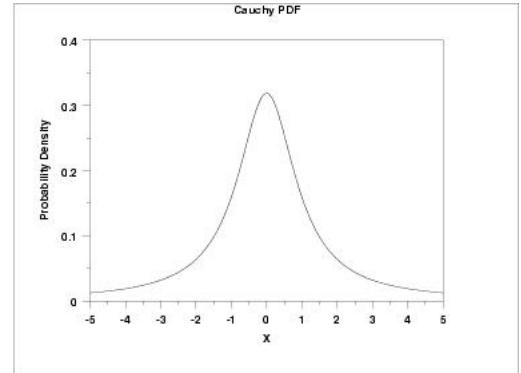
$$d_h - d_{h-1} = \frac{n!(n-2h-1)}{(h+1)!(n-h)!} - \frac{n!(n-2h+1)}{h!(n-h+1)!} = \frac{n!}{(h+1)!(n-h+1)!} \Delta$$

odvisen samo od predznaka razlike

$$\Delta := (n-2h-1)(n-h+1) - (n-2h+1)(h+1) = n^2 - 4nh + 4h^2 - n - 2 = (n-2h)^2 - (n+2). \quad \square$$

Zopet smo pred novim presenečenjem: **vse popravljenje krivulje so si sedaj na moč podobne!** Za katero obliko/krivuljo pa gre? Če smo prej našli konstanto c , potem se moramo sedaj nekako podviziati 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 + x²)** prepustimo bralcu (tudi ta ni prava, a je matematična utemeljitev nekoliko bolj zapečetena). V resnici naša funkcija ni niti racionalna (kvocient dveh polinomov).



Slika 4.8: Graf funkcije $1/(\pi(1 + x^2))$, kjer smo pod ulomkovo črto dodali še π , da bo ploščina pod krivuljo enaka 1.

Morda pa je usoda hotela, da na ta način spoznamo novo funkcijo. Zaenkrat jo poimenujmo kar **usojena** funkcija (da ne bomo pretiravali s tujkami - **fatalka**).

Sledimo Abrahamu de Moivru (The Doctrine of Chance, 1733), ki predлага, da vpeljemo

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

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

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

Z nekaj računske spremnosti in uporabo vrste za $\ln(1 + x)$ (glej škatlo) pridemo do

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

Izpeljali smo naslednjo trditev.

Izrek 4.3. (De Moivrov točkovni obrazec)

Za velike n velja

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



De Moivrov točkovni obrazec je poseben primer **Laplaceovega točkovnega obrazca**. Slednjega smemo uporabljati, ko je n velik in p blizu $1/2$:

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



Skica dokaza De Moivrovega točkovnega obrazca

Vpeljimo zamenjavo $n = 4m^2$ (da se znebimo korenov in ulomkov) in si približo oglejmo 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 preoblikujmo 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 z $m \rightarrow \infty$ 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 $(N + 1)/2$. □

Poglavlje 5

Slučajne spremenljivke in porazdelitve



Najbolj značilne oblike histogramov opišemo z grafi funkcij.

5.1 Slučajne spremenljivke

Denimo, da imamo poskus, katerega izidi so števila (npr. pri metu kocke so izidi števila pik). Se pravi, da je poskusom pripojena 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**. Vpeljemo/predstavimo jo tako, da povemo

1. kakšne vrednosti more imeti (*zaloga vrednosti*, oznaka \mathcal{Z}) 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 (če $x_i \notin Z$ gre za nemogoč dogodek).

Primer: (i) Trikrat vržemo pošten kovanec. X je število padlih grbov.

(ii) Streljamo v tarčo, Y pa predstavlja razdaljo zadetka od središča tarče.

(iii) Z je število klikov na določeni spletni strani v časovni enoti. \diamond

Porazdelitveni zakon slučajne spremenljivke X je poznan, če je mogoče za vsako realno število x določiti verjetnost

$$F(x) = P(X \leq x).$$

Predpis $F(x)$ ali bolj natančno $F_X(x)$ imenujemo **porazdelitvena funkcija** (tudi kumulativna 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 (bodisi končnega ali neskončnega).

V nekaterih primerih je slučajna spremenljivka lahko tudi kombinirana, tj. na nekem območju diskretna in drugje zvezna.

Izrek 5.1. [Lastnosti porazdelitvene funkcije F]

1. Funkcija F je definirana na vsem \mathbb{R} in zanjo velja $0 \leq F(x) \leq 1$ za vsak $x \in \mathbb{R}$.
2. Funkcija F je nepadajoča: $x_1 < x_2 \implies F(x_1) \leq F(x_2)$.
3. $F(-\infty) := \lim_{x \rightarrow -\infty} F(x) = 0$ in $F(\infty) := \lim_{x \rightarrow \infty} F(x) = 1$.
4. Funkcija je v vsaki točki z desne zvezna: $F(x+) := \lim_{0 \leq h \rightarrow 0} F(x+h) = F(x)$.
5. Funkcija ima lahko v nekaterih točkah skok. Vseh skokov je največ števno mnogo.
6. $P(x_1 < X \leq x_2) = F(x_2) - F(x_1)$.
7. $P(x_1 \leq X \leq x_2) = F(x_2) - F(x_1+)$.
8. $P(X > x) = 1 - F(x)$.
9. $P(X = x) = F(x) - F(x-)$. □

5.2 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

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

sestavlja popoln sistem dogodkov. Označimo verjetnost posameznega dogodka s p_i , tj.

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

Vsota verjetnosti vseh dogodkov je enaka 1:

$$p_1 + p_2 + \dots + p_m + \dots = 1. \quad (5.1)$$

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:

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

Porazdelitveno funkcijo v diskretnem primeru lahko izrazimo na naslednji način

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

Je torej vsota vseh verjetnosti, ki pripadajo vrednostim, do neke določene vrednosti (kumulativa). Pare $(x_1, p_1), (x_2, p_2), \dots, (x_m, p_m), \dots$ lahko predstavimo v ravnini s točkami in se tako približamo histogramom ter razmišljamo, kakšne oblike je katera porazdelitev.

Pričakovana vrednost (končne) diskretne spremenljivke X , oznaka $\mathbb{E}(X)$ oz. μ_X , je posplošitev povprečne vrednosti. Začnemo z **uteženim povprečjem**, kjer je $k_1 + \dots + k_m = N$ in upoštevamo $f_i = k_i/N$:

$$\bar{X} = \frac{x_1 k_1 + \dots + x_m k_m}{N} = x_1 f_1 + \dots + x_m f_m$$

na osnovi česar vpeljemo

$$\mathbb{E}(X) = x_1 p_1 + x_2 p_2 + \dots + x_m p_m + \dots$$

Primer: Recimo, da želimo organizirati igro, pa nas zanima, pod kakšnimi pogoji nam bo na dolgi rok prinašala dobiček, igralci pa bodo še vedno zainteresirani za igranje. \diamond

V nadaljevanju si poglejmo nekaj primerov diskretnih slučajnih spremenljivk in za kakšno od njih izračunajmo še njeno pričakovano vrednost. Slednje sicer ne bo vedno enostavno opravilo, sicer pa to niti ni čudno, saj je pričakovana vrednost lahko tudi neskončna (kot bomo videli nekaj poglavij kasneje), a za končne zaloge vrednosti si nam s tem ni potrebno beliti glave.

5.2.1 Enakomerna diskretna porazdelitev – $U(n)$

Diskretna slučajna spremenljivka se porazdeljuje **enakomerno** (angl. uniform, zato oznaka U), če so vse njene vrednosti enako verjetne. Slučajna spremenljivka, ki je porazdeljena enakomerno, mora imeti vedno končno zалогу vrednosti, tj. $n < \infty$, pri čemer je $n := |\mathcal{Z}|$. V primeru enakomerne diskretne porazdelitve se pričakovana vrednost slučajne spremenljivke ujema s povprečjem njene zaloge vrednosti. **Kakšna bi izgledala njena grafična predstavitev?** Narišimo jo v primeru $n = 6$. Lahko bi rekli, da izgleda precej vodoravno (za večje n bi bilo to

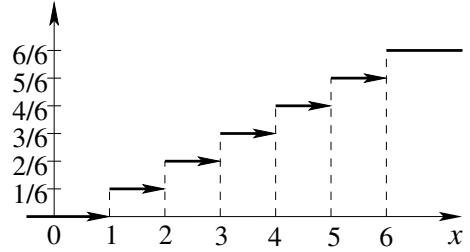


Primer: Naj bo X slučajna spremenljivka, ki beleži število pik pri metu kocke. Če je kocka poštena, potem je

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

in smo dobili primer enakomerno porazdeljene slučajne spremenljivke. Na desni je predstavljen graf njene porazdelitvene funkcije.

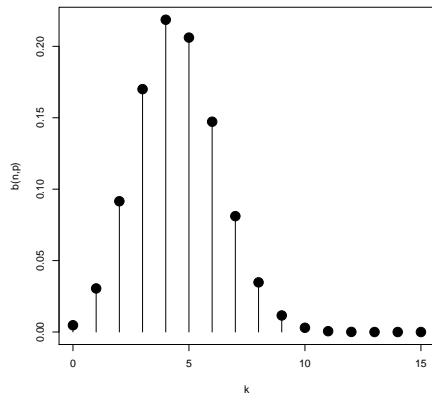
(Še bolj enostaven primer bi bil met poštenega kovanca. Naredite ga sami!) \diamond



5.2.2 Binomska porazdelitev – $B(n, p)$

Spomnimo se Bernoullijevega zaporedja neodvisnih poskusov. Poskus ponovimo n -krat in opazujemo, kolikokrat se je pri tem zgodil dogodek A , katerega verjetnost je p .

Primer: Kovanec vržemo 10 krat. Kolikšne so verjetnosti, da pade 'cifra' 0-krat, 1-krat, 2-krat, vsaj 3-krat itd.? \diamond



```
> 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 (4.1):

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

$k = 0, 1, 2, \dots, n$, glej sliko 4.3. Binomska porazdelitev je natanko določena z dvema podatkom – parametromi: $n \in \mathbb{N}$ in $p \in [0, 1]$. Če se slučajna spremenljivka X porazdeljuje binomsko s parametromi n in p , zapišemo:

$$X \sim B(n, p).$$

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:

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

Na primer

$$P(X = 2) = P_4(2) = \binom{4}{2} \left(\frac{1}{2}\right)^2 \left(1 - \frac{1}{2}\right)^{4-2} = \frac{6}{16} = \frac{3}{8}.$$

Porazdelitev obravnavane slučajne spremenljivke X je v tem primeru simetrična (glede na premico $x = 2$). \diamond

Pokazati se da, da je binomska porazdelitev simetrična, le, če je $p = 0.5$, sicer je asimetrična. Zanima nas, koliko je njena pričakovana vrednost. Pri tem nam pride prav naslednji izrek.

Izrek 5.2. *Delovanje pričakovane vrednosti na (diskretne) slučajne spremenljivke je linearne, tj. $\mathbb{E}(aX + bY) = a\mathbb{E}(X) + b\mathbb{E}(Y)$.*

Dokaz. *Homogenost*, tj. $\mathbb{E}(aX) = a\mathbb{E}(X)$, sledi iz homogenosti seštevanja (tu gre v bistvu za izpostavljanje). Preostane nam še preverjanje *aditivnosti*, tj. $\mathbb{E}(X) + \mathbb{E}(Y) = \mathbb{E}(X + Y)$. Za $Z = X + Y$ velja $\mathbb{E}(Z) = \sum_h z_h P(Z = z_h)$. Zveza $z_h = x_i + y_j$ lahko velja pri fiksni h za različne i in j , a ker seštejemo po vseh možnostih za i in j , dobimo

$$\mathbb{E}(Z) = \sum_h z_h P(Z = z_h) = \sum_i \sum_j (x_i + y_j) P(\{X = x_i\} \cdot \{Y = y_j\}).$$

Ker vrstni red seštevanja ne vpliva na končni rezultat¹, velja nadalje

$$\mathbb{E}(Z) = \sum_i x_i \sum_j P(\{X = x_i\} \cdot \{Y = y_j\}) + \sum_j y_j \sum_i P(\{X = x_i\} \cdot \{Y = y_j\}).$$

Naj bo $q_j = P(Y = y_j)$ za $j = 1, 2, \dots$. Če vemo, da sta slučajni spremenljivki X in Y neodvisni, potem velja $P(\{X = x_i\} \cdot \{Y = y_j\}) = p_i \cdot q_j$ ter zaradi $\sum_i p_i = 1 = \sum_j q_j$ tudi

$$\mathbb{E}(Z) = \sum_i x_i p_i \sum_j q_j + \sum_j q_j y_j \sum_i p_i = \mathbb{E}(X) + \mathbb{E}(Y).$$

V primeru, da slučajni spremenljivki X in Y nista neodvisni, pa pridemo do istega zaključka z uporabo obrazca za popolno verjetnost, glej izrek 3.3 za $A = \{X = x_i\}$ (oz. $\{Y = y_j\}$) in popoln sistem dogodkov $\{Y = y_1\}, \{Y = y_2\}, \dots$ (oz. $\{X = x_1\}, \{X = x_2\}, \dots$):

$$\mathbb{E}(Z) = \sum_i x_i P(X = x_i) + \sum_j y_j P(Y = y_j) = \mathbb{E}(X) + \mathbb{E}(Y). \quad \square$$

Za $n = 1$ pravimo slučajni spremenljivki X , ki je porazdeljenja binomsko, **Bernoullijeva** (oz. indikatorska) slučajna spremenljivka. Naj bo X_i ($1 \leq i \leq n$) Bernoullijeva slučajna spremenljivka za i -to ponovitev poskusa, tj. X_i zavzame vrednost 1, če se je pri i -ti ponovitvi poskusa izbrani dogodek zgodil, in 0 sicer. Potem je

$$X = X_1 + \dots + X_n \quad \text{in} \quad \mathbb{E}(X_i) = 0 \cdot (1 - p) + 1 \cdot p = p,$$

po zgornjem izreku pa sledi tudi $\mathbb{E}(X) = np$.

Primer: Pri pošteni kocki je verjetnost, da pade šestica, enaka $1/6$. Po n metih pričakujemo, da bo šestica padla približno $(n/6)$ -krat, povsem pričakovani rezultat torej. \diamond

¹ Pri neskončnih vsotah se moramo na tem mestu sklicati na absolutno konvergenco, tj. upoštevati $\sum_{i=1}^{\infty} |x_i| p_i, \sum_{j=1}^{\infty} |y_j| q_j < \infty$ in zaradi trikotniške neenakosti še $\sum_{h=1}^{\infty} |z_h| P(Z = z_h) < \infty$.

5.2.3 Poissonova porazdelitev – $P(\lambda)$

Strokovnjaki za promet opazujejo neko križišče. Radi bi napovedali, kakšna je verjetnost, da v danem časovnem intervalu skozi križišče zapelje npr. 100 avtomobilov. Definirajo slučajno spremenljivko X = število avtomobilov, ki pridejo v križišče na uro, in v modelu privzamejo še dve predpostavki:

- vsaka ura je enaka kot vsaka druga ura (čeprav vemo, da to ne drži, saj je v konici več prometa kot npr. ponoči).
- če v enem časovnem intervalu pride veliko avtomobilov, to še ne pomeni, da bo podobno tudi v naslednjem časovnem intervalu (takšni dogodki so med seboj neodvisni).

Za začetek si s štetjem prometa pridobijo grobo oceno za pričakovano število avtomobilov na časovno enoto: $E(X) = \lambda$, npr. 9 avtomobilov na uro, in poiskusijo v modelu uporabiti binomsko porazdelitev $B(n, p)$. Zanjo velja $E(X) = np$, število n , ki predstavlja število ponovitev poskusa, je v našem primeru enako številu manjših časovnih enot, verjetnost posameznega dogodka p pa je enaka verjetnosti, da je v dani manjši časovni enoti prišel mimo vsaj en avto. Potem lahko sklepajo na naslednji način:

$$\lambda \text{ avtomobilov/uro} = (60 \text{ min./uro}) \times ((\lambda/60) \text{ avtomobilov/min.}),$$

formula za verjetnost v primeru binomske porazdelitve pa je enaka

$$P(X = k) = P_{60}(k) = \binom{60}{k} \left(\frac{\lambda}{60}\right)^k \left(\frac{1-\lambda}{60}\right)^{60-k}.$$

Vendar tu naletijo na problem, da gre v eni manjši časovni enoti mimo lahko več kot en avtomobil na minuto. V tem primeru štejejo, da se je dogodek zgodil enkrat, četudi gre v tisti minuti mimo npr. 5 avtomobilov. Ta problem rešijo tako, da zmanjšajo časovno enoto:

$$P(X = k) = P_{3600}(k) = \binom{3600}{k} \left(\frac{\lambda}{3600}\right)^k \left(\frac{1-\lambda}{3600}\right)^{3600-k}.$$

Če pa tudi to ni dovolj, zamenjajo 3600 še z večjim številom. Pa poglejmo, kaj se zgodi, če gre to število proti neskončno. Iz analize se spomnimo naslednje zveze, glej (A.2):

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

Limito verjetnosti $P_n(k)$ lahko izračunamo za velike n in $\lambda = np$ na naslednji način:

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

Izpeljali² smo:

Izrek 5.3. (Poissonov obrazec) Za velike n in majhne verjetnosti, tj. p blizu 0 velja

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

□

Dobljena limita predstavlja verjetnost pri porazdelitvi, ki jo želimo vpeljati v tem razdelku. **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. Uporabimo jo lahko tudi za število dogodkov v drugih intervalih, npr. razdalja, prostornina, ... Ima zalogo vrednosti $\{0, 1, 2, \dots\}$, njena verjetnostna funkcija pa je

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

kjer je $\lambda > 0$ dani parameter. 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 presegati števila Bernoullijevih poskusov n . Poissonov obrazec lahko sedaj zapišemo tudi v naslednji obliki: $B(n, p) \approx P(np)$.

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 $\lambda = 15 \times 3 = 45$. ◇

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).

² Za dokaz $\lim_{n \rightarrow \infty} F = e^0 = 1$, kjer je $F := n(n-1)\cdots(n-k+1)/n^k$, lahko uporabimo Stirlingovo aproksimacijo (A.4) na izrazu $\ln F = \ln(n!) - \ln[(n-k)!] - k \ln n$ in za velike n res dobimo:

$$\ln F \approx [n \ln n - n] - [(n-k) \ln(n-k) - (n-k)] - [k \ln n] = -\underbrace{\left(1 - \frac{k}{n}\right)}_{\rightarrow -1} \ln \underbrace{\left(1 - \frac{k}{n}\right)^n}_{\rightarrow -k} - k \approx k - k = 0.$$

5.2.4 Negativna binomska oziroma Pascalova porazdelitev – $P(m, p)$

Primer: Kolikokrat moramo vreči kocko, da z verjetnostjo vsaj 0,99 pričakujemo, da bo padla vsaj ena šestica? \diamond

Negativna binomska $\text{NegBin}(m, p)$, ki je znana tudi pod imenom **Pascalova porazdelitev** $P(m, p)$, ima zalogo vrednosti $\{m, m + 1, m + 2, \dots\}$, njena verjetnostna funkcija pa je

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

kjer je $0 < p < 1$ dani parameter za verjetnost dogodka A v posameznem poskusu. Opisuje porazdelitev potrebnega števila poskusov, da se dogodek A zgodi m -krat.



Če številu poskusov sledimo s slučajno spremenljivko X , potem verjetnost $P(X = k)$, da se bo pri k ponovitvah poskusa dogodek A zgodil v zadnjem poskusu ravno m -tič, izračunamo po zgornji formuli za p_k . Za $m = 1$ dobimo **geometrijsko porazdelitev $G(p)$** . Le-ta opisuje porazdelitev števila poskusov, da se dogodek A v zadnji ponovitvi poskusa 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. \diamond

Izračunajmo pričakovano vrednost za geometrijsko porazdelitev:

$$\mathbb{E}(X) = \sum_{i=1}^{\infty} i p q^{i-1} = p \sum_{i=1}^{\infty} (q^i)' = p \left(\sum_{i=1}^{\infty} q^i \right)' = p \left(\frac{1}{1-q} \right)' = \frac{1}{p}.$$

Primeri: **(P1)** Marko zbir kartice svojih športnih junakov. Celotno zbirko sestavlja n različnih kartic. **Poisci verjetnost, da bo moral kupiti vsaj t kartic, da bo sestavil celotno zbirko.** Lahko se vprašamo tudi, koliko je pričakovano število kartic, ki jih bo Marko kupil, da bo zbral vse kartice. Izkušnje kažejo, da je npr. za $n = 50$ potrebno kupiti 225 kartic, da zberemo vseh 50 različnih kartic. Glej tudi <http://www.presek.si/5/5-3-Batagelj.pdf>.

(P2) Zbiranje Shrek magnetkov: V trgovinah Spar je pred časom potekalo zbiranje magnetkov z junaki iz risanke *Shrek za vedno*, glej, <http://www.sparsi/spar/aktualnozakupce/Shrek.htm>. Ob nakupu nad 15 EUR je kupec prejel enega od 24 Shrek magnetkov. Magneti so v ovitkih, tako da ni možno izbiranje magnetkov, ki jih kupec še nima.



(P3) Zbiranje nalepk: V trgovinah lahko kupimo *Kraševe* čokoladice *Kraljestvo živali*. Vsaka čokoladica ima priloženo nalepko določene živali. Vseh nalepk je kar 250, posamezna pa stane okoli 0,40 EUR

(P4) Kolikokrat je potrebno povleči karto iz kupa z 52 kartami, da dobimo vsako karto enkrat? Eksperimentirajmo:³ 303, 253, 281, 218, 202, 144, 190, 294, 302, 194, 183, 209, 244, 177, 280, 221, 312, 234, 249, 245, 201, 240, 110, 224, 215, 201, 221, 203, 197, 391. Mimogrede, zares neverjetno, koliko časa se porabi za zadnje manjkajoče karte: npr. 35 za 26, 90 za 45 kart in vse ostalo za preostalih 7, 140 za 50, Nekdo mi je prišepnil, da bo potrebno potegniti karto **($n \ln n$)-krat**.⁴ Pa poglejmo, kaj pravi na to naš eksperiment. V povprečju dobimo

$$(303 + 253 + 281 + 218 + 202 + 144 + 190 + 294 + 302 + 194 + 183 + 209 + 244 + 177 + 280 + 221 + 312 + 234 + 249 + 245 + 201 + 240 + 110 + 224 + 215 + 201 + 221 + 203 + 197 + 391) / 30$$

ozziroma **231,267**, kar je nekaj več kot $52 \times \ln 52 \doteq 205$.

Naj bo T čas, da zberemo vseh n kartic in za $i = 1, 2, \dots, n$ naj bo T_i čas da zberemo i -to kartico, potem ko smo že zbrali $i - 1$ kartic. Potem so T in T_i za $i = 1, 2, \dots, n$, slučajne spremenljivke. Verjetnost, da bomo izbrali novo kartico, če smo jih izbrali že $i - 1$, je enaka $p_i = (n - i + 1)/n$. Torej je T_i slučajna spremenljivka z *geometrijsko porazdelitvijo* in pričakovano vrednostjo $1/p_i$. Velja pa tudi $T = T_1 + T_2 + \dots + T_n$. Zaradi linearnosti pričakovane vrednosti je

$$E(T) = E(T_1) + E(T_2) + \dots + E(T_n) = n \cdot \left(\frac{1}{1} + \frac{1}{2} + \dots + \frac{1}{n} \right) = n \cdot H_n,$$

kjer je **H_n harmonično število**. Če upoštevamo še asimptotično vedenje harmoničnih števil, dobimo

$$E(T) = n \cdot H_n = n \ln n + \gamma n + \frac{1}{2} + o(1), \quad \text{za } n \rightarrow \infty,$$

kjer je $\gamma \doteq 0,5772156649$ *Euler-Mascheronijeva konstanta*. Za $n = 52$ dobimo po tej formuli in zaokrožanju 236, kar je enako, kot če bi seštevali ulomke:

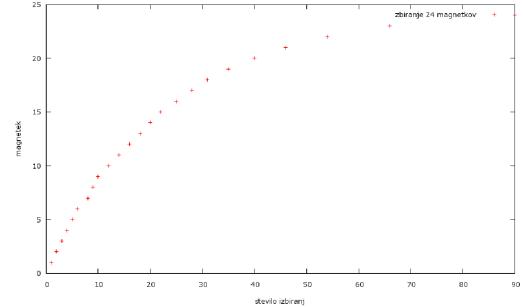
$$52 \times \left(1 + \frac{1}{2} + \frac{1}{3} + \dots + \frac{1}{52} \right) \doteq 236.$$

Analiza P2: Graf zbiranja magnetkov. Na x -osi je število potrebnih izbiranj za pridobitev magnetka, na y -osi pa števila magnetkov. Za pridobitev vrednosti je bil uporabljen simulator s 105 ponovitvami poskusa, napisan v Javi, naloga B11-63100014. V primeru, da bi ob



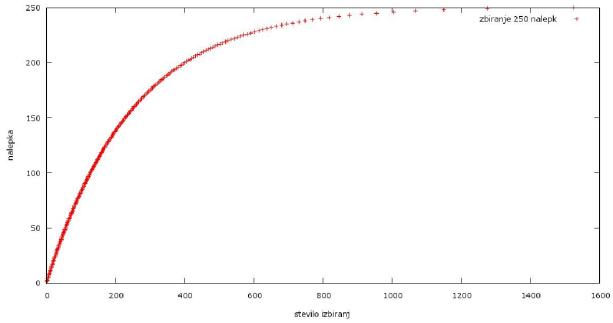
⁴ Kakšen nejeveren Tomaž bi se znal vprašati, od kje kar naenkrat logaritem.

vsakem nakupu porabili točno 15 EUR in bi ob 91. nakupu dobili zadnji magnetek, bi ob tem porabili **1365 EUR** (velika verjetnost je, da bi dejansko porabili več denarja, saj bi dostikrat ob nakupu porabili več kot 15 EUR in manj kot 30 EUR, tako da bi dobili samo en magnetek; če bi izbrali zelo pesimističen scenarij, pri katerem bi pri vsakem nakupu porabili 29 EUR in pri tem dobili en magnetek, bi za 91 nakupov porabili kar 2639 EUR).



Analiza P3: Na x -osi je število potrebnih izbiranj za pridobitev nalepke, na y -osi pa števila nalepk. V primeru, da ne bomo menjavali sličic ter bomo kupovali čokoladice v trgovini po

$0 \cdot 40$ EUR, je pričakovana vrednost porabljenih evrov **610 EUR**. Za primerjavo cena knjige **Živali – velika ilustrirana enciklopedija** znaša **120 EUR**. *Zbiranje nalepk se očitno ne splača.* Pa vendar ljudi zbiranje različnih predmetov veseli, saj je izziv (B11-63100014).



Nauk: Če se torej še vseeno raje odločite za zbiranje, pa je zelo smiselno, da proti koncu zbiranja nalepke izmenjujete z ostalimi. \diamond

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

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

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

za $\max(0, n - (N - M)) \leq k \leq \min(M, n)$ in $n \leq N$. 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. Zanima nas primer, ko je velikost n vzorca občutna glede na velikost populacije N , tj. $n/N > 0.05$. Sicer bi za velike N, M in $N - M$ lahko uporabljali kar binomsko porazdelitev, pri čemer bi vzeli $p = M/N$. Slednje pomeni, da je pri veliki seriji praktično vseeno ali izbiramo vzorec z vračanjem ali brez.

Če na neki množici z N elementi opazujemo neko lastnost A , se množica razdeli na 2 dela: na tiste elemente, ki to lastnost imajo, in tiste, ki je nimajo (npr., med serijo izdelkov se nahajajo neoporečni in oporečni izdelki).

Primer: Dnevna produkcija nekega obrata je 850 izdelkov, od katerih je 50 defektnih. Naključno izberemo 2 izdelka *brez vračanja*. Naj bo A dogodek, da je prvi izdelek defekten, B pa, da je drugi izdelek defekten. Potem velja $P(A) = 50/850$ in $P(B|A) = 49/849$. Informacija, da je prvi defekten, pomeni da je manj verjetno, da bo defekten tudi drugi, torej v tem primeru poskusi *niso* neodvisni (če pa bi prvi izdelek *vrnili*, preden bi izbrali drugega, bi bila dogodka neodvisna in bi šlo za binomsko porazdelitev). Naj bo X število defektnih izdelkov. Potem velja:

$$\begin{aligned} P(X = 0) &= (800/850)(799/849) = 0.886, \\ P(X = 1) &= (800/850)(50/849) + (50/850)(800/849) = 0.111, \\ P(X = 2) &= (50/850)(49/849) = 0.003. \end{aligned}$$

◇

Primer: Pošiljka rezervnih delov vsebuje 100 kosov od domačega dobavitelja in 200 kosov iz tujine.

- (a) Če izberemo naključno 4 rezervne dele, kakšna je verjetnost, da bodo vsi deli narejeni doma? $P(X = 4) = 0.0119$.
- (b) Kakšna je verjetnost, da sta dva ali več delov v vzorcu narejena doma?

$$P(X \geq 2) = P(X = 2) + P(X = 3) + P(X = 4) = 0.298 + 0.098 + 0.0119 = 0.408.$$

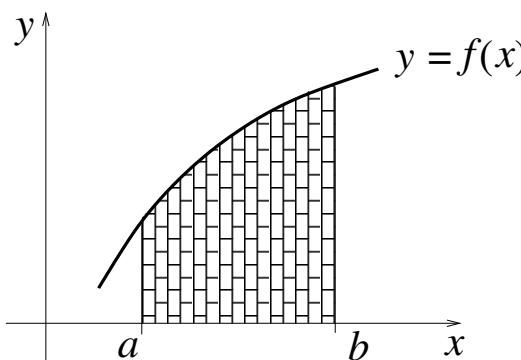
- (c) Kakšna je verjetnost, da je vsaj en del narejen doma? $P(X \geq 1) = 1 - P(X = 0) = 0.804$.

Primer: Mešalec dobro premeša 52 standardnih kart in nam jih dodeli pet.

- Kakšna je verjetnost, da dobimo med prejetimi kartami štiri ase?
- Kolikšno je pričakovano število asov, ki jih prejmemo?

◇

5.3 Ponovitev: integrali



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

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

Trditev 5.4. [Lastnosti določenega integrala]

(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. **A brez preplaha!** Iščemo le celotno ploščino pod krivuljo, od enega konca do drugega, le da konca pravzaprav sploh ni.



5.4 Zvezne slučajne spremenljivke

Slučajna spremenljivka X je **zvezno porazdeljena**, če obstaja taka integrabilna funkcija p ,⁵ imenovana **gostota verjetnosti**, da za vsak $x \in \mathbb{R}$ velja:

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

⁵V zahodni literaturi se pogosto namesto črke p uporablja kar črka f , kar je v skladu s prejšnjim razdelkom.

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 $p(x)$. Spomnimo se še Newton-Liebnitzove formule iz analize in zapišimo vse omenjene zveze:

$$\int_{-\infty}^{\infty} p(x) dx = 1, \quad P(x_1 < X \leq x_2) = \int_{x_1}^{x_2} p(t) dt = F(x_2) - F(x_1) \quad \text{ter} \quad p(x) = F'(x).$$

5.4.1 Enakomerna zvezna porazdelitev – $U(a, b)$

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

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

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

5.4.2 Normalna ali Gaussova porazdelitev – $N(\mu, \sigma)$



Leta 1738 je Abraham De Moivre (1667-1754) objavil aproksimacijo binomske porazdelitve, ki je normalna krivulja, glej prejšnje poglavje, konkretno De Moivrov točkovni obrazec in nje-govo posplošitev, Laplaceov točkovni obrazec. Leta 1809 je Carl Frederich Gauss (1777-1855) raziskoval matematično ozadje planetarnih orbit, ko je prišel do normalne porazdelitvene funkcije.

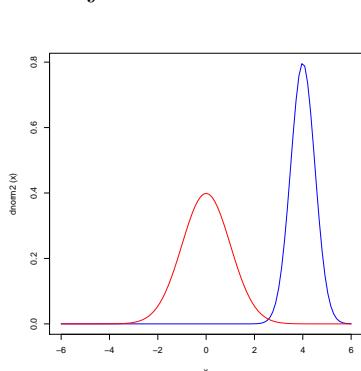
Zaloga vrednosti **normalno porazdeljene** slučajne spremenljivke so vsa realna števila, gostota verjetnosti pa je:

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

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

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

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



Na sliki je vrh krivulje predstavljen s točko $(\mu, 1/(\sigma\sqrt{2\pi}))$, v primeru rdeče krivulje, ko je $\mu = 0$ in $\sigma = 1$, pa je višina vrha enaka $1/\sqrt{2\pi} = 0.3989 \approx 0.4$. Mimogrede, tangenta na standardizirano krivuljo najbolj strma v točkah -1 in 1 , ko seče x os pod kotom približno 13.6° (enoti za x in y os sta torej precej različni).

Če je krivulja, ki ustreza gostoti porazdelitve, simetrična glede na premico $x = a$ za nek $a \in R$, potem velja $p(a-x) = p(a+x)$ za vsak $x \in \mathbb{R}$ in rečemo, da je porazdelitev **simetrična**. Ni se težko prepričati, da je pričakovana vrednost poljubne simetrične porazdelitve enaka mediani:

$$(\mu - x) p(\mu - x) + (\mu + x) p(\mu + x) = (\mu - x + \mu + x) p(\mu + x) = 2\mu p(\mu + x),$$

če gre pa za zvonasto porazdelitev, potem je enaka tudi modusu.

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

Spremenljivko $X \sim N(\mu, \sigma)$ pretvorimo s transformacijo $z = \frac{x - \mu}{\sigma}$ v standardizirano spremenljivko $Z \sim N(0, 1)$.

Laplaceov intervalski obrazec

Iz Laplaceovega točkovnega obrazca izhaja, da za p blizu $1/2$ in velike n velja:

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

Sedaj pa nas zanima še, 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,

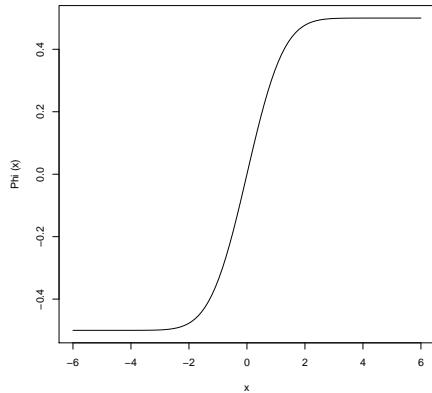
$$P_n(k_1, k_2) = \sum_{k=k_1}^{k_2-1} P_n(k) = \frac{1}{\sqrt{2\pi}} \sum_{k=k_1}^{k_2-1} e^{-\frac{1}{2}x_k^2} \Delta x_k.$$

Za (zelo) velike n lahko vsoto zamenjamo z integralom

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

Funkcijo napake imenujemo funkcijo

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



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

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 programih⁶, glej na učilnico ali pa npr. <http://www.mathsisfun.com/data/standard-normal-distribution-table.html>

```
> 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
```

Porazdelitveno funkcijo slučajne spremenljivke $X \sim N(\mu, \sigma)$ izrazimo s funkcijo napake na naslednji način:

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

ozioroma velja

$$P(x_1 \leq X < x_2) = \Phi\left(\frac{x_2-\mu}{\sigma}\right) - \Phi\left(\frac{x_1-\mu}{\sigma}\right). \quad (5.2)$$



C. F. Gauss

Izpeljimo še izrek, ki opravičuje statistično definicijo verjetnosti.

Izrek 5.5 (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*

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

⁶Bodite pozorni, da se v literaturi pojavlja tudi $\text{Erf}(x) = (2/\sqrt{\pi}) \int_0^x e^{-t^2} dt$, kar pomeni $\text{Erf}(x) = 2\Phi(x\sqrt{2})$. Da bo zmeda še večja, včasih najdemo tabelirano porazdelitveno funkcijo F standardizirane normalne porazdelitve, ki pa jo avtorji označijo kar s Φ , čeprav v tem primeru velja $F(x) = 0.5 + \Phi(x)$. Je pa tudi res, da ta razlika ne vpliva na formulo (5.2).

Bolj natančno, za dovolj velika naravna števila n velja

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

Dokaz. Ker je n naravno število, lahko oba izraza v neenakost iz zgornje verjetnosti pomnožimo z n , nato pa še odpravimo absolutno vrednost in z upoštevanjem, da je tudi k celo število med 0 in n , dobimo oceno:

$$P\left(\left|\frac{k}{n} - p\right| \leq \varepsilon\right) = P(np - n\varepsilon \leq k \leq np + n\varepsilon) = P_n(k_1) + P_n(k_1 + 1) + \cdots + P_n(k_2),$$

kjer so $k_1 < k_1 + 1 < \cdots < k_2$ vsa cela števila na intervalu $[np - n\varepsilon, np + n\varepsilon]$, tj. intervalu s središčem v točki np in radijem $n\varepsilon$. Dobljeno vsoto smo označili s $P(k_1 - 1, k_2)$ in jo ocenili s funkcijo napake, kar nam da:

$$P(k_1 - 1, k_2) \approx \Phi\left(\frac{k_2 - np}{\sqrt{npq}}\right) - \Phi\left(\frac{(k_1 - 1) - np}{\sqrt{npq}}\right) \approx 2\Phi\left(\frac{n\varepsilon}{\sqrt{npq}}\right),$$

saj sta razliki $k_2 - np$ in $(np - 1) - k_1$ za velike n približno enaki omenjenemu radiju, Φ pa je liha funkcija. Za konec upoštevamo še $\Phi(\infty) = 1/2$ ter da so števila ε , p in $q = 1 - p$ vnaprej izbrane konstante, število n pa pošljemo proti neskončno in zgornja limita je izračunana. \square

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 $p = 1/2$, $n = 3600$, $\varepsilon = 0·01$, tako da iz zgornje formule dobimo

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

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

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

(c) Zoran Dragić zadane na treningih 70% metov za tri. *Kaj je bolj verjetno, da zadane 10 metov od desetih ali 80 metov od stotih?* \diamond

5.4.3 Eksponentna porazdelitev – $\text{Exp}(\lambda)$

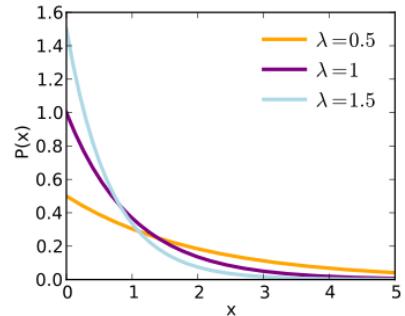
Eksponentna porazdelitev opisuje čas med dvema zaporednima dogodkoma v Poissonovemu procesu, tj. procesu, pri katerem se dogodki pojavljajo zvezno in neodvisno, pri konstantni povprečni hitrosti pojavljanja. Hkrati je tudi zvezni analog geometrijske porazdelitve.

Gostota **eksponentne porazdelitve**, ki ji pravimo tudi porazdelitev Poissonovega toka, je enaka

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

Potem za porazdelitveno funkcijo velja

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



Primeri: (a) Študenti prihajajo v klub porazdeljeni približno po Poissonu, s povprečno stopnjo 30 študentov na uro. **Kakšna je verjetnost, da bo vratar čakal več kot 3 minute na naslednjega študenta?** Če označimo z X število študentov, je pričakovana vrednost Poissonove porazdelitve λ enaka 30 študentov na uro oz. $1/2$ študenta na minuto. Naj bo W čas med zaporednima študentoma. Potem je pričakovani čas čakanja na naslednjega študenta $\theta = 1/\lambda = 2$ minuti. Ker je slučajna spremenljivka W po predpostavki eksponentno porazdeljena, s pričakovano vrednostjo $\theta = 2$, je njena gostota verjetnosti $p(w) = 0.5 e^{-w/2}$, za $w \geq 0$. Ploščina pod to krivuljo, za w vsaj 3 je enaka $e^{-3/2} = 0.223$.

(b) Število kilometrov, ki jih naredi določen avtomobil predno se izrabi akumulator, je eksponentno porazdeljeno s pričakovano vrednostjo 10.000 km. Lastnik avtomobila gre na 5000km dolgo pot. **Kakšna je verjetnost, da bo lahko dokončal potovanje brez zamenjave akumulatorja?** Na prvi pogled se zdi, da manjka nek bistven del informacije. Ali ne rabimo vedeti, koliko kilometrov smo že prevozili z danim akumulatorjem, preden lahko odgovorimo na vprašanje? Naj X označuje število km, ki jih lahko prevozi avtomobil z danim akumulatorjem. Predvidevamo, da velja naslednje

$$P(X > x + y \mid X > x) = P(X > y).$$

Če to drži, to pomeni, da za verjetnost, da se akumulator izprazni v več kot $y = 5000$ km, ni važno, če smo akumulator že uporabljali za $x = 0$ km, za $x = 1000$ km ali za $x = 15000$ km. Vemo, da je slučajna spremenljivka X eksponentno porazdeljena. Izkaže se, da zgornja izjava drži za eksponentno porazdelitev (dokažite to za domačo nalogu)! Zaradi tega rečemo, da je eksponentna porazdelitev “brez spomina”. Prav tako drži (ali želite pokazati tudi to?), da če je X eksponentno porazdeljena slučajna spremenljivka s pričakovano vrednostjo θ , potem je $P(X > k) = e^{-k/\theta}$ in $P(X > 5000) = e^{-5000/10000} = 1/\sqrt{e} = 0.604$.⁷ ◇

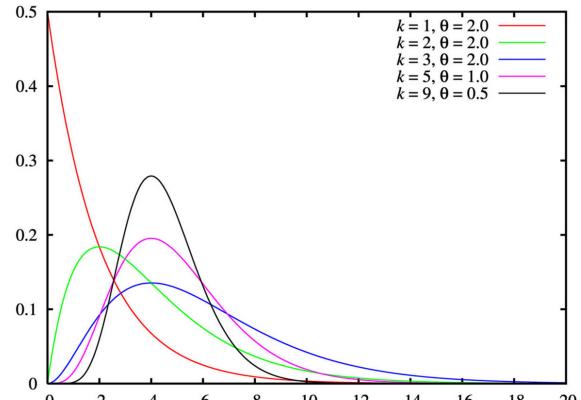
⁷Lastnik naj se sedaj sam odloči, ali je ta verjetnost dovolj velika, da bo raje že prej zamenjal akumulator in ne bo tvegal, da ga avto pusti na cedilu na kakšni zapuščeni cesti.

5.4.4 Gama porazdelitev – $\Gamma(k, \lambda)$

EkspONENTNO porazdelitev še posplošimo: tokrat pri Poissonovem procesu merimo čas, da se zgodi k dogodkov. Naj bosta $k, \lambda > 0$. Tedaj ima **Gama porazdelitev**, oznaka $\Gamma(k, \lambda)$, gostoto:

$$p(x) = \begin{cases} \frac{\lambda^k}{\Gamma(k)} x^{k-1} e^{-\lambda x}, & \text{za } x > 0, \\ 0, & \text{sicer,} \end{cases}$$

kjer je k parameter oblike, λ pa parameter raztega.



Na sliki je $\theta = 1/\lambda$. Za $k = 1$ seveda dobimo eksponentno porazdelitev. **Funkcijo Gama**⁸ 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 zaradi (A.5) tudi $\Gamma(1/2) = \sqrt{\pi}$.

Z uvedbo nove spremenljivke dobimo še

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

x	$\Gamma(x)$
1.1	0.951351
1.2	0.918168
1.3	0.897471
1.4	0.887264
1.5	0.886225
1.6	0.893514
1.7	0.908642
1.8	0.931384
1.9	0.961767

Integracija po delih (po realnem argumentu) nam da

(za $v = t^n$ in $du = e^{-t} dx$ velja $dv = nt^{n-1} dt$ in $u = -e^{-t}$):

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

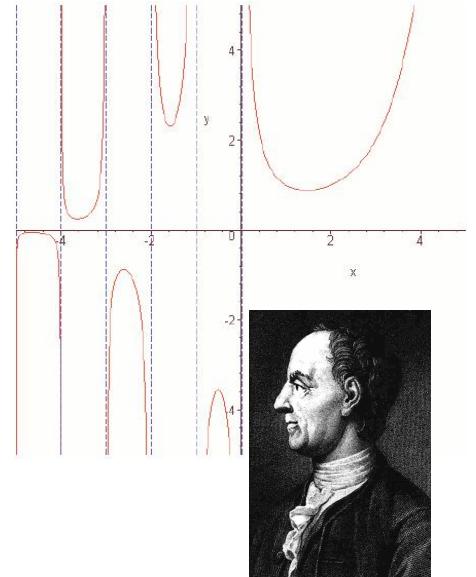
Na levi je še majhna tabela za funkcijo Gama ter njen graf.

Za naravno število x ($n = 1, 2, 3, \dots$), dobimo

$$\begin{aligned} \Gamma(n) &= (n-1)\Gamma(n-1) = (n-1)(n-2)\Gamma(n-2) \\ &= (n-1)(n-2)\dots 1 = (n-1)!, \end{aligned}$$

torej se Γ funkcija zreducira v ‘faktorjel’

(in jo znamo izračunati za vsako naravno število).



Leonhard Euler (1707-1783)⁹

⁸ <http://functions.wolfram.com/webMathematica/FunctionEvaluation.jsp?name=Gamma>

⁹ Eulerju se lahko zahvalimo za zapis $f(x)$ za funkcijo (1734), e za osnovo naravnega logaritma (1727), i

5.4.5 Hi-kvadrat porazdelitev – $\chi^2(n)$

Hi-kvadrat porazdelitev je poseben primer Gama porazdelitve:

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

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

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

Leta 1863 jo je prvi izpeljal nemški fizik Ernst Abbe, ko je preučeval kvadratov napak. Če je $n = 2k$ za $k \in \mathbb{N}$, je $\Gamma(n/2) = (k-1)!$, če pa je

$$\Gamma(n/2) = 1 \cdot 3 \cdot 5 \cdot \dots \cdot (2k-1)\sqrt{\pi}/2^n = (2k-1)!!\sqrt{\pi}/2^n.$$

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.



Ernst Abbe
(1840-1905)



χ^2 -test za ugotavljanje razlike kategoričnih spremenljivk (npr. religija, politične opcije,...) uporabimo v dveh podobnih (a različnih primerih):

- kako dobro se opazovana/izmerjena porazdelitev prilega pričakovani porazdelitvi – **kvaliteta-prilagoditve** (angl. goodness-of-fit test)
- za ocenjevanje ali sta naključni spremenljivki neodvisni.

χ^2 -test. Po dveh Presekovih člankih Tomaža Pisanskega. Pri mnogih družabnih igrah uporabljamo kocko. Če je kocka poštena, pade šestica z verjetnostjo 1/6. To pomeni, da lahko v povprečju pričakujemo na vsakih šest metov eno šestico. Če vržemo kocko 60-krat, lahko pričakujemo približno 10 šestic. Včasih se nam zazdi, da šestica noče pasti. **Ali je kocka morda obtežena, torej nepoštena?** Statistika pozna metodo, s katero lahko dokaj zanesljivo preverimo, ali je kocka obtežena ali ne. Metodi rečemo **test hi-kvadrat**. Sama uporaba testa hi-kvadrat je zelo preprosta. Za pomoč sem prosil svoja otroka. Pobrskali smo malo in našli tri kocke. Vsako smo vrgli 60-krat. Rezultate smo zabeležili v tabeli.

število pik	1	2	3	4	5	6	skupaj
rdeča kocka	5	16	8	14	12	5	60
črna kocka	7	8	8	11	10	16	60
bela kocka	6	10	10	9	14	11	60
teoretično	10	10	10	10	10	10	60

Prikazane so absolutne frekvence (število pojavitev posameznih izidov). Najbolj sumljivo se vede **rdeča** kocka, ker dejanske frekvence najbolj odstopajo od teoretičnih. **Ali je obtežena?** Dodali smo vrstico s teoretičnimi absolutnimi frekvencami. Dejanske absolutne frekvence

za kvadratni koren števila -1 (1777), π za konstanto 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>

odstopajo od teoretičnih, kar je popolnoma razumljivo, četudi so kocke poštene. Če dejanske frekvence malo odstopajo od teoretičnih, lahko z veliko verjetnostjo sklepamo, da je kocka poštena. Če pa dejanske frekvence močno odstopajo od teoretičnih, tedaj je malo verjetno, da gre za pošteno, neobteženo kocko. Test hi-kvadrat napravi dvoje. Dejanskim in teoretičnim frekvencam priredi število, s katerim merimo odstopanje frekvenc. Čim večje je dobljeno število, tem večje je odstopanje. Za odstopanje dopuščamo dve razlagi:

- lahko, da gre za slučajno odstopanje, ali pa
- gre (poleg slučajnega) še za sistematično odstopanje, torej teoretične frekvence ne ustrezajo dejanski porazdelitvi.

V našem primeru pomeni prva domneva, da gre za slučajno odstopanje poštene kocke, druga pa, da imamo opravka z obteženo kocko.

Čas je, da si pogledamo vso stvar čisto splošno. Denimo, da ima poskus n izidov. Denimo, da poskus ponovimo N -krat. Naj bodo E_1, E_2, \dots, E_n teoretične absolutne frekvence, O_1, O_2, \dots, O_n dejanske absolutne frekvence. To pomeni, da se je pri N ponovitvah poskusa izid i dogodil O_i -krat, medtem ko smo pričakovali, da se zgodi E_i -krat. Izraz

$$\chi^2(n-1) = \frac{(E_1 - O_1)^2}{E_1} + \frac{(E_2 - O_2)^2}{E_2} + \dots + \frac{(E_n - O_n)^2}{E_n}$$

imenujemo **hi-kvadrat z $(n-1)$ -prostostnimi stopnjami**. V statističnih priročnikih najdemo preglednice za hi-kvadrat. Za naše namene bo zadostovala priložena tabela.

Preden si še pogledamo, kaj pravi hi-kvadrat za naše kocke, še pomembno opozorilo: če je teoretična absolutna frekvenca kakega dogodka premajhna, je tudi zanesljivost testa hi-kvadrat vprašljiva. Običajno zahtevamo, da je vrednost vsakega E_i vsaj 5. Denimo, da ima poskus 6 izidov in je vrednost hi-kvadrat enaka 12·7. Število prostostnih stopenj je 5. Pogledamo v vrstico s petimi prostostnimi stopnjami in vidimo, da leži 12·7 med 11·1 in 15·1. Verjetnost, da so odstopanja med dejanskimi in teoretičnimi frekvencami zgolj slučajna, je manj kot 5% in več kot 1%. Pesimist bo verjetno domnevo o slučajnem odstopanju zavrnil.

št. prost. stopenj	$P=.1$	$P=.05$	$P=.01$	$P=0\cdot001$
2	4·6	6·0	9·2	13·8
3	6·3	7·8	11·3	16·3
4	7·8	9·5	13·3	18·5
5	9·2	11·1	15·1	20·5
6	10·6	12·6	16·8	22·5
7	12·0	14·1	18·5	24·3
8	13·4	15·5	20·1	26·1
9	14·7	16·9	21·7	27·9
10	16·0	18·3	23·2	29·6
12	18·6	21·0	26·2	32·9
14	21·1	23·7	29·1	36·1
16	23·5	26·3	32·0	39·3
18	26·0	28·9	34·8	42·3
20	28·4	31·4	37·6	45·3
25	34·4	37·6	44·3	52·6
30	40·3	43·8	50·9	59·7
40	51·8	55·8	63·7	73·4
60	74·4	79·1	88·4	99·6
80	96·6	101·9	112·3	124·8
100	118·5	124·3	135·8	149·5

Če bi bila vrednost hi-kvadrat pri istih pogojih 18·5, pa lahko domnevo mirno zavrnemo, saj je verjetnost manjša od enega odstotka. Zelo verjetno gre za resnično neujemanje med teoretičnimi in dejanskimi frekvencami. Običajno se vnaprej dogovorimo, katero mejo vzamemo za ločilo med sprejetjem oziroma zavrnitvijo domneve. Ta meja je običajno 5% ali 1%. Če pa je število prostostnih stopenj tako, da ga ni v naši preglednici, si pri oceni pomagamo z dvema vrsticama, tisto, ki je neposredno pred, in tisto, ki je za manjkajočo vrstico.

Najbolje je, da se vrnemo k našim trem kockam. Pokazali bomo, kako lahko uporabimo tabelo za χ^2 . Najbolj sumljiva je rdeča kocka, saj sta pri 60 metih enica in šestica samo po 5-krat:

$$\frac{(10-5)^2}{10} + \frac{(10-16)^2}{10} + \frac{(10-16)^2}{10} + \frac{(10-8)^2}{10} + \frac{(10-14)^2}{10} + \frac{(10-12)^2}{10} + \frac{(10-5)^2}{10} = 11.$$

Iz tabele za χ^2 razberemo, da je vsaj v petih odstotkih mogoče pričakovati tako odstopanje, če gre zgolj za slučajnost. Zato ne moremo sklepati, da je rdeča kocka nepoštena.

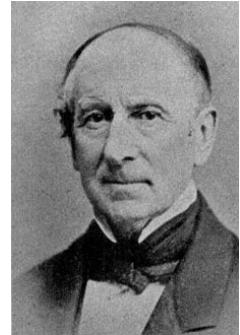
5.4.6 Cauchyeva porazdelitev

Dobimo jo s slučajno spremenljivko $X = \operatorname{tg} U$, kjer je U zvezno enakomerno porazdeljena na intervalu $[-\pi/2, \pi/2]$. Zato lahko rečemo, da je Cauchyjeva porazdelitev normalizirana intenzivnost svetlobe na premici iz točkovnega izvora. To je porazdelitev z vsemi realnimi števili kot zalogo vrednosti in z gostoto

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

Njena porazdelitvena funkcija je enaka

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



Ali se da izračunati pričakovano vrednost za to porazdelitev? Do uporabe Cauchyjeve porazdelitve pridemo tudi pri testiranju programske opreme, kjer moramo uporabiti podatke, ki vsebujejo nekaj ekstremnih vrednosti za potencialno povzročanje negativnih reakcij.

Porazdelitve v R-ju

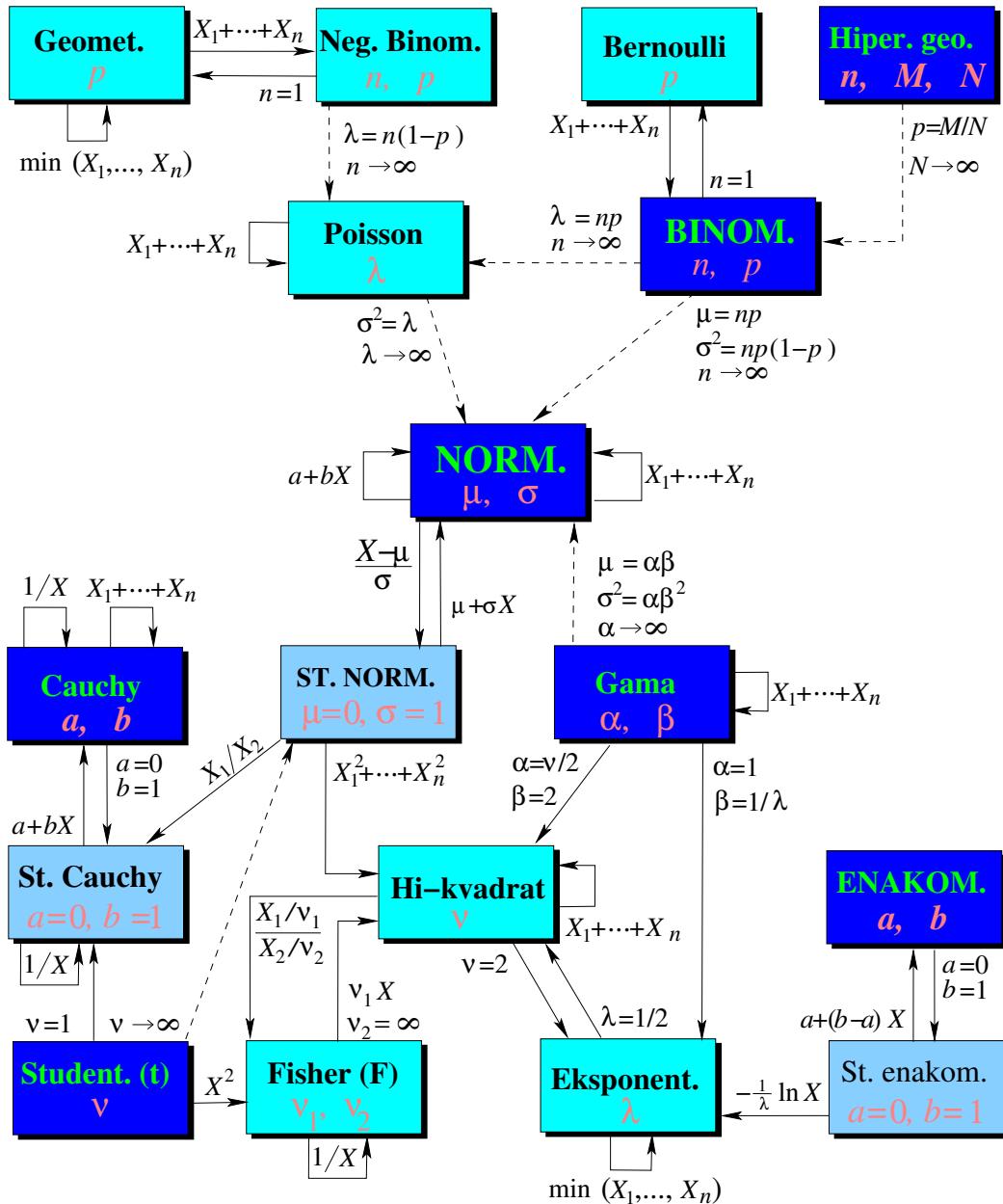
V R-ju so za delo s pomembnejšimi porazdelitvami na voljo funkcije:

- **dime** – gostota porazdelitve ime $p_{ime}(x)$
- **pime** – porazdelitvena funkcija ime $F_{ime}(q)$
- **qime** – obratna funkcija: $q = F_{ime}(p)$
- **rime** – slučajno zaporedje iz dane porazdelitve

Za ime lahko postavimo:

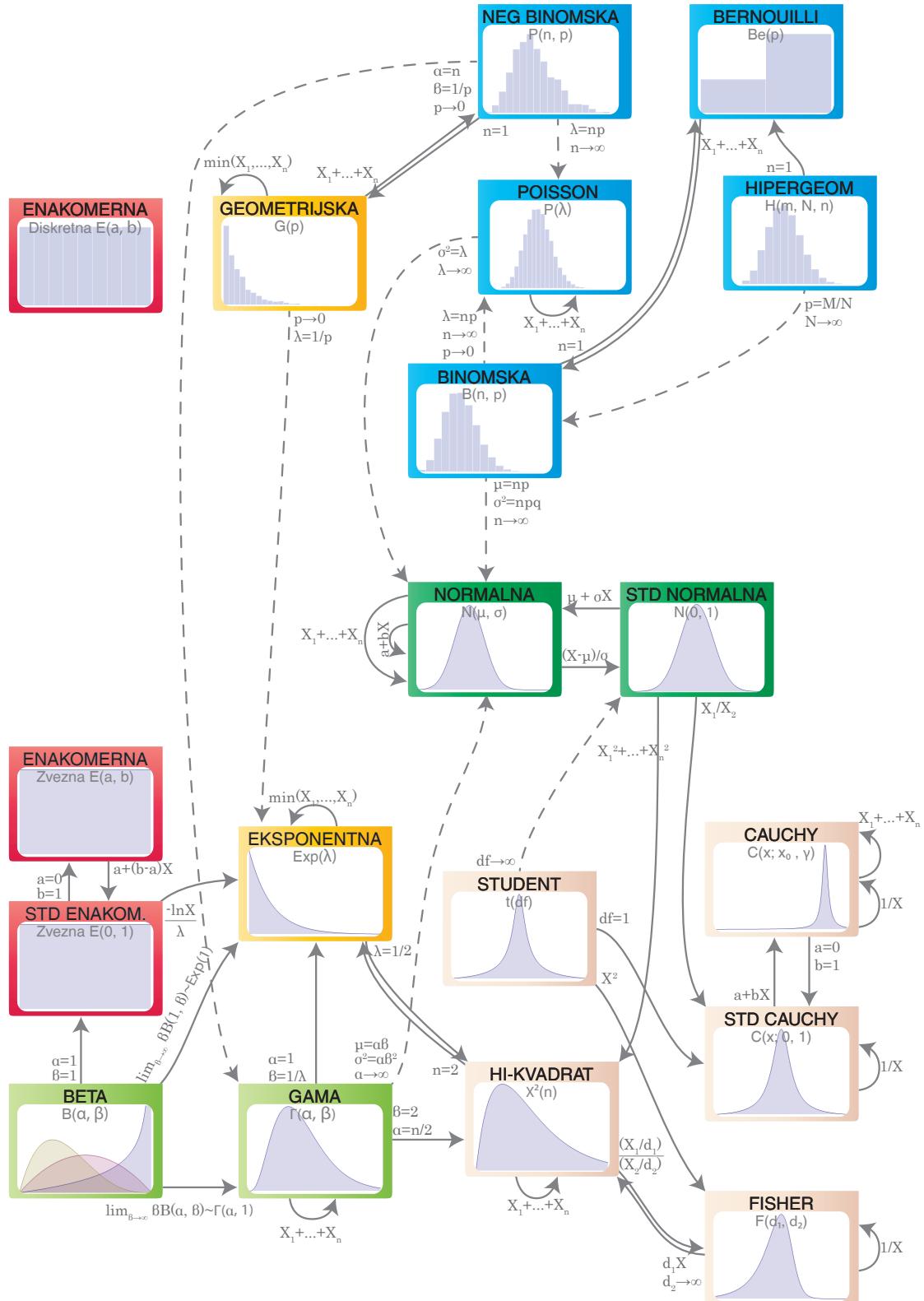
- **unif** – zvezna enakomerna,
- **binom** – binomska,
- **norm** – normalna,
- **exp** – eksponentna,
- **lnorm** – logaritmičnonormalna,
- **chisq** – porazdelitev χ^2 , ...

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



Slika 5.16: Povezave med znanimi porazdelitvami. Iz normalne na standardno normalno pridemo s substitucijo $(X - \mu)/\sigma$. Iz binomske na normalno pridemo tako, da pošljemo $n \rightarrow \infty$ in dobimo, da je $\mu = np$ in $\sigma^2 = np(1 - p)$, do Poissonove pa tako, da je $\lambda = np$...

Primer: Če je slučajna spremenljivka X porazdeljena Cauchyjevo, kako je porazdeljena slučajna spremenljivka $1/X$. Pravtako Cauchyjevo, saj dobimo Cauchijevo porazdelitev kot kvocient dveh neodvisnih standardnih normalnih, npr. Z_1 in Z_2 , tj. $X = Z_1/Z_2$, in od tod $1/X = Z_2/Z_1$. \diamond



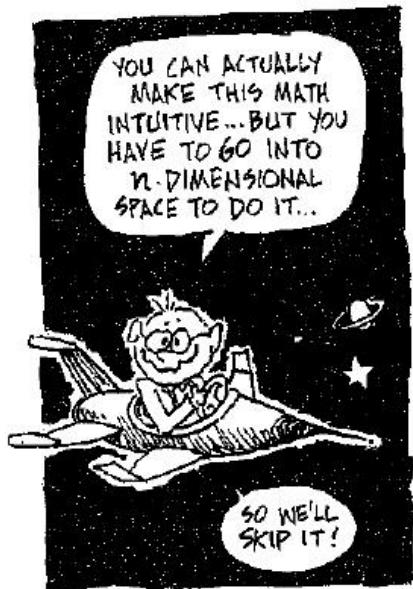
Slika 5.17: Porazdelitve in oblike podatkov.

Poglavlje 6

Slučajni vektorji

Pristop z eno slučajno spremenljivko lahko posplošimo na več slučajnih spremenljivk in v tem primeru rečemo, da gre za slučajne vektorje. Pri tem bo porazdelitvena funkcija ene spremenljivke prešla v funkcijo več spremenljivk, za katere verjetno mislite, da jih ne poznate dobro, vendar temu ni čisto tako. Ko želite npr. izračunati ploščino trikotnika, zapišete $f(a, v) = a \cdot v/2$, kjer je v višina nad stranico a . Prostornina kvadra je enaka $g(a, b, c) = a \cdot b \cdot c, \dots$

Istočasno opazovanje več slučajnih spremenljiv, tj. študija slučajnih vektorjev je pomembna že zato, ker bi radi razumeli, kdaj so slučajne spremenljivke odvisne oz. neodvisne, zato je ne moremo spustiti.



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

Pošči 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

$$X \sim \begin{pmatrix} 1 & 2 & 3 & 4 \\ 0.08 & 0.37 & 0.40 & 0.15 \end{pmatrix}.$$

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

Ne nista, saj velja npr.: $P(X = 4, Y = 3) = 0 \neq 0.15 \times 0.15 = P(X = 4) \cdot P(Y = 3)$. \diamond

Primer: Oglejmo si slučajni spremenljivki X in Y , ki spremeljata, koliko lihih¹ pik pokažeta bela in črna kocka. Zanima nas, ali sta slučajni spremenljivki $U := XY$ in $V = X + Y$ odvisni. Če je npr. $V = 5$, potem imamo za U eno samo možnost, tj. $U = 6$, čeprav jo lahko dobimo na dva načina: $X = 1, Y = 5$ in $X = 5, Y = 1$, kar nam sugerira odvisnost. Pa vendar, pojdimo lepo počasi. Brez težav zapišemo

$$V \sim \begin{pmatrix} 2 & 4 & 6 & 8 & 10 \\ \frac{1}{9} & \frac{2}{9} & \frac{3}{9} & \frac{2}{9} & \frac{1}{9} \end{pmatrix} \quad \text{in} \quad U \sim \begin{pmatrix} 1 & 3 & 5 & 9 & 15 & 25 \\ \frac{1}{9} & \frac{2}{9} & \frac{2}{9} & \frac{1}{9} & \frac{2}{9} & \frac{1}{9} \end{pmatrix}$$

(upoštevali smo, da elementi zaloge vrednosti slučajne spremenljivke U niso deljivi s praštevilom večjim od 5). Sedaj pa premislimo še verjetnosti v tabeli²:

$+ \setminus \cdot$	1	3	5	9	15	25	Σ
2	$\frac{1}{9}$	0	0	0	0	0	$\frac{1}{9}$
4	0	$\frac{2}{9}$	0	0	0	0	$\frac{2}{9}$
6	0	0	$\frac{2}{9}$	$\frac{1}{9}$	0	0	$\frac{3}{9}$
8	0	0	0	0	$\frac{2}{9}$	0	$\frac{2}{9}$
10	0	0	0	0	0	$\frac{1}{9}$	$\frac{1}{9}$
Σ	$\frac{1}{9}$	$\frac{2}{9}$	$\frac{2}{9}$	$\frac{1}{9}$	$\frac{2}{9}$	$\frac{1}{9}$	1

Enako kot v prejšnjem primeru lahko zaključimo, da sta slučajni spremenljivki V in U odvisni, saj bi v nasprotnem primeru bile *vse* verjetnosti v tabeli neničelne. \diamond

Primer: Študent ima $66 \cdot \bar{6}$ možnost, da naredi izpit. Na voljo ima tri roke in nobenega ne izpusti, če izpita še ni opravil. Naj bo X slučajna spremenljivka, ki šteje število polaganj, Y pa slučajna spremenljivka, ki je zavzame vrednost 1, če študent uspešno opravi izpit in 0

¹Brez te predpostavke bi bil primer nekoliko daljši.

² Nekaj računanja si lahko prihranimo z naslednjimi ugotovitvami (posebej če delamo brez predpostavke, da sledimo na kockah le liha števila):

- V vrstici, ki se prične s številom $i \leq 7$, so na mestu verjetnosti ničle vse do stolpca, ki ima na vrhu $i - 1$, nato pa do stolpca, ki ima na vrhu $k(i - k)$, pri čemer je $i - k = 6$.
- Iz $A_2 \geq G_2$ oziroma $(a + b)/2 \geq \sqrt{ab}$ za $a, b \geq 0$, kar je po odpravi korena in ulomka ekvivalentno $(a - b)^2 \geq 0$, lahko zaključimo, da je produk dveh števil na zgoraj omejen s kvadratom polovične vsote teh dveh števil. Slednja neenakost nam zagotovi precej ničel tudi na koncu začetnih vrstic.
- Vsota verjetnosti v vsakem stolpcu oziroma vrstici osrednjega dela tabele je enaka zadnjemu elementu v stolpcu oziroma vrstici, ki smo ga izračunali na začetku.

sicer. [Ali sta slučajni spremenljivki \$X\$ in \$Y\$ neodvisni?](#)

$$X \sim \begin{pmatrix} 1 & 2 & 3 \\ 2/3 & 2/9 & 1/9 \end{pmatrix}, \quad Y \sim \begin{pmatrix} 0 & 1 \\ 1/27 & 26/27 \end{pmatrix}$$

in $E(X) = 13/9 = 1\bar{4}\dots$, $E(Y) = 26/27$. Nadomestimo $2/3$ s p in $1 - p$ s q . Potem velja:

$$\begin{array}{c|ccc|c} Y \setminus X & 1 & 2 & 3 & \Sigma \\ \hline 0 & 0 & 0 & q^3 & q^3 \\ 1 & p & pq & pq^2 & 1 - q^3 \\ \hline \Sigma & p & pq & q^2 & 1 \end{array} \quad \begin{array}{l} \text{in } E(X) = 1 + q + q^2, \quad E(Y) = 1 - q^3, \\ XY \sim \begin{pmatrix} 0 & 1 & 2 & 3 \\ q^3 & p & pq & pq^2 \end{pmatrix}, \quad E(XY) = p(1 + 2q + 3q^2) \\ \text{Končno je } E(XY) - E(X)E(Y) = -pq^3(q + 2). \quad \diamond \end{array}$$

Slučajni vektor je n -terica slučajnih spremenljivk $\mathbf{X} = (X_1, \dots, X_n)$. Tudi za slučajni vektor opišemo porazdelitveni zakon s porazdelitveno funkcijo ($x_i \in \mathbb{R}$):

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

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

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

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

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

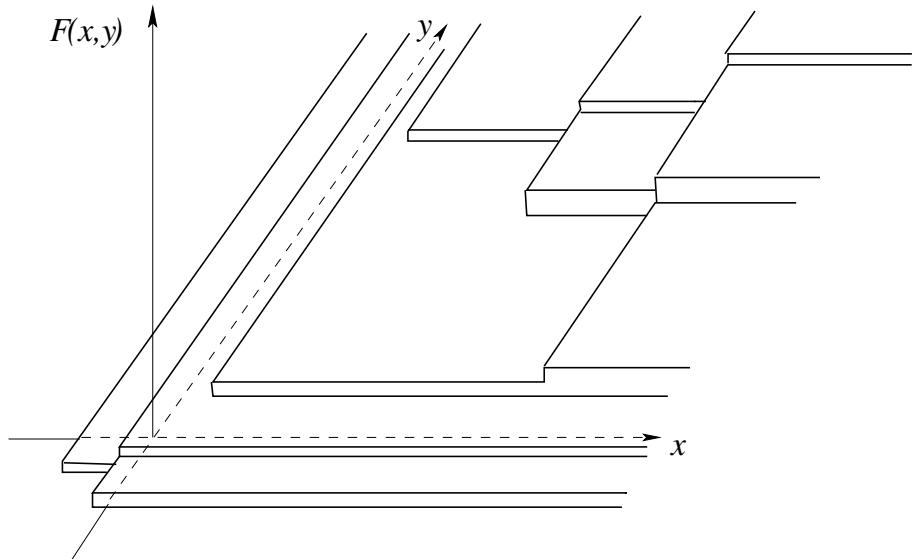
Tedaj je verjetnost, da je slučajna točka (X, Y) v pravokotniku $(a, b] \times (c, d]$ enaka

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

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 $n = 2$, $X : \{x_1, \dots, x_k\}$, $Y : \{y_1, \dots, y_m\}$ 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 $P(n; p_1, \dots, p_r)$, $\sum p_i = 1$, $\sum k_i = n$ je določena s predpisom

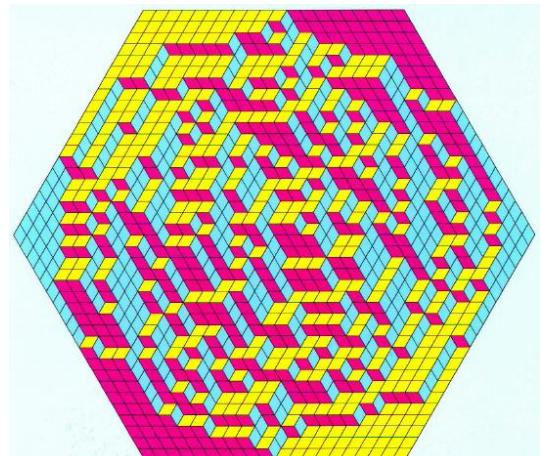
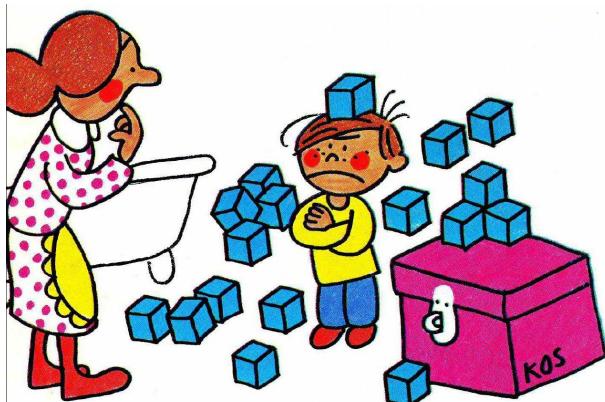
$$P(X_1 = k_1, \dots, X_r = k_r) = \frac{n!}{k_1! \cdots k_r!} p_1^{k_1} \cdots p_r^{k_r}.$$

Koeficient šteje permutacije s ponavljanjem.³ Za $r = 2$ dobimo binomsko porazdelitev, tj. $B(n, p) = P(n; p, q)$.

Primer: Iz kupa igralnih kart (52) na slepo izberemo eno karto in jo nato vrnemo nazaj. Postopek ponovimo 5-krat. **Kakšna je verjetnost, da bomo videli dvakrat srce, po enkrat pa pika, križa in karo?** Število razredov r je 4, $n = 5$ in $p_1 = p_2 = p_3 = p_4 = 1/4$. Po zgornji formuli dobimo

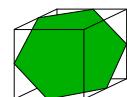
$$P(X_1 = 1, X_2 = 2, X_3 = 1, X_4 = 1) = \frac{5!}{1! \cdot 2! \cdot 1! \cdot 1!} 0 \cdot 25^1 0 \cdot 25^2 0 \cdot 25^1 0 \cdot 25^1 = 0 \cdot 05859. \quad \diamond$$

6.2 Ponovitev: dvojni integral



Slika: **Kako bi najlažje prešteli vse kocke?** Leva slika: 13 jih leži po tleh, ena je na fantovi glavi, na vijoličasti škatli pa so še 4 (čeprav četrte v resnici ne vidimo), torej jih je skupaj 18.

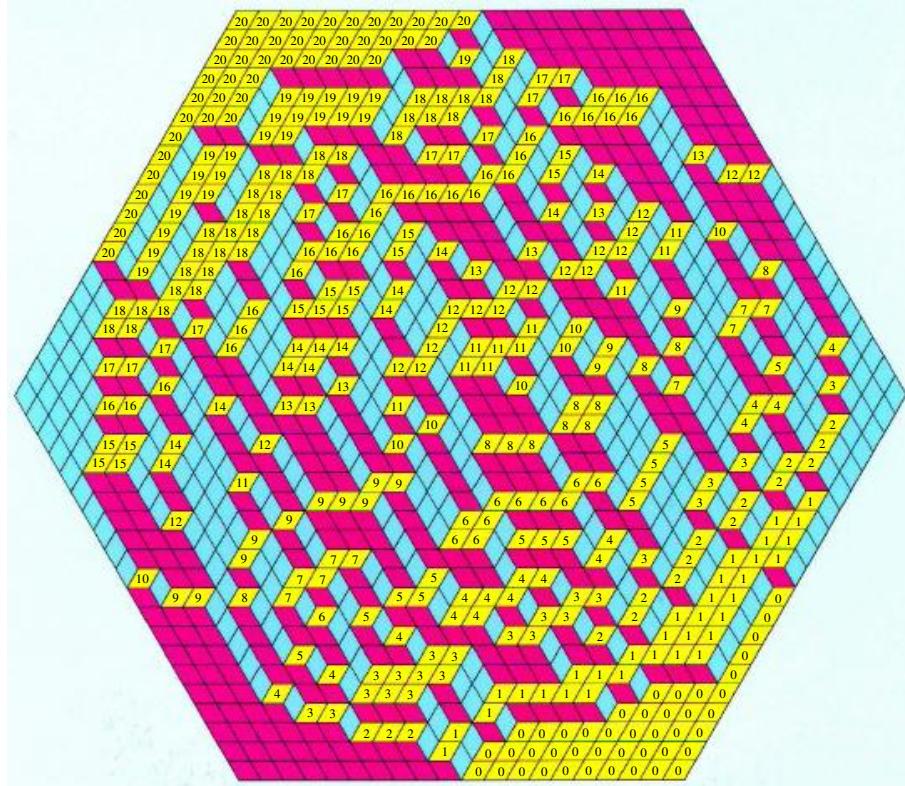
Sedaj pa poglejmo še desno sliko, kjer se kocke nahajajo znotraj večje kocke $20 \times 20 \times 20$, torej jih ne more biti več kot 8000, a izgleda kot da jih je ravno polovica. Po občutku jih je torej približno 4000.



Če pa hočemo biti do kockice natančni, potem nam ne preostane drugega, kot da jih začnemo šteti. Problem pa ni samo v natančnosti, pač pa se moramo vprašati tudi katera ‘polovica’ nas zanima, kajti če sliko obrnemo na glavo, nam npr. zgornje ploskve doslej najvišjih kock postanejo tla.

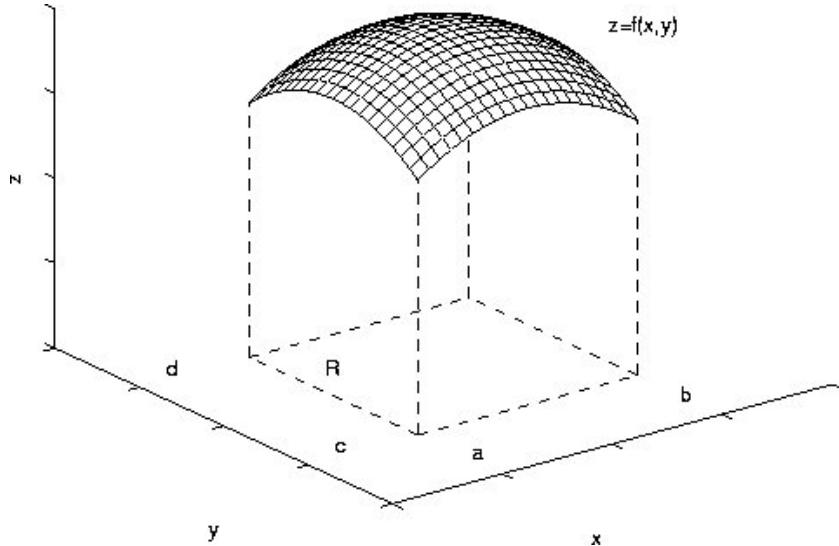
³Glej npr. http://en.wikipedia.org/wiki/Multinomial_distribution.

Pa začnimo šteti po nivojih/slojih, ki so vzporedni s ploskvami kocke. Tako vidimo, da imamo tri možnosti: v rumeni smeri (naj bo to z os), v modri smeri (naj bo to y os) in v vijolični smeri (ki naj bo x os). Mi se odločimo za rumeno in začnemo pri vrhu (te kocke namreč vse vidimo). V prvem stolpcu jih je 11, nato 10, pa 8, trikrat po 3 in sedemkrat po 1, skupaj torej $11 + 10 + 8 + 3 \times 3 + 7 \times 1 = 45$ (ki smo jih prešeli skoraj tako kot zidake pri integriranju). Ker pa so te kocke v 20 ‘nadstropju’, pomnožimo njihovo število z 20 in dobimo 900. Nadaljujemo en nivo nižje. Kocke, ki jih ne vidimo, smo že šteli, tako da preštejemo še $0 + 1 + 0 + 2 \times 5 + 4 \times 2 + 4 \times 1 = 23$ kock in jih pomnožimo z 19, tako da dobimo 437. Tako nadaljujemo vse do konca:



Vsota vseh števil v vseh oklepajih je ravno 400, kajti od zgoraj se vidjo vsa polja v mreži 20×20 . Zadnje vrstice sicer ne potrebujemo, a je dobrodošla zaradi preiskusa. Načeloma so nas zanimala samo vsote števil v oklepajih, a tudi ta razčlenitev je v resnici koristna, saj se morajo števila v oklepajih po stolpcih sešteviti v 20. Velja pa pripomniti še, da nam zgornja informacija zadostuje tudi za rekonstrukcijo celotne slike.

Končni odgovor je 4098. Če bi si sliko izpisali in bi jo po nesreči obrnili na glavo, bi dobili $8000 - 4098 = 3902 = 0 \cdot 45 + 1 \cdot 23 + 2 \cdot 34 + 3 \cdot 12 + 4 \cdot 29 + 5 \cdot 13 + 6 \cdot 13 + 7 \cdot 7 + 8 \cdot 19 + 9 \cdot 11 + 10 \cdot 7 + 11 \cdot 13 + 12 \cdot 11 + 13 \cdot 9 + 14 \cdot 11 + 15 \cdot 13 + 16 \cdot 16 + 17 \cdot 22 + 18 \cdot 16 + 19 \cdot 33 + 20 \cdot 44$. V resnici smo računali dvojni integral, to pa smo počeli z uporabo običajnega integrala.



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 $[a, b] \times [c, d]$). Prostornina 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 z uporabo 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

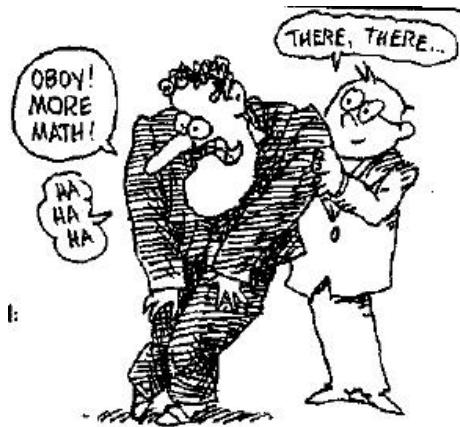
$$\iint_R f(x, y) \, dx \, dy \leq \iint_R g(x, y) \, dx \, dy. \quad \square$$

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 $\mathbf{X} = (X_1, X_2, \dots, X_n)$ 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

$$p_X(x) = F'_X(x) = \int_{-\infty}^{\infty} p(x, y) dy, \quad \text{in} \quad p_Y(y) = F'_Y(y) = \int_{-\infty}^{\infty} p(x, y) dx.$$

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:

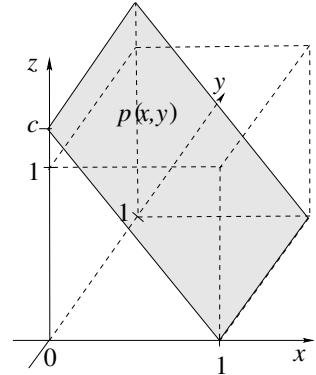
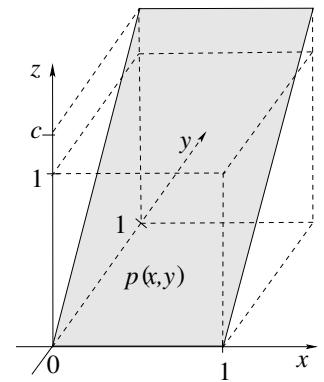
$$p_X(x) = \int_0^1 2y \, dy = y^2 \Big|_{y=0}^1 = 1.$$

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

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

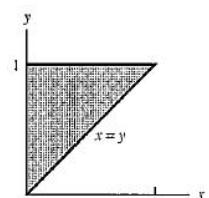
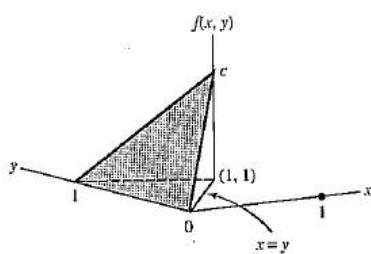
Za vajo poskusi odgovoriti na ista vprašanja še za

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



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

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

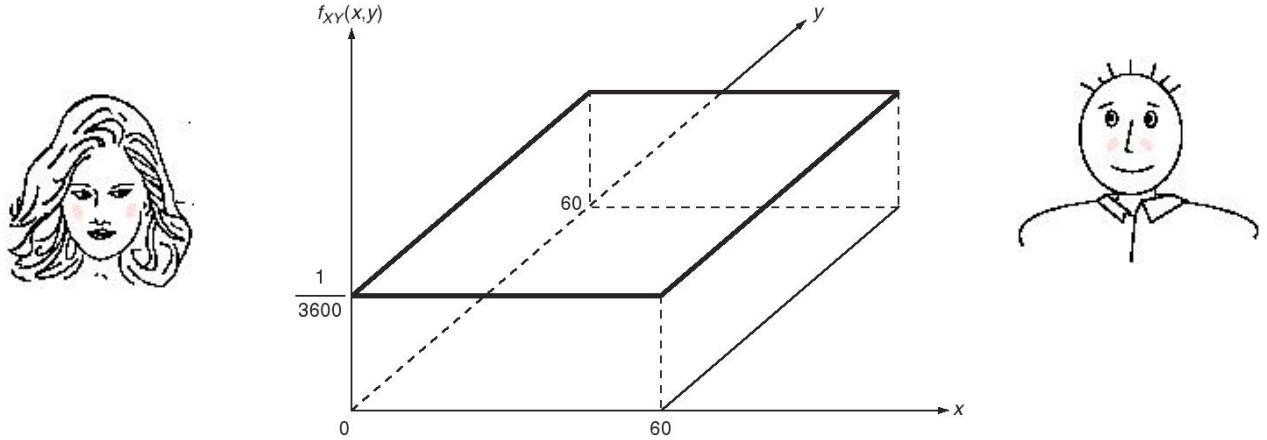


Izberi vrednost konstante c med:

- (a) 6, (b) 2, (c) 3, (d) 1. ◊

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

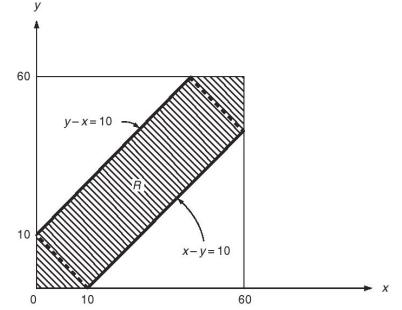
se bosta srečala. Naj bo čas prihoda fanta X minut po 9-i, pravtako pa naj bo čas prihoda dekleta Y minut po 9-i.



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:

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



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 $p_Y(y) = 0$, torej sta X in Y obe enakomerno porazdeljeni slučajni spremenljivki na intervalu $[0, 60]$. \diamond

Večrazsežna normalna porazdelitev

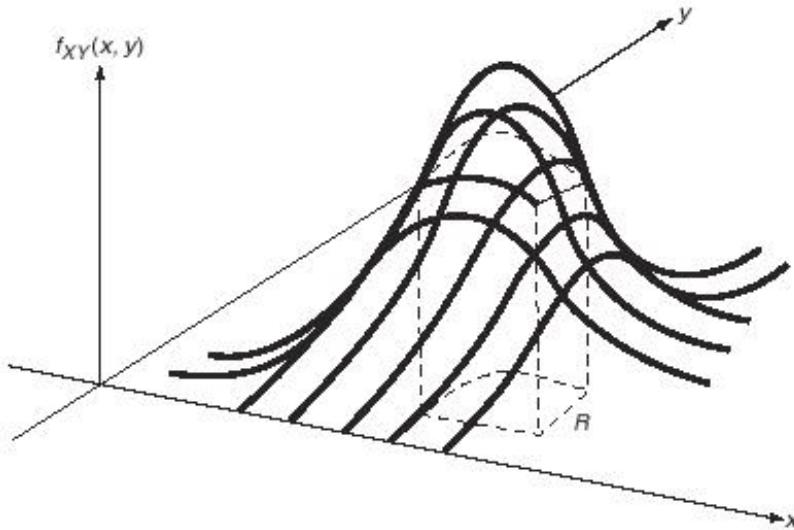
V dveh razsežnostih označimo normalno porazdelitev z $N(\mu_X, \mu_Y, \sigma_X, \sigma_Y, \rho)$ in ima gostoto

$$p(x, y) = \frac{1}{2\pi\sigma_X\sigma_Y\sqrt{1-\rho^2}} e^{-\frac{1}{2(1-\rho^2)}\left(\left(\frac{x-\mu_X}{\sigma_X}\right)^2 - 2\rho\frac{x-\mu_X}{\sigma_X}\frac{y-\mu_Y}{\sigma_Y} + \left(\frac{y-\mu_Y}{\sigma_Y}\right)^2\right)}.$$

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, $\boldsymbol{\mu} = (\mu_1, \dots, \mu_n)^T$, 'eksponent' T pa pomeni transponiranje (tj. zrcaljenje matrike preko glavne diagonale). Obe robni porazdelitvi sta normalni.



Primer: Pri študiju upora Y strurnega 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

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

Izračunajmo prvi integral. Upoštevamo $a dx = d(ax) = -d(-ax - by)$:

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

Za vajo izračunate tudi drugi dvojni integral (vemo, da sta enaka).

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*

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

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:

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

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:

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

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

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

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:

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

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

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

kar nam da $P(X = x_i, Y \leq y_j) = p_i \cdot P(Y \leq y_j)$. Potem seveda velja tudi $P(X = x_i, Y \leq y_{j+1}) = p_i \cdot P(Y \leq y_{j+1})$, 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 žeeli dokazati. \square

Podobno dokažemo tudi naslednjo trditev za zvezno porazdeljeni slučajni spremenljivki (le da v prvem 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

$$p(x, y) = \frac{1}{2\pi\sigma_X\sigma_Y} e^{-\frac{1}{2}((\frac{x-\mu_X}{\sigma_X})^2 + (\frac{y-\mu_Y}{\sigma_Y})^2)} = p_X(x) \cdot p_Y(y).$$

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

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

Poglavlje 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

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

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

$$F_Y(y) = P(Y \leq y).$$

Če je funkcija f linear, 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 števni 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 njena inverzna funkcija f^{-1} in velja

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

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

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

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

$$F_Y(y) = \int_{-\infty}^{f^{-1}(y)} p(x) dx \quad \text{in, če je } f \text{ odvedljiva, še } p_Y(y) = p(f^{-1}(y))f^{-1}(y)'.$$

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

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

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

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 gostoto verjetnosti porazdelitve hi-kvadrat $\chi^2(1)$. \diamond

Primer: Naj bo X porazdeljena normalno $N(\mu, \sigma)$. Za slučajno spremenljivko $Y = e^X$ izračunaj njeni porazdelitveni funkcijo in gostoto verjetnosti. Za $\mu = 0$ in $\sigma^2 = 1$ nariši graf gostote verjetnosti $p_Y(x)$.

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

in $f(x) = e^x$. Torej $f^{-1}(y) = \ln y$. Potem pa še

$$(f^{-1}(y))' = (\ln y)' = 1/y$$

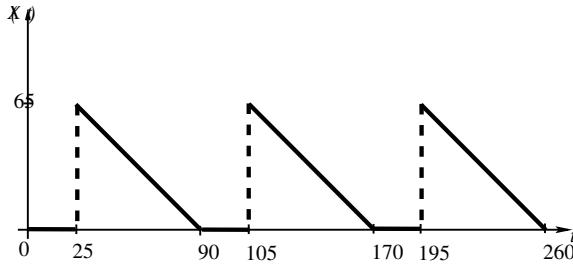
in

$$p(f^{-1}(y)) = \frac{1}{\sqrt{2\pi}\sigma} e^{(\ln y - \mu)^2/(2\sigma^2)}.$$

\diamond

Primer: Na nekem pešprehodu zelena luč gori 25 sekund (vključno z utripanjem - če do tega sploh prihaja), nato pa sledi rdeča luc 65 sekund. Jure pride naključno na ta prehod. Naj bo X čas, ki ga mora čakati, da lahko prečka cesto.

- (a) Kakšna je porazdelitev za slučajno spremenljivko X (diskretna, zvezna, mešana), nariši sliko?
- (b) Kakšna je verjetnost, da bo Jure čakal več kot 20 sekund?
- (c) Kakšna je pogojna verjetnost, da bo Jure čakal vsaj dodatnih 20 sekund, če vemo, da je že čakal 20 sekund?
- (d) Izračunaj 10%, 25%, 50% in 90% centil slučajne spremenljivke X .
- (e) Izračunaj pričakovano vrednost $E(X)$.
 - (a) Verjetnost, da Jure ne čaka, tj. $P(X = 0)$ je enaka $25/(25 + 65) = 5/18$, kar pomeni, da spremenljivka ni zvezna (verjetnost, da zvezna slučajna spremenljivka zavzame neko konkretno vrednost iz zaloge vrednosti, je namreč vedno 0), po drugi strani pa zavzame vse vrednosti na intervalu $[0, 65]$, kar pomeni, da gre za mešano slučajno spremenljivko.



Čeprav je v primeru časa zaloga vrednosti enaka $[0, \infty]$, se lahko za preostala vprašanja omejimo na interval $[0, 90]$, glede na to, da se obnašanje semaforja periodično ponavlja. Hkrati lahko s slike vidimo, da je povprečen čas čakanja v času rdeče luči enak $65/2 = 32.5$ sekund, kar pomeni, da je pričakovana vrednost spremenljivke X enaka $0 \cdot \frac{5}{18} + 32.5 \cdot \frac{13}{18} = 23.47$. S tem smo odgovorili tudi na vprašanje iz točke (e).

- (b) S slike lahko odčitamo tudi, da je graf funkcije $X(t)$ na intervalu $[0, 90]$ nad premico $y = 20$ na podintervalu $[25, 70]$, kar pomeni, da je verjetnost iz vprašanja (b) enaka $(70 - 25)/90 = 1/2$.
- (c) Pogoj, da je Jure čakal 20 sekund nam pove, da je prišel znotraj podintervala $[25, 70]$ in je sedaj znotraj intervala $[45, 90]$. Na tem intervalu pa je graf funkcije $X(t)$ nad premico $y = 20$ na podintervalu $[45, 70]$ in je izkana verjetnost enaka $(70 - 45)/45 = 5/9$.
- (d) Iz razmisleka v prvi točki, sta 10. in 25. centil enaka 0, saj je $5/18 = 0.2\bar{7}$ in je $90/4 < 25$. Končno je 50. centil enak $(65 - (90/2 - 25)) = 45$ sekund, 90. centil pa $(65 - (81 - 25)) = 9$ sekund. \diamond

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

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

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

Y funkcija slučajnega vektorja \mathbf{X} , tj. $\mathbf{Y} = f(\mathbf{X})$. Porazdelitve komponent dobimo na običajen način

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

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

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

in

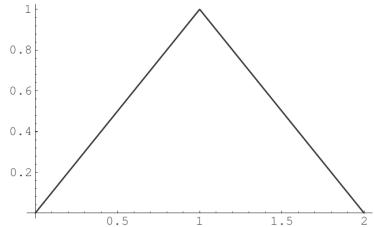
$$p_Z(z) = F'_Z(z) = \int_{-\infty}^{\infty} p(x, z-x) dx = \int_{-\infty}^{\infty} p(z-y, y) dy.$$

Če sta spremenljivki X in Y neodvisni, dobimo zvezo

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

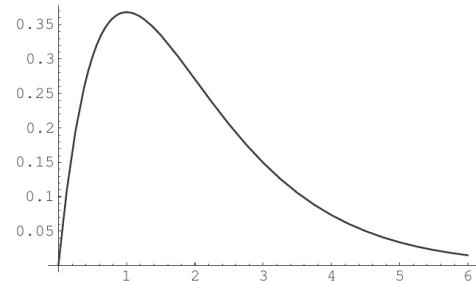
Gostota $p_Z = p_X * p_Y$ je **konvolucija** funkcij p_X in p_Y .

Primer: Enakomerno porazdeljeni slučajni spremenljivki. Predpostavimo, da izberemo neodvisno naključni števili na intervalu $[0, 1]$ z enkomerno porazdelitvijo. Gostota porazdelitve njune vsote je prikazana na sliki. \diamond



Primer: Eksponentno porazdeljeni slučajni spremenljivki. Predpostavimo, da izberemo neodvisno naključni števili na intervalu $[0, \infty]$ z eksponentno porazdelitvijo $\text{Exp}(\lambda)$. Torej je za $z > 0$:

$$\begin{aligned} p_Z(z) &= \int_{-\infty}^{\infty} p_X(x) p_Y(z-x) dx \\ &= \int_0^z \lambda e^{-\lambda x} \lambda e^{-\lambda(z-x)} dx \\ &= \int_0^z \lambda^2 e^{-\lambda z} dx = \lambda^2 z e^{-\lambda z}, \end{aligned}$$



medtem ko za $z < 0$ dobimo $p_Z(z) = 0$. Torej smo dobili Gama porazdelitev $\Gamma(2, \lambda)$. \diamond

Primer: Če je $(X, Y) \sim N(\mu_X, \mu_Y, \sigma_X, \sigma_Y, \rho)$, je vsota $Z = X + Y$ zopet normalno porazdeljena: $Z \sim N(\mu_X + \mu_Y, \sqrt{\sigma_X^2 + 2\rho\sigma_X\sigma_Y + \sigma_Y^2})$.

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

Še bolj splošno, če sta $X \sim \Gamma(k_1, \lambda)$ in $Y \sim \Gamma(k_2, \lambda)$ neodvisni slučajni spremenljivki, je tudi njuna vsota $Z = X + Y$ porazdeljena po tej porazdelitvi $Z \sim \Gamma(k_1 + k_2, \lambda)$. \diamond

Dosedanje ugotovitve povezane s hi-kvadrat lahko združimo v naslednjo trditev.

Trditev 7.2. Če so X_1, X_2, \dots, X_n neodvisne standardizirane 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

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

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 funkcija f bijektivna z zveznimi parcialnimi odvodi, lahko nadaljujemo

$$F_{U,V}(u, v) = \iint_{A(u, v)} p(x(u, v), y(u, v)) \left| \det \begin{pmatrix} \frac{\partial u}{\partial x} & \frac{\partial u}{\partial y} \\ \frac{\partial v}{\partial x} & \frac{\partial v}{\partial y} \end{pmatrix} \right| du dv.$$

Zgornja determinanta je poznana pod imenom **Jacobijeva determinanta**, oznaka $J(u, v)$ (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: Za $\Omega = \{(x, y) \mid 0 < x \leq 1, 0 < y \leq 1\}$, 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

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

Jacobijeva determinanta je torej enaka $r(-2)\pi/(rx) = -2\pi/x$ 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. \diamond

7.4 Pogojne porazdelitve

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

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

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

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

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

Primer: Nadalujmo primer s katerim smo pričeli poglavje 6. **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:
$$X|y=2 \sim \begin{pmatrix} 1 & 2 & 3 & 4 \\ 0.25 & 0.50 & 0.25 & 0 \end{pmatrix}. \diamond$$

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

Poisci robno porazdelitev za Y .

$$Y \sim \begin{pmatrix} 0 & 1 & 2 & 3 \\ 0.11 & 0.25 & 0.40 & 0.24 \end{pmatrix}.$$

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

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

Zvezne pogojne porazdelitve

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

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

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

Primer: Za 2-razsežno normalno porazdelitev dobimo

$$p_X(x|y) \sim N\left(\mu_X + \rho \frac{\sigma_X}{\sigma_Y}(y - \mu_Y), \sigma_X \sqrt{1 - \rho^2}\right). \quad \diamond$$

Poglavlje 8

Momenti in kovarianca



V tem poglavju bomo videli, da ima tudi verjetnost svoje momente.

8.1 Pričakovana vrednost

Kot posplošitev povprečne vrednosti smo že vpeljali **pričakovano vrednost** $E(X)$ (oz. matematično upanje) za končno diskretno slučajno spremenljivko X . Splošna diskretna slučajna spremenljivka X z verjetnostno funkcijo p_k pa ima pričakovano vrednost

$$E(X) = \sum_{i=1}^{\infty} x_i p_i, \quad \text{če je} \quad \sum_{i=1}^{\infty} |x_i| p_i < \infty.$$

Zvezna slučajna spremenljivka X z gostoto verjetnosti $p(x)$ ima pričakovano vrednost

$$E(X) = \int_{-\infty}^{\infty} x p(x) dx, \quad \text{če je} \quad \int_{-\infty}^{\infty} |x| p(x) dx < \infty.$$

Primer: Poišči vsaj dve slučajni spremenljivki, za katere pričakovana vrednost ne obstaja.

V diskretnem primeru vzemimo: $x_k = (-1)^{k+1} 2^k / k$ in $p_k = 2^{-k}$, v zveznem pa Cauchyjevo porazdelitev, tj. $p_X(x) = 1/(\pi(1+x^2))$. V prvem primeru bi morala biti končna naslednja vsota: $S = \frac{1}{1} + \frac{1}{2} + \frac{1}{3} + \dots$. Opazimo: $\frac{1}{3} + \frac{1}{4} > \frac{1}{4} + \frac{1}{4} = \frac{1}{2}$ in $\frac{1}{5} + \frac{1}{6} + \frac{1}{7} + \frac{1}{8} > \frac{1}{8} + \frac{1}{8} + \frac{1}{8} + \frac{1}{8} = \frac{1}{2}$ ter v splošnem

$$\frac{1}{2^n + 1} + \dots + \frac{1}{2^{n+1}} > \frac{2^n}{2^{n+1}} = \frac{1}{2}, \quad \text{torej velja} \quad S > \sum_{n=1}^{\infty} \frac{1}{2},$$

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

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

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 pričakovane vrednosti

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

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

Velja splošno: pričakovana vrednost poljubne funkcije $f(x)$ slučajne spremenljivke X , tj. $\mathbb{E}(f(X))$, obstaja in je v diskretnem primeru enaka $\sum_{i=1}^{\infty} f(x_i) p_i$, v zveznem pa $\int_{-\infty}^{\infty} f(x) p(x) dx$, če ustrezeni izraz absolutno konvergira.

Primer: Za diskretno slučajno spremenljivko X je $\mathbb{E}(X^2) = x_1^2 p_1 + x_2^2 p_2 + \dots$ \diamond

Naj bo a realna konstanta. Če ima slučajna spremenljivka X pričakovano vrednost, potem jo ima tudi spremenljivka aX in velja $\mathbb{E}(aX) = a\mathbb{E}(X)$. Če imata slučajni spremenljivki X in Y pričakovano vrednost, jo ima tudi njuna vsota $X + Y$ in velja $\mathbb{E}(X + Y) = \mathbb{E}(X) + \mathbb{E}(Y)$. 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 pričakovano vrednost:

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

Sedaj pa še zvezo

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

Torej deluje pričakovana vrednost na diskretne in zvezne slučajne spremenljivke linearno, tj.

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

Z matematično indukcijo posplošimo to na poljubno končno število členov

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

V Dodatku A.9 izpeljemo še naslednjo trditev.

Trditev 8.1. Če obstajata pričakovani vrednosti $\mathbb{E}(X^2)$ in $\mathbb{E}(Y^2)$, obstaja tudi pričakovana vrednost produkta $\mathbb{E}(XY)$ ter velja ocena $\mathbb{E}(|XY|) \leq \sqrt{\mathbb{E}(X^2)\mathbb{E}(Y^2)}$.

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



Primer: Življenska doba varovalk (merjena v stotinah ur), ki jih uporabljam 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** $D(X)$ slučajne spremenljivke, ki ima pričakovano vrednost, je določena z izrazom

$$D(X) = \mathbb{E}(X - \mathbb{E}(X))^2 = D(X) = \mathbb{E}(X^2) - (\mathbb{E}(X))^2.$$

Disperzija je vedno nenegativna, $D(X) \geq 0$, je pa lahko tudi neskončna, tj. ne obstaja. Naj bo a realna konstanta. Če je $P(X = a) = 1$, je $D(X) = 0$. Iz linearnosti pričakovane vrednosti sledi tudi $D(aX) = a^2D(X)$ in $D(X + a) = D(X)$.

Trditev 8.2. Če obstaja $D(X)$ in je a realna konstanta, obstaja tudi $E(X - a)^2$ in velja

$$E(X - a)^2 \geq D(X).$$

Enakost velja natanko za $a = E(X)$.

Količino $\sigma_X = \sqrt{D(X)}$ imenujemo **standardna deviacija** ali **standardni odklon**. Za poljubno realno število a in slučajno spremenljivko X, za katero obstaja disperzija, velja $\sigma_{aX} = a\sigma_X$ za $a \geq 0$

8.3 Standardizirane spremenljivke

Slučajno spremenljivko X **standardiziramo** s transformacijo

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

kjer sta $\mu = E(X)$ in $\sigma = \sqrt{D(X)}$. Za X_S velja $E(X_S) = 0$ in $D(X_S) = 1$, saj je

$$E(X_S) = E\left(\frac{X - \mu}{\sigma}\right) = \frac{E(X - \mu)}{\sigma} = \frac{\mu - \mu}{\sigma} = 0,$$

kjer smo upoštevali linearnost pričakovane vrednosti in da je pričakovanata vrednost konstante kar ta konstanta, ter

$$D(X_S) = D\left(\frac{X - \mu}{\sigma}\right) = \frac{D(X - \mu)}{\sigma^2} = \frac{\sigma^2(+0)}{\sigma^2} = 1,$$

kjer naj bralec sam premisli, kaj je potrebno upoštevati.

Pričakovane vrednosti in disperzije nekaterih porazdelitev

Izračunajte si jih za vajo.

porazdelitev	$E(X)$	$D(X)$
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
Hipergeometrijska $H(n; M, N)$	nM/N	$\frac{M(N - M)n(N - n)}{N^2(N - 1)}$
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

$$\text{Cov}(X, Y) = \mathbb{E}((X - \mathbb{E}(X)) \cdot (Y - \mathbb{E}(Y))).$$

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

$$\text{Cov}(X, Y) = \mathbb{E}(XY) - \mathbb{E}(X)\mathbb{E}(Y).$$

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

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

V Dodatku A.9 izpeljemo še naslednjo trditev.

Trditev 8.3. Če obstajata $\mathbb{D}(X)$ in $\mathbb{D}(Y)$, obstaja tudi $\text{Cov}(X, Y)$ in velja

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

Enakost velja natanko takrat, ko je $Y - \mathbb{E}(Y) = \pm \frac{\sigma_Y}{\sigma_X} (X - \mathbb{E}(X))$ z verjetnostjo 1. \square

Trditev 8.4. Če sta slučajni spremenljivki, ki imata pričakovano vrednost, neodvisni, obstaja tudi pričakovana vrednost njunega produkta in velja $\mathbb{E}(XY) = \mathbb{E}(X) \cdot \mathbb{E}(Y)$, tj. $\text{Cov}(X, Y) = 0$.

\square

Spremenljivki, za kateri velja $\mathbb{E}(XY) \neq \mathbb{E}(X) \cdot \mathbb{E}(Y)$, tj. $\text{Cov}(X, Y) \neq 0$, imenujemo **korelirani**, (sicer rečemo, da sta **nekorelirani**).

Obstajajo tudi odvisne spremenljivke, ki so nekorelirane. Z drugimi besedami, nekoreliranost ni zadosten pogoj za neodvisnost.

Primer: Npr., če je X zvezna enakomerno porazdeljena slučajna spremenljivka na intervalu $[-1, 1]$ (lahko bi bila tudi diskretna z zalogo vrednostmi i/n za $i = -n, \dots, 0, \dots, n$) in je $Y = X^2$, potem sta X in Y nekorelirani, čeprav spremenljivka X natanko določa spremenljivko Y , kar pomeni, da sta očitno odvisni. \diamond

Primer: Naj bo slučajna spremenljivka X porazdeljena standardizirano normalno. Potem je $Y = X^2$ porazdeljena po $\chi^2(1)$. Velja tudi $\mathbb{E}(X) = 0$, $\mathbb{E}(XY) = \mathbb{E}(X^3) = 0$ in zato $\mathbb{E}(XY) = 0 = \mathbb{E}(X) \cdot \mathbb{E}(Y)$. 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

Če imata spremenljivki X in Y končni disperziji, jo ima tudi njuna vsota $X + Y$ in velja

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

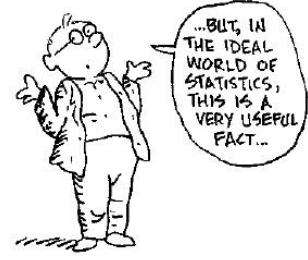
Če pa sta spremenljivki nekorelirani, je enostavno $D(X + Y) = D(X) + D(Y)$,
za odklon pa velja nekakšen ‘Pitagorov izrek’: $(\sigma_{X+Y})^2 = (\sigma_X)^2 + (\sigma_Y)^2$.

Zvezo lahko posplošimo na

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

in za paroma nekorelirane spremenljivke v aditivno pravilo

$$D\left(\sum_{i=1}^n X_i\right) = \sum_{i=1}^n D(X_i) \text{ oz. } (\sigma_{\sum_{i=1}^n X_i})^2 = \sum_{i=1}^n (\sigma_{X_i})^2.$$



Koreacijski koeficient slučajnih spremenljivk X in Y je definiran z izrazom

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

Primer: Za $(X, Y) \sim N(\mu_X, \mu_Y, \sigma_X, \sigma_Y, \rho)$ je $r(X, Y) = \rho$. Torej sta normalno porazdeljeni slučajni spremenljivki X in Y neodvisni natanko takrat, ko sta nekorelirani. \diamond

V splošnem velja:

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

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

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

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

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

8.5 Kovariančna matrika

Pričakovana vrednost slučajnega vektorja $\mathbf{X} = (X_1, \dots, X_n)$ je tudi vektor, in sicer $\mathbb{E}(\mathbf{X}) = (\mathbb{E}(X_1), \dots, \mathbb{E}(X_n))$.

Primer: Za $(X, Y) \sim N(\mu_X, \mu_Y, \sigma_X, \sigma_Y, \rho)$ je $\mathbb{E}(X, Y) = (\mu_X, \mu_Y)$. \diamond

Naj slučajna spremenljivka Y prestavlja linearno kombinacijo spremenljivk X_1, \dots, X_n . Spomnimo se, da je njena pričakovana vrednost

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

Za disperzijo spremenljivke Y pa dobimo

$$\begin{aligned} D(Y) &= E(Y - E(Y))^2 = E\left(\sum_{i=1}^n \sum_{j=1}^n a_i a_j (X_i - E(X_i))(X_j - E(X_j))\right) \\ &= \sum_{i=1}^n \sum_{j=1}^n a_i a_j \text{Cov}(X_i, X_j) = \mathbf{a}^T \mathbf{K} \mathbf{a}, \end{aligned}$$

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

Kovariančna matrika $\mathbf{K} = [K_{ij}]$ je *simetrična*, tj. $K_{ij} = K_{ji}$. Diagonalne vrednosti so disperzije spremenljivk: $K_{ii} = D(X_i)$. Ker je $\mathbf{a}^T \mathbf{K} \mathbf{a} = D(Y) \geq 0$, 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 $0 \leq D(Y) = \mathbf{a}^T \mathbf{K} \mathbf{a} = \lambda$, 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 $(X, Y) \sim N(\mu_X, \mu_Y, \sigma_X, \sigma_Y, \rho)$ 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 pričakovane vrednosti in disperzije. **Moment reda** $k \in \mathbb{N}$ glede na točko $a \in \mathbb{R}$ imenujemo količino

$$\mathbf{m}_k(\mathbf{a}) = E((X - a)^k).$$

Moment obstaja, če obstaja pričakovana vrednost $E(|X - a|^k) < \infty$. Za $a = 0$ dobimo **začetni moment** $\mathbf{z}_k = m_k(0)$; za $a = E(X)$ pa **centralni moment** $\mathbf{m}_k = m_k(E(X))$.

Primer: $E(X) = z_1$ in $D(X) = m_2$. \diamond

Trditev 8.5. 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 $a \in \mathbb{R}$ ter velja

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

Posebej za centralni moment velja $m_0 = 1$, $m_1 = 0$, $m_2 = z_2 - z_1^2$, $m_3 = z_3 - 3z_2 z_1 + 2z_1^3$, ...

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

Asimetrija spremenljivke X imenujemo količino $A(X) = m_3/\sigma^3$, **sloščenost** pa količino $K(X) = m_4/\sigma^4 - 3$, kjer je $\sigma = \sqrt{m_2}$. Obe sta meri za obliko.

Trditev 8.6. Za simetrično glede na $z_1 = \mathbb{E}(X)$ porazdeljene spremenljivke so vsi lihi centralni momenti enaki 0. Če sta spremenljivki X in Y neodvisni, je $m_3(X + Y) = m_3(X) + m_3(Y)$.

Primeri: Za $X \sim N(\mu, \sigma)$ so $m_{2k+1} = 0$ in $m_{2k} = (2k-1)!!\sigma^{2k}$. Zato sta tudi $A(X) = 0$ in $K(X) = 0$. Za $X \sim B(n, p)$ je $m_3(X) = npq(q-p)$ in $A(X) = (q-p)/\sqrt{npq}$. \diamond

Kadar spremenljivka nima momentov, uporabljamo kvantile. **Kvantil reda** $p \in (0, 1)$ je vsaka vrednost $x \in \mathbb{R}$, za katero velja

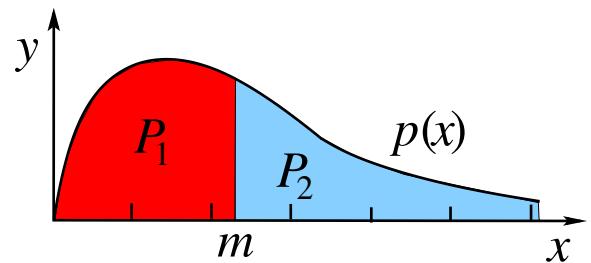
$$P(X \leq x) \geq p \text{ in } P(X \geq x) \geq 1 - p$$

ozziroma $F(x-) \leq p \leq F(x)$. Kvantil reda p označimo z x_p . Za zvezno spremenljivko je $F(x_p) = p$.

Kvantil $x_{\frac{1}{2}}$ imenujemo **mediana** (glej sliko).

Vrednosti $x_{\frac{i}{4}}$, $i = 0, 1, 2, 3, 4$ so **kvartili**. Kot nadomestek za standardni odklon uporabljamo **kvartilni razmik**

$$\frac{x_{\frac{3}{4}} - x_{\frac{1}{4}}}{2}$$



Slika: Navpična črta pri vrednosti m razdeli ploščino med grafom gostote verjetnosti $p(x)$ in x -osjo na dva dela: P_1 in P_2 . Ker je ploščina pod krivuljo enaka 1, velja $P_1 + P_2 = 1$. Vemo tudi, da je $P(X < m) = P_1$ in $P(X \geq m) = P_2$. Če je torej m mediana, potem je $P_1 = P_2$. Z drugimi besedami, **mediana** je tista vrednost ozziroma tisto realno (predstavljeno na osi x), pri katerem ustrezna napičnica razdeli ploščino grafom $y = p(x)$ in x -osjo na dva enaka dela. Če je krivulja $y = p(x)$ simetrična, potem je mediana ravno na sredini.

Poglavlje 9

Karakteristične funkcije in limitni izreki



Reprodukcijska lastnost normalne porazdelitve

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

Izrek 9.1. Če so slučajne spremenljivke X_1, \dots, X_n neodvisne in normalno porazdeljene $N(\mu_i, \sigma_i)$, potem je njihova vsota X tudi normalno porazdeljena:

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

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

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

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

9.1 Karakteristična funkcija

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

$$E(Z) = E(X) + iE(Y) \quad \text{in} \quad D(Z) = E(|Z - E(Z)|^2) = D(X) + D(Y).$$

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 $\varphi_X(t) = \mathbb{E}(e^{itX})$. 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.

9.2 Limitni izreki

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

$$\lim_{n \rightarrow \infty} P(|X_n - X| \geq \varepsilon) = 0, \quad \text{tj.} \quad \lim_{n \rightarrow \infty} P(|X_n - X| < \varepsilon) = 1.$$

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

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

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

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

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

Naj bo $\mathbf{X}_1, \dots, \mathbf{X}_n$ zaporedje spremenljivk, ki imajo pričakovano vrednost. Označimo $\mathbf{S}_n = \sum_{k=1}^n X_k$ in

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

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
$$\lim_{n \rightarrow \infty} P\left(\left|\frac{S_n - \mathbb{E}(S_n)}{n}\right| < \varepsilon\right) = 1;$$
- **krepki zakon velikih števil**, če gre skoraj gotovo $Y_n \rightarrow 0$, tj., če velja
$$P\left(\lim_{n \rightarrow \infty} \frac{S_n - \mathbb{E}(S_n)}{n} = 0\right) = 1.$$

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

Izrek 9.2. [Neenakost Čebiševa] Če ima slučajna spremenljivka X končno disperzijo, tj. $D(X) < \infty$, potem za vsak $\varepsilon > 0$ velja

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



Dokaz: Pokažimo jo za zvezne spremenljivke

$$\begin{aligned} P(|X - E(X)| \geq \varepsilon) &= \int_{|x - E(X)| \geq \varepsilon} p(x) dx = \frac{1}{\varepsilon^2} \int_{|x - E(X)| \geq \varepsilon} \varepsilon^2 p(x) dx \\ &\leq \frac{1}{\varepsilon^2} \int_{-\infty}^{\infty} (x - E(X))^2 p(x) dx = \frac{D(X)}{\varepsilon^2}. \quad \square \end{aligned}$$

Posledica 9.3. [Markov] Če gre za zaporedje slučajnih spremenljivk X_i izraz $\frac{D(S_n)}{n^2} \rightarrow 0$, ko gre $n \rightarrow \infty$, velja za zaporedje šibki zakon velikih števil.

Posledica 9.4. [Čebišev] Če so slučajne spremenljivke X_i paroma nekorelirane in so vse njihove disperzije omejene z isto konstanto C , tj. $D(X_i) < C$ za vsak i , velja za zaporedje šibki zakon velikih števil.

Za Bernoullijevo zaporedje X_i so spremenljivke paroma neodvisne, $D(X_i) = pq$, $S_n = k$. Pogoji izreka Čebiševa so izpolnjeni in zopet smo prišli do Bernoullijevega zakona velikih števil, tj. izreka 5.5.

Izrek 9.5. [Hinčin] Če so neodvisne slučajne spremenljivke X_i enako porazdeljene in imajo pričakovano vrednost $E(X_i) = a$ 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 9.6. [Kolmogorov] Naj bodo slučajne spremenljivke X_i neodvisne.

- Če imajo končno disperzijo in velja $\sum_{n=1}^{\infty} \frac{D(S_n)}{n^2} < \infty$, potem velja zanje

$$\text{krepki zakon velikih števil: } P\left(\lim_{n \rightarrow \infty} \frac{S_n - E(S_n)}{n} = 0\right) = 1.$$

- Če so enako porazdeljene in imajo pričakovano vrednost $E(X_i) = \mu$, potem velja zanje

$$\text{krepki zakon velikih števil: } P\left(\lim_{n \rightarrow \infty} \frac{1}{n} \sum_{i=1}^n X_i = \mu\right) = 1.$$

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.



(Osnovni CLI) Če so slučajne spremenljivke X_i neodvisne, enako porazdeljene s končno pričakovano vrednostjo in končno disperzijo ter je $S_n = X_1 + \dots + X_n$, potem zanje velja **centralni limitni zakon**, tj. porazdelitvene funkcije za standardizirane spremenljivke

$$Z_n = \frac{S_n - \mathbb{E}(S_n)}{\sigma(S_n)}$$

gredo proti porazdelitveni funkciji standardizirane normalne porazdelitve oziroma simbolično za vsak $x \in \mathbb{R}$ velja

$$\lim_{n \rightarrow \infty} P\left(\frac{S_n - \mathbb{E}(S_n)}{\sigma(S_n)} < x\right) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^x e^{-\frac{t^2}{2}} dt.$$

Opomba: v praksi uporabimo $n > 30$.

Skica dokaza 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 $k = E(Z_i^3)$. Nadalje velja

$$\log M_n(t) = n \log \left(1 - \frac{t^2}{2n} + \frac{t^3}{3!n^{3/2}} k + \dots \right).$$

Od tod za $x = \left(-\frac{t^2}{2n} + \frac{t^3}{3!n^{3/2}} k + \dots \right)$ dobimo $\log M_n(t) = n \log(1 + x) = n(x - \frac{x^2}{2} + \dots)$ oziroma

$$\log M_n(t) = n \left[\left(-\frac{t^2}{2n} + \frac{t^3}{3!n^{3/2}} k + \dots \right) - \frac{1}{2} \left(-\frac{t^2}{2n} + \frac{t^3}{3!n^{3/2}} k + \dots \right)^2 + \dots \right]$$

in končno še

$$\lim_{n \rightarrow \infty} \log M_n(t) = -\frac{t^2}{2} \quad \text{oziora} \quad \lim_{n \rightarrow \infty} M_n(t) = e^{-t^2/2}.$$

Iz konvergencije karakterističnih funkcij φ_{Y_n} proti karakteristični funkciji standardizirano normalne porazdelitve lahko sklepamo po obratnem konvergenčnem izreku, da tudi porazdelitvene funkcije za Y_n konvergirajo proti porazdelitveni funkciji standardizirano normalne porazdelitve. Torej velja centralni limitni zakon. □

Del II

STATISTIKA



Skozi življenje se prebijamo z odločitvami,
ki jih naredimo na osnovi nepopolnih informacij ...

Pridobivanje podatkov

Novice so polne številk. 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 številk 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 pogоворov, analize DNK iz vzorcev krvi in podobno. Pri statističnem sklepanju nam dokaze priskrbijo podatki. Po domače statistično sklepanje velikokrat temelji na grafični predstavitev 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 diplomo, ki je dolga 100 strani, kolega pa ima za 20 strani daljšo diplomo. Če za trenutek pustimo ob strani samo vsebino (kvaliteto), je še vedno vprašanje ali je bil res boljši od nas v kvantiteti. Uporabljaj je drugačen font, njegov rob je nekoliko večji,... Kako lahko na hitro ocenimo dejansko stanje (brez da bi primerjali sami datoteki)?

Nadalujmo 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 številk? 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 modeliranja igrata 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 programsko opremo, 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”, ”Zares 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 kocic 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 številk 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 vsaj 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 žeeli 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 kocic (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 2116 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 17 nezgodami s smrtnim izidom pri enakem številu prepotovanih km z avtom. Ta članek vsebuje veliko številk in statistike, toda kaj vse to sploh pomeni?

Ljudje običajno besedo *statistika* povezujejo z zbiranjem in urejanjem podatkov o nekem pojavu, izračunom raznih značilnosti iz teh podatkov, njih predstavljivjo 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. Kot smo že v uvodu omenili, pravimo Tej veji statistike *opisna statistika*. Druga veja, *inferenčna statistika*, poskuša spoznanja iz zbranih podatkov posplošiti (razširiti, podaljšati, napovedati, ...) in oceniti kakovost teh pospošitev. Statistiko lahko razdelimo tudi na *uporabno* in *teoretično* (računalniško in matematično) statistiko.



Statistika je veda, ki proučuje množične pojave.

(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 ustreerne 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.

Poglavlje 10

Opisna statistika

10.1 Vrste slučajnih spremenljivk ozziroma podatkov

Glede na vrsto vrednosti delimo slučajne spremenljivke na:

1. *številske* (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.

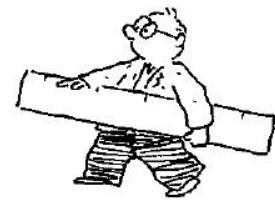


Če so kategorije brez odgovarjajočega vrstnega reda/urejenosti, rečemo, da so podatki *imenski* (nominalni). Kakor hitro imajo kategorije neko urejenost pa rečemo, da so podatki *ordinalni* (ali urejenostni).

Glede na vrsto merske lestvice delimo slučajne spremenljivke na:

1. *imenske* (ali nominalne) spremenljivke – vrednosti lahko le razlikujemo med seboj: dve vrednosti sta enaki ali različni (npr. spol: moški/ženski, kraj bivanja, narodnost, GSM številka, inventarna številka); smiselno je računati število enot, pogostost oz. frekvenco, modus, deleže,...

2. *urejenostne* (ali ordinalne) spremenljivke – vrednosti lahko uredimo od najmanjše do največje, a ni mogoče določiti za koliko je ena vrednost večja od druge (npr. uspeh: odličen, pravdober, dober, zadosten, nezadosten); smiselno je računati še kvantile (mediana, kvartile, decile,...), ne pa aritmetične sredine;
3. *razmične* (ali intervalne) spremenljivke – lahko primerjamo razlike med vrednostima dvojic enot (npr. temperatura, čas po koledarju, merjenje inteligenčnega kvocienta,...); smiselno je računati še aritmetično sredino, standarni odklon,...
4. *razmernostne* spremenljivke – lahko primerjamo razmerja med vrednostima dvojic enot (npr. starost, vrednost povzročene škode, število prometnih nesreč); smiselno je računati še geometično sredino, harmonično sredino,...
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

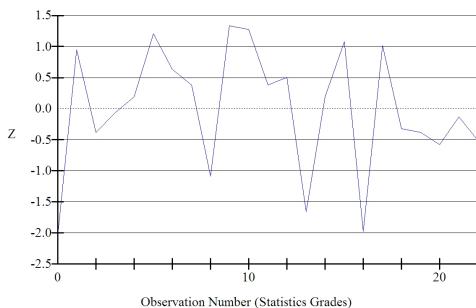
Vrste spremenljivk so urejene od tistih z najslabšimi merskimi lastnostmi do tistih z najboljšimi. Urejenostne spremenljivke zadoščajo lastnostim, ki jih imajo imenske spremenljivke, podobno razmernostne spremenljivke zadoščajo lastnostim, ki jih imajo razmične, urejenostne in imenske spremenljivke, itd.:

$$\text{absolutna} \subset \text{razmernostna} \subset \text{razmična} \subset \text{urejenostna} \subset \text{imenska}$$

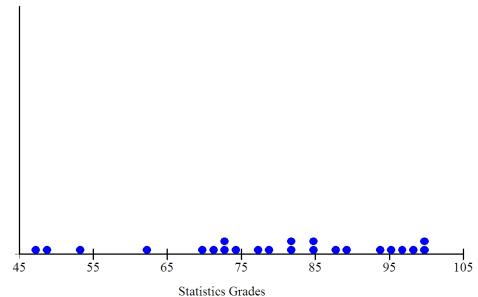
Merske lestvice	kategorije	urejenost	interval	dejanska ničla
razmična	x	x	x	x
razmernostna	x	x	x	
urejenostna	x	x		
imenska	x			

Posamezne statistične metode predpostavljajo določeno vrsto spremenljivk. Največ učinkovitih statističnih metod je razvith za številske spremenljivke. V teoriji merjenja pravimo, da je nek stavek *smiseln*, če ohranja resničnost/lažnost pri zamenjavi meritev z enakovrednimi (glede na dovoljene transformacije) meritvami.

10.2 Grafična predstavitev kvantitativnih podatkov



(a) grafični prikaz podatkov (angl. runs chart/plot),



(b) diagram frnikul (angl. dot plot)

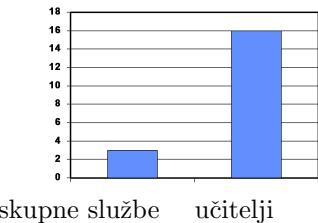
Primer: Oddelek sistemskih inženirjev

kategorija	frekvenca	relativna frekvenca
vrsta	število	
zaposlenih	zaposlenih	delež
učitelji	16	0.8421
skupne službe	3	0.1579
skupaj	19	1.0000

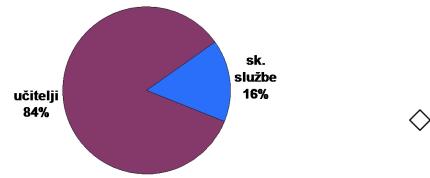


Stolpčni prikaz (tudi stolpčni graf, poligonski diagram):

Na eni osi prikažemo (urejene) razrede. Nad vsakim na redimo stolpec/črto višine sorazmerne frekvenci razreda.



Krožni prikaz (tudi strukturni krog, pogača, kolač): Vsakemu razredu priredimo krožni izsek s kotom $\alpha_i = \frac{k_i}{n} 360$ stopinj.



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.

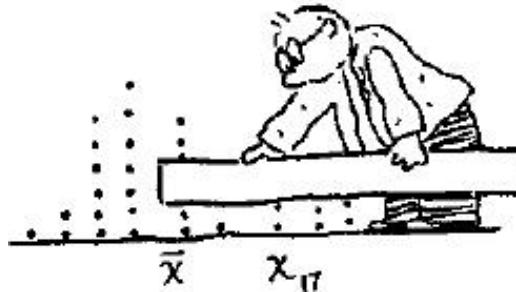
Meje razredov so določene tako, da velja $x_{i,max} = x_{i+1,min}$.

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. *Relativna frekvenca* rel. $f_i = f_i/N$ in *relativna kumulativa* rel. $F_i = F_i/N$, kjer je N število vseh podatkov.

Poligon: v koordinatnem sistemu zaznamujemo točke (x_i, f_i) , kjer je x_i sredina i -tega razreda in f_i njegova frekvenca (tj. koliko podatkov je padlo v i -ti razred). 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.



Primer: Oglejmo si zaporedje podatkov

- (a) Konstruiraj urejeno zaporedje.
- (b) Nariši steblo-list diagram.
- (c) Naredi histogram.

88	103	113	122	132
92	108	114	124	133
95	109	116	124	133
97	109	116	124	135
97	111	117	128	136
97	111	118	128	138
98	112	119	128	138
98	112	120	131	142
100	112	120	131	146
100	113	122	131	150

◇

Koraki za konstrukcijo steblo-list predstavitve

1. Razdeli vsako opazovanje-podatke na dva dela: **steba** (angl. stem) in **listi** (angl. leaf).
2. Naštej steba 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 frekvence za vsako steblo.

Steblo-list diagram

stebla	listi	f_i	$rel.f_i$
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%



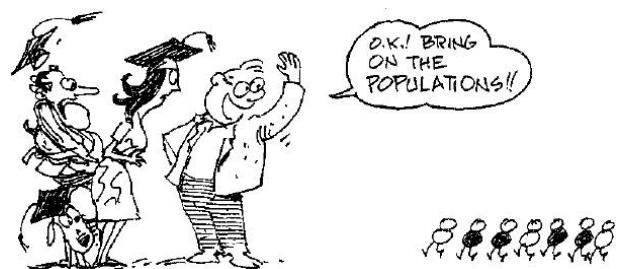
Histogrami

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 $(x_{i,min}, F_i)$.

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**.



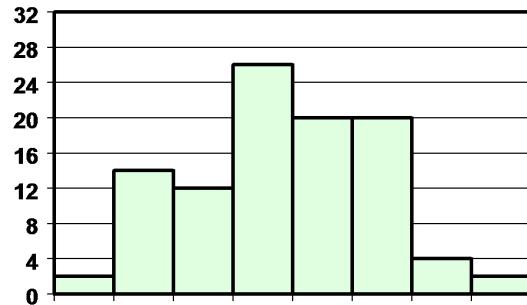
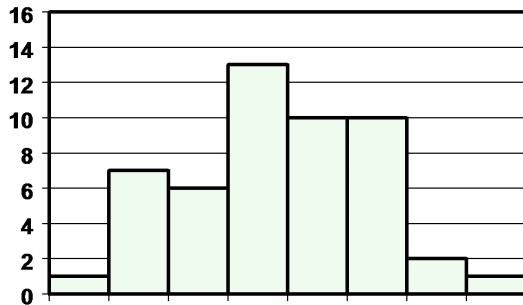
Pravilo za določanje števila razredov v histogramu

število vzorcev v množici podatkov	število razredov
manj kot 25	5 ali 6
25 – 50	7 – 14
več kot 50	15 – 20

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] [,9]
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="s",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="s")
> curve(pnorm(x,mean=175,sd=30),add=T,col="red")

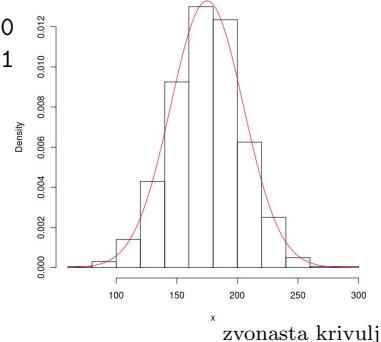
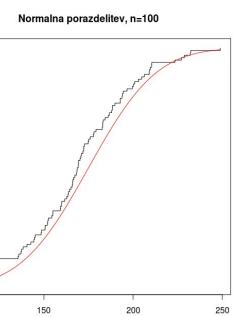
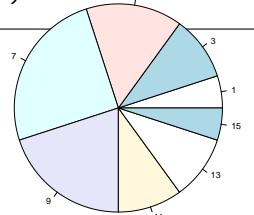
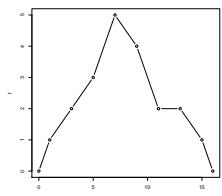
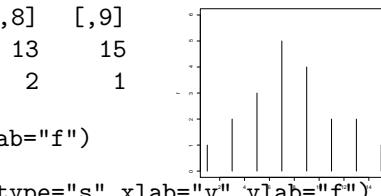
```

```

> 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%      40%      50%
 92.09012 135.83928 148.33908 158.53864 166.96955 174.4801
 60%      70%      80%      90%     100%
182.08577 191.29261 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")

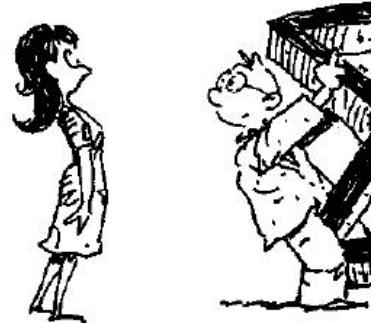
```

Zaporedje (20ih) števil spravimo v vektor X. Za X lahko izračunamo njegovo dolžino (število komponent) in frekvenco vseh elementov na razponu od najmanjšega do največjega (tabulate). Sledi seznam elementov, katerih frekvence so pozitivne in grafični prikazi...

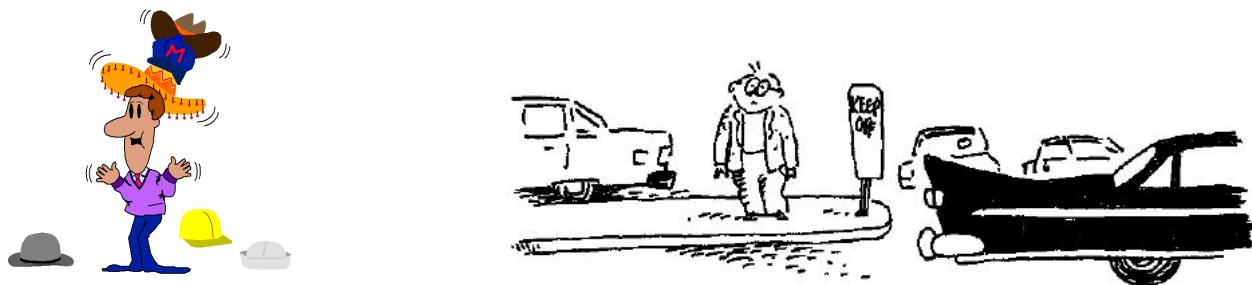


10.3 Mere za lokacijo in razpršenost

- srednje vrednosti
- razpon (min/max)
- centili, kvartili
- varianca
- standardni odklon
- Z -vrednosti

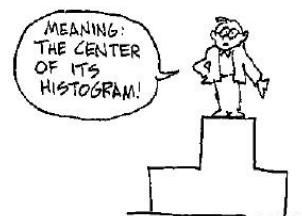


Modus (oznaka M_0) množice podatkov je tista vrednost, ki se pojavi z največjo frekvenco.



Da bi prišli do **medianе** (oznaka M_e) za neko množico podatkov, naredimo naslednje:

1. Podatke uredimo po velikosti v naraščajočem 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.



Povrečje populacije:
$$\mu = \frac{x_1 + \dots + x_N}{N} = \frac{\sum_{i=1}^N x_i}{N}.$$

Povrečje vzorca:
$$\bar{x} = \frac{x_1 + \dots + x_n}{n} = \frac{\sum_{i=1}^n x_i}{n}.$$

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**.

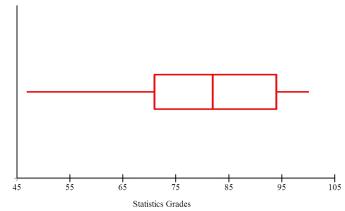
50. centil se imenuje **2. kvartil** ali **mediana**.

75. centil se imenuje tudi **3. kvartil**.



Še nekaj ukazov v R-ju: škatle in Q-Q-prikazi

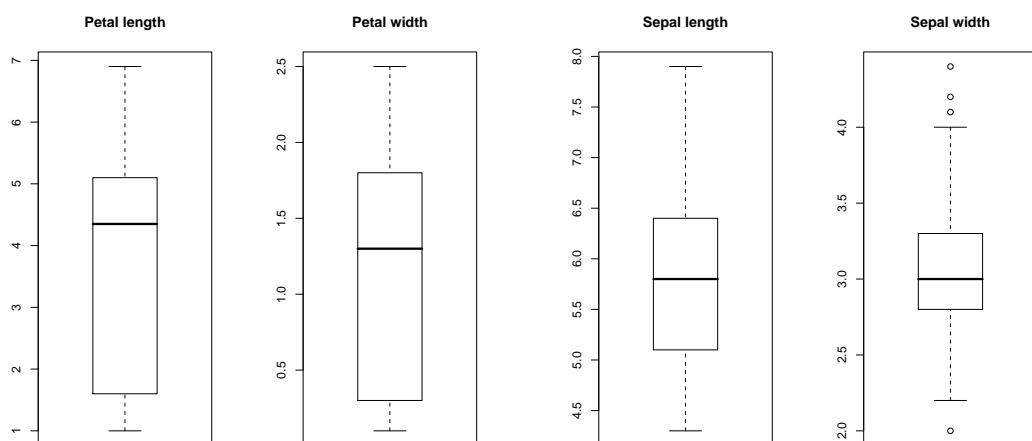
Škatle z brki (angl. box plot) (box-and-whiskers plot; grafikon kvantilov) **boxplot**: škatla prikazuje notranja kvartila razdeljena z mediansko črto. Daljici – brka vodita do robnih podatkov, ki sta največ za 1.5 dolžine škatle oddaljena od nje. Ostali podatki so prikazani posamično.



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.

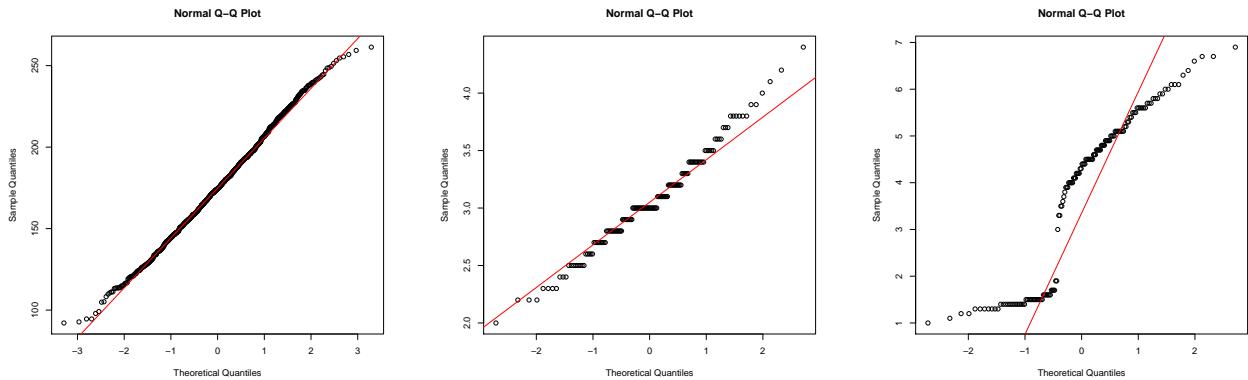
Obstaja tudi splošnejši ukaz `qqplot`, ki omogoča prikaz povezanosti poljubnega para porazdelitev. S parametrom `datax=T` zamenjamo vlogo koordinatnih osi.

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



```
> qnorm(x)
> qqline(x,col="red")
> qnorm(iris$Sepal.Width)
> qqline(iris$Sepal.Width,col="red")
> qnorm(iris$Petal.Length)
> qqline(iris$Petal.Length,col="red")
```

Mere razpršenosti: varianca in standardni odklon

Varianca populacije je povprečje kvadratov odklonov od pričakovane vrednosti

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

Varianca vzorca (prodobljenega z n meritvami), v tem primeru vsoto kvadratov odklonov delimo s stopnjo prostosti (vzorca), tj. $n - 1$ namesto n :

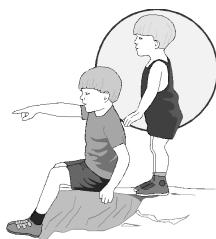
$$s^2 = \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n - 1} = \frac{\sum_{i=1}^n x_i^2 - \frac{\left(\sum_{i=1}^n x_i\right)^2}{n}}{n - 1}.$$

Standardni odklon (deviacija) je pozitivno predznačen kvadratni koren variance, **koeficient variacije** pa standardni odklon deljen s povprečjem.

populacija vzorec

varianca	σ^2	s^2
D, V		

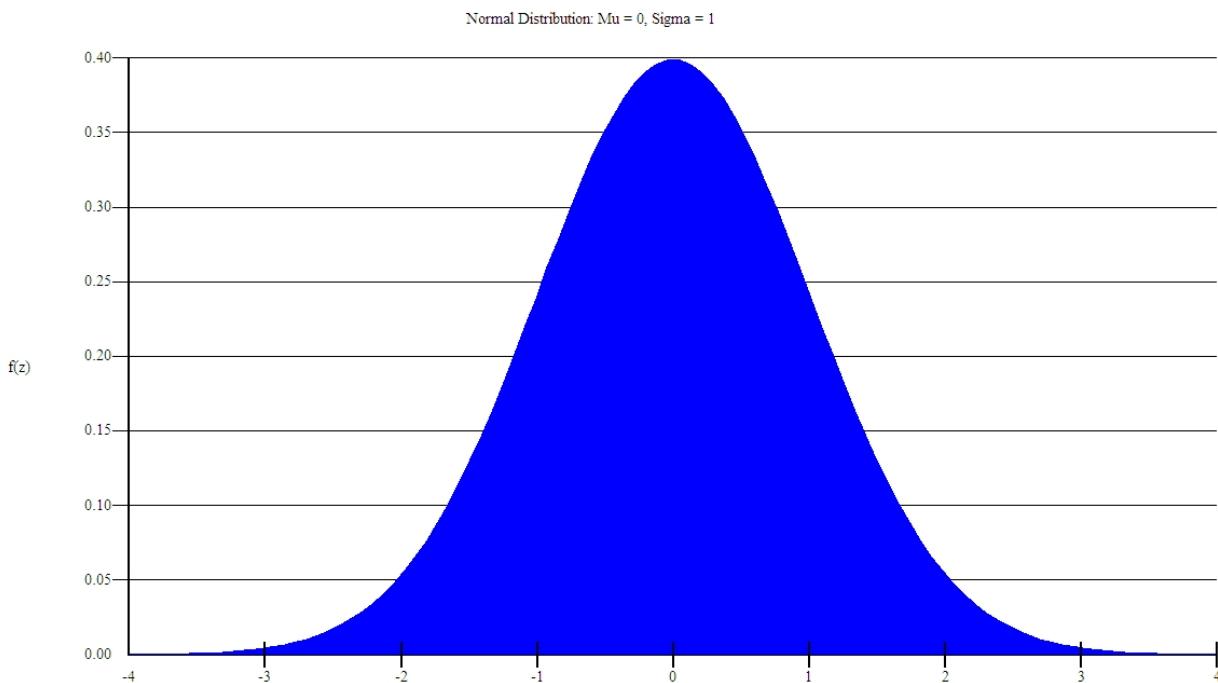
standardni odklon	σ	s
----------------------	----------	-----



Za vzorec vzamemo npr. osebje FRI in zabeležimo naslednje število otrok: 1, 2, 2, 1, 2, 5, 1, 2. V tem primeru je $n = 8$, $\bar{x} = 2$ in $s^2 = 12/7 = 1.71$, če pa bi delili z 8 bi dobili $s_0^2 = 3/2$.

Normalna porazdelitev

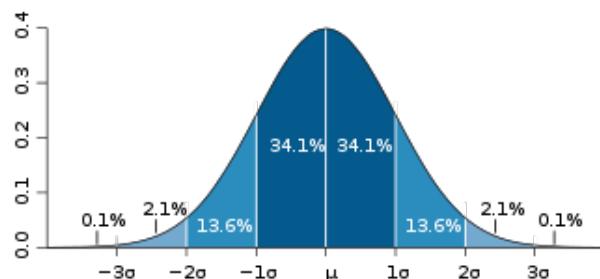
Veliko podatkovnih množic ima porazdelitev približno **zvonaste oblike** (unimodalna oblika - ima en sam vrh):



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 intervalu $[\mu - \sigma, \mu + \sigma]$,
2. približno **95·4%** meritev leži na intervalu $[\mu - 2\sigma, \mu + 2\sigma]$,
3. skoraj vse meritve (**99·7%**) ležijo na intervalu $[\mu - 3\sigma, \mu + 3\sigma]$.



Mere za obliko

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*).

Za $\ell \in \mathbb{N}$ je **ℓ -ti centralni moment** enak

$$m_\ell = \frac{(x_1 - \bar{x})^\ell + \cdots + (x_n - \bar{x})^\ell}{n}.$$

$$m_1 = 0, m_2 = \sigma^2, \dots$$

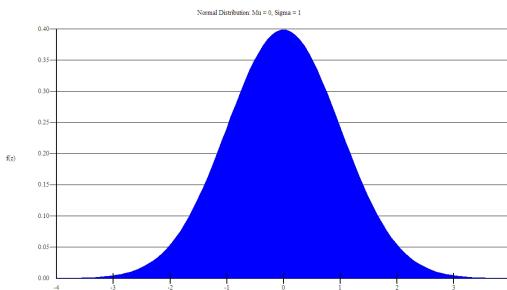
Koeficient asimetrije (s centralnimi momenti): $g_1 = m_3/m_2^{3/2}$. Mere asimetrije dobim tako, da opazujemo razlike med srednjimi vrednostimi. Le-te so tem večje čim bolj je porazdelitev asimetrična:

$$KA_{M_0} = (\bar{x} - M_0)/s, \quad KA_{M_e} = 3(\bar{x} - M_e)/s.$$

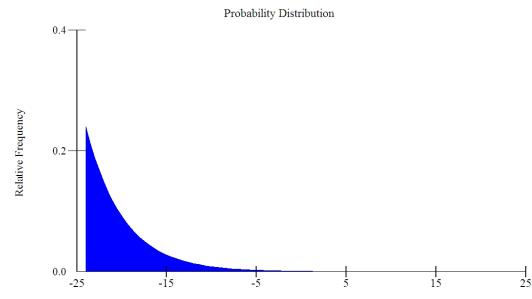
Koeficient sploščenosti (kurtosis) (s centralnimi momenti): $K = g_2 = m_4/m_2^2 - 3$

- $K = 3$ (ali 0) normalna porazdelitev zvonaste-oblike (*mesokurtic*),
- $K < 3$ (ali < 0) bolj kopasta kot normalna porazdelitev, s krajsimi repi (*platykurtic*),
- $K > 3$ (ali > 0) bolj špičasta kot normalna porazdelitev, z daljšimi repi (*leptokurtic*).

Normalna in asimetrična porazdelitev

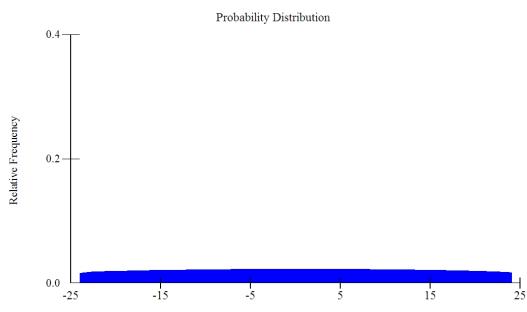


asim.= 0, sploščenost= 3 (mesokurtic).

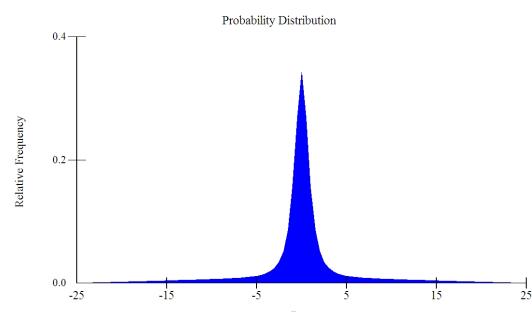


asim.= 1.99, sploščenost= 8.85.

Kopasta in špičasta porazdelitev



asim.= 0, sploščenost= 1.86 (platykurtic).



asim.= -1.99, sploščenost= 8.85 (leptokurtic).

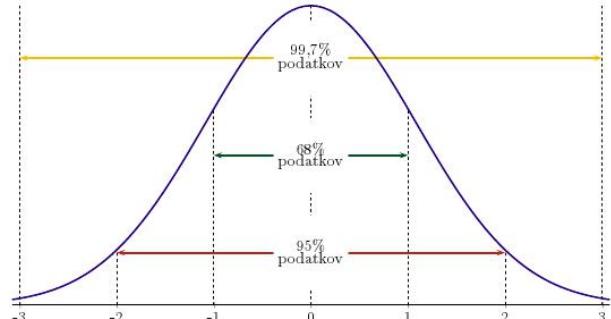
10.4 Standardizacija

Vsaki vrednosti x_i spremenljivke X odštejemo njeni pričakovano vrednosti μ in delimo z njenim standardnim odklonom σ :

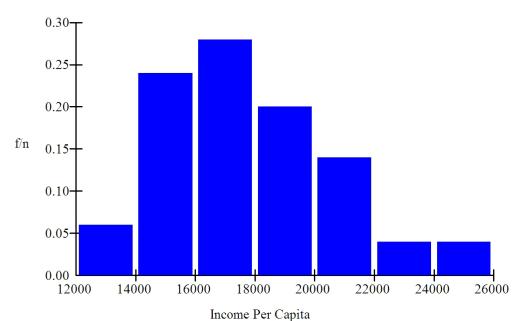
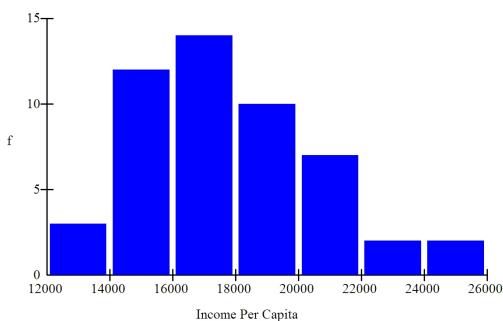
$$z_i = \frac{x_i - \mu}{\sigma}.$$

Za novo spremenljivko Z bomo rekli, da je **standardizirana**, z_i pa je **standardizirana vrednost**.

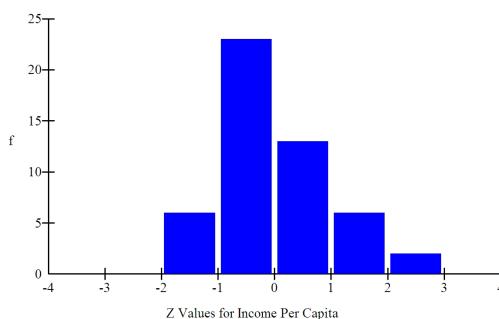
Potem je $\mu(Z) = 0$ in $\sigma(Z) = 1$.



Frekvenčni in relativni frekvenčni histogram

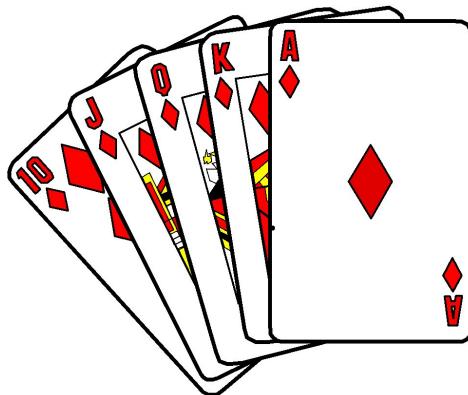


Histogram standardiziranih Z -vrednosti



Poglavlje 11

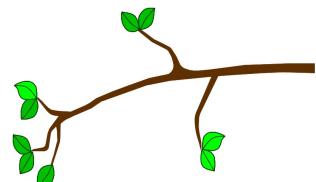
Vzorčenje



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.

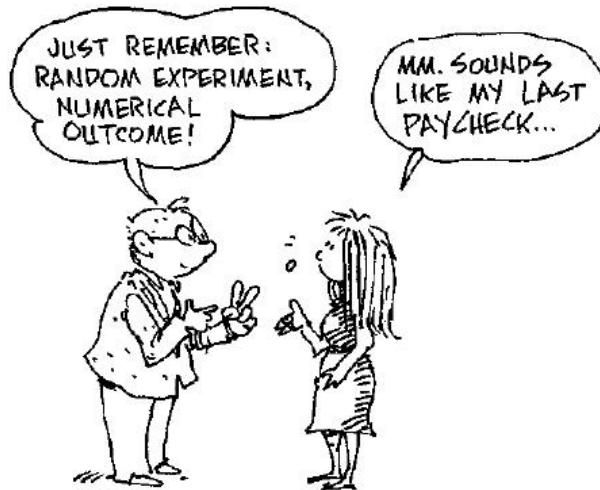
Kdaj vzorec dobro predstavlja celo populacijo? Preprost odgovor je:

- 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

$$(X_1, \dots, X_n),$$

vrednostim meritev (x_1, \dots, x_n) pa 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.

11.1 Osnovni izrek statistike

Spremenljivka X ima na populaciji porazdelitev $F(x) = P(X \leq x)$. Toda tudi vsakemu vzorcu ustreza neka porazdelitev. Za realizacijo vzorca (x_1, \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

$$Y_i(x) \sim \begin{pmatrix} 1 & 0 \\ F(x) & 1 - F(x) \end{pmatrix},$$

velja

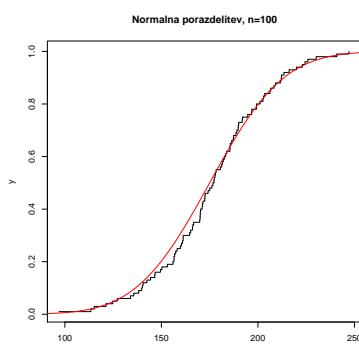
$$V_n(x) = \frac{1}{n} \sum_{i=1}^n Y_i(x).$$

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 $n = 1, 2, \dots$. Oznaka sup je okrajšava za supremum, ki predstavlja najmanjšo zgornjo mejo. V primeru končne množice razlik, gre za največjo razliko. Zanjo 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 bolje prilega funkciji $F(x)$ – vse bolje povzema razmere na celotni populaciji.

11.2 Vzorčne ocene

Najpogostejsa parametra, ki bi ju radi ocenili sta: *sredina populacije* μ glede na izbrano lastnost – pričakovana vrednost 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 (x_1, x_2, \dots, x_n) . Zato jim tudi rečemo *vzorčne ocene*.

Sredinske mere

Kot sredinske mere se pogosto uporablja:

Vzorčni modus – najpogosteša 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): $G(x) = \sqrt[n]{\prod_{i=1}^n x_i}$.
(definiran za pozitivne x_i).

Mere razpršenosti

Za oceno populacijskega odklona uporabljam *mere razpršenosti*.

Vzorčni razmah $= \max_i x_i - \min_i x_i$.

Vzorčna disperzija $s_0^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2$.

Popravljena vzorčna disperzija $s^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2$.

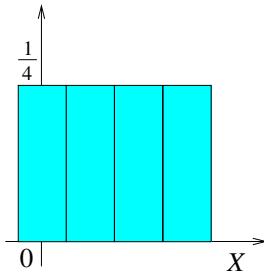
ter ustrezna *vzorčna odklona* s_0 in s .

11.3 Porazdelitve vzorčnih povprečij

Denimo, da je v populaciji N enot in da iz te populacije slučajno izbiramo n enot v enostavnih slučajnih 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 . 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. in izračunamo populacijsko pričakovano vrednost in varianco:



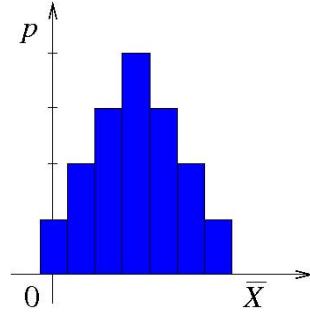
$$\mu = \frac{1}{N} \sum_{i=1}^N x_i = \frac{3}{2}, \quad \sigma^2 = \frac{1}{N} \sum_{i=1}^N (x_i - \mu)^2 = \frac{5}{4}.$$

Slika 5. Histogram porazdelitve spremenljivke X .

Sedaj pa tvorimo vse možne vzorce velikosti $n = 2$ s ponavljanjem, in na vsakem izračunajmo vzorčno povprečje \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



Slika 6. Histogram porazdelitve spremenljivke \bar{X} .

Zapišimo verjetnostno shemo slučajne spremenljivke vzorčno povprečje \bar{X} (glej sliko 6):

$$\bar{X} : \begin{pmatrix} 0 & 0.5 & 1 & 1.5 & 2 & 2.5 & 3 \\ 1/16 & 2/16 & 3/16 & 4/16 & 3/16 & 2/16 & 1/16 \end{pmatrix}.$$

Izračunajmo pričakovano vrednost ter varianco vzorčnega povprečja:

$$E(\bar{X}) = \sum_{i=1}^m \bar{x}_i p_i = \frac{0 + 1 + 3 + 6 + 6 + 5 + 3}{16} = \frac{3}{2},$$

$$D(\bar{X}) = \sum_{i=1}^m \left(\bar{x}_i - E(\bar{X}) \right)^2 p_i = \frac{5}{8}. \quad \diamond$$

Na primeru smo spoznali, da je statistika ‘vzorčno povprečje’ slučajna spremenljivka s svojo porazdelitvijo. Sedaj pa izračunajmo pričakovano vrednost in razpršenost za slučajno spremenljivko, ki spremišča vzorčno povprečje in je določena z naslednjo zvezo

$$\bar{X} = \frac{X_1 + \dots + X_n}{n} = \frac{1}{n} \sum_{i=1}^n X_i.$$

Zaradi linearnosti pričakovane vrednosti $E(X)$, homogenosti odklona σ_X , ‘Pitagorovega izreka’ za σ_X oziroma aditivnost variance $D(X)$ (za paroma nekorelirane spremenljivke X_i) in

dejstva, da je $E(X_i) = E(X)$ in $(\sigma_{X_i})^2 = \sigma^2$ za $i = 1, \dots, n$, velja

$$E(\bar{X}) = \frac{1}{n} \sum_{i=1}^n E(X_i) = \frac{nE(X)}{n} = \mu, \quad \text{in} \quad (\sigma_{\bar{X}})^2 = \frac{1}{n^2} \sum_{i=1}^n (\sigma_{X_i})^2 = \frac{1}{n^2} \sum_{i=1}^n \sigma^2 = \frac{\sigma^2}{n}.$$

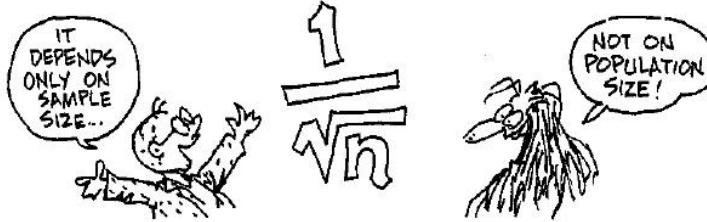
Zapišimo dobljeni ugotovitvi še v obliki trditve.

Trditve 11.1. *Na neki populaciji nas zanima lastnost, ki jo spremišča slučajna spremenljivka X s končnima parametroma $E(X) = \mu$ in $D(X) = \sigma^2$. Iz te populacije zberemo naključni slučajni vzorec (x_1, \dots, x_n) . Če so X_1, \dots, X_n slučajne spremenljivke, ki spremiščajo elemente tega vzorca po koordinatah, potem so neodvisne in zanje velja*

$$E(X_i) \approx \mu \quad \text{in} \quad D(X_i) \approx \sigma^2 \quad (1 \leq i \leq n).$$

Nadalje lahko pričakovano vrednost in standardni odklon vzorčnega povprečja \bar{X} ocenimo z

$$\mu_{\bar{X}} \approx \mu \quad \text{in} \quad \sigma_{\bar{X}} \approx \frac{\sigma}{\sqrt{n}}.$$



Iz druge zveze vidimo, da pada standardni odklon $\sigma_{\bar{X}}$ proti 0 z naraščanjem velikosti vzorca, tj. $\bar{X} \rightarrow \mu$ (enako nam zagotavlja tudi krepki zakon velikih števil).

Hitrost centralne tendence pri CLI

Spomnimo se **Centralnega limitnega izreka** in ga priredimo za primer vzorčnega povprečja:

Če je naključni vzorec velikosti n izbran iz populacije s

- končno pričakovano vrednostjo μ in končno varianco σ^2

ter če je

- n dovolj velik (npr. $n > 30$),

potem je porazdelitev standardiziranega vzorčnega povprečja \bar{X} ,
tj. $(\bar{X} - \mu_{\bar{X}})/\sigma_{\bar{X}} = (\bar{X} - \mu)/\sqrt{n}/\sigma$, aproksimirana z $N(0, 1)$.

Dokaz CLI je precej tehničen, kljub temu pa nam ne da občutka, kako velik mora biti n , da se porazdelitev slučajne spremenljivke

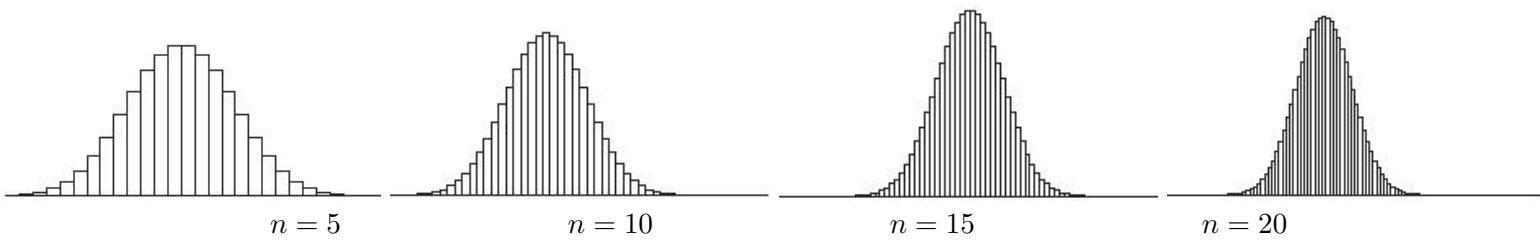
$$X_1 + \dots + X_n$$

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

$$p_1 = 1/6, \quad p_2 = 1/6, \quad p_3 = 1/6, \quad p_4 = 1/6, \quad p_5 = 1/6, \quad p_6 = 1/6.$$

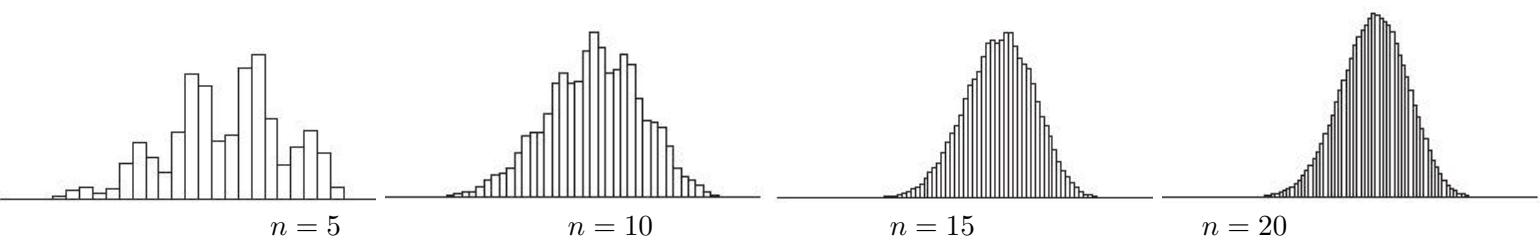
in slučajno spremenljivko $X_1 + X_2 + \dots + X_n$:



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 $X_1 + X_2 + \dots + X_n$:



11.4 Vzorčna statistika

Vzorčna statistika je poljubna simetrična funkcija (tj. njena vrednost je neodvisna od permutacije argumentov) vzorca

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

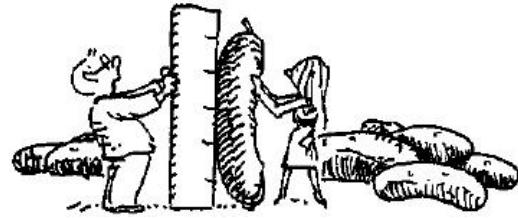
Tudi vzorčna statistika je slučajna spremenljivka, za katero lahko določimo porazdelitev iz porazdelitve vzorca. Najzanimivejši sta značilni vrednosti

- pričakovana vrednost $E(Y)$, za katero uporabimo vzorčno povprečje, in
- standardni odklon σ_Y , ki mu pravimo tudi *standardna napaka* statistike Y (angl. standard error – zato oznaka $SE(Y)$), za katerega uporabimo vzorčni odklon.



11.4.1 (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 (iskana/neznana) pričakovana vrednost dolžin, in je σ standardni odklon porazdelitve dolžin kumar, potem velja

$$\mathbb{E}(X_i) = \mu, \quad \text{in} \quad \mathbb{D}(X_i) = \sigma^2,$$

za vsak i , ker bi X_i bila lahko dolžina katerekoli kumare.



Denimo, da se spremenljivka X na populaciji porazdeljuje normalno $N(\mu, \sigma)$. Na vsakem vzorcu (s ponavljanjem) izračunamo vzorčno povprečje \bar{X} . Po reprodukcijski lastnosti normalne porazdelitve je *porazdelitev vzorčnih povprečij normalna, kjer je*

- pričakovana vrednost vzorčnih povprečij enaka pričakovani vrednosti spremenljivke na populaciji, tj.

$$\mathbb{E}(\bar{X}) = \mu,$$

- standardni odklon vzorčnih povprečij

$$\text{SE}(\bar{X}) = \frac{\sigma}{\sqrt{n}}.$$

Če tvorimo vzorce iz končne populacije brez vračanja, je standardni odklon vzorčnih povprečij

$$\text{SE}(\bar{X}) = \frac{\sigma}{\sqrt{n}} \sqrt{\frac{N-n}{N-1}}.$$

Za dovolj velike vzorce ($n > 30$) je porazdelitev vzorčnih povprečij približno normalna, tudi če spremenljivka X ni normalno porazdeljena. Če se spremenljivka X porazdeljuje vsaj približno normalno s standardno napako $SE(X)$, potem se

$$Z = \frac{X - E(X)}{SE(X)}$$

porazdeljuje standardizirano normalno.

Vzorčno povprečje in normalna porazdelitev

Naj bo $X \sim N(\mu, \sigma)$. Tedaj je $\sum_{i=1}^n X_i \sim N(n\mu, \sigma\sqrt{n})$ in dalje $\bar{X} \sim N(\mu, \sigma/\sqrt{n})$. Torej je vzorčna statistika

$$Z = \frac{\bar{X} - \mu}{\sigma} \sqrt{n} \sim N(0, 1).$$

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 $N(0 + \mu, 1 \cdot \sigma/\sqrt{n}) = N(\mu, \sigma/\sqrt{n})$, kar je enako kot v primeru, ko smo začeli z normalno porazdeljenimi slučajnimi spremenljivkami.

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. Zanjo je $\mu = 3.5$ in standardni odklon $\sigma = 1.7$. Vseh 36 ponovitev meta lahko obravnavamo kot slučajni vzorec velikost 36. Tedaj je

$$P(\bar{X} \geq 4) = P\left(Z \geq (4 - \mu)\sqrt{n}/\sigma\right) = P(Z \geq 1.75) \doteq 0.04.$$

```
> 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 1.756662 0.03949129
```

◇

11.4.2 (B) Vzorčna disperzija

Naj bo slučajna spremenljivka X na neki populaciji porazdeljena normalno, tj. $N(\mu, \sigma)$. Kako bi določili porazdelitev za vzorčno disperzijo ali popravljeno vzorčno disperzijo, tj.

$$S_0^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2 \quad \text{oziroma} \quad S^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2?$$

Raje izračunamo porazdelitev za naslednjo vzorčno statistiko

$$\chi^2 = \frac{nS_0^2}{\sigma^2} = \frac{(n-1)S^2}{\sigma^2} = \frac{1}{\sigma^2} \sum_{i=1}^n (X_i - \bar{X})^2.$$

Statistiko χ^2 lahko preoblikujemo takole (pri čemer opozorimo, da vzorčno povprečje \bar{X} in populacijsko povprečje μ običajno nista enaka):

$$\begin{aligned} \chi^2 &= \frac{1}{\sigma^2} \sum_{i=1}^n (X_i - \mu + \mu - \bar{X})^2 \\ &= \frac{1}{\sigma^2} \sum_{i=1}^n (X_i - \mu)^2 + \frac{2}{\sigma^2} (\mu - \bar{X}) \sum_{i=1}^n (X_i - \mu) + \frac{n}{\sigma^2} (\mu - \bar{X})^2 \end{aligned}$$

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 za katere velja $Y_i = (X_i - \mu)/\sigma$. Porazdelitvena funkcija za izbrano vzorčno statistiko χ^2 je

$$F_{\chi^2} = P(\chi^2 < z) = \iint \cdots \int_{\sum_{i=1}^n y_i^2 - \frac{1}{n}(\sum_{i=1}^n y_i)^2 < z} e^{-(y_1^2 + y_2^2 + \dots + y_n^2)/2} dy_n \dots dy_1,$$

z ustreznou ortogonalno transformacijo v nove spremenljivke z_1, z_2, \dots, z_n dobimo po nekaj računanja (glej Hladnik)

$$F_{\chi^2} = \frac{1}{(2\pi)^{(n-1)/2}} \iint \cdots \int_{\sum_{i=1}^{n-1} z_i^2 < z} e^{-(z_1^2 + z_2^2 + \dots + z_{n-1}^2)/2} dz_{n-1} \dots dz_1.$$

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

$$E(S_0^2) = E\left(\frac{\sigma^2 \chi^2}{n}\right) = \frac{(n-1)\sigma^2}{n}, \quad E(S^2) = E\left(\frac{\sigma^2 \chi^2}{n-1}\right) = \sigma^2$$

in

$$\mathsf{D}(S_0^2) = \mathsf{D}\left(\frac{\sigma^2 \chi^2}{n}\right) = \frac{2(n-1)\sigma^4}{n^2} \quad \mathsf{D}(S^2) = \mathsf{D}\left(\frac{\sigma^2 \chi^2}{n-1}\right) = \frac{2\sigma^4}{n-1}$$

Za dovolj velike n je

- statistika χ^2 porazdeljena približno normalno in sicer po zakonu $N(n-1, \sqrt{2(n-1)})$,
- vzorčna disperzija S_0^2 približno po $N\left(\frac{(n-1)\sigma^2}{n}, \frac{\sigma^2\sqrt{2(n-1)}}{n}\right)$ in
- popravljena vzorčna disperzija S^2 približno po $N\left(\sigma^2, \sigma^2 \sqrt{\frac{2}{n-1}}\right)$.

11.5 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

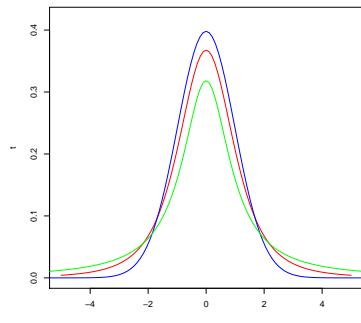
$$T = \frac{\bar{X} - \mu}{S} \sqrt{n} \text{ ?}$$



'Student' in 1908

11.5.1 Studentova t -porazdelitev

Leta 1908 je W.S. Gosset (1876-1937) pod psevdonimom 'Student' objavil članek, v katerem je pokazal, da ima statistika T porazdelitev z gostoto



$$p(x) = \frac{\left(1 + \frac{x^2}{n-1}\right)^{-n/2}}{\sqrt{n-1} B\left(\frac{n-1}{2}, \frac{1}{2}\right)}.$$

Tej porazdelitvi pravimo **Studentova t -porazdelitev** z **$n - 1$ prostostnimi stopnjami** in jo označimo na kratko s **$t(n - 1)$** .

Pri tem smo uporabili Beta funkcijo, ki jo lahko vpeljemo z Gama funkcijo na naslednji način:

$$B(x, y) = \frac{\Gamma(x)\Gamma(y)}{\Gamma(x+y)} = B(y, x).$$

Posebne vrednosti:

$$B\left(\frac{1}{2}, \frac{1}{2}\right) = \pi, \quad B(m, n) = \frac{(m-1)!(n-1)!}{(m+n-1)!} \text{ za } m, n \in \mathbb{N}.$$

```
> 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 $t(1)$ dobimo Cauchyevu porazdelitev z gostoto

$$p(x) = \frac{1}{\pi(1+x^2)}$$

Za $n \rightarrow \infty$ pa gre

$$\frac{1}{\sqrt{n-1} B\left(\frac{n-1}{2}, \frac{1}{2}\right)} \rightarrow \sqrt{2\pi} \quad \text{in} \quad \left(1 + \frac{x^2}{n-1}\right)^{-n/2} \rightarrow e^{-x^2/2}.$$

Torej ima limitna porazdelitev gostoto

$$p(x) = \frac{1}{\sqrt{2\pi}} e^{-x^2/2}$$

standardizirane normalne porazdelitve.



Če zadnji sliki dodamo

```
> curve(dnorm(x),col="magenta",add=T)
```

ta pokrije modro krivuljo.

11.5.2 Fisherjeva porazdelitev

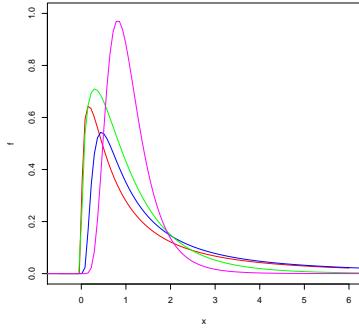
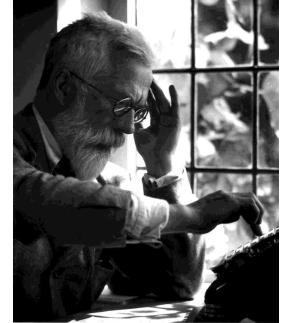
Poskusimo najti še porazdelitev kvocienta $Z = \frac{U}{V}$,

kjer sta $U \sim \chi^2(m)$ in $V \sim \chi^2(n)$ 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 **Snedenecorjeva 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 $U \sim F(m, n)$, je $1/U \sim F(n, m)$,

če je $U \sim t(n)$, je $U^2 \sim F(1, n)$.

11.6 Cenilke



11.6.1 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 rastocim n zaporedje C_n verjetnostno konvergira k parametru ζ , tj. za vsak $\varepsilon > 0$ velja

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

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 $z_k = \mathbb{E}(X^k)$, če le-ti obstajajo.

Primer: Vzorčna mediana \tilde{X} je dosledna cenilka za populacijsko mediano. \diamond

Trditev 11.2. Če pri pogoju $n \rightarrow \infty$ velja $\mathbb{E}(C_n) \rightarrow \zeta$ in $\mathbb{D}(C_n) \rightarrow 0$, je C_n dosledna cenilka parametra ζ .

Dokaz. To sprevidimo takole:

$$1 - P(|C_n - \zeta| < \varepsilon) = P(|C_n - \zeta| \geq \varepsilon) \leq P(|C_n - \mathbb{E}(C_n)| + |\mathbb{E}(C_n) - \zeta| \geq \varepsilon),$$

upoštevajmo še, da za dovolj velike n velja $|\mathbb{E}(C_n) - \zeta| < \varepsilon/2$, in uporabimo neenakost Čebiševa

$$P(|C_n - \mathbb{E}(C_n)| \geq \varepsilon/2) \leq \frac{4D(C_n)}{\varepsilon^2} \rightarrow 0. \quad \square$$

Primer: Naj bo $X \sim N(\mu, \sigma^2)$. Ker za $n \rightarrow \infty$ velja

$$\mathbb{E}(s_0^2) = \frac{(n-1)\sigma^2}{n} \rightarrow \sigma^2 \quad \text{in} \quad D(s_0^2) = \frac{2(n-1)\sigma^4}{n^2} \rightarrow 0,$$

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 $\mathbb{E}(C_n) = \zeta$ (za vsak n); in je **asimptotično nepristranska**, če je $\lim_{n \rightarrow \infty} \mathbb{E}(C_n) = \zeta$. Količino $B(C_n) = \mathbb{E}(C_n) - \zeta$ imenujemo **pristransost** (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.

Izrek 11.3. Če je razred cenilk parametra ζ konveksen (vsebuje tudi njihove konveksne kombinacije), obstaja v razredu nepristranskih cenilk istega parametra ζ v bistvu ena sama cenilka z najmanjšo disperzijo.

Dokaz. 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

$$D\left(\frac{C + C'}{2}\right) = \frac{D(C) + D(C') + 2\text{Cov}(C, C')}{4} \leq \left(\frac{\sqrt{D(C)} + \sqrt{D(C')}}{2}\right)^2 = \sigma^2.$$

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 $D(C) = D(C')$ izhaja $a = 1$, iz $\mathbb{E}(C) = \mathbb{E}(C')$ pa še $b = 0$. \square

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*

$$q(C) = \mathbb{E}(C - \zeta)^2.$$

Ker velja $q(C) = \mathbb{E}(C - \mathbb{E}(C) + \mathbb{E}(C) - \zeta)^2 = \mathbb{E}(C - \mathbb{E}(C))^2 + (\mathbb{E}(C) - \zeta)^2$, jo lahko zapišemo tudi v obliki

$$q(C) = \mathbb{D}(C) + B(C)^2.$$

Za nepristranske cenilke je $B(C) = 0$ in zato $q(C) = \mathbb{D}(C)$. Če pa je disperzija cenilke skoraj 0, je $q(C) \approx B(C)^2$.

11.6.2 Rao-Cramérjeva ocena*

Naj bo p gostotna ali verjetnostna funkcija slučajne spremenljivke X in naj bo odvisna še od parametra a , tako da je $p(x; a)$ 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*)

$$L(x_1, \dots, x_n; a) = p(x_1; a) \cdots p(x_n; a).$$

Velja

$$\int \int \cdots \int L(x_1, \dots, x_n; a) dx_1 \cdots dx_n = 1. \quad (11.1)$$

$L(X_1, \dots, X_n)$ je funkcija vzorca – torej slučajna spremenljivka. Privzemimo, da je funkcija L vsaj dvakrat zvezno odvedljiva po a na nekem intervalu I in naj na tem intervalu tudi integral odvoda L po a enakomerno konvergira. Odvajajmo enakost (11.1) po a in upoštevajmo $\frac{\partial \ln L}{\partial a} = \frac{1}{L} \frac{\partial L}{\partial a}$ pa dobimo

$$\int \int \cdots \int \frac{\partial \ln L}{\partial a} L dx_1 dx_2 \cdots dx_n = 0,$$

kar lahko tolmačimo kot

$$\mathbb{E}\left(\frac{\partial \ln L}{\partial a}\right) = 0.$$

Naj bo sedaj C nepristranska cenilka parametra a , torej $\mathbb{E}(C) = a$, oziroma zapisano z integrali $\int \int \cdots \int C L dx_1 dx_2 \cdots dx_n = a$. Ker C ni odvisna od a , dobimo z odvajanjem po parametru a :

$$\int \int \cdots \int C \frac{\partial \ln L}{\partial a} L dx_1 dx_2 \cdots dx_n = 1$$

kar z drugimi besedami pomeni

$$\mathbb{E}\left(C \frac{\partial \ln L}{\partial a}\right) = 1.$$

Če to enakost združimo s prejšnjo (pomnoženo s a), dobimo:

$$\mathbb{E}\left((C-a)\frac{\partial \ln L}{\partial a}\right) = 1.$$

Od tu po $(\mathbb{E}(XY))^2 \leq \mathbb{E}(X^2)\mathbb{E}(Y^2)$ izhaja naprej

$$1 = \left(\mathbb{E}\left((C-a)\frac{\partial \ln L}{\partial a}\right)\right)^2 \leq \mathbb{E}(C-a)^2 \mathbb{E}\left(\frac{\partial \ln L}{\partial a}\right)^2 = \mathbb{D}(C) \mathbb{E}\left(\frac{\partial \ln L}{\partial a}\right)^2,$$

kar da *Rao-Cramérjevo oceno*

$$\mathbb{D}(C) \geq \left(\mathbb{E}\left(\frac{\partial \ln L}{\partial a}\right)^2\right)^{-1} = \left(-\mathbb{E}\frac{\partial^2 \ln L}{\partial a^2}\right)^{-1} = \left(n\mathbb{E}\left(\frac{\partial \ln p}{\partial a}\right)^2\right)^{-1}.$$

11.6.3 Učinkovitost cenilk*

Rao-Cramérjeva ocena da absolutno spodnjo mejo disperzije za vse nepristranske cenilke parametra a (v dovolj gladkih porazdelitvah). Ta meja ni nujno dosežena. Cenilka, ki jo doseže, se imenuje *najučinkivitejša cenilka* parametra a in je ena sama (z verjetnostjo 1).

Kdaj pa je ta spodnja meja dosežena?

V neenakosti $(\mathbb{E}(XY))^2 \leq \mathbb{E}(X^2)\mathbb{E}(Y^2)$, 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

$$\frac{\partial \ln L}{\partial a} = A(a)(C-a),$$

kjer je $A(a)$ konstanta, odvisna od a in neodvisna od vzorca. Zato je tudi

$$(\mathbb{D}(C))^{-1} = \mathbb{E}\left(\frac{\partial \ln L}{\partial a}\right)^2 = A(a)^2 \mathbb{E}(C-a)^2 = A(a)^2 \mathbb{D}(C)$$

ozziroma končno

$$\mathbb{D}(C) = |A(a)|^{-1}.$$

Primer: Naj bo $X \sim N(\mu, \sigma)$. Tedaj je

$$L = \frac{1}{(2\pi)^{n/2}\sigma^n} e^{-\left(\frac{(X_1-\mu)^2}{\sigma^2} + \dots + \frac{(X_n-\mu)^2}{\sigma^2}\right)/2}$$

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

$$\frac{\partial \ln L}{\partial \sigma^2} = -\frac{n}{2\sigma^2} + \frac{1}{2\sigma^4}((X_1 - \mu)^2 + \cdots + (X_n - \mu)^2) = \frac{n}{2\sigma^4}(S_\mu^2 - \sigma^2).$$

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 $D(S_\mu^2) = 2\sigma^4/n$. \diamond

Primer: Za Poissonovo porazdelitev $P(\lambda)$ s parametrom λ , tj. $p_k = \lambda^k e^{-\lambda} / k!$, je

$$L = e^{-n\lambda} \frac{\lambda^{x_1 + \cdots + x_n}}{x_1! \cdots x_n!}$$

in dalje

$$\ln L = -n\lambda + (x_1 + \cdots + x_n) \ln \lambda - \ln(x_1! \cdots x_n!)$$

ter končno

$$\frac{\partial \ln L}{\partial \lambda} = -n + \frac{x_1 + \cdots + x_n}{\lambda} = \frac{n}{\lambda}(\bar{X} - \lambda). \quad (11.2)$$

Najučinkovitejša cenilka za parameter λ je \bar{X} z disperzijo $D(\bar{X}) = \lambda/n$. \diamond

Naj bo C_0 najučinkovitejša cenilka parametra a in C kaka druga nepristranska cenilka. Tedaj je *učinkovitost* cenilke C določena s predpisom

$$e(C) = \frac{D(C_0)}{D(C)}.$$

Učinkovitost najučinkovitejše cenilke je $e(C_0) = 1$. Če najučinkovitejša cenilka ne obstaja, vzamemo za vrednost $D(C_0)$ desno stran v Rao-Cramérjevi oceni.

Primer: Naj bo $X \sim N(\mu, \sigma)$. Pri velikih n -jih je vzorčna mediana \tilde{X} – ocena za μ , porazdeljena približno po $N(\mu, \sigma \sqrt{\pi/2n})$. Torej je

$$e(\tilde{X}) = \frac{D(\bar{X})}{D(\tilde{X})} = \frac{\sigma^2/n}{\pi\sigma^2/(2n)} = \frac{2}{\pi} \doteq 0.64. \quad \diamond$$

Primer: Naj bo $X \sim N(\mu, \sigma)$. Če poznamo μ , je najučinkovitejša cenilka za σ^2 statistika S_μ^2 z disperzijo $D(S_\mu^2) = 2\sigma^4/n$. Popravljena vzorčna disperzija S^2 pa je nepristranska cenilka istega parametra z disperzijo $D(S^2) = 2\sigma^4/(n-1)$. Tokrat je

$$e(S^2) = \frac{D(S_\mu^2)}{D(S^2)} = \frac{2\sigma^4/n}{2\sigma^4/(n-1)} = \frac{n-1}{n}.$$

Iz tega vidimo, da $e(S^2) \rightarrow 1$, ko $n \rightarrow \infty$. Zato je cenilka S^2 asimptotično najučinkovitejša cenilka za σ^2 . \diamond

11.6.4 Metoda momentov

Parametre populacije (ki jih ne poznamo) določimo tako, da so momenti slučajne spremenljivke X (npr. $\bar{X} = \hat{\mu}_X$, $s_X^{*2} = \hat{\sigma}_X^2$) enaki ocenam teh momentov, ki jih izračunamo iz vzorca.

Primer: $Y \sim B(n, p)$: $\mu_Y = np$, $\sigma_Y^2 = np(1-p)$. \diamond

Če iščemo en parameter (tj. a) in velja $\mu_X = f(a)$ izračunamo: $a = f^{-1}(\mu_X)$ oz. $\hat{a} = f^{-1}(\hat{\mu}_X) = f^{-1}(\bar{X})$, za dva parametra (tj. a_1, a_2) pa velja $\mu_X = f_1(a_1, a_2)$, $\sigma_X^2 = f_2(a_1, a_2)$, torej rešimo sistem enačb: $f_1(a_1, a_2) = \bar{X}$, $f_2(a_1, a_2) = s_X^{*2}$.

Primer: Za vzorec s 30-imi elementi velja $\bar{X} = 20$. Če je

- $X \sim B(80, p)$: $\mu_X = np$ oz. $p = \mu_X/n$ in $\hat{p} = \hat{X}/n = 20/80 = 0.25$.
- $X \sim \text{Exp}(\lambda)$: $\mu_X = 1/\lambda$ oz. $\lambda = 1/\mu_X$ in $\hat{\lambda} = 1/\bar{X} = 1/20 = 0.05$. \diamond

Primer: Za vzorec s 30-imi elementi velja $\bar{X} = 30$, $s_X^{*2} = 5$. Če je

- $X \sim U(a, b)$, iz $\mu_X = (a+b)/2$ in $\sigma_X = (b-a)/\sqrt{12}$ dobimo,
 $b = \mu_X + \sigma_X\sqrt{3} = 30 + 5\sqrt{3} = 38\bar{6}$ in $a = \mu_X - \sigma_X\sqrt{3} = 30 - 5\sqrt{3} = 21\bar{3}$.
- $X \sim B(n, p)$: iz $\mu_X = np$, $\sigma_X = \sqrt{np(1-p)}$ je $p = 1 - \frac{\sigma_X^2}{\mu_X} = 0.16$, $n = \frac{\mu_X}{p} = 180$.
- $X \sim N(\mu, \sigma)$: $\mu_X = \mu_X$ in $\sigma_X = \sigma_X$, sledi $\hat{\mu}_X = 30$ in $\hat{\sigma}_X = 5$. \diamond

Recimo, da je za zvezno slučajno spremenljivko X njena gostota p odvisna od m parametrov $p(x; a_1, \dots, a_m)$ in naj obstajajo momenti

$$z_k = z_k(a_1, \dots, a_m) = \int_{-\infty}^{\infty} x^k p(x; a_1, \dots, a_m) dx \quad \text{za } k = 1, \dots, m.$$

Če se dajo iz teh enačb enolično izračunati parametri a_1, \dots, a_m kot funkcije momentov z_1, \dots, z_m :

$$a_k = \varphi_k(z_1, \dots, z_m),$$

potem so

$$C_k = \varphi_k(Z_1, \dots, Z_m)$$

cenilke parametrov a_k po *metodi momentov*. k -ti vzorčni začetni moment $Z_k = \frac{1}{n} \sum_{i=1}^n X_i^k$ je cenilka za ustrezni populacijski moment z_k . *Cenilke, ki jih dobimo po metodi momentov so dosledne.*

Primer: Naj bo $X \sim N(\mu, \sigma)$. Tedaj je $z_1 = \mu$ in $z_2 = \sigma^2 + \mu^2$. Od tu dobimo $\mu = z_1$ in $\sigma^2 = z_2 - z_1^2$. Ustrezni cenilki sta

$$Z_1 = \bar{X} \quad \text{za } \mu \quad \text{in} \quad Z_2 - Z_1^2 = \bar{X}^2 - \bar{X}^2 = s_0^2 \quad \text{za } \sigma^2,$$

torej vzorčno povprečje in disperzija. \diamond

11.6.5 Metoda največjega verjetja

Naj bo $V = (X_1, \dots, X_n)$ slučajni vzorec. Odvisen je od porazdelitve in njenih parameterov a_1, \dots, a_m . Želimo določiti ocene: $\hat{a}_1, \dots, \hat{a}_m$ tako, da je verjetnost, da se je zgodil vzorec V največja. Velja

$$P(X_1 = x_1, \dots, X_n = x_n) = P(X_1 = x_1) \cdots P(X_n = x_n) = p_X(x_1) \cdots p_X(x_n).$$

Naj bo $m = 1$ in $a = a_1$. Funkcija verjetja $L(a)$, je verjetnost, da se je zgodil nek vzorec, tj.

$$L(a) = P(\text{vzorec}) = \prod_{i=1}^n p_X(x_i).$$

Oceno parametra \hat{a} določimo tako, da bo imela funkcija verjetja $L(a)$ največjo vrednost. Če je L vsaj dvakrat zvezno odvedljiva, mora veljati

$$\frac{\partial L}{\partial a} = 0 \quad \text{in} \quad \frac{\partial^2 L}{\partial a^2} < 0.$$

Največja vrednost parametra je še odvisna od x_1, \dots, x_n : $a_{\max} = \varphi(x_1, \dots, x_n)$. Tedaj je cenilka za parameter a 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: $Y \sim B(40, p)$: $(Y_1, Y_2, Y_3) = (5, 4, 7)$ (velikost vzorca je 3)

$$\begin{aligned} L(p) &= P(\text{vzorec}) = P(Y_1 = 5)P(Y_2 = 4)P(Y_3 = 7) \\ &= \prod_{i=1}^n \binom{n}{y_i} p^{y_i} q^{n-y_i} = \binom{n}{y_1} \binom{n}{y_2} \binom{n}{y_3} p^{\sum y_i} q^{3n - \sum y_i}. \end{aligned}$$

Veljati mora $L'(p) = 0$, a je lažje izraz za $L(p)$ najprej logaritmirati

$$\ln L(p) = \ln \left(\dots (y_1) (y_2) (y_3) \dots \right) + \sum y_i \ln p + + (3n - \sum y_i) \ln q$$

in nato odvajati

$$\frac{d \ln L(p)}{dp} = \frac{\sum y_i}{p} - \frac{3n - \sum y_i}{1-p} = 0$$

Torej velja $(1-p) \sum y_i - p(3n - \sum y_i) = 0$ oz. $\sum y_i = 3np$ in končno $p = (\sum y_i)/(3n) = (16/3)/40 = 0.13$. \diamond

Primer: Naj bo $X \sim B(1, p)$. 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 $E(P) = E(X)/n = p$. Ker gre $D(P) = D(X)/n^2 = p(1-p)/n \rightarrow 0$ 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: Nadalujmo primer Poissonove porazdelitve. Odvod (11.2) 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

11.7 Vzorčna statistika (nadaljevanje)

11.7.1 (C) Vzorčna povprečja

Primer: Denimo, da se spremenljivka inteligenčni kvocient¹ na populaciji porazdeljuje normalno s pričakovano vrednostjo $\mu = 100$ in standardnim odklonom $\sigma = 15$, tj.

$$X \sim N(100, 15)$$

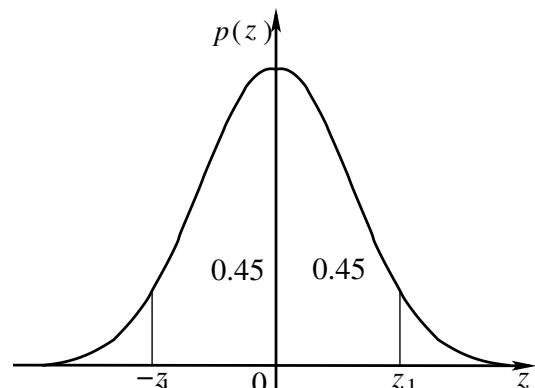
Denimo, da imamo vzorec velikosti $n = 225$. Tedaj se vzorčna povprečja porazdeljujejo normalno $\bar{X} \sim N(100, 15/\sqrt{225}) = N(100, 1)$. Izračunajmo, kolikšna vzorčna povprečja ima 90% vzorcev (simetrično na pričakovano vrednost).

90% vzorčnih povprečij se nahaja na intervalu:

$$P(\bar{X}_1 < \bar{X} < \bar{X}_2) = 0.90$$

$$P(-z_1 < z < z_1) = 0.90 \implies 2\Phi(z_1) = 0.90$$

$$\Phi(z_1) = 0.45 \implies z_1 = 1.65.$$



¹ Inteligenčnemu kvocientu (IQ) lahko med vsemi dejavniki, ki narekujejo uspeh v človekovem življenju, pripisujemo le 20% deleža. 80% uspeha se skriva v drugih dejavnikih, ki jih združujemo pod pojmom čustvena inteligenca. DANIEL GOLEMAN Emotional Intelligence Bloomsbury, London, Anglija

Potem se vzorčne povprečja nahajajo v intervalu

$$P\left(\mu - z_1 \frac{\sigma}{\sqrt{n}} < \bar{X} < \mu + z_1 \frac{\sigma}{\sqrt{n}}\right) = 0.90 \quad \text{ozioroma} \quad P\left(100 - 1.65 < \bar{X} < 100 + 1.65\right) = 0.90.$$

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

$$P\left(100 - 1.65 \frac{15}{\sqrt{2500}} < \bar{X} < 100 + 1.65 \frac{15}{\sqrt{2500}}\right) = 0.90 \quad \text{ozioroma} \quad (99.5, 100.5). \quad \diamond$$

11.7.2 (D) Vzorčni deleži

Prihajajoče obdobje blaginje

1930: pričakovana življenska doba: 60 let - 525,000 ur

Spanje 175,000 ur, Delo 100,000 ur,

Otroštvo, šolanje, zabava, konjički, šport 250,000 ur

2000: pričakovana življenska doba: 75 let - 657,000 ur

Spanje 219,000 ur, Delo 50,000 ur,

Otroštvo, šolanje, zabava, konjički, šport 388,000 ur

Denimo, da želimo na populaciji oceniti delež enot π z določeno lastnostjo. V ta namen 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



- *s pričakovano vrednostjo vzorčnih deležev enako deležu na populaciji $E(p) = \pi$, in*
- *standardnim odklonom vzorčnih deležev $SE(p) = \sqrt{\frac{\pi(1-\pi)}{n}}$.*

Za manjše vzorce se vzorčni deleži porazdeljujejo *binomsko*. Mimogrede, cenilka populacij skega 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:

$$p : N\left(0.5, \sqrt{\frac{\pi(1-\pi)}{n}}\right) = N\left(0.5, \sqrt{\frac{0.5 \times 0.5}{25}}\right) = N(0.5, 0.1).$$

Zato je

$$P(p > 0.55) = P\left(Z > \frac{0.55 - 0.5}{0.1}\right) = P(Z > 0.5) = 0.5 - \Phi(0.5) = 0.5 - 0.1915 = 0.3085.$$

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

11.7.3 (E) Razlika vzorčnih povprečij

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 povprečje \bar{X}_1 in podobno na vsakem vzorcu druge populacije \bar{X}_2 . Po reprodukcijski lastnosti normalne porazdelitve *je porazdelitev razlik vzorčnih povprečij normalna, kjer je*

- pričakovana vrednost razlik vzorčnih povprečij enako

$$\mathbb{E}(\bar{X}_1 - \bar{X}_2) = \mathbb{E}\bar{X}_1 - \mathbb{E}\bar{X}_2 = \mu_1 - \mu_2,$$

- disperzija razlik vzorčnih povprečij enaka

$$\mathbb{D}(\bar{X}_1 - \bar{X}_2) = \mathbb{D}(\bar{X}_1) + \mathbb{D}(\bar{X}_2) = \frac{\sigma^2}{n_1} + \frac{\sigma^2}{n_2} = \sigma^2 \cdot \frac{n_1 + n_2}{n_1 n_2}.$$

Primer: Populacijama študentov na neki univerzi (tehnikom in družboslovcem) so izmerili neko sposobnost s pričakovanima vrednostima $\mu_t = 70$ in $\mu_d = 80$ točk in standardnim odklonom, ki je na obeh populacijah enak, $\sigma = 7$ točk.

Kolikšna je verjetnost, da je pri naključnih vzorcih vzorčno povprečje družboslovcev ($n_d = 36$) večje za več kot 12 točk od vzorčnega povprečja tehnikov ($n_t = 64$)? Zanima nas torej verjetnost:

$$P(\bar{X}_d - \bar{X}_t > 12) = P\left(Z > \frac{12 - 10}{7\sqrt{\frac{36+64}{36 \cdot 64}}}\right) = P(Z > 1.37) = 0.5 - \Phi(1.37) = 0.0853.$$

Torej, približno 8.5% parov vzorcev je takih, da je povprečje družboslovcev večje od povprečja tehnikov za 12 točk. \diamond

11.7.4 (F) Razlika vzorčnih deležev

Podobno kot pri porazdelitvi razlik vzorčnih povprečij 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*

- pričakovano vrednostjo razlik vzorčnih deležev $E(p_1 - p_2) = E(p_1) - E(p_2) = \pi_1 - \pi_2$,
- disperzijo razlik vzorčnih deležev $D(p_1 - p_2) = D(p_1) + D(p_2) = \frac{\pi_1(1 - \pi_1)}{n_1} + \frac{\pi_2(1 - \pi_2)}{n_2}$.



Poglavlje 12

Intervali zaupanja

Denimo, da s slučajnim vzorcem ocenjujemo parameter γ . Poskušamo najti statistiko g , ki je nepristranska, tj. $E(g) = \gamma$, 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

$$Z = \frac{g - E(g)}{SE(g)} = \frac{g - \gamma}{SE(g)}$$

porazdeljuje standardizirano normalno $N(0, 1)$. Tveganje α porazdelimo simetrično polovico na levo in polovico na desno na konci normalne porazdelitve. Naj bodo ti konci določeni s številoma $\pm z_{\alpha/2}$ oziroma sta ustrezena poltraka enaka $(-\infty, -z_{\alpha/2})$ in $(z_{\alpha/2}, \infty)$. Potem je $P(Z < -z_{\alpha/2}) = \alpha/2 = P(Z > z_{\alpha/2})$ in lahko zapišemo

$$P\left(-z_{\alpha/2} < \frac{g - \gamma}{SE(g)} < z_{\alpha/2}\right) = 1 - \alpha.$$

Po ustreznji preureditvi lahko izpeljemo naslednji interval zaupanja za parameter γ

$$P\left(g - z_{\alpha/2} SE(g) < \gamma < g + z_{\alpha/2} SE(g)\right) = 1 - \alpha.$$

Vrednosti $z_{\alpha/2}$ lahko razberemo iz tabele za verjetnosti za standardizirano normalno porazdelitev, ker velja $\Phi(z_{\alpha/2}) = 0.5 - \alpha/2$. Podajmo vrednost $z_{\alpha/2}$ za nekaj najbolj standardnih tveganj:

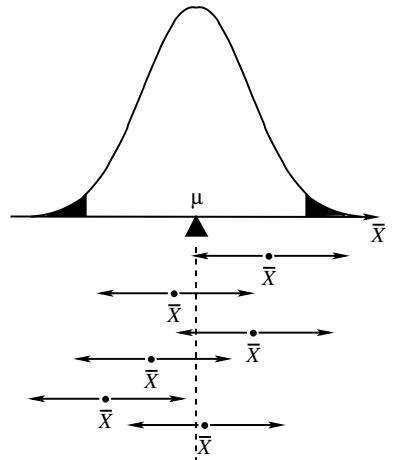
- $\alpha = 0.10$, $z_{\alpha/2} = 1.65$,
- $\alpha = 0.05$, $z_{\alpha/2} = 1.96$,
- $\alpha = 0.01$, $z_{\alpha/2} = 2.58$.

Naj bo x_p kvantil p -tega reda za porazdelitveno funkcijo F , tj. tak x_p , da velja $F(x_p) = p$.¹ Potem velja $z_{\alpha/2} = x_{1-\alpha/2}$ in $-z_{\alpha/2} = x_{\alpha/2}$ (slednje je zaradi simetričnosti normalne porazdelitve enako tudi $-x_{1-\alpha/2}$).

12.1 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 (glej sliko). V tem primeru ocenjujemo parameter pričakovano vrednost inteligenčnega kvocienta. Kot vemo, se vzorčna povprečja \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 $\alpha = 0.05$. Predstavitev več intervalov zaupanja za pričakovano vrednost μ pri 5% stopnji tveganja: približno 95% intervalov pokrije parameter μ .

12.2 Intervalsko ocenjevanje parametrov

Naj bo X slučajna spremenljivka na populaciji G z gostoto verjetnosti odvisno od parametra a . Slučajna množica $M \subset \mathbb{R}$, ki je odvisna le od slučajnega vzorca, ne pa od parametra a , se imenuje *množica zaupanja* za parameter a , če obstaja tako število α , $0 < \alpha < 1$, da velja $P(a \in M) = 1 - \alpha$. Število $1 - \alpha$ imenujemo tedaj *stopnja zaupanja*; število α pa *stopnja tveganja*. Stopnja zaupanja je običajno 95% ali 99%, tj. $\alpha = 0.05$ ali $\alpha = 0.01$. Pove nam, kakšna je verjetnost, da M vsebuje vrednost parametra a ne glede na to, kakšna je njegova dejanska vrednost. Če je množica M interval $M = [A, B]$, ji rečemo *interval zaupanja* (za parameter a). Njegovi krajišči sta funkciji slučajnega vzorca – torej statistiki.

¹Se pravi, da lahko na x_p gledamo kot na inverzno funkcijo od F , seveda bi bilo v tem primeru lepše pisati $x(p)$ namesto x_p in bi prejšnjo zvezo napisali v naslednji obliki: $F(x(p)) = p$.

Naj bo $X \sim N(\mu, \sigma)$ 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 $Z = (\bar{X} - \mu) \sqrt{n} / \sigma$. 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 $\Phi(b) = (1 - \alpha)/2$. Slednje pomeni, da je $b = z_{\alpha/2}$. Iskani interval je torej določen s točkama

$$A = \bar{X} - z_{\alpha/2} \frac{\sigma}{\sqrt{n}} \text{ in } B = \bar{X} + z_{\alpha/2} \frac{\sigma}{\sqrt{n}},$$

tj., z verjetnostjo $1 - \alpha$ je $|\bar{X} - \mu| < z_{\alpha/2} \sigma / \sqrt{n}$.

Od tu zaključimo, *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 $X \sim N(\mu, \sigma)$ tudi parameter σ ni znan, ga nadomestimo s cenilko s in moramo zato uporabiti Studentovo statistiko $T = (\bar{X} - \mu) \sqrt{n} / s$. 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 $a^2 p(a) = b^2 p(b)$ in seveda $\int_a^b p(t) dt = 1 - \alpha$.

Teoretična interpretacija koeficiente 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 = $1 - \text{stopnja zaupanja}$

12.2.1 Pričakovana vrednost μ z znanim odklonom σ

Točki

$$\bar{x} \pm z_{\alpha/2} \frac{\sigma}{\sqrt{n}}$$

prestavlja 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{x} je vrednost vzorčnega povprečja.

12.2.2 Velik vzorec za pričakovano vrednost μ

$$\bar{x} \pm z_{\alpha/2} \frac{s}{\sqrt{n}}, \quad \text{kjer je } s \text{ standardni odklon vzorca.}$$

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 $\bar{x} = 40\cdot4$ let in standardni odklon $s = 10\cdot2$ let. Pri 5% tveganju želimo z intervalom zaupanja oceniti povprečno starost podjetnikov v majhnih podjetjih v Sloveniji:

$$P\left(\bar{x} - z_{\alpha/2} \frac{s}{\sqrt{n}} < \mu < \bar{x} + z_{\alpha/2} \frac{s}{\sqrt{n}}\right) = 1 - \alpha.$$

oziroma

$$40\cdot4 - \frac{1\cdot96 \times 10\cdot2}{\sqrt{151}} < \mu < 40\cdot4 + \frac{1\cdot96 \times 10\cdot2}{\sqrt{151}}$$

in končno $40\cdot4 - 1\cdot6 < \mu < 40\cdot4 + 1\cdot6$. Torej je 95% interval zaupanja za povprečno starost podjetnikov v majhnih podjetjih v Sloveniji med 38·8 in 42·0 leti. \diamond

12.2.3 Majhen vzorec za pričakovano vrednost μ

$$\bar{x} \pm t_{\alpha/2} \frac{s}{\sqrt{n}},$$

kjer je porazdelitev spremenljivke y vzeta na osnovi $(n - 1)$ prostostnih stopenj.

Privzeli smo, da je vzorec iz približno **normalno porazdeljene** populacije.

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 pričakovano vrednost 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

$$\bar{x} = \frac{\sum_{i=1}^n x_i}{n} = \frac{49}{7} = 7,$$

$$s^2 = \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n-1} = \frac{22}{6} = 3.67.$$

Iz tabele za t -porazdelitev preberemo, da je $t_{\alpha/2}(n-1) = t_{0.05}(6) = 1.943$ in interval zaupanja je

$$7 - 1.943 \times \frac{1.9}{\sqrt{7}} < \mu < 7 + 1.943 \times \frac{1.9}{\sqrt{7}} \quad \text{ozziroma} \quad 7 - 1.4 < \mu < 7 + 1.4. \quad \diamond$$

12.2.4 Razlika pričakovanih vrednosti $\mu_1 - \mu_2$ z znanima odklonoma σ_1 in σ_2

$$\bar{x}_1 - \bar{x}_2 \pm z_{\alpha/2} \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}.$$

12.2.5 Velika vzorca za razliko pričakovanih vrednosti $\mu_1 - \mu_2$

$$\bar{x}_1 - \bar{x}_2 \pm z_{\alpha/2} \sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}.$$

12.2.6 Majhna vzorca za razliko pričakovanih vrednosti $\mu_1 - \mu_2$ z neznanima $\sigma_1 = \sigma_2$

$$\bar{x}_1 - \bar{x}_2 \pm t_{\alpha/2, n_1+n_2-2} \sqrt{s_p^2 \left(\frac{1}{n_1} + \frac{1}{n_2} \right)}, \quad \text{kjer je} \quad s_p^2 = \frac{(n_1-1)s_1^2 + (n_2-1)s_2^2}{n_1 + n_2 - 2}.$$

Privzeli smo, da sta obe populaciji porazdeljeni **približno normalni**, da sta varianci **enaki** in da so naključni vzorci izbrani **neodvisno**.

12.2.7 Majhna vzorca za razliko pričakovanih vrednosti $\mu_1 - \mu_2$

$$\bar{x}_1 - \bar{x}_2 \pm t_{\alpha/2, \nu} \sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}, \quad \text{kjer je} \quad \nu = \left\lfloor \frac{\left(s_1^2/n_1 + s_2^2/n_2 \right)^2}{(s_1^2/n_1)^2/(n_1-1) + (s_2^2/n_2)^2/(n_2-1)} \right\rfloor.$$

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. 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 ni=):

T = -0.96 P = 0.37 DF = 7



◇

12.2.8 Velik vzorec za razliko $\mu_d = \mu_1 - \mu_2$ ujemajočih se parov

$$\bar{d} \pm z_{\alpha/2} \frac{s_d}{\sqrt{n}}, \quad \text{kjer je } n \text{ število parov.}$$

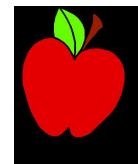
12.2.9 Majhen vzorec za razliko $\mu_d = \mu_1 - \mu_2$ ujemajočih se parov

$$\bar{d} \pm t_{\alpha/2, n-1} \frac{s_d}{\sqrt{n}}, \quad \text{kjer je } n \text{ število parov.}$$

Privzeli smo, da je populacija razlik parov 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 dan	Residue noč	razlika
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, 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 $(-58.9, -18.9)$. \diamond

Za π – delež populacije, p – delež vzorca, kjer je $p = x/n$ in je x število uspehov v n poskusih.

12.2.10 Delež populacije π z znanim odklonom σ

$$p \pm z_{\alpha/2} \sigma.$$

12.2.11 Velik vzorec za delež populacije π

$$p \pm z_{\alpha/2} \sqrt{\frac{pq}{n}}.$$

Privzeli smo, da je velikost vzorca n dovolj velika, da je aproksimacija veljavna.

Izkušnje kažejo, da je za izpolnitev pogoja ‐dovolj velik vzorec‐ priporočljivo privzeti (angl. rule of thumb): $np \geq 4$ in $nq \geq 4$.

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.** Po zgornji formuli dobimo

$$0.50 - 1.96 \sqrt{\frac{0.50 \times 0.50}{151}} < \pi < 0.50 + 1.96 \sqrt{\frac{0.50 \times 0.50}{151}}$$

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. \diamond

12.2.12 Razlika deležev $\pi_1 - \pi_2$ z znanim odklonom $\sigma_{\pi_1 - \pi_2}$

$$(p_1 - p_2) \pm z_{\alpha/2} \sigma_{\pi_1 - \pi_2}.$$

12.2.13 Velik vzorec za razliko deležev $\pi_1 - \pi_2$

$$(p_1 - p_2) \pm z_{\alpha/2} \sqrt{\frac{p_1 q_1}{n_1} + \frac{p_2 q_2}{n_2}}.$$

Privzeli smo, da je velikost vzorca n dovolj velika, da je aproksimacija veljavna.

Izkušnje kažejo, da je za izpolnitev pogoja ‐dovolj velik vzorec‐ priporočljivo privzeti: $n_1 p_1 \geq 4$, $n_1 q_1 \geq 4$, $n_2 p_2 \geq 4$ in $n_2 q_2 \geq 4$.

12.2.14 Majhen vzorec za varianco σ^2

$$\frac{(n-1)s^2}{\chi^2_{1-\alpha/2}(n-1)} \leq \sigma^2 \leq \frac{(n-1)s^2}{\chi^2_{\alpha/2}(n-1)}.$$

Privzeli smo, da je vzorec iz približno **normalno porazdeljene** populacije.

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

$$\chi^2_{1-\alpha/2}(n-1) = \chi^2_{0.95}(6) = 12.6, \quad \chi^2_{\alpha/2}(n-1) = \chi^2_{0.05}(6) = 1.64.$$

90% interval zaupanja za varianco je tedaj

$$\frac{6 \times 3.67}{12.6} < \sigma^2 < \frac{6 \times 3.67}{1.64} \quad \text{oziroma} \quad 1.75 < \sigma^2 < 13.43. \quad \diamond$$

12.2.15 Majhen vzorec za kvocient varianc σ_1^2/σ_2^2

$$\frac{s_1^2}{s_2^2} \cdot \frac{1}{F_{1-\alpha/2}(n_1-1, n_2-1)} \leq \frac{\sigma_1^2}{\sigma_2^2} \leq \frac{s_1^2}{s_2^2} \cdot \frac{1}{F_{\alpha/2}(n_2-1, n_1-1)}.$$

Privzeli smo, da sta obe populaciji iz katerih izbiramo vzorce, imata **približno normalni porazdelitvi** relativnih frekvenc in da so naključni vzorci izbrani **neodvisno** iz obeh populacij.

12.3 Izbira velikosti vzorca

V tem razdelku bomo izbirali velikosti vzorcev za oceno različnih parametrov, ki je pravilna znotraj ε enot z verjetnostjo $(1 - \alpha)$.

- Populacijsko povprečje μ : $n = \left(\frac{z_{\alpha/2} \sigma}{\varepsilon}\right)^2$. Populacijski odklon je običajno aproksimiran.

- Razlika med parom populacijskih povprečij, tj.

$$\mu_1 - \mu_2: n_1 = n_2 = \left(\frac{z_{\alpha/2}}{\varepsilon}\right)^2 (\sigma_1^2 + \sigma_2^2).$$

- Populacijski delež π : $n = \left(\frac{z_{\alpha/2}}{\varepsilon}\right)^2 pq$.

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.

$$\pi_1 - \pi_2: n_1 = n_2 = \left(\frac{z_{\alpha/2}}{\varepsilon}\right)^2 (p_1 q_1 + p_2 q_2).$$



Poglavlje 13

Preverjanje statističnih domnev

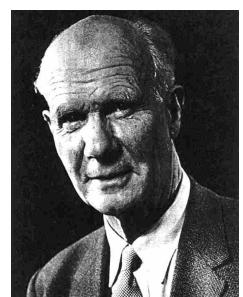


Cilj

- postaviti domnevo (trditev) o populaciji,
- izberati vzorec, s katerim bomo preverili domnevo,
- zavrniti ali sprejeti domnevo.

Domneva je testirana z določanjem verjetja, da dobimo določen rezultat, kadar jemljemo vzorce iz populacije s predpostavljenimi vrednostimi parametrov.

Teorijo preverjanja 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 $p(x; a)$ in postavljamo/raziskujemo domnevo o parametru a , govorimo o *parametrični domnevi*. Če pa je vprašljiva tudi sama

vrsta porazdelitve, je domneva *neparametrična*. Domneva je *enostavna*, če natačno določa porazdelitev (njeno vrsto in točno vrednost parametra); sicer je *sestavljen*a.

Primer: Naj bo $X \sim N(\mu, \sigma)$. Če poznamo σ , je domneva $H : \mu = 0$ enostavna; Če pa parametra σ ne poznamo, je sestavljen. Primer sestavljenih domnev je tudi $H : \mu > 0$. ◇

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)**

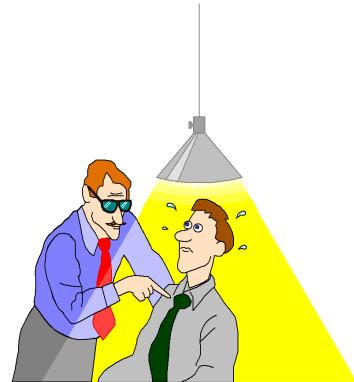
- je trditev o lastnosti populacije za katero predpostavimo, da drži (oziroma za katero verjamemo, da je resnična),
- je trditev, ki jo test želi ovreči.

- **Alternativna (nasprotna) domneva (H_a)**

- je trditev, ki ni združljiva z ničelno domnevo,
- je trditev, ki jo s testiranjem skušamo dokazati.

Okvirni postopek preverjanja domneve

- postavimo ničelno in alternativno domnevo,
- izberemo testno statistiko,
- določimo zavrnitveni kriterij,
- izberemo naključni vzorec,
- izračunamo vrednost na osnovi testne statistike,
- sprejmemo odločitev,
- naredimo ustrezni zaključek.



13.1 Ilustrativni primeri

Primer: Oglejmo si ameriški sodni sistem.

- H_0 : obtoženec je nedolžen (ničelna domneva),
- H_a : obtoženec je kriv (alternativna domneva).

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 preverjanja domneve

		<i>odločitev</i>	
		nedolžen	kriv
	nedolžen	pravilna <i>odločitev</i>	napaka 1. vrste (α)
	kriv	napaka 2. vrste (β)	moč testa ($1 - \beta$)

- 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, tj. 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. \diamond

Primer: Postavimo domnevo o vrednosti parametra, npr. π – delež enot z določeno lastnostjo na populaciji, tj. $H_0 : \pi = \pi_0$, kjer je $\pi_0 = 0.36$. 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

$$N\left(\pi_0, \sqrt{\frac{\pi_0(1 - \pi_0)}{n}}\right).$$

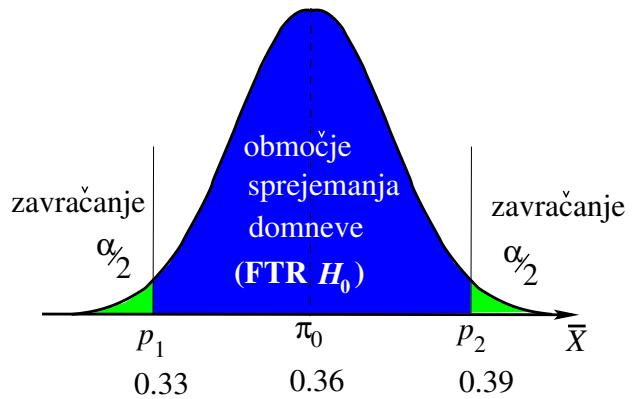
Vzemimo en slučajni vzorec z vzorčnim deležem p . Ta se lahko bolj ali manj razlikuje od π_0 . Če se zelo razlikuje, lahko podvomimo o resničnosti domneve $\pi = \pi_0$. Zato okoli π_0 naredimo

območje sprejemanja domneve in izven tega 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:

$$p_{1,2} = \pi_0 \pm z_{\alpha/2} \sqrt{\frac{\pi_0(1 - \pi_0)}{n}},$$

$$p_{1,2} = 0,36 \pm 1,96 \sqrt{\frac{0,36 \times 0,64}{900}}$$

$$= 0,36 \pm 0,03.$$



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 $\pi = \pi_0$ zavrnemo. Pri tem pa vemo, da ob resnični domnevi $\pi = \pi_0$ obstajajo vzorci, ki imajo vrednosti v območju zavračanja. Število α je verjetnost, da vzorčna vrednost pade v območje zavračanja ob prostopavki, da je domneva resnična. Zato je α verjetnost, da zavrnemo pravilno domnevo – **napaka 1. vrste**. Ta napaka je merljiva in jo lahko poljubno manjšamo.

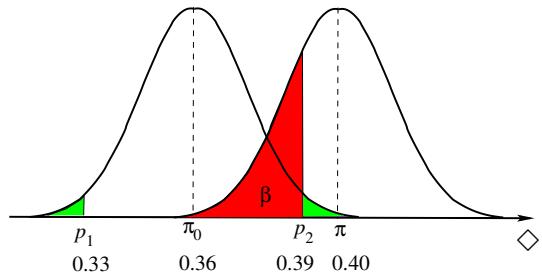
Napaka 2. vrste (β): Vzorčna vrednost lahko pade v območje sprejemanja, čeprav je domnevna vrednost parametra napačna. V primeru, ki ga obravnavamo, naj bo prava vrednost deleža na populaciji $\pi = 0,40$. Tedaj je porazdelitev vzorčnih deležev

$$N\left(\pi, \sqrt{\frac{\pi(1 - \pi)}{n}}\right) = N(0,40, 0,0163).$$

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.



13.2 Alternativna domneva in definicije

Za začetek si oglejmo parametrične statistične domneve:

- ničelna domneva $H_0 : q = q_0$
- alternativna domneva
 - $H_a : q \neq q_0$
 - $H_a : q > q_0$
 - $H_a : q < q_0$



Nekaj konkretnih primerov ničelnih domnev:

$$H_0 : \mu = 9\text{mm} \text{ (premer 9 milimetrskega kroga),}$$

$$H_0 : \mu = 600\text{km} \text{ (doseg novih vozil),}$$

$$H_0 : \mu = 3 \text{ dnevi}$$

Neusmerjena H_a

Merjenje 9 mm kroga:

$$H_0 : \mu = 9\text{mm},$$

$$H_a : \mu \neq 9\text{mm}.$$

“Manj kot” H_a

Proizvalajec trdi, da je to doseg novih vozil:

$$H_0 : \mu = 600 \text{ km,}$$

$$H_a : \mu < 600 \text{ km.}$$

“Več kot” H_a

Čas odsotnosti določenega artikla pri neposredni podpori:

$$H_0 : \mu = 3 \text{ dnevi,}$$

$$H_a : \mu > 3 \text{ dnevi.}$$

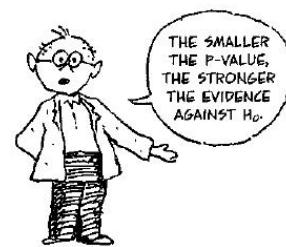
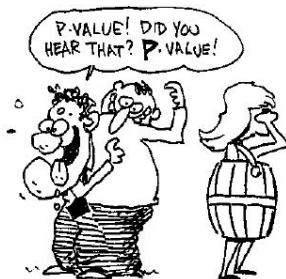
Definicije

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. **P-vrednost** oziroma **(bistvena) stopnja značilnosti testa** (tudi **signifikantnosti**) je največja vrednost parametra α , ki jo je vodja eksperimenta pripravljen sprejeti (zgornja meja za napako 1. vrste) glede na dan vzorec.
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.

		odločitev					
		FTR H_0	zavrni H_0	velikost vzorca	napaka 1. vrste	napaka 2. vrste	moč
dejansko stanje	H_0 je pravilna		(α)	n	α	β	$1 - \beta$
	H_0 je napacna	(β)	($1 - \beta$)	konst.	↑	↓	↑
				konst.	↓	↑	↓
				povečanje	↓	↓	↑
				zamnjišanje	↑	↑	↓

Še enkrat z drugimi besedami: P -vrednost (ali ugotovljena bistvena stopnja značilnosti 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 zavrni H_0 .



Primer: 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 preveriti domnevo, da je Pascal bolj učinkovit jezik kot Agol, pri katerem je 20% spremenljivk tabelaričnih.

Postavitev statistične domneve

- $H_0 : \pi = 0.20$ (ničelna domneva),
- $H_a : \pi > 0.20$ (alternativna domneva),

Naj bo p verjetnost, da izberemo tabelarično spremenljivko na vsakem posameznem poskusu.

(a) Določi α za območje zavrnitve $y > 8$. Izračunati želimo verjetnost, da se bo zgodila napaka 1. vrste, torej da bomo zavrnili 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 zavrnili, 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 (a):

$$\begin{aligned}\alpha &= P(Y \geq 5) = 1 - P(Y \leq 4) = 1 - \sum_{i=0}^4 P(Y = i) \\ &= 1 - \sum_{i=0}^4 \binom{20}{i} 0.2^i 0.2^{20-i} = 1 - 0.6296 = 0.3704 = 37.04\%.\end{aligned}$$

(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.

$$\beta = P(y \leq 7) = \sum_{i=0}^7 \binom{20}{i} 0.5^i 0.5^{20-i} = 0.1316 = 13.16\%.$$

(d) Določi β za območje zavrnitve $y \geq 5$, če je $p = 0.5$. Do rezultata pridemo na enak način kot v (c):

$$\beta = P(y \leq 4) = \sum_{i=0}^4 \binom{20}{i} 0.5^i 0.5^{20-i} = 0.0059 = 0.59\%.$$

(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 izkušenje 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

$$\beta = P(y \leq 8) = \sum_{i=0}^8 \binom{20}{i} 0.4^i 0.6^{20-i} = 0.5956 = 59.56\%.$$

Moč testa znaša 0.4044.

(h) Za območje zavrnitev določeno v točki (f) določi moč testa, če je v resnici $p = 0.7$.

Tokrat velja $Y \sim B(20, 0.7)$ in

$$\beta = P(Y \leq 8) = \sum_{i=0}^8 \binom{20}{i} 0.7^i 0.3^{20-i} = 0.0051 = 0.51\%.$$

Moč testa znaša 0.995. \diamond

13.3 Predznačni test

- $H_0 : \tau = \tau_0$ (mediana populacije)

Predpostavke: naključno vzorčenje iz zvezne porazdelitve.

Primer: Ali je dejanska mediana (τ) pH iz določene regije 6.0? Da bi odgovorili na to vprašanje 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

- $H_0 : \tau = 6.0$ (ničelna domneva),
- $H_a : \tau < 6.0$ (alternativna domneva).

Izbira testne statistike (T.S.)

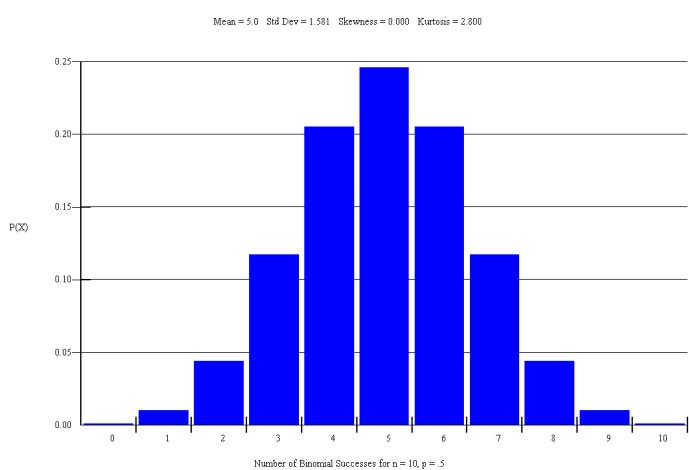
- $S_+ =$ število vzorcev, ki so večji od mediane τ_0 iz domneve.
- $S_- =$ število vzorcev, ki so manjši od mediane τ_0 iz domneve.

Porazdelitev testne statistike

- vsak poskus je bodisi uspeh ali neuspeh,
- fiksni vzorec, velikosti n ,
- naključni vzorci
 - neodvisni poskusi,
 - konstantna verjetnost uspeha.

Gre za binomsko porazdelitev:

$S_+ \sim B(10, 0.5)$ (ker gre za mediano!),
torej je $E(S_+) = np = 5$.



Določimo zavrnitveni kriterij

- 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.

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

Ali empirični podatki podpirajo trditev, da je dejanska mediana manjša ali enaka 6.0?

pH	predznak
5.93	—
6.08	+
5.86	—
5.91	—
6.12	+
5.90	—
5.95	—
5.89	—
5.98	—
5.96	—

Izračunajmo P -vrednost iz testne statistike:

$$S_+ = 2, P\text{-vrednost} = P(S_+ \geq 2 \mid \tau = 6.0) = 0.05469.$$

$$S_- = 8, P\text{-vrednost} = P(S_- \geq 8 \mid \tau = 6.0) = 0.05469.$$

Odločitev in zaključek: P -vrednost $= 0.05469 > 0.01074 = \alpha$, zato zaključimo, da ni osnove za zavrnitev ničelne domneve (angl. fail to reject, kratica FTR), tj. da nimamo dovolj osnov, da bi dokazali, da velja alternativna trditev. Torej privzemimo, da je pH enaka 6.0 v tej konkretni regiji, saj imamo premalo podatkov, da bi pokazali, da je dejanska mediana pH manjša od 6.0. \diamond

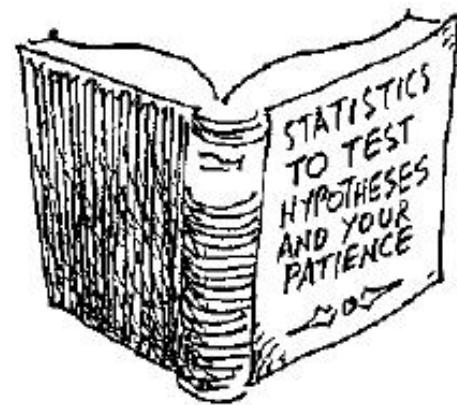
konkretni regiji, saj imamo premalo podatkov, da bi pokazali, da je dejanska mediana pH manjša od 6.0. \diamond

13.4 Wilcoxonov predznačni-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 (medianu).

Testna statistika

- $S_+ =$ vsota rangov, ki ustreza pozitivnim številom.
- $S_- =$ vsota rangov, ki ustreza 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,
- vsota vseh rangov je $= n(n+1)/2$, pričakovana vrednost za S^+ oz. S^- pa je $= n(n+1)/4$.

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	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, tj. najmanjši α pri katerem zavrnemo domnevo H_0 :

- Če je P -vrednost $> \alpha$, potem FTR H_0 .
- Če je P -vrednost $< \alpha$, potem zavrni H_0 .
- Če je P -vrednost $= 2 P(Z > 1,278) = 2 \times 0,1003 = 0,2006$.

Odločitev in zaključek:

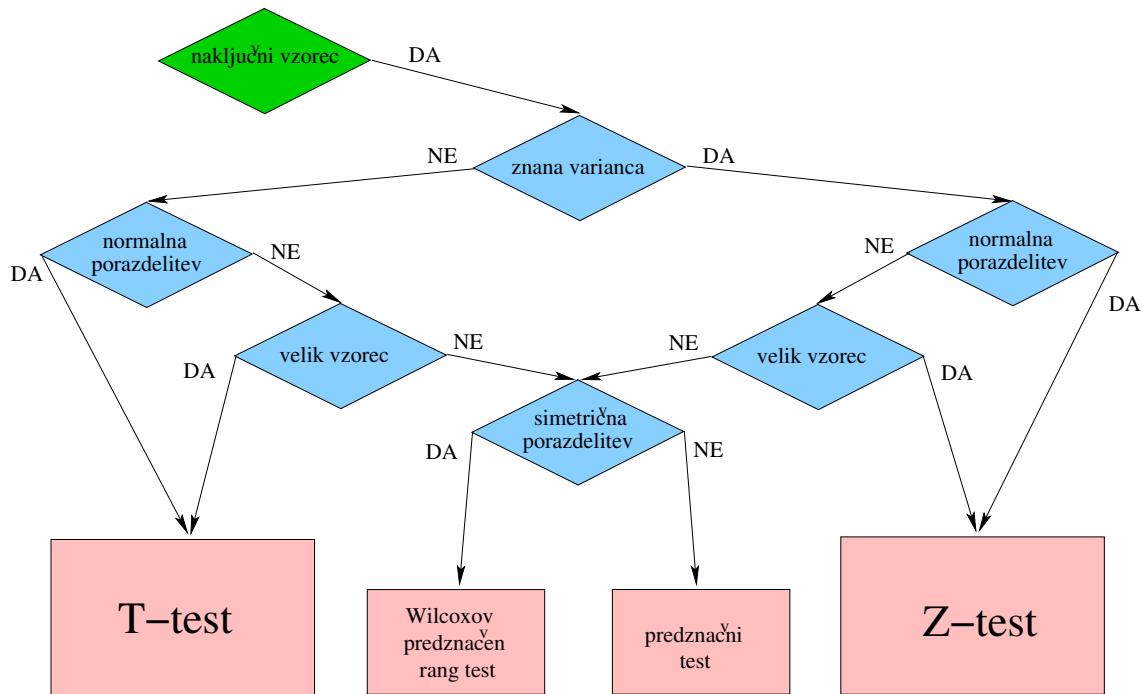
$$P\text{-vrednost} = P(S_+ = 3) = P(S_- = 12) = 0,15625 > 0,1 = \alpha \implies \text{FTR } H_0.$$

Privzemimo $\tau = 500$, saj ni osnov, da bi pokazali $\tau > 500$

◇

13.5 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 njen 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 razilike med parametrom in vzorčno oceno.



13.6 Domneve za pričakovano vrednost $H_0 : \mu = \mu_0$

Delovne *predpostavke* v tem rezdelku so:

- naključno vzorčenje
- izbiramo vzorce iz normalne porazdelitve in/ali imamo vzorec pri katerem je n velik.

13.6.1 Znan odklon σ

Če poznamo odklon populacije σ , potem na osnovi predpostavk na začetku tega razdelka in diagrama iz razdelka 13.5

$$\text{T.S.} = \frac{\bar{X} - \mu_0}{\frac{\sigma}{\sqrt{n}}} \quad \text{sledi } z\text{-porazdelitev.}$$

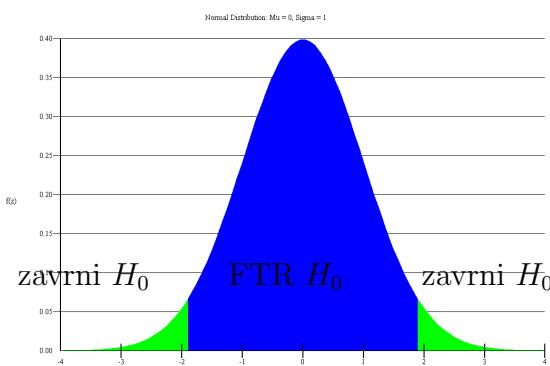
Primer: Proizvajalec omake za špagete da v vsako posodo 28 unče omake za špagete. Količina omake, ki je v vsaki posodi, je porazdeljena normalno s standardnim odklonom 0.05 unče. Podjetje ustavi proizvodni trak in popravi napravo za polnenje, če so posode bodoši premalo napolnjene (to razjezi kupce), ali preveč napolnjene (kar seveda pomeni manjši profit). **Ali naj na osnovi vzorca iz 15ih posod ustavijo proizvodno linijo?** Uporabi stopnjo značilnosti 0.05. Postavimo domnevi

- $H_0 : \mu = 28$ (ničelna domneva),
- $H_a : \mu \neq 28$ (alternativna domneva).

Dodatna predpostavka: poznamo varianco populacije.

Izberemo testno statistiko: Z-Test: $H_0 : \mu = \mu_0$ (povprečje populacije)

Določimo zavnitveni kriterij (glej sliko).

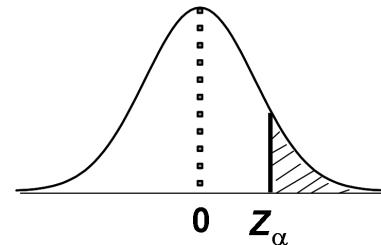


Rezultati testiranja

- Za naključni vzorec 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 = $2 P(Z > 1.278) = 2 \times 0.1003 = 0.2006$.
 - Najmanjši α pri katerem zavnemo domnevo H_0
 - P -vrednost $> \alpha$, zato FTR H_0 .
- Za $H_a : \mu > \mu_0$ je **odločitveno pravilo**: zavnri H_0 , če je **T.S. $\geq z_\alpha$** .



Za $H_a : \mu < \mu_0$ **odločitveno pravilo**: zavnri H_0 , če je **T.S. $\leq -z_\alpha$** .

Za $H_a : \mu \neq \mu_0$ **odločitveno pravilo**: zavnri H_0 če je **T.S. $\leq -z_{\alpha/2}$** ali **če je T.S. $\geq z_{\alpha/2}$** . \diamond

13.6.2 Neznan odklon σ in velik vzorec

Če ne poznamo odklona σ in je $n \geq 30$, potem

$$\text{T.S.} = \frac{\bar{x} - \mu_0}{\sqrt{s}} \quad \text{sledi } t\text{-porazdelitev z } n-1 \text{ prostostnimi stopnjami.}$$

(Velja omeniti še, da se pri tako velikem n z - in t -porazdelitev tako ne razlikujeta kaj dosti.)

13.6.3 Neznan odklon σ , normalna populacija in majhen vzorec

Če ne poznamo odklona σ , populacija je normalna in je $n < 30$, potem

$$\text{T.S.} = \frac{\bar{x} - \mu_0}{\sqrt{s}} \quad \text{sledi } t\text{-porazdelitev z } n-1 \text{ prostostnimi stopnjami.}$$

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

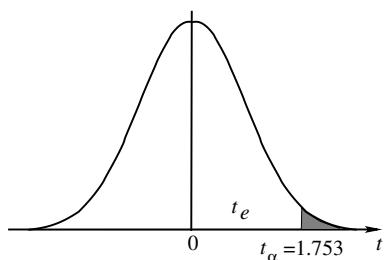
populaciji porazdeljuje normalno. **Ali lahko sprejmemo domnevo, da imajo odrasli Slovenci v povprečju več kot osemletko pri 5% stopnji značilnosti?** Postavimo najprej ničelno in osnovno domnevo:

$$H_0 : \mu = 8 \quad \text{in} \quad H_1 : \mu > 8.$$

Ustrezna testna statistika

$$t = \frac{\bar{x} - \mu_H}{s} \sqrt{n}$$

je porazdeljena po t -porazdelitvi s 15imi 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:

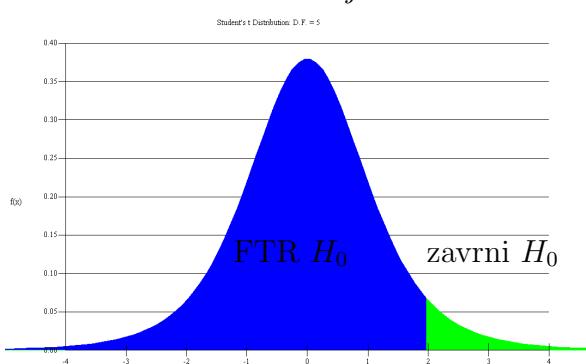


$$t_e = \frac{9 - 8}{3} \sqrt{16} = 1.3.$$

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 bežigrajske 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 bežigrajske gimnazije statistično večje od 240**. *Dodatna predpostavka:* ne poznamo varianco populacije. Postavimo domnevi: $H_0 : \mu = 240$ in $H_a : \mu > 240$ in zaradi dodatne predpostavke izberemo T -test za testno statistiko.

Določimo zavrnitveni kriterij



Rezultati testiranja

- 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.

$$P\text{-vrednost} = P(T > 0.9413) = 0.1949 > \alpha, \text{ zato FTR } H_0.$$



\diamond

Primer: Specifikacije za mestni vodovod zahtevajo, da cevi zdržijo pritisk 2500 enot. Proizvajalec cevi testira svoje cevi in dobi naslednje podatke: 2610, 2750, 2420, 2510, 2540, 2490,

2680. [Ali je to dovolj, da nadzornik zaključi izpolnjevanje pogojev?](#) Uporabimo stopnjo zaupanja $\alpha = 0.10$. Vstavimo podatke v Minitab:

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	T	p-VALUE
C1 7	2571.4	115.1	43.5	1.64	0.076

Razlaga P-vrednosti

1. Izberi največjo vrednost za α , ki smo jo pripravljeni tolerirati.

2. Če je P -vrednost manjša kot maksimalna vrednost parametra α , potem zavrnji H_0 . \diamond

13.7 Domneve za razliko povprečij $H_0 : \mu_1 - \mu_2 = D_0$

13.7.1 Znana odklona σ_1 in σ_2

Vzorce jemljemo neodvisno, zato



$$\text{T.S.} = \left((\bar{x}_1 - \bar{x}_2) - D_0 \right) / \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}} \quad \text{sledi } z\text{-porazdelitev.}$$

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:

$$\bar{x}_F = 7.0, \quad s_F = 1 \quad \text{in} \quad \bar{x}_D = 7.2, \quad s_D = 1$$

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 razliko pričakovanih vrednosti lastnosti dveh populacij na vzorcih računamo razliko vzorčnih povprečij, ki se za dovolj velike vzorce porazdeljuje normalno

$$\bar{x}_D - \bar{x}_F : N \left(\mu_D - \mu_F, \sqrt{\frac{s_D^2}{n_D} + \frac{s_F^2}{n_F}} \right).$$

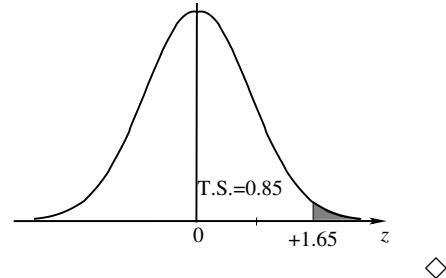
oziroma statistika

$$z = \left(\bar{x}_D - \bar{x}_F - (\mu_D - \mu_F)_H \right) / \sqrt{\frac{s_D^2}{n_D} + \frac{s_F^2}{n_F}}$$

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. Vrednost testne statistike je

$$\text{T.S.} = (7.2 - 7 - 0) / \sqrt{\frac{1}{36} + \frac{1}{36}} = 0.852.$$

Torej ne pade v kritično območje. Ničelne domneve ne moremo zavrniti. Povprečna uspešnost deklet in fantov ni statistično značilno različna. \diamond



13.7.2 Neznana odklona σ_1 in/ali σ_2 , $n_1 \geq 30$ in/ali $n_2 \geq 30$

Ker vzorce jemljemo neodvisno, velja

$$\text{T.S.} = ((\bar{x}_1 - \bar{x}_2) - D_0) / \sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}} \quad \text{sledi } z\text{-porazdelitev.}$$

13.7.3 Neznana σ_1 in/ali σ_2 , norm. pop., $\sigma_1 = \sigma_2$, $n_1 < 30$ ali $n_2 < 30$

Ker vzorce jemljemo neodvisno, velja:

$$\text{T.S.} = \frac{(\bar{x}_1 - \bar{x}_2) - D_0}{\sqrt{\frac{(n_1 - 1)s_1^2 + (n_2 - 1)s_2^2}{n_1 + n_2 - 2} \left(\frac{1}{n_1} + \frac{1}{n_2} \right)}}$$

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.

13.7.4 Neznana σ_1 in/ali σ_2 , norm. pop., $\sigma_1 \neq \sigma_2$, $n_1 < 30$ ali $n_2 < 30$

Ker jemljemo vzorce neodvisno, velja

$$\text{T.S.} = \frac{(\bar{x}_1 - \bar{x}_2) - D_0}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}} \quad \text{sledi } t\text{-porazdelitev z } \nu \text{ prostostnimi stopnjami,}$$

kjer je

(Če ν ni naravno število, zaokroži ν navzdol do najbližjega naravnega števila za uporabo t -tabele.)

$$\nu = \frac{\left(\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2} \right)^2}{\frac{\left(\frac{s_1^2}{n_1} \right)^2}{n_1 - 1} + \frac{\left(\frac{s_2^2}{n_2} \right)^2}{n_2 - 1}}$$

13.8 Domneve za povprečje razlik $H_0 : \mu_d = D_0$

13.8.1 Velik vzorec

Ker vzorce jemljemo neodvisno velja

$$\text{T.S.} = \frac{(\bar{d} - D_0)}{s_d} \sqrt{n} \quad \text{sledi } z\text{-porazdelitev.}$$

13.8.2 Normalna populacija razlik in majhen vzorec

Če je populacija razlik normalno porazdeljena in če je $n \leq 30$, potem velja

$$\text{T.S.} = \frac{\bar{d} - D_0}{s_d} \sqrt{n} \quad \text{sledi } t\text{-porazdelitev z } n - 1 \text{ prostostnimi stopnjami.}$$

naloga	človek. 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

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. Parjenj 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 $= 0.034$.

13.9 Domneve za delež $H_0 : \pi = \pi_0$

13.9.1 Velik vzorec

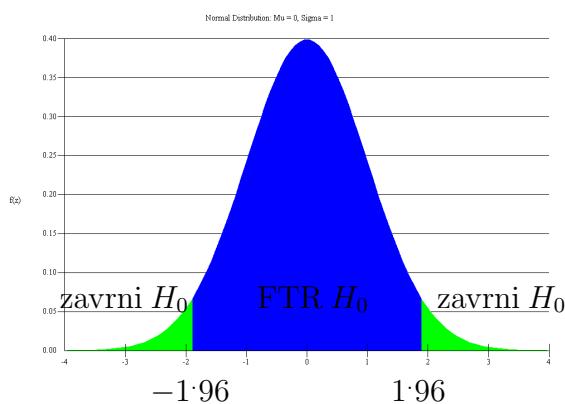
$$\text{T.S.} = \frac{p - \pi_0}{\sqrt{\frac{\pi_0(1 - \pi_0)}{n}}} \quad \text{sledi } z\text{-porazdelitev.}$$

Kot splošno pravilo bomo zahtevali, da velja $np \geq 4$ in $nq \geq 4$.

Primer: Državni zapisи 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 $\alpha = 0.05$.

- Ničelna domneva $H_0 : \pi = 0.7$,
- alternativna domneva $H_1 : \pi \neq 0.7$.

Določimo zavrnitveni kriterij



Rezultati testiranja

- Iz vzorca dobimo $p = 156/200 = 0.78$.
 - Izračunaj vrednost testne statistike:
- $$Z = (0.78 - 0.7) / 0.0324 = 2.4688.$$
- Naredi odločitev: zavrnji domnevo H_0
 - Zaključek: pokrajna ima drugačen kriterij.

P -vrednost (tj. najmanjši α pri katerem zavrnemo domnevo H_0). Gre za sprejemljivost domneve H_0 na osnovi vzorca (možnost za opazovanje vzorca ali bolj ekstremno podatkov, če je domneva H_0 pravilna): P -vrednost $= 2P(Z > 2.469) = 2 \times 0.0068 = 0.0136 < \alpha$, zato zavrnji domnevo H_0 . \diamond

13.10 Razlika deležev dveh populacij $H_0 : \pi_1 - \pi_2 = D_0$

Če želimo, da je vzorec za testiranje domneve o $\pi_1 - \pi_2$ velik,



bomo kot splošno pravilo zahtevali

$$n_1 p_1 \geq 4, \quad n_1 q_1 \geq 4,$$

$$n_2 p_2 \geq 4 \quad \text{in} \quad n_2 q_2 \geq 4.$$

13.10.1 Velik vzorec in $D_0 = 0$

$$\text{T.S.} = \frac{p_1 - p_2}{\sqrt{pq\left(\frac{1}{n_1} + \frac{1}{n_2}\right)}} \quad \text{sledi } z\text{-porazdelitev, kjer je } p = \frac{n_1 p_1 + n_2 p_2}{n_1 + n_2}.$$

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 značinosti.

$$H_0 : \pi_1 = \pi_2 \quad \text{ozziroma} \quad \pi_1 - \pi_2 = 0,$$

$$H_1 : \pi_1 \neq \pi_2 \quad \text{ozziroma} \quad \pi_1 - \pi_2 \neq 0.$$

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 pričakovana vrednost 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 uteženim povprečjem vzorčnih deležev p_1 in p_2

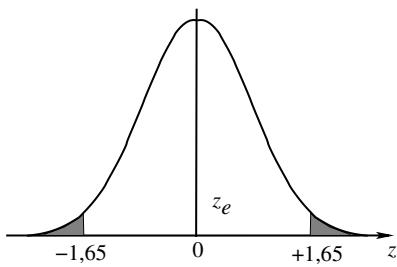
$$p = \frac{n_1 p_1 + n_2 p_2}{n_1 + n_2} = \frac{k_1 + k_2}{n_1 + n_2}.$$

Vrnimo se na primer. Vzorčna deleža sta: $p_1 = 90/300 = 0.30$ in $p_2 = 50/200 = 0.25$. Ocena populacijskega deleža je $p = (50 + 90)/(200 + 300) = 0.28$. Kot smo že omenili, se testna statistika

$$z = \left(p_1 - p_2 - (\pi_1 - \pi_2)_H\right) / \sqrt{p(1 - p)\left(\frac{1}{n_1} + \frac{1}{n_2}\right)}.$$

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

$$z_e = (0.30 - 0.25 - 0) / \sqrt{0.28(1 - 0.28)\left(\frac{1}{300} + \frac{1}{200}\right)} = 1.22.$$



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. ◇

13.10.2 Velik vzorec in $D_0 \neq 0$.

$$\text{T.S.} = \left((p_1 - p_2) - D_0 \right) / \sqrt{\frac{p_1 q_1}{n_1} + \frac{p_2 q_2}{n_2}} \quad \text{sledi } z\text{-porazdelitev.} \quad \text{🚬}$$

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 tveganja, da bo prodaja znamke A boljša od prodaje znamke B za 10% proti alternativni domnevi, da bo razlika manj kot 10% (slednje pomeni le, da je v tem primeru alternativna domneva negacija ničelne domneve). ◇

13.11 Analiza variance

Če opravljamo isti poskus v nespremenjenih pogojih, kljub temu v rezultatu poskusa opažamo 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 spreminja, seveda dobimo dodatna odstopanja od povprečja. Analiza tega, ali so odstopanja zaradi sprememb različnih faktorjev ali pa zgolj slučajna, in kateri faktorji vplivajo na varianco, se imenuje *analiza variance*.

Zgleda:

- 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 $H_0 : \mu_1 = \mu_2$
- 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 spreminja, ločimo t.i. *enojno klasifikacijo* ali *enofaktorski eksperiment*, *dvojno klasifikacijo* ali *dovfaktorski eksperiment*, itd. V izrazu

$$Q_v^2 = \sum_{i=1}^r \sum_{k=1}^{n_i} (X_{ik} - \bar{X}_i)^2 = \sum_{i=1}^r (n_i - 1) s_i^2 \quad \text{je} \quad s_i^2 = \frac{1}{n_i - 1} \sum_{k=1}^{n_i} (X_{ik} - \bar{X}_i)^2$$

nepristranska cenilka za disperzijo v i -ti skupini; neodvisna od s_v^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 $\mathbb{E}\left(\frac{Q_v^2}{\sigma^2}\right) = n - r$, 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 pa spredimo, da je testna 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

$$F = \frac{s_m^2}{s_v^2} = \frac{Q_m^2/(r - 1)}{Q_v^2/(n - r)} = \frac{\frac{Q_m^2}{\sigma^2}/(r - 1)}{\frac{Q_v^2}{\sigma^2}/(n - r)}$$

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	$n - 1$		

vidimo, da gre za Fisherjevo (Snedecorjevo) porazdelitev $F(r - 1, n - r)$, glej *tabelo analize variance*.

13.11.1 Domneva o varianci $\sigma^2 = \sigma_0^2$

$$\text{T.S.} = \frac{(n - 1)s^2}{\sigma_0^2} \text{ sledi } \chi^2\text{-porazdelitev}$$



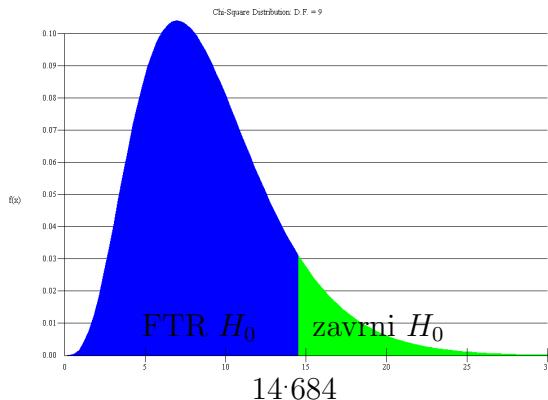
Če je

- $H_a : \sigma^2 > \sigma_0^2$, potem je **odločitveno pravilo**: zavrni ničelno domnevo, če je testna statistika večja ali enaka $\chi_{\alpha}^2(n - 1)$,
- $H_a : \sigma^2 < \sigma_0^2$ potem je **odločitveno pravilo**: zavrni ničelno domnevo, če je testna statistika manjša ali enaka $\chi_{1-\alpha}^2(n - 1)$,
- $H_a : \sigma^2 \neq \sigma_0^2$, potem je **odločitveno pravilo**: zavrni ničelno domnevo, če je testna statistika manjša ali enaka $\chi_{1-\alpha/2}^2(n - 1)$ ali če je testna statistika večja ali enaka $\chi_{\alpha/2}^2(n - 1)$.

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 zavrnjeno tekočino. Če je razpršenost zavrnjene količine prevelika, potem mora naprava na servis. **Ali naj jo odpeljejo na servis?** Uporabi $\alpha = 0,1$.

- Ničelna domneva $H_0 : \sigma^2 = 0\cdot01$,
- Alternativna domneva $H_a : \sigma^2 > 0\cdot01$,
- Predpostavke: naključni vzorec in vzorčenje iz normalne porazdelitve.
- Testna statistika $\chi^2 = (n - 1)s^2/\sigma_0^2$.

Določimo zavrnitveni kriterij



Rezultati testiranja

- naredi naključni vzorec izračunamo naslednjo varianco vzorca: $0\cdot02041$,
- izračunaj vrednost testne statistike $\chi^2 = (0\cdot02041)(9)/(0\cdot01) = 18\cdot369$,
- naredi odločitev: zavrnji H_0 ,
- zaključek popravi napravo.

P -vrednost $= P(\chi^2 > 18\cdot369) = 0\cdot0311 < \alpha$, zato zavrnji domnevo H_0 . \diamond

13.11.2 Domneva o kvocientu varianc $\sigma_1^2/\sigma_2^2 = 1$

Če velja $H_a : \sigma_1^2/\sigma_2^2 > 1$, potem je **testna statistika** enak s_1^2/s_2^2 , **odločitveno pravilo** pa je: zavrnji ničelno domnevo, če velja $T.S. \geq F_\alpha(n_1 - 1, n_2 - 1)$.

Če velja $H_a : \sigma_1^2/\sigma_2^2 < 1$, potem je **testna statistika** enaka

$$\frac{\text{varianca večjega vzorca}}{\text{varianca manjšega vzorca}},$$

odločitveno pravilo pa je: zavrnji ničelno domnevo, če velja $s_1^2 > s_2^2$ in $T.S. \geq F_\alpha(n_1 - 1, n_2 - 1)$ oziroma zavrnji ničelno domnevo, če velja $s_1^2 < s_2^2$ in $T.S. \geq F_\alpha(n_2 - 1, n_1 - 1)$.

13.12 Domneve 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 predpostavljenim porazdelitvijo.

13.12.1 Domneva o enakomerni porazdelitvi

Za primer vzemimo met kocke in za spremenljivko število pik pri metu kocke. Preverimo 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.

Testna 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

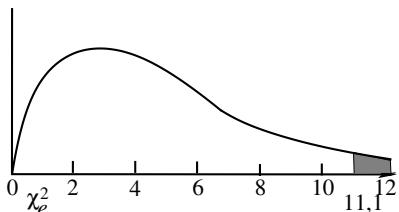
$$H_0 : \chi^2 = 0 \quad \text{in} \quad H_1 : \chi^2 > 0.$$

Domnevo preverimo pri stopnji značilnosti $\alpha = 5\%$. Ker gre za enostranski test, je kritična vrednost enaka

$$\chi^2_{1-\alpha}(k-1) = \chi^2_{0.95}(5) = 11.1.$$

Eksperimentalna vrednost statistike pa je

$$\begin{aligned} \chi^2_e &= \frac{(20-20)^2}{20} + \frac{(22-20)^2}{20} + \frac{(17-20)^2}{20} + \frac{(18-20)^2}{20} + \frac{(19-20)^2}{20} + \frac{(24-20)^2}{20} \\ &= \frac{4+9+4+1+16}{20} = \frac{34}{20} = 1.7. \end{aligned}$$



Ker eksperimentalna vrednost statistike ne pade v kritično območje, ničelne domneve ne moremo zavrniti. Empirične in teoretične frekvence niso statistično značilno različne med seboj.

13.12.2 Domneva 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: Preverimo 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 priloženi frekvenčni porazdelitvi. Ničelna in osnovna domneva sta tedaj

$$H_0 : \chi^2 = 0 \quad \text{in} \quad H_1 : \chi^2 \neq 0.$$

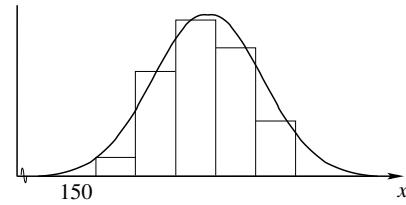
Za test uporabimo statistiko

$$\chi^2 = \sum_{i=1}^k \frac{(f_i - f'_i)^2}{f'_i},$$

ki se porazdeljuje po χ^2 porazdelitvi z $m = n - 1$ prostostnimi stopnjami.

	f_i
nad 150-160	2
nad 160-170	20
nad 170-180	40
nad 180-190	30
nad 190-200	8
	100

V našem primeru je $n = 5$ in $m = 4$. Kritična vrednost je $\chi^2_{0.95}(4) = 9.49$. 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 (glej sliko). Tako je na primer verjetnost, da je višina med 150 in 160 cm:



$$\begin{aligned} P(150 < X < 160) &= P\left(\frac{150 - 177}{10} < Z < \frac{160 - 177}{10}\right) \\ &= P(-2.7 < Z < -1.7) = \Phi(2.7) - \Phi(1.7) = 0.4965 - 0.4554 \\ &= 0.0411. \end{aligned}$$

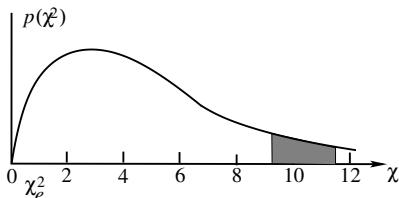
Podobno lahko izračunamo ostale verjetnosti.

Teoretične frekvence so $f'_i = n \times p_i$. Izračunane verjetnosti p_i in teoretične frekvence f'_i postavimo v tabelo.

Eksperimentalna vrednost statistike je tedaj

	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

$$\chi_e^2 = \frac{(2 - 4.11)^2}{4.11} + \frac{(20 - 19.74)^2}{19.74} + \frac{(40 - 37.59)^2}{37.59} + \frac{(30 - 28.53)^2}{28.53} + \frac{(8 - 8.61)^2}{8.61} \doteq 1$$



Ker eksperimentalna vrednost ne pade v kritično območje, ne moremo zavrniti ničelne domneve, da je spremenljivka normalno porazdeljena. ◇

Obstajajo tudi drugi testi za preverjanje porazdelitve spremenljivke, npr. Kolmogorov-Smirnov test.

Poglavlje 14

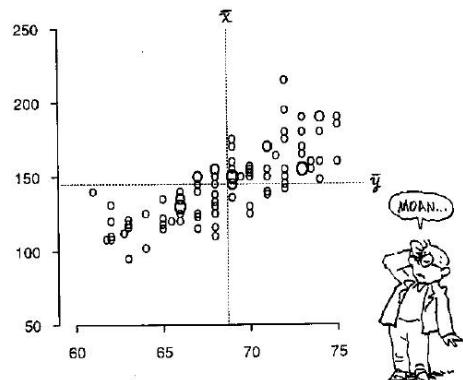
Bivariatna analiza in regresija

Bivariatna analiza

$X \longleftrightarrow Y$ povezanost

$X \rightarrow Y$ odvisnost

Mere povezanosti ločimo glede na tip spremenljivk:



1. IMENSKI/NOMINALNI tip para spremenljivk (ena od spremenljivk je imenska/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.

14.1 Povezanost dveh imenskih (nominalnih) spremenljivk

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. Denimo, da so podatki za vzorec urejeni v naslednji kontingenčni tabeli. Ker nas zanima ali študentke drugače stanujejo v času študija kot

	starši	št. dom	zasebno	skupaj
moški	16	40	24	80
ženske	48	36	36	120
skupaj	64	76	60	200

študentje, moramo porazdelitev stanovanja študentk primerjati s porazdelitvijo študentov.

Ker je število študentk različno od števila študentov, moramo zaradi primerjave izračunati relativne frekvence. Če med spoloma ne bi bilo razlik, bi bili obe porazdelitvi (za

	starši	št. dom	zasebno	skupaj
moški	20	50	30	100
ženske	40	30	30	100
skupaj	32	38	30	100

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. Kontingenčna tabela kaže podatke za slučajni vzorec. Zato nas zanima, ali so

	starši	št. dom	zasebno	skupaj
moški	25	56·6	40	40
ženske	75	43·4	60	60
skupaj	100	100	100	100

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 imenskima (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:

$$P(M) = \frac{80}{200} = 0.40, \quad P(S) = \frac{64}{200} = 0.32, \quad P(M S) = P(M) \cdot P(S) = \frac{80}{200} \times \frac{64}{200} = 0.128.$$

Teoretična frekvenca je verjetnost $P(M S)$ pomnožena s številom enot v vzorcu:

$$f'(M S) = n \cdot P(M S) = 200 \times \frac{80}{200} \times \frac{64}{200} = 25.6.$$

Podobno izračunamo teoretične frekvence f'_i tudi za druge celice kontingenčne tabele. Spomnimo se tabel empiričnih (dejanskih) frekvenc f_i :

	starši	št. dom	zasebno	skupaj
moški	26	30	24	80
ženske	38	46	36	120
skupaj	64	76	60	200

Testna statistika χ^2 , 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 (razredov) v kontingenčni tabeli. Testna 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 s število stolpcev. Ničelna in osnovna domneva sta v primeru tega testa

$$\begin{aligned} H_0: \quad \chi^2 = 0 & \quad (\text{spremenljivki nista povezani}) \\ H_1: \quad \chi^2 > 0 & \quad (\text{spremenljivki sta povezani}) \end{aligned}$$

Iz tabele za porazdelitev χ^2 lahko razberemo kritične vrednosti te statistike pri 5% stopnji značilnosti:

$$\chi^2_{1-\alpha}[(s-1)(v-1)] = \chi^2_{0.95}(2) = 5.99.$$

Ekperimentalna vrednost statistike χ^2 pa je:

$$\chi^2_e = \frac{(16-26)^2}{26} + \frac{(40-30)^2}{30} + \frac{(24-24)^2}{24} + \frac{(48-38)^2}{38} + \frac{(36-46)^2}{46} + \frac{(36-36)^2}{36} = 12.$$

Ker je ekperimentalna vrednost večja od kritične vrednosti, pomeni, da pade v kritično območje. To pomeni, da ničelno domnevo zavrnemo. Pri 5% stopnji značilnosti lahko sprejmemo osnovno domnevo, da sta spremenljivki statistično značilno povezani med seboj. ◇

Testna statistika χ^2 je lahko le pozitivna. Zavzame lahko vrednosti v intervalu $[0, \chi^2_{\max}]$, kjer je $\chi^2_{\max} = n(k-1)$, če je $k = \min(v, s)$. Statistika χ^2 v splošnem ni primerljiva. Zato je definiranih več **kontingenčnih koeficientov**, ki so bolj ali manj primerni. Omenimo naslednje:

1. **Pearsonov koeficient:** $\Phi = \chi^2/n$, 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:** $C = \sqrt{\chi^2/(\chi^2 + n)}$, ki je definiran na intervalu $[0, C_{\max}]$, kjer je $C_{\max} = \sqrt{k/(k-1)}$.

14.2 Koeficienti asociacije

Denimo, da imamo dve imenski (nominalni) spremenljivki, ki imata le po dve vrednosti (sta dihotomni). Povezanost med njima lahko računamo poleg kontingenčnih koeficientov s **koeficienti asociacije** na osnovi frekvenc iz kontingenčne tabele,

$Y \setminus X$	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 definirnih več koeficientov asociacije. Omenimo najpomembnejše:

- **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:** $J = \frac{a}{a + b + c} \in [0, 1]$,

Primer: povezanost med kaznivimi dejanji in alkoholizmom. Tabela kaže podatke za $N = 10.750$ ljudi. Izračunajmo koeficiente asociacije:

alk. \ kaz. d.	DA	NE	skupaj
DA	50	500	550
NE	200	10.000	10.200
skupaj	250	10.500	10.750

$$Q = \frac{50 \times 10000 - 200 \times 500}{50 \times 10000 + 200 \times 500} = 0.67,$$

$$S = \frac{10050}{10750} = 0.93 \text{ in } J = \frac{50}{50 + 500 + 200} = 0.066.$$

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). \diamond

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.

14.3 Povezanost dveh ordinalnih spremenljivk

V tem primeru gre za študij povezanosti med dvema spremenljivkama, ki sta vsaj ordinalnega značaja. Povezanost med spremenljivkama lahko merimo s **koeficientom korelacije rangov r_s (Spearman)**, ki je definiran takole:

$$r_s := 1 - \frac{6 \cdot \sum_{i=1}^n d_i^2}{n \cdot (n^2 - 1)}, \quad \text{kjer je } d_i \text{ razlika med rangoma v } i\text{-ti enoti.}^1$$

Če je ena od obeh spremenljivk številska, moramo vrednosti pred izračunom d_i rangirati. Če so kakšne vrednosti enake, zanje izračunamo povprečne pripadajoče range.

Koeficient korelacije rangov lahko zavzame vrednosti na 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

¹ Mimogrede: ni se težko prepričati o naslednji identiteti: $1^2 + 2^2 + 3^2 + \dots + n^2 = n(n+1)(2n+1)/6$ (npr. s popolno indukcijo, ali pa neposredno iz vsote enačb $(1+k)^3 = 1 + 3k + 3k^2 + k^3$ za $k = 0, \dots, n$).

spremenljivki rangi po drugi spremenljivki manjšajo, gre za *negativno povezanost*. Koeficient je tedaj negativen in blizu -1 . V naslednjem (preprostem) primeru gre negativno povezanost. Če ne gre za pozitivno in ne za negativno povezanost, rečemo, da spremenljivki nista povezani.

Postavimo ničelno in osnovno domnevo:

$$H_0: \rho_s = 0 \text{ (spremenljivki nista povezani)}$$

$$H_1: \rho_s \neq 0 \text{ (spremenljivki sta povezani)}$$

kjer korelacijski koeficient (populacije) označimo s ρ_s . Pokazati se da, da se **testna statistika**

$$t = \frac{r_s \cdot \sqrt{n-2}}{\sqrt{1-r_s^2}}$$

porazdeljuje približno po *t*-porazdelitvi z $m = (n-2)$ prostostnimi stopnjami.

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 (glej tabelo).

poklic	R_0	R_N
A	1	6
B	2	4
C	3	5
D	4	2
E	5	3
F	6	1

Izračunajmo koeficient korelacije rangov za primer šestih poklicev:

poklic	R_0	R_N	d_i	d_i^2
A	1	6	-5	25
B	2	4	-2	4
C	3	5	-2	4
D	4	2	2	4
E	5	3	2	4
F	6	1	5	25
vsota			0	66

$$r_s = 1 - \frac{6 \times 66}{6 \times 35} = 1 - 1.88 = -0.88.$$

Res je koeficient blizu -1 , 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. 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

$$t_e = \frac{-0.88 \times 2}{\sqrt{1 - (-0.88)^2}} = \frac{-1.76}{0.475} = -3.71.$$

Ekperimentalna vrednost pade v kritično območje. Pri 5% stopnji značilnosti lahko rečemo, da sta odgovornost in fizična napornost (negativno) povezani med seboj. \diamond

14.4 Povezanost dveh številskih spremenljivk

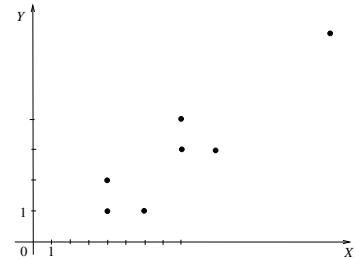
Vzemimo primer dveh številskih spremenljivk:

X - izobrazba (število priznanih let šole)

Y - število ur branja dnevnih časopisov na teden

Podatki za 8 slučajno izbranih oseb so:

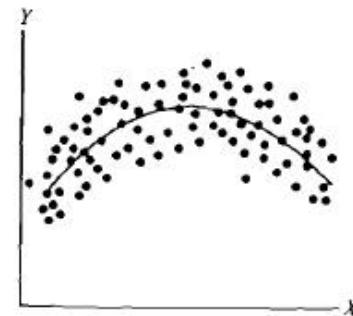
X	10	8	16	8	6	4	8	4
Y	3	4	7	3	1	2	3	1



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.

Tipi povezanosti:

- **funkcijska** povezanost: vse točke ležijo na krivulji,
- **korelacijska** (stohastična) povezanost: točke so od krivulje bolj ali manj odklanjajo (manjša ali večja povezanost).

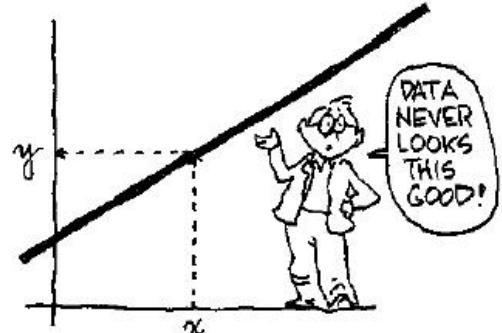


Primer nelinearne povezanosti spremenljivk.

(Vzorčna) kovarianca

$$k(X, Y) = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x}) \cdot (y_i - \bar{y})$$

meri povezanost med spremenljivkama.



(Pearsonov) koeficient korelacije je definiran s formulo

$$r_{XY} = \frac{k(X, Y)}{s_X s_Y} = \frac{\sum_{i=1}^n (x_i - \bar{x}) \cdot (y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2 \cdot \sum_{i=1}^n (y_i - \bar{y})^2}}.$$

Koeficient korelacije lahko zavzame vrednosti v intervalu $[-1, 1]$. Če se z večanjem vrednosti prve spremenljivke večajo tudi vrednosti druge spremenljivke, gre za **pozitivno linearne povezanost**. Tedaj je koeficient povezanosti blizu 1. Če pa se z večanjem vrednosti prve spremenljivke vrednosti druge spremenljivke manjšajo, gre za **negativno linearne povezanosti**.

zanost. 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:

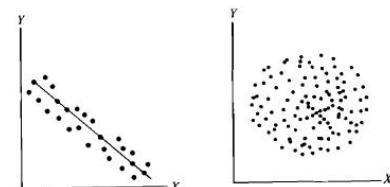
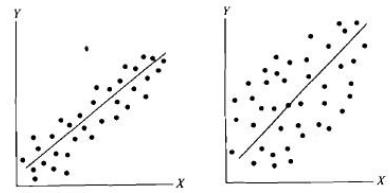
$H_0: \rho = 0$ (spremenljivki nista linearno povezani)

$H_1: \rho \neq 0$ (spremenljivki sta linearno povezani)

Pokaže se, da se **testna statistika**

$$t = \frac{r \cdot \sqrt{n-2}}{\sqrt{1-r^2}}$$

porazdeljuje po t porazdelitvi z $m = (n-2)$ prostostnimi stopnjami.



Tipični primeri linearne povezanosti spremenljivk

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	$x_i - \bar{x}$	$y_i - \bar{y}$	$(x_i - \bar{x})^2$	$(y_i - \bar{y})^2$	$(x_i - \bar{x})(y_i - \bar{y})$
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

$$r = \frac{48}{\sqrt{104 \times 26}} = 0.92.$$

Eksperimentalna vrednost statistike je:

$$t_e = \frac{0.92 \sqrt{8-2}}{\sqrt{1 - 0.92^2}} = 2.66.$$

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 pade v kritično območje. *Zaključek:* ob 5% stopnji značilnosti lahko rečemo, da je izobrazba linearno povezana z branjem dnevnih časopisov. \diamond

14.5 Parcialna korelacija

Včasih je potrebno meriti zvezo med dvema spremenljivkama in odstraniti vpliv vseh ostalih spremenljivk. To zvezo dobimo s koeficientom 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.

	X	Z	U	Y
X	1	0.51	0.41	0.36
Z		1	0.29	0.60
U			1	0.49
Y				1

Izračunani koeficienti korelacije so podani v zgornji tabeli. 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 poskusimo odstraniti vpliv revščine iz zveze: “nebelo prebivalstvo : kazniva dejanja”:**

$$r_{XY,Z} = \frac{0.36 - 0.51 \times 0.60}{\sqrt{1 - 0.51^2} \sqrt{1 - 0.60^2}} = 0.08.$$

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

14.6 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

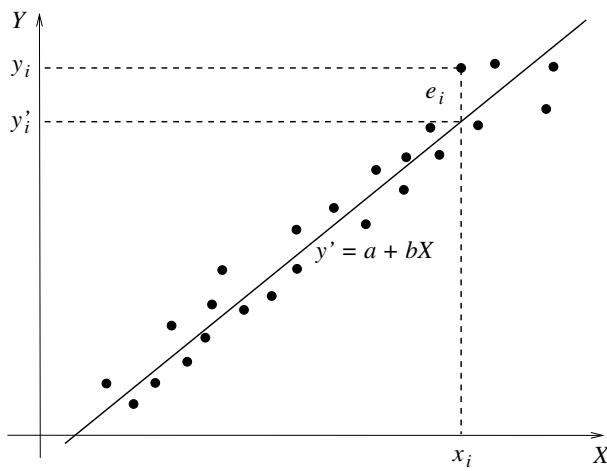
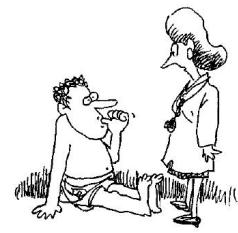
$$Y = Y' + E = f(X) + E,$$

kjer X imenujemo neodvisna spremenljivka, Y odvisna spremenljivka in E člen napake (ali motnja, disturbanca). Če je regresijska funkcija linearna, tj. $Y' = f(X) = a + bX$,

je regresijska odvisnost $Y = Y' + E = a + bX + E$ ozziroma za i to enoto

$$y_i = y'_i + e_i = a + b x_i + e_i.$$

Regresijsko odvisnost si lahko zelo nazorno predstavimo v razsevnem grafikonu.



Regresijsko funkcijo lahko v splošnem zapišemo kot

$$Y' = f(X, a, b, \dots),$$

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 razsevnem grafikonu.

Pri dvorazsežno normalno porazdeljenem slučajnem vektorju $(X, Y) : N(\mu_X, \mu_Y, \sigma_X, \sigma_Y, \rho)$ je, kot vemo

$$\mathbb{E}(Y|x) = \mu_Y + \rho \frac{\sigma_Y}{\sigma_X} (x - \mu_X).$$

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 linearja, regresijska premica pa gre skozi točko (μ_X, μ_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 $y = \alpha + \beta x$.

Preverjanje 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 porazdeljeni cenilki B in A ?

$$B = \frac{C_{xy}}{C_x^2} = \sum_{i=1}^n \frac{X_i - \bar{X}}{C_x^2} (Y_i - \bar{Y}) = \sum_{i=1}^n \frac{X_i - \bar{X}}{C_x^2} Y_i.$$

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 $Y_i \sim N(\alpha + \beta X_i, \sigma)$, 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 smo upoštevali, 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 $B \sim N\left(\beta, \frac{\sigma}{C_x}\right)$, oziroma $\frac{B - \beta}{\sigma} C_x \sim N(0, 1)$. 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 računanju se izkaže, da velja $\mathbb{E}(Q^2/\sigma^2) = n - 2$. Torej je $S^2 = Q^2/(n - 2) = \frac{\sigma^2}{n - 2} \frac{Q^2}{\sigma^2}$ nepristranska cenilka za σ^2 . S^2 je neodvisna od A in B . Testni statistiki za A in B sta tedaj

$$T_A = \frac{A - \mathbb{E}(A)}{\sqrt{\mathbb{D}(A)}} = \frac{A - \alpha}{S} \sqrt{\frac{nC_x^2}{C_x^2 + n\bar{X}^2}} = \frac{A - \alpha}{S} C_x \sqrt{\frac{n}{\sum_{i=1}^n X_i^2}}, \quad T_B = \frac{B - \mathbb{E}(B)}{\sqrt{\mathbb{D}(B)}} = \frac{B - \beta}{S} C_x,$$

ki obe sledita Studentovo t -porazdelitev z $n - 2$ prostostnimi stopnjami. Statistika za σ^2 pa je spremenljivka $Q^2/\sigma^2 = (n - 2)S^2/\sigma^2$, 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 .

14.7 Linearni model

Pri proučevanju pojavov pogosto teorija postavi določeno funkcionalno zvezo med obravnavanimi spremenljivkami. Oglejmo si primer *linearnega 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 ustreznnejši *verjetnostni linearni model*

$$Y = \alpha + \beta X + U,$$

kjer so X , Y in U slučajne spremenljivke in $\mathbb{E}(U) = 0$ – 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 s pričakovano vrednostjo 0 in disperzijo σ^2 . Torej je:

$$\mathbb{E}(U_i) = 0, \quad \mathbb{D}(U_i) = \sigma^2 \quad \text{in} \quad \mathbb{E}(U_i U_j) = 0, \quad \text{za } i \neq j.$$

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

$$\frac{\partial f}{\partial a} = \sum_{i=1}^n 2(y_i - (bx_i + a)) = 0, \quad \frac{\partial f}{\partial b} = \sum_{i=1}^n 2(y_i - (bx_i + a))x_i = 0,$$

z rešitvijo

$$b = \frac{n \sum xy - \sum x \sum y}{n \sum x^2 - (\sum x)^2}, \quad a = \frac{1}{n} \left(\sum y - b \sum x \right).$$

ozioroma, če vpeljemo oznako $\bar{z} = \frac{1}{n} \sum z$:

$$b = \frac{\bar{xy} - \bar{x}\bar{y}}{\bar{x^2} - \bar{x}^2}, \quad a = \bar{y} - b\bar{x}.$$

Poglejmo še matriko dvojnih parcialnih odvodov:

$$H = \begin{bmatrix} \frac{\partial^2 f}{\partial a^2} & \frac{\partial^2 f}{\partial a \partial b} \\ \frac{\partial^2 f}{\partial b \partial a} & \frac{\partial^2 f}{\partial b^2} \end{bmatrix} = 2 \begin{bmatrix} n & \sum x \\ \sum x & \sum x^2 \end{bmatrix}.$$

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 f 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 koeficiente 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

$$\mathbb{E}(A) = \alpha \text{ in } D(A) = \sigma^2 \left(\frac{1}{n} + \frac{\bar{X}^2}{C_x^2} \right), \quad \mathbb{E}(B) = \beta \text{ in } D(B) = \frac{\sigma^2}{C_x^2}, \quad \text{Cov}(A, B) = -\sigma^2 \frac{\bar{X}}{C_x^2}.$$

Cenilki za A in B sta najboljši linearni nepristranski cenilki za α in β .

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.

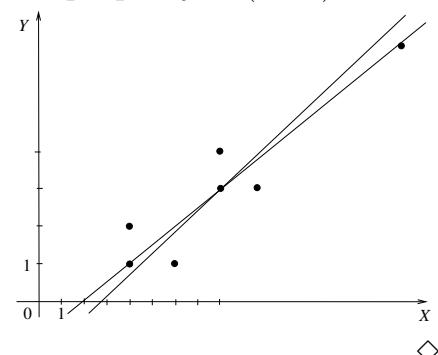
X	10	8	16	8	6	4	8	4
Y	3	4	7	3	1	2	3	1

Zanje izračunajmo obe regresijski premici in ju vrišimo v razsevni grafikon. Ko smo računali koeficient korelacije, smo že izračunali vzorčni povprečji $\bar{x} = 64/8 = 8$ in $\bar{y} = 24/8 = 3$, vsoti kvadratov odklonov od vzorčnega povprečja za obe spremenljivki $(n-1)s_X^2 = 104$ in $(n-1)s_Y^2 = 26$ ter vsoto produktov odklonov od obeh vzorčnih povprečij $n k(X, Y) = 48$:

$$y = 3 + \frac{48}{104} (x - 8) = 0.46x - 0.69,$$

$$x = 8 + \frac{48}{26} (y - 3) = 1.85y + 2.46, \text{ tj. } y = 0.54x + 1.33.$$

Obe regresijski premici lahko vrišemo v razsevni grafikon in preverimo, če se res najbolje prilegata točкам v grafikonu.



Regresijski premici se sečeta v točki, določeni s pričakovanima vrednostima spremenljivk X in Y . **Dokažite, da se premici vedno sečeta v tej točki.**

14.7.1 Statistično sklepanje o regresijskem koeficientu

Vpeljimo naslednje oznake:

$Y = \alpha + \beta X$ regresijska premica na populaciji, $y = a + b x$ regresijska premica na vzorcu.

Denimo, da želimo preveriti domnevo o regresijskem koeficientu β . Postavimo ničelno in osnovno domnevo: $H_0: \beta = \beta_0$, $H_1: \beta \neq \beta_0$.

Nepristranska cenilka za regresijski koeficient β je $b = \text{Cov}(X, Y) / s_X^2$. Njena pričakovana vrednost in standardna napaka sta:

$$\mathbb{E}(b) = \beta; \quad \text{SE}(b) = \frac{s_Y \sqrt{1 - r^2}}{s_X \sqrt{n - 2}}.$$

Testna statistika za zgornjo ničelno domnevo je:

$$t = \frac{s_Y \sqrt{n - 2}}{s_X \sqrt{1 - r^2}} (b - \beta_0)$$

in 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. **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:

$$\pm t_{\alpha/2}(n - 2) = \pm t_{0.025}(6) = \pm 2.447, \quad \text{TS} = \sqrt{\frac{104 \times (8 - 2)}{26 \times (1 - 0.92^2)}} \cdot (0.46 - 0) = 5.8.$$

Regresijski koeficient je statistično značilno različen od 0. ◇

14.7.2 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 $R = \sigma_{Y'}^2 / \sigma_Y^2$. Imenujemo ga **determinacijski koeficient** in je definiran na intervalu $[0, 1]$. Pokazati se da, da je v primeru linearne regresijske odvisnosti determinacijski koeficient $R = \rho^2$, kjer je ρ koeficient korelacije. Kvadratni koren iz nepojasnjene variance σ_e^2 imenujemo **standardna napaka regresijske ocene**. Meri razpršenost točk okoli regresijske krivulje oziroma kakovost ocenjevanja vrednosti odvisne spremenljivke z regresijsko funkcijo. V primeru linearne regresijske odvisnosti je standardna napaka enaka: $\sigma_e = \sigma_Y \sqrt{1 - \rho^2}$, saj velja: $\sigma_e^2 = \sigma_Y^2 - \sigma_{Y'}^2 R^2 = \sigma_Y^2 - \sigma_Y^2 (\sigma_{Y'}^2 / \sigma_Y^2) = \sigma_Y^2 - \sigma_{Y'}^2$.

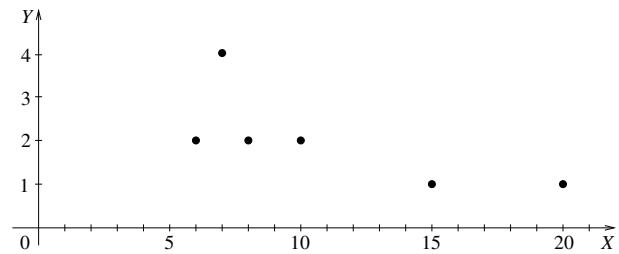
Primer: Vzemimo spremenljivki

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



(a) Z linearno regresijsko funkcijo ocenimo, kolikokrat bo šla oseba v kino na mesec, če gleda 18 ur na teden televizijo. (b) Kolikšna je standardna napaka? (c) Kolikšen delež variance obiska kinopredstav lahko pojasnimo z gledanjem televizije?

x_i	y_i	$x_i - \bar{x}$	$y_i - \bar{y}$	$(x_i - \bar{x})^2$	$(y_i - \bar{y})^2$	$(x_i - \bar{x}) \cdot (y_i - \bar{y})$
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

$$y' = 2 - \frac{21}{148} (x - 11) \\ = 3.54 - 0.14x,$$

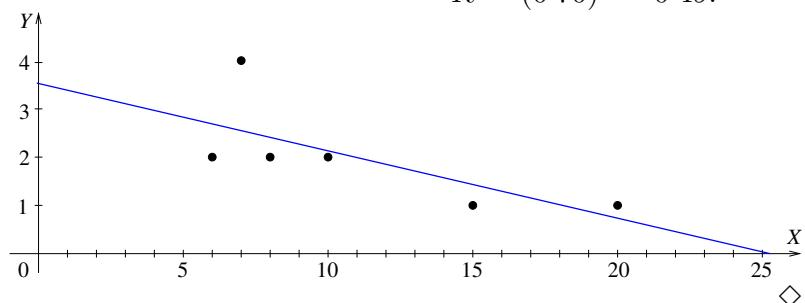
$$y'(18) = 3.54 - 0.14 \cdot 18 = 1.02,$$

$$r = \frac{-21}{\sqrt{148 \times 6}} = -0.70,$$

$$\sigma_e^2 = \frac{6}{6} \sqrt{1 - (-0.70)^2} = 0.71,$$

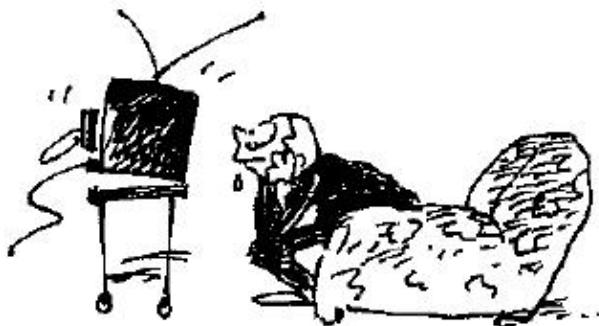
$$R = (0.70)^2 = 0.49.$$

Če oseba gleda TV 18 ur na teden, lahko pričakujemo, da bo 1-krat na mesec šla v kino, pri čemer je standardna napaka 0.71. 49% variance obiska kino predstav lahko pojasnimo z gledanjem TV.



Poglavlje 15

Časovne vrste in trendi



Družbeno-ekonomski pojavi so časovno spremenljivi. Spremembe so rezultat delovanja naj-različnejših dejavnikov, ki tako ali dugače vplivajo na pojave. Sliko dinamike pojavov dobimo s časovimi 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 spremenjanje vrednosti členov v časovih vrstah in iščemo zakonitosti tega spremenjanja. 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). ◇

15.1 Primerljivost členov v časovni vrsti

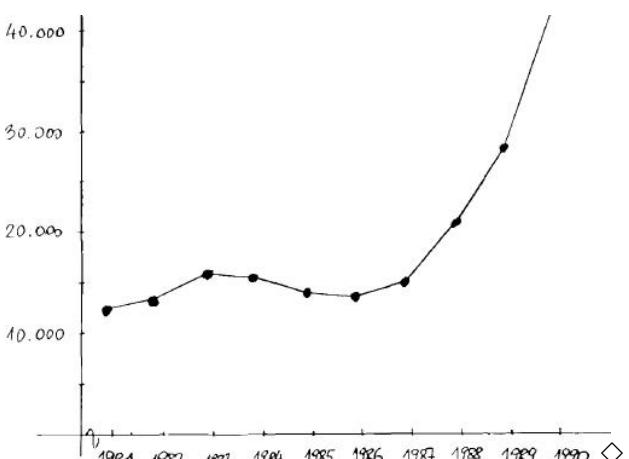
Kljub temu, da so členi v isti časovni vrsti istovrstne količine, dostikrat niso med seboj neposredno primerljivi. Osnovni pogoj za primerljivost členov v isti časovni vrst je pravilna in nedvoumna opredelitev pojave, ki ga časovna vrsta prikazuje. Ta opredelitev mora biti vso dobo opazovanja enaka in se ne sme spremenjati. 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.

15.2 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.



15.3 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 : $I_{k/0} = (X_k/X_0) \cdot 100$.
- **Verižni indeksi.** Za dano časovno vrsto računamo vrsto verižnih indeksov tako, da za vsak člen vzamemo za osnovo predhodni člen: $I_k = (X_k/X_{k-1}) \cdot 100$.

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): $T_k = ((X_k - X_{k-1})/X_{k-1}) \cdot 100 = I_k - 100$.

Interpretacija indeksov

indeks\pojav	raste	stagnira	pada
stalna osnova	$I_{k+1/0} > I_{k/0}$	$I_{k+1/0} = I_{k/0}$	$I_{k+1/0} < I_{k/0}$
verižni	$I_k > 100$	$I_k = 100$	$I_k < 100$
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).

15.4 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:

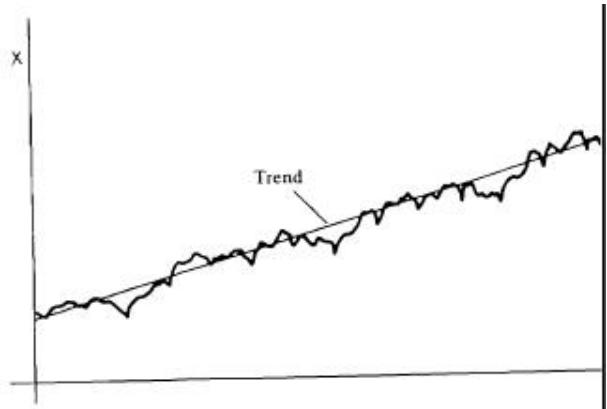
1. Dolgoročno gibanje ali trend - X_T podaja dolgoročno smer razvoja. Običajno ga je mogoče izraziti s preprostimi rahlo ukrivljenimi krivuljami.
2. Ciklična gibanja - X_C , so oscilacije okoli trenda. Poriode so ponavadi daljše od enega leta in so lahko različno dolge.
3. Sezonske oscilacije - X_S so posledice vzrokov, ki se pojavljajo na stalno razdobje. Periode 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 sestavnimi je mogoče prikazati z nekaj osnovnim modeli. Npr.:

$$X = X_T + X_C + X_S + X_E$$

$$\text{ali } X = X_T \cdot X_C \cdot X_S \cdot X_F;$$

$$\text{ali } X = X_T \cdot X_C \cdot X_S + X_E.$$



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:

$$H_0: \rho_e = 0 \quad \text{trend ne obstaja}$$

$$H_1: \rho_e \neq 0 \quad \text{trend obstaja}$$

Ustrezna testna statistika

$$t = \frac{r_s \sqrt{n-2}}{\sqrt{1-r_s^2}}$$

se porazdeluje približno po t porazdelitvi z $(n-2)$ prostostnimi stopnjami.

Metode določanja trenda

- Prostoročno
- 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 ne moremo izračunati.

- Včasih se uporablja utežena 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 metodo najmanjših kvadratov $\sum_{i=1}^n (X_i - X_{iT})^2 = \min$. V primeru linearnega trenda

$$X_T = a + bT, \quad \sum_{i=1}^n (X_i - a - bT_i)^2 = \min.$$

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 $\bar{T} = 0$. 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

$$\sigma_e = \sqrt{\frac{1}{n} \sum_{i=1}^n (X_i - X_{iT})^2}, \quad \text{kjer je } X_{iT} \text{ enak } X_T \text{ v času } T_i.$$

Primer: Kot primer si oglejmo število doktoratov znanosti v Sloveniji v razdobju od leta 1986 do 1990. (a) **Z linearnim trendom ocenimo koliko doktorjev znanosti je v letu 1991.** (b) **Izračunajmo tudi standardno napako ocene.** Izračunajmo najprej trend:

t	x_i	t_i	$x_i - \bar{x}$	$(x_i - \bar{x})t_i$	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	0	80	10

$$\bar{x} = \frac{544}{5} = 108.8,$$

$$x_t = 108.8 + \frac{80}{10} t = 108.8 + 8t,$$

$$x_t(1991) = 108.8 + 8 \times 3 = 132.8.$$

(a) Ocena za leto 1991 je približno 133 doktorjev znanosti.

(b) Za vsako leto je potrebno najprej izračunati x_{it} iz regresijske premice (trenda).

t	x_i	x_{it}	$x_i - x_{it}$	$(x_i - x_{it})^2$
1986	89	92.8	-3.8	14.14
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

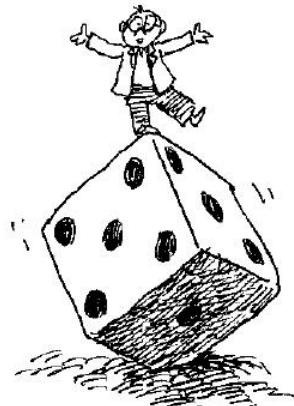
$$\sigma_e = \sqrt{\frac{114.8}{5}} = 4.8.$$

◇

Del III

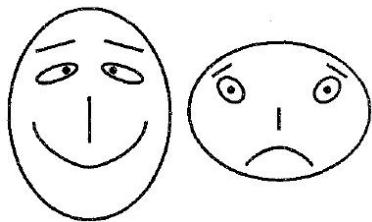
KAM NAPREJ

Osnovne principe in orodja, ki smo jih spoznali pri OVS/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:

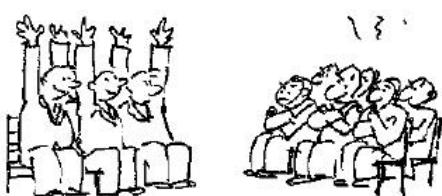
- X =naklon obrvi,
- Y =velikost oči,
- Z =dolžina nosu,
- T =dolžina ust,
- U =višino obraza, itd.



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.



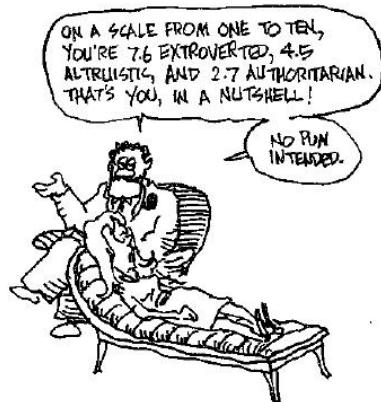
Diskriminacijska analiza

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 prispevkvi) ali pa ne bodo uspešni (gre delati dobro po svetu in šola nikoli več ne sliši zanj(o)).



Analiza faktorjev

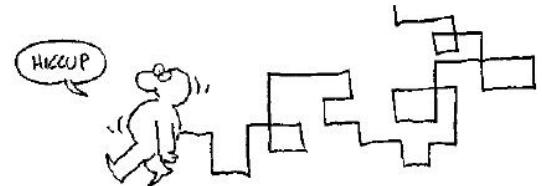
išče poenostavljen razlagi večrazsežnih podatkov z manjšo skupino spremenljivk. Npr. Psihiater lahko postavi 100 vprašanj, skrivoma pa pričakuje, da so odgovori odvisni samo od nekaterih faktorjev: ekstravertiranost, avtoritativnost, alutarizem, itd. Rezultate testa lahko potem povzamemo le z nekaterimi sestavljenimi rezultati v ustreznih dimenzijah.



Naključni sprehodi

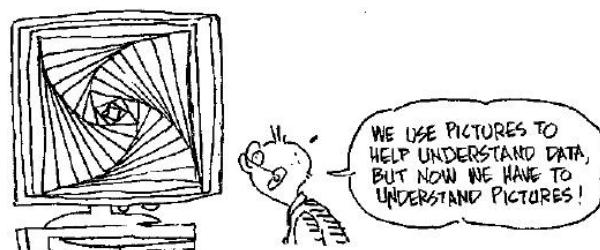
Pričnejo se z metom kovanca, recimo, da se pomaknemo korak nazaj, če pade grb, in korak naprej, če pade cifra. (z dvema kovancema se lahko gibljemo v 2-razsežnemu prostoru - tj.

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 pikslov, 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.



Moderno 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

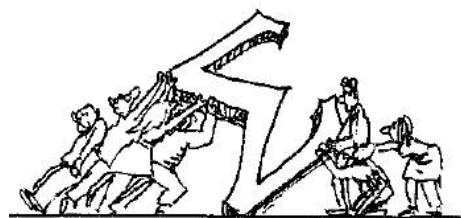
Še tako uspešna in bistroumna analiza je zelo malo vredna, če je ne znamo jasno predstaviti, vključujoč stopnjo statistične značilnosti.



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.



Kako skleniti predavanja

V knjigi Revolucija učenja (RU) sta kot primer navedena dva od mnogih načinov, ki prikazujeta, kako je treba na pozitiven način skleniti poučevanje. Vpletena naj bo dobršna mera humorja, obenem pa na hitro povzamite najpomembnejše točke:

1. (a) Vsak študent(ka) naj na list papirja zapiše stavek, s katerim bo povedal, kaj se je naučil(a).
- (c) Združitev posameznikov v pare z nalogo, da drug drugega prepričate, kaj je bistvo snovi. Na voljo imate 45 sekund.
- (e) Po dva para se združita v skupine po štiri. Naloga je enaka kot prej.
- (f) Združevanje se nadaljuje v skupine po osem in tako dalje do združitve v dve veliki skupini.¹

Za vsako stopnjo je na voljo le približno dve minuti časa. Ko postaneta skupini veliki, čas nekoliko podaljšamo.²

2. (a) Udeležencem odmerite pet minut in jim dajte nalogo, da v kratkih stavkih opišejo glavne vsebine, ki so se jih naučili. Vsaka vsebina naj bo napisana na svoj list papirja.
- (c) Liste potem pritrdijo na veliko oglasno desko.
- (d) Začne se postopek razvrščanja podobnih vsebin v skupine. Razvrščanje spremišča debata o razlogih, zakaj je nek listek uvrščen v eno skupino, drug pa v drugo.

Velikost posameznih skupin pomaga skupini, da naredi primerne sklepe. Učitelj pa lahko doda še končne pripombe.³

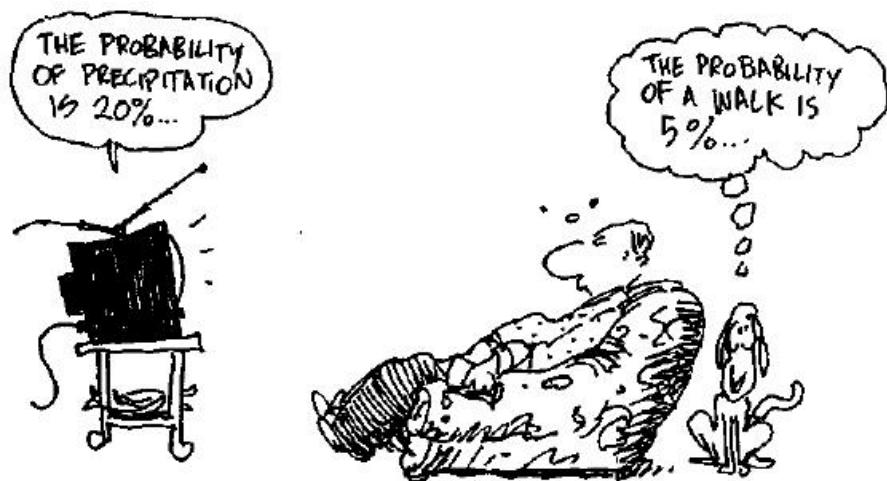
¹ Ko postanejo skupine velike, določite predstavnika, ki zastopa vaše mnenje.

² Skupina 300 ljudi lahko zaključi debato v dvajsetih minutah.

³ Oba primera sta bila uporabljena v delavnicah, ki jih pripravlja združenje International Alliance for Learning Conferences in America.

Poglavlje 16

Nekaj primerov uporabe verjetnosti



16.1 Zamenjalna šifra

Tomaž Pisanski, Skrivnostno sporočilo, *Presek V/1, 1977/78*, str. 40-42.

YHW?HD+CVODHVTHVO-! JVG: 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-?+CV3J0-YC

(črko Č smo zamenjali s C, črko Č pa z D). Imamo $26! = 40329146112665635584000000$ možnosti z direktnim preverjanjem, zato v članku dobimo 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 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 tekstrom:

YHW?HD+C ODH TH 0-!J G:CDCYJ(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 3J0-YC

- (3) Kandidati za samoglasnike e,a,i,o so znaki z visokimi frekvancami. 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 3

- (5) Preučimo besede z dvema črkama:

Samoglasnik na koncu

- 1) da ga na pa ta za (ha ja la)
- 2) če je le me ne se še te ve že (he)
- 3) bi ji ki mi ni si ti vi
- 4) bo do (ho) jo ko no po so to
- 5) ju mu tu (bu)
- 6) rž rt

Samoglasnik na začetku

- 1) ar as (ah aj au)
- 2) en ep (ej eh)
- 3) in iz ig
- 4) on ob od os on (oh oj)
- 5) uk up uš ud um ur (uh ut)

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 preveriti 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 zmogel problema zaradi premajhne dolžine tajnopisa, vsekakor pa bi bilo problem mogoče rešiti z uporabo 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 angleščini¹ jeziku veliko členov THE, A in AN, vendar pa zato običajno najprej izpustimo presledke iz teksta, ki ga želimo spraviti v tajnopsis. 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.

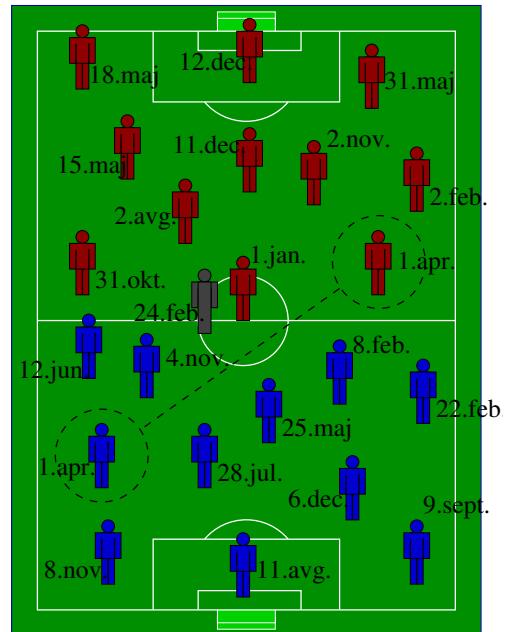
¹ Preprost začetek učenja angleščine: Angleščina ima 550.000 besed, vendar v 90% govora uporablja le 2.000 besed, 400 besed sestavlja 65% večine zapisanih tekstov. Angleščina ima 26 črk in 44 glasov. Obstaja le 70 glavnih kombinacij izgovorjave. Polovica besed je zapisana fonetično, polovica ni. GORDON DRYDEN, Povzeto po prosojnici, predstavljeni na konferenci The Peoples, Netvork Mastermind, Dallas, Texas, junija 1996

16.2 Kakšno naključje!!! Mar res?

Na nogometni tekmi sta na igrišču dve enajsterici in sodnik, skupaj **23 oseb.**

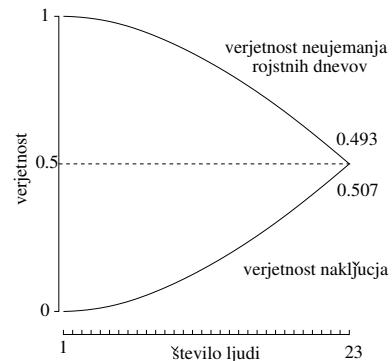
Kakšna je verjetnost, da imata **dve osebi** isti rojstni dan?

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:

$$\frac{365}{365} \times \frac{364}{365} \times \frac{363}{365} \times \cdots \times \frac{365 - k + 1}{365} = \begin{cases} 0.493; & \text{če je } k = 22 \\ 0.507; & \text{če je } k = 23 \end{cases}$$



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 paradoksom 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

$$e^{-x} = 1 - x + \frac{x^2}{2!} - \frac{x^3}{3!} + \dots$$

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 $1 - e^{-k(k-1)/(2n)}$. Potem velja $e^{-k(k-1)/(2n)} \approx 1 - \varepsilon$ oziroma $-k(k-1)/(2n) \approx \log(1 - \varepsilon)$, tj. $k^2 - k \approx -2n \log(1 - \varepsilon)$ in če ignoriramo $-k$, dobimo končno

$$k \doteq \sqrt{2n \log \frac{1}{1 - \varepsilon}}.$$

Za $\varepsilon = 0,5$ je

$$k \doteq 1.17\sqrt{n},$$

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 paradoksom rojstnih dnevov s tem določi spodnjo mejo za velikost zaloge vrednosti zgoščevalnih funkcij, ki jih uporabljam 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.

16.3 Ramseyjeva teorija

- intuitivna ideja
- Ramseyjev izrek
- Erdősev 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 prostop: 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.

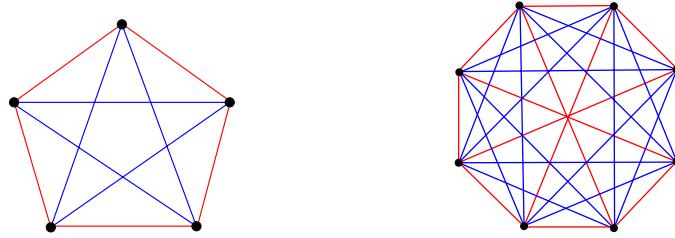
kar pomeni, da $G[T]$ vsebuje ali $(k - 1)$ -kliko ali ℓ -antikliko.

Od tod sledi, da G vsebuje bodisi k -kliko bodisi ℓ -antikliko. Naj bosta $r(k, \ell - 1)$ in $r(k - 1, \ell)$ sudi 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 $35 < r(3, 9) = 36$.

Erdős 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 nalogu:

Med petimi točkami v ravnini, od katerih nobene tri niso kolinearne (ležijo na premici), lahko vedno izberemo štiri, ki določajo konveksen četverokotnik.

Rešitev: Vpeljemo pojem **konveksne ogrinjače** ...

Če je konveksna ogrinjača teh petih točk

- (a) **petkotnik**, potem vsake 4 točke med njimi sestavljajo konveksen četverokotnik,
- (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 četverokotnik $ABCD$ konveksen. □

Nalogo lahko posplošimo na 9 točk in iskanje konveksnega petkotnika ter počasi pridemo do Erdőseve 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 Erdős postavil svojo domnevo.

Erdőseva probabilistična metoda (1947)

34 točk določa 561 premic. Da se to zgodi v eni barvi, je verjetnost

$$2^{-561} \doteq 2 \cdot 6 \times 10^{-169}.$$

Velja tudi $\binom{1.000.000}{34} = 3 \cdot 4 \times 10^{165}$. Torej lahko pričakujemo $\binom{10^6}{34} = 3 \cdot 4 \times 10^{165} \doteq 0 \cdot 01$ 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 $\mathcal{G}_n = \{G \mid \overline{G} \in \mathcal{G}_n\}$ 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.

$$\text{PODVOJ}(x) = 2x, \quad \text{EKSPONENT}(x) = 2^x, \quad \text{STOLP}(x) = 2^{2^{2^{\dots^2}}} \quad (x \text{ dvojk})$$

$$\text{UAU}(1) = \text{STOLP}(1) = 2. \quad \text{UAU}(2) = \text{STOLP}(2) = 4. \quad \text{UAU}(3) = \text{STOLP}(4) = 65,536$$

$$\text{UAU}(4) = \text{prevelik za vse knjige, za vse računalnike} \dots, \quad \text{UAU}(x) = \dots$$

Zaporedje 1, 2, ..., $\text{ACKERMANN}(k)$ pobarvamo z dvema barvama.

Potem obstaja monokromatično (enobarvno) aritmetično podzaporedje s k členi.

16.4 Teorija kodiranja

Claude Shannon je postavil teoretične osnove **teorije informacij** in zanesljivega prenosa digitalnih podatkov kmalu po koncu druge svetovne vojne.



Glavni mejnički 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).

Poglavlje 17

Uporaba statistike

(Testiranje PRNG, Teorija informacij in entropija)

17.1 Načrtovanje eksperimentov

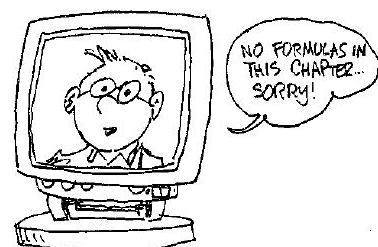


Statistično raziskovanje po Šadl [6]:

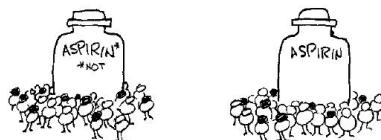
1. **načrtovanje statističnega raziskovanja** (vsebina, org./teh. vpr.),
2. **statistično opazovanje** (namen, populacija, spremenljivka),
3. **obdelava in urejanje podatkov** (časovne, krajevne in stvarne vrste),
4. **prikazovanje podatkov** (enostavne ter 2-razsežne tabele in grafikoni),
5. **analiza podatkov proučevanega pojava** (vse to smo počeli do sedaj).

Za podrobnejšo razlago priporočam pravkar omenjeno referenco.

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 programske opreme.



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:

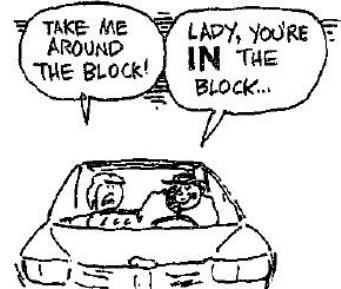
- optimizacija industrijskih procesov,
- medicina,
- sociologija.



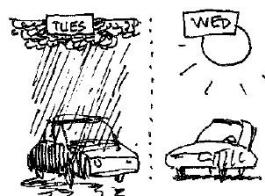
Na primeru bomo predstavili tri osnovne principe načrtovanja eksperimentov:

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 taxi 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.

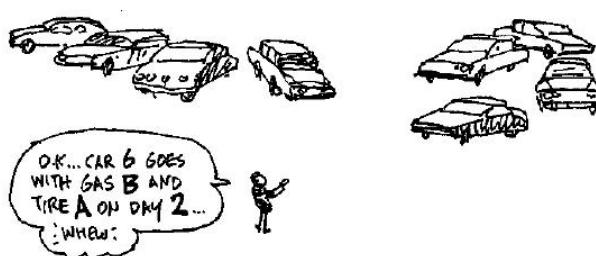


Latinski kvadrati

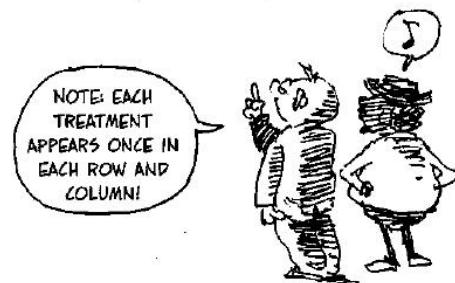
Latinski kvadrat reda v je $v \times v$ -razsežna matrika, v kateri vsi simboli iz množice

$$\{1, \dots, v\}$$

nastopajo v vsaki vrstici in vsakem stolpcu.



		DAY			
		1	2	3	4
CAB	1	a	b	c	d
		b	c	d	a
2	3	c	d	a	b
3	4	d	a	b	c



Trije paroma ortogonalni latinski kvadrati reda 4, tj. vsak par znak-črka ali črka-barva ali barva-znak se pojavi natanko enkrat.

Projektivni prostor $PG(d, q)$ (razsešnosti d nad q) dobimo iz vektorskega prostora $[GF(q)]^{d+1}$, tako da naredimo kvocient po 1-razsežnih podprostорih.

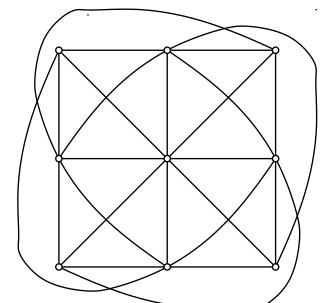
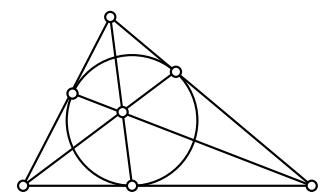
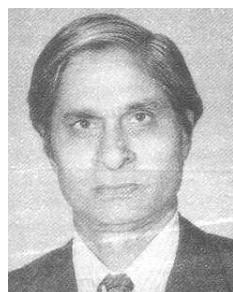
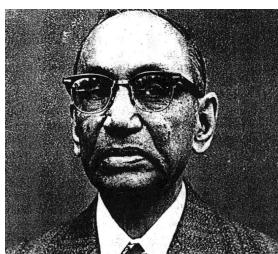
Projektivna ravnina $PG(2, q)$ je incidenčna struktura z 1- in 2-dim. podprostori prostora $[GF(q)]^3$ kot **točkami** in **premicami**, kjer je “ \subset ” incidenčna relacija za katero velja:

- $v = q^2 + q + 1$ je število točk (in število premic b),
- vsaka premlica ima $k = q + 1$ točk (in skozi vsako točko gre $r = q + 1$ premic),
- vsak par točk leži na $\lambda = 1$ premicah (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.

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).
- [6] M. Šadl, *Statistika* (Srednješolski program, ekonomski tehnik), Mohorjeva Hermagoras, 2004:.

Obstaja obilna literatura na spletu in v knjižnicah.

Gradiva bodo dosegljiva preko internetne učilnice (moodle).

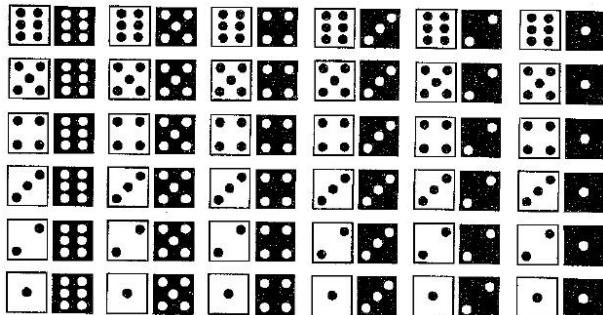
Pri delu z dejanskimi podatki se bomo v glavnem naslonili na prosti statistični program R. Program je prosto dostopen na: <http://www.r-project.org/>
Proti koncu semestra pa morda tudi Minitab.

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 Računalna nove dobe

Ste že kdaj razmišljali o računanju (aritmetiki), ki ga uporabljam 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 z osnovnošolskim računanjem približati računalom, ki jih preko številnih naprav, kot so osebni računalniki, diskmani in pametne kartice, uporabljam 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 odštevati. 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 *nevtralni 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 $b + c$.

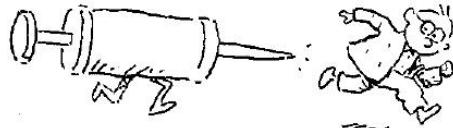
Primer: Za cilj postavimo iskanje obsega s končno mnogo elementi, v katerem bo računanje v nekem smislu še udobnejše kot v obsegih, ki jih srečamo v osnovni ali srednji šoli (racionalna števila \mathbb{Q} , realna števila \mathbb{R} ali celo kompleksna števila \mathbb{C}).

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 množikativno grupo: $1 * 1 = 1$ zadovolji aksiom (O2). Ne pozabite, da operaciji $+$ in $*$ ne predstavlja (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

Naj bo n poljubno naravno število. **Permutacija** elementov $1, \dots, n$ je bijekcija, ki slika iz množice $\{1, \dots, n\}$ v množico $\{1, \dots, n\}$. Lahko tudi rečemo, da gre za razvrstitev n -tih različnih elementov v *vrsto* (mesto elementa i v vrsti je ravno slika tega elementa).

Npr. permutacija kart je običajno premešanje kart – ko jih vrenem na kup (spremeni se vrstni red v kupu, karte pa ostanejo iste – nobene nismo ne dodali ne odvzeli). Permutacijo torej predstavimo z (novim) vrstnim redom kart.

Število permutacij n elementov je enako n -faktorijel (oz. n -fakulteta), tj.

$$n! := 1 \cdot 2 \cdot \dots \cdot n$$

(ozziroma 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 končna množica. Permutacije množice A med seboj *množimo* (oz. komponiramo) po naslednjem pravilu: produkt $\pi = \pi_2 \circ \pi_1$ je permutacija množice A , ki preslika vsak element $a \in A$ v $\pi_2(\pi_1(a))$. Ni se težko prepričati, da *permutacije množice A tvorijo grupo* (premisli kaj je enota za to grupo in katera permutacija je inverzna poljubni permutaciji).

Primer: Za permutaciji

$$\pi_1 = \begin{pmatrix} 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 & 10 & 11 \\ 3 & 4 & 5 & 10 & 2 & 1 & 7 & 9 & 11 & 6 & 8 \end{pmatrix}$$

in

$$\pi_2 = \begin{pmatrix} 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 & 10 & 11 \\ 8 & 2 & 1 & 3 & 10 & 9 & 4 & 5 & 7 & 6 & 11 \end{pmatrix}$$

je njun produkt enak

$$\pi = \pi_2 \circ \pi_1 = \begin{pmatrix} 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 & 10 & 11 \\ 1 & 3 & 10 & 6 & 2 & 8 & 4 & 7 & 11 & 9 & 5 \end{pmatrix}. \quad \diamond$$

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 (ozziroma tujih) ciklov.* \square

Primer: Za permutacije π_1, π_2 in π iz prejšnjega primera 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:

$$(a_1 a_2 a_3 \dots a_r) = (a_{r-1} a_r) \circ \dots \circ (a_2 a_3) \circ (a_1 a_2),$$

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 poglejmo v katerih razlikah se pojavit:

$x_1 - x_a, \dots, x_{a-1} - x_a,$		$x_a - x_{a+1}, \dots, x_a - x_{b-1},$	$x_a - x_b,$	$x_a - x_{b+1}, \dots, x_a - x_n,$
$x_1 - x_b, \dots, x_{a-1} - x_b,$	$x_a - x_b,$	$x_{a+1} - x_b, \dots, x_{b-1} - x_b,$		$x_b - x_{b+1}, \dots, x_b - x_n.$

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 predznače (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. Več o uporabnosti permutacij ter parnosti (npr. pri igri *Petnajst*) pa si lahko preberete v <http://www.presek.si/16/930-Klavzar.pdf>. Ne pozabite pa tudi na **Rubikovo kocko**, ki je prav tako povezana s permutacijami, glej npr. http://sl.wikipedia.org/wiki/Rubikova_kocka.

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: $P_8^{2,2,2,2} = 8 \times 7 \times 6 \times 5 \times 4 \times 3 \times 2/2^4 = 7 \times 6 \times 5 \times 4 \times 3 = 42 \times 60 = 2520$. 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 sošolki 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 $3^5 = 3^2 \times 3^2 \times 3 = 9 \times 9 \times 3 = 81 \times 3 = 243$), 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 nobeno dvojko, 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 $S(5, 3) = 3(3^4 - 2^5 + 2 - 1) = 3(81 - 32 + 2 - 1) = 3 \times 50 = 150$.

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	$3^4 - 3(2^4 - 2) - 3 = 36$	$= 3 \times 4!/2!$	$= 3 \cdot P_4^{112}$
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	$3^7 - 3(2^7 - 2) - 3 = 1806$	$= \frac{3 \times 7!}{5!} + \frac{6 \times 7!}{2!4!} + \frac{3 \times 7!}{3!3!} + \frac{3 \times 7!}{2!2!3!}$	$= 3P_7^{115} + 6P_7^{124} + 3P_7^{133} + 3P_7^{223}$

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 $P_3^3 = 1$.

00000*	01000x	02000y	10000x	11000x	12000	20000y	21000	22000y
00001x	01001x	02001	10001x	11001x	12001	20001	21001	22001
00002y	01002	02002y	10002	11002	12002	20002y	21002	22002y
00010x	01010x	02010	10010x	11010x	12010	20010	21010	22010
00011x	01011x	02011	10011x	11011x	12011	20011	21011	22011
00012y	01012	02012	10012	11012	12012	20012	21012	22012
00020	01020	02020y	10020	11020	12020	20020y	21020	22020y
00021	01021	02021	10021	11021	12021	20021	21021	22021
00022y	01022	02022y	10022	11022	12022	20022y	21022	22022y
00100x	01100x	02100	10100x	11100x	12100	20100	21100	22100
00101x	01101x	02101	10101x	11101x	12101	20101	21101	22101
00102	01102	02102	10102	11102	12102	20102	21102	22102
00110x	01110x	02110	10110x	11110x	12110	20110	21110	22110
00111x	01111x	02111	10111x	11111*	12111z	20111	21111z	22111z
00112	01112	02112	10112	11112z	12112z	20112	21112z	22112z
00120	01120	02120	10120	11120	12120	20120	21120	22120
00121	01121	02121	10121	11121z	12121z	20121	21121z	22121z
00122	01122	02122	10122	11122z	12122z	20122	21122z	22122z
00200y	01200	02200y	10200	11200	12200	20200y	21200	22200y
00201	01201	02201	10201	11201	12201	20201	21201	22201
00202y	01202	02202y	10202	11202	12202	20202y	21202	22202y
00210	01210	02210	10210	11210	12210	20210	21210	22210
00211	01211	02211	10211	11211z	12211z	20211	21211z	22211z
00212	01212	02212	10212	11212z	12212z	20212	21212z	22212z
00220y	01220	02220y	10220	11220	12220	20220y	21220	22220y
00221	01221	02221	10221	11221z	12221z	20221	21221z	22221z
00222y	01222	02222y	10222	11222z	12222z	20222y	21222z	22222z
7x 7y	8x	8y	8x	7x 7z	8z	8y	8z	7y 7z
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 šesteric označenih z z ni potrebno več odštevati, dve šesterici, ki sta označeni z zvezdico pa bi morali označiti z x 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 $S(n, 1)$:

$$S(n, 4) = 4^n - 4 \left(3^n - 3(2^n - 2) - 3 \right) - 6(2^n - 2) - 4 = 4^n - 4 \cdot 3^n + 6 \cdot 2^n - 4.$$

Bralcu prepuščamo, da pravkar dobljeno formulo testira (npr. bodisi s permutacijami s ponavljanjem ali pa kar štetjem izračuna $S(5, 4)$ (glej sp. 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 števkami 0, 1, 2, 3.

A.4 Kombinacije

Binomski simbol oz. število **kombinacij**, tj. število m -elementnih podmnožic množice moči n , za $0 \leq m \leq n$, je enako

$$\binom{n}{m} = \frac{n \cdot (n-1) \cdots (n-m+1)}{1 \cdot 2 \cdots m} = \frac{n!}{m!(n-m)!},$$

saj lahko prvi element izberemo na n načinov, drugi na $n-1$ načinov, ..., zadnji na $n-m+1$ načinov, ker pa vrstni red izbranih elementov ni pomemben, dobljeno število še delimo s številom permutacij (dobljeno število se seveda ujama s številom permutacij s ponavljanjem z dvema skupinama, tj. $P_n^{k,n-k}$, glej prejšnji razdelek).

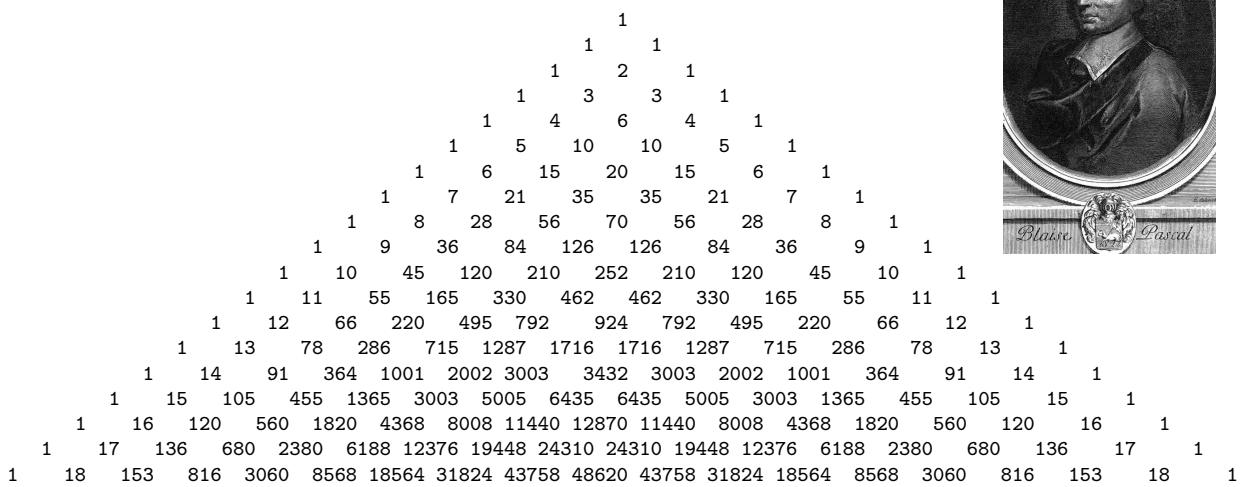
Trditev A.4. Za binomske simbole velja $\binom{n}{0} = 1$, $\binom{n}{m+1} = \binom{n}{m} \frac{n-m}{m+1}$ za $1 \leq m+1 \leq n$, $\binom{n}{m} = \binom{n}{n-m}$ za $0 \leq m \leq n$,

$$(a+b)^n = \sum_{k=0}^n \binom{n}{k} a^{n-k} b^k \quad \text{in} \quad \binom{n}{m} + \binom{n}{m+1} = \binom{n+1}{m+1}.$$

Zgornji enakosti na levi pravimo **binomski obrazec**, na desni pa **pravilo Pascalovega trikotnika** (glej sliko 1).

Dokaz. Prve tri relacije sta očitni že iz same definicije binomskega simbola, četrta pa sledi neposredno iz dejstva, da binomski simbol enak številu vseh kombinacij ustreznega reda. Po definiciji binomskega simbola in množenju obeh strani z $(m+1)!(n-m)!/n!$ je zadnja relacija ekvivalentna $m+1+n-m=n+1$, kar očitno drži. \square

Pascalov trikotnik



Slika 1: To so v obliki trikotnika zapisani binomski simboli, vsaka vrstica pa ustreza enemu binoskemu obrazcu.

Primer: Na kvizu je danih 6 vprašanj. Pri vsakem so možni 4 odgovori (npr. (a), (b), (c) in (d)), od katerih je natanko en pravilen. **Kakšna je verjetnost, da bomo na vsaj polovico vprašanj odgovorili pravilno, če odgovore izbiramo naključno?** Vseh možnih odgovorov je $4^6 = 2^{12} = 4096$, kar je vsekakor preveč, da bi jih vse držali v glavi, tudi na papir bi jih bilo zoprno pisati, s kopiraj in prilepi pa že ni več tako naporno. Privzameš lahko, da je pravilna rešitev kar aaaaaa. Ni pa nič narobe, če nalogo najprej nekoliko poenostavimo in si najprej pogledamo kviz z dvemi, tremi ali celo štirimi vprašanji:

aa DA	ba DA	ca DA	da DA		aaaa DA	acaa DA	baaa DA	bcaa DA	caaa DA	ccaa DA	daaa DA	dcaa DA	
ab DA	bb	cb	db		aaab DA	acab DA	baab DA	bcab	caab DA	ccab	daab DA	dcab	
ac DA	bc	cc	dc		aaac DA	acac DA	baac DA	bcac	caac DA	ccac	daac DA	dcac	
ad DA	bd	cd	dd		aaad DA	acad DA	baad DA	bcad	caad DA	ccad	daad DA	dcad	
$n=2$					aaba DA	acba DA	baba DA	bcba	caba DA	ccba	daba DA	dcba	
					aabb DA	acbb	babb	bcbb	cabb	ccbb	dabb	dcbb	
					aabc DA	acbc	babc	bcbc	cabc	ccbc	dabc	dcbc	
					aabd DA	acbd	babd	bcbd	cabd	ccbd	dabd	dcbd	
aaa DA	baa DA	caa DA	daa DA		aaca DA	acca DA	baca DA	bcca	caca DA	ccca	daca DA	dcca	
aab DA	bab	cab	dab		aacb DA	accb	bacb	bccb	cacb	cccb	dacb	dccb	
aac DA	bac	cac	dac		aacc DA	accc	bacc	bccc	cacc	cccc	dacc	dccc	
aad DA	bad	cad	dad		aacd DA	accd	bacd	bcdd	cacd	cccd	dacd	dcdd	
aba DA	bba	cba	dba		aada DA	acda DA	bada DA	bcda	cada DA	ccda	dada DA	dcca	
abb	bbb	cbb	dbb		aadb DA	acdb	badb	bcdb	cadb	ccdb	dadb	dcdb	
abc	bcb	cbc	dbc		aadc DA	acdc	badc	bcdc	cadc	ccdc	dadc	dcdc	
abd	bbd	cbd	dbd		aadd DA	acdd	badd	bcdd	cadd	ccdd	dadd	dcdd	
aca DA	bca	cca	dca		abaa DA	adaa DA	bbaa DA	bdaa DA	cbaa DA	cdaa DA	dbaa DA	ddaa DA	
acb	bcb	ccb	dcb		abab DA	adab DA	bbab	bdab	cbab	cdab	dbab	ddab	
acc	bcc	ccc	dcc		abac DA	adac DA	bbac	bdac	cbac	cdac	dbac	ddac	
acd	bcd	ccd	dcd		abad DA	adad DA	bbad	bdad	cbad	cdad	dbad	ddad	
ada DA	bda	cda	dda		abba DA	adba DA	bbba	bdba	cbba	cdba	dbba	ddba	
adb	bdb	cdb	ddb		abbb DA	adbb	bbbb	bdbb	cbbb	cdbb	dbbb	ddbb	
adc	bdc	cdc	ddc		abbc DA	adbc	bbbc	bdbc	cbbc	cdbc	dbbc	ddbc	
add	bdd	cdd	ddd		abbd DA	adbd	bbbb	bdbd	cbbd	cdbd	dbbd	ddbd	
$n=3$					abca DA	adca DA	bbca	bdca	cbca	cdca	dbca	ddca	
					abcb DA	adcb	bbcb	bdcb	cbcb	cdcb	dbc	ddcb	
					abcc DA	adcc	bbcc	bdcc	cbcc	cdcc	dbcc	ddcc	
					abcd DA	adcd	bbcd	bdcd	cbcd	cdcd	dcad	ddcd	
					abda DA	adda DA	bbda	bdda	cbda	cdda	dbda	ddda	
					abdb DA	addb	bbdb	bddb	cbdb	cddb	dbdb	dddb	
					abdc DA	addc	bbdc	bddc	cbdc	cddc	dbdc	dddc	
					abdd DA	addd	bbdd	bddd	cbdd	cddd	dbdd	dddd	
$n=4$													

in dobimo zaporedoma $7/16$, $10/64$ in $67/256$. Študent računalništva zna seveda napisati kratek programček, ki opravi delo namesto njega. Nalogo lahko tudi malo posplošimo. Namesto 6ih vprašanj bi lahko vzeli, da jih kviz vsebuje n , kjer je $n \in \mathbb{N}$. Tudi število možnih odgovorov je lahko namesto 4 kar $m \in \mathbb{N}$, $m \geq 2$. Vpeljimo še en parameter, označimo ga s k , in si postavimo dodatno vprašanje. **Kakšna je verjetnost, da smo s poskušanjem odgovorili pravilno na natanko k vprašanj, pri čemer je $0 \leq k \leq n$?** Če to verjetnost označimo s $P_{n,m,k}$, potem je iskana verjetnost enaka $P_{6,4,3} + P_{6,4,4} + P_{6,4,5} + P_{6,4,6}$ oziroma $1 - P_{6,4,0} - P_{6,4,1} - P_{6,4,2}$. Z uporabo binomskih koeficientov dobimo

$$P_{n,m,0} = \frac{(m-1)^n}{m^n}, \quad P_{n,m,1} = \frac{n(m-1)^{n-1}}{m^n}, \dots, \quad P_{n,m,k} = \binom{n}{k} \frac{(m-1)^{n-k}}{m^n}, \dots,$$

glede na to, da najprej med n nalogami izberemo k tistih s pravilnimi odgovori, ostale naloge pa nimajo pravilnih odgovorov, zato imamo pri vsakem odgovoru le $m-1$ možnosti. Iskana

verjetnost je torej enaka

$$1 - \frac{1}{m^n} \sum_{i=0}^{\lfloor (n-1)/2 \rfloor} \binom{n}{i} (m-1)^{n-i} \quad (\text{A.1})$$

(če bi seštevali po vseh $i \in \{0, 1, \dots, n\}$ bi po binomskem obrazcu dobili 0) oziroma za

- $n = 2$ in $m = 4$: $(4^2 - 3^2)/4^2 = 7/16 = 0.438$,
- $n = 3$ in $m = 4$: $(4^3 - 3^3 - 3 \times 3^2)/4^3 = 10/64 = 0.156$,
- $n = 4$ in $m = 4$: $(4^4 - 3^4 - 4 \times 3^3)/4^4 = 67/256 = 0.262$,
- $n = 5$ in $m = 4$: $(4^5 - 3^5 - 5 \times 3^4 - 10 \times 3^3)/4^5 = 106/1024 = 0.104$, in končno
- $n = 6$ in $m = 4$: $(4^6 - 3^6 - 6 \times 3^5 - 15 \times 3^4)/4^6 = 694/4096 = 0.169$.

Rešitvi (A.1) nekoliko bolj zaupamo, ker se odgovori za $n = 2, 3, 4$ ujemajo z verjetnostmi, ki smo jih dobili s preštevanjem. Zaupanje pa bi se še povečalo, če bi preverili še $n = 6$ in $m = 3$ ali vsaj $n = 6$ in $m = 2$. \diamond

Primer: Na polico bi radi postavili 4 matematične, 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 skupaj¹,
- (c) če morajo matematične knjige stati na začetku.

\diamond

Primer: Na koliko načinov lahko sestavimo iz 7ih soglasnikov in 5ih samoglasnikov besedo, ki ima 4 soglasnike in 3 samoglasnike?² \diamond

¹Namig: če želimo ponagajati prijateljici, dajmo tisti dve kemijski knjigi namenoma skupaj.

² Namig: nalogo lahko rešimo tako s kombinacijami in permutacijami, kot tudi z variacijami in permutacijami s ponavljanjem.

A.5 Vrsta za e

Število e , ki predstavlja osnovo za naravni logaritem, pogosto definiramo s formulo

$$e = \lim_{n \rightarrow \infty} \left(1 + \frac{1}{n}\right)^n. \quad (\text{A.2})$$

Leta 1683 je Jacob Bernoulli poskušal izračunati limito $(1 + 1/n)^n$, ko gre n proti neskončno, do katere je prišel pri študiju obrestnih obresti. Na računu pričnemo z 100€ in 100 procenčnimi letnimi obrestmi. Če se obresti izplačajo enkrat letno, bo vrednost na računu ob koncu leta znašala 200€. Kaj se zgodi, če se obresti izplačajo večkrat v letu?

Uporabil je binomski obrazec:

$$\left(1 + \frac{1}{n}\right)^n = 1 + \binom{n}{1} \frac{1}{n} + \binom{n}{2} \frac{1}{n^2} + \binom{n}{3} \frac{1}{n^3} + \cdots + \binom{n}{n} \frac{1}{n^n}.$$

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.

$$e = 1 + \frac{1}{1!} + \frac{1}{2!} + \frac{1}{3!} + \cdots. \quad (\text{A.3})$$

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 logaritmom³. 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.3), hkrati pa izračunal prvih 18 decimalk števila e : $e \doteq 2.718281828459045235$

Še splošnejšo vrsto pa dobimo z uporabo Taylorjeve vrste:

Trditev A.5.

$$e^x = 1 + \frac{x}{1!} + \frac{x^2}{2!} + \frac{x^3}{3!} + \cdots.$$

Za konec omenimo še dobro znano dejstvo, da sta eksponentna funkcija $f(x) = e^x$ in logaritemska funkcija $g(x) = \ln x$ inverzni, tj. $f(g(x)) = x = g(f(x))$ za vsak $x \in \mathbb{R}^+$. To pomeni, da sta njuna grafa simetrična glede na simetralo lihih kvadrantov, tj. premico $y = x$.

³Glej <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):

$$n! \approx \sqrt{2\pi n} \left(\frac{n}{e}\right)^n, \quad \text{ko gre } n \rightarrow \infty. \quad (\text{A.4})$$

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 naročnemu številu k , ki je zelo blizu $\log n$, pri čemer za logaritemsko osnovo vzamemo številsko osnovo, v kateri zapišemo število n). V čem smo torej na boljšem, ko računamo izraz na desni strani (A.4)? Računanje potence lahko izvedemo tako, da najprej izračunamo naslednje potence $(n/e)^{2^0}, (n/e)^{2^1}, (n/e)^{2^2}, \dots, (n/e)^{2^k}$ z zaporednim kvadriranjem, nato pa zmnožimo med seboj tiste potence katerih eksponenti ustrezajo mestom enic v binarni predstavitevi števila n , kar pomeni da smo opravili največ $2k$ množenj.

Primer: Namesto, da bi izračunali $P = a^{21}$ z dvajsetimi množenji ($P = a$, dvajsetkrat ponavlja $P := P * a$), 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. ◇

Formulo (A.4) je prvi odkril Abraham de Moivre v naslednji obliki

$$n! \approx [\text{konstanta}] \cdot n^{n+1/2} e^{-n}, \quad \text{ko gre } n \rightarrow \infty,$$

pri čemer je konstanto izrazil s hyperboličnim logaritmom. James Stirlingov prispevek pa je bil, da je konstanta v resnici enaka $\sqrt{2\pi}$. Formulo tipično uporabljamo v aplikacijah v obliki

$$\ln n! \approx n \ln n - n, \quad \text{ko gre } n \rightarrow \infty.$$

V zgornji verziji manjka faktor $\frac{1}{2} \ln(2\pi n)$, ki ga lahko za velike n zanemarimo v primerjavi z drugimi sumandi. 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 Bernoullijskih števil. Glej npr. http://en.wikipedia.org/wiki/Stirling's_approximation.

A.7 Normalna krivulja v prostoru

Če normalno krivuljo zarotiramo okoli osi simetrije ($x = \mu$) dobimo zvonasto ploskev. Želimo pokazati naslednjo identiteto.⁴

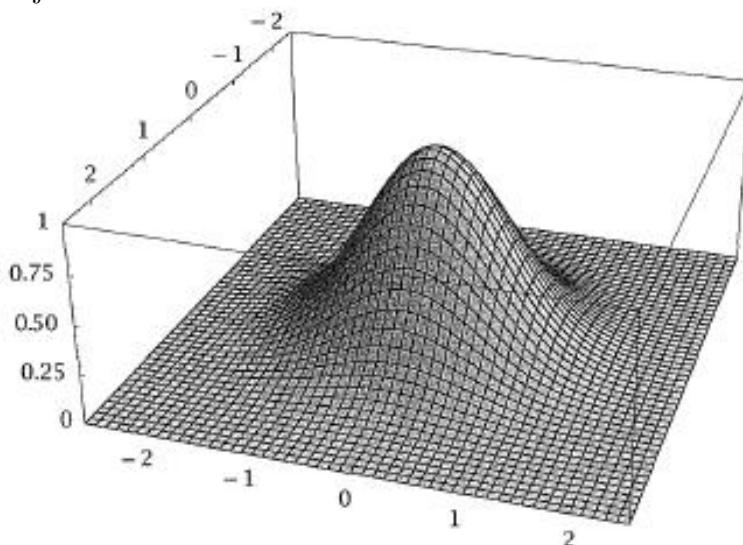
Izrek A.6.

$$\int_{-\infty}^{\infty} e^{-t^2} dt = \sqrt{\pi}. \quad (\text{A.5})$$

Dokaz. Označimo vrednost integrala na levi strani (A.5) 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} :

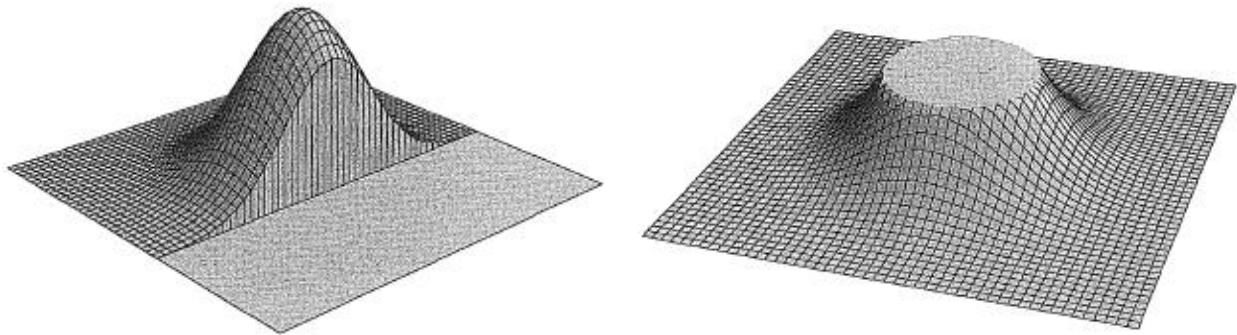
$$\int_{-\infty}^{\infty} e^{-1} e^{-t^2} dt = e^{-1} I.$$

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

$$V = \int_{-\infty}^{\infty} e^{-s^2} I ds = I^2.$$

Preostane nam le še, da dokažemo, da je $V = \pi$. Tokrat presekajmo normalno ploskev z ravnino $z = h$.

⁴Stara zgodba pravi, da je Lord Kelvin nekoč dejal, da je matematik nekdo, za katerega je ta identiteta očitna.



Slika: Vertikalni in horizontalni prerez normalne ploskve.

Potem za točke preseka velja

$$e^{-s^2-t^2} \geq h \quad \text{ozziroma} \quad s^2 + t^2 \leq -\ln h.$$

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

$$\int_0^1 \ln x dx = x \ln x \Big|_0^1 - \int_0^1 \frac{x}{x} dx = -1$$

in od tod $V = \pi$. □

A.8 Sredine nenegativnih števil a_1, \dots, a_n



Aritmetična:

$$A_n = \frac{a_1 + \dots + a_n}{n}$$

Geometrična:

$$G_n = \sqrt[n]{a_1 \cdot \dots \cdot a_n}$$

Harmonična:

$$H_n = \frac{n}{\frac{1}{a_1} + \dots + \frac{1}{a_n}}$$

Kvadratna:

$$K_n = \sqrt{\frac{a_1^2 + \dots + a_n^2}{n}}$$

Sredine dveh števil: $H_2 \leq G_2 \leq A_2 \leq K_2$

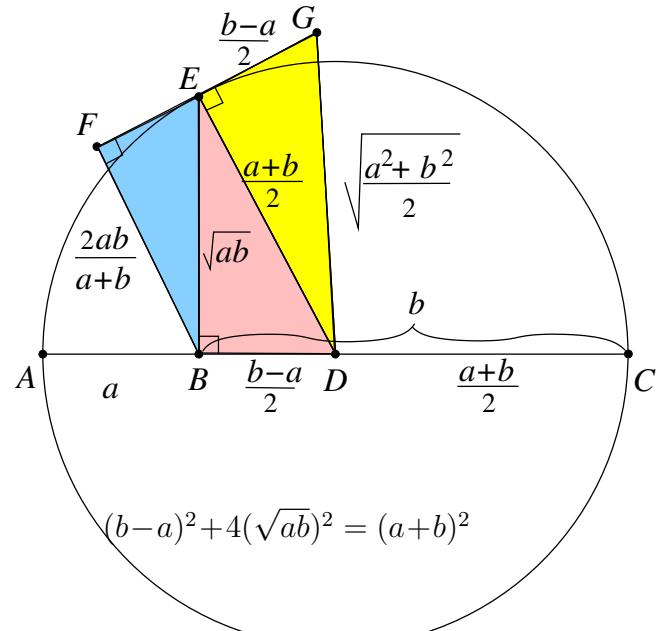
$$a, b \geq 0$$

$$H_2 = \frac{2}{\frac{1}{a} + \frac{1}{b}}$$

$$G_2 = \sqrt{ab}$$

$$A_2 = \frac{a+b}{2}$$

$$K_2 = \sqrt{\frac{a^2 + b^2}{2}}$$



Sidney H. Kung

(iz R.B. Nelsenove knjige "Dokazi brez besed")

Potenčna (stopnje k):

$$P_{n,k} = \sqrt[k]{\frac{a_1^k + \dots + a_n^k}{n}}$$

Velja:

$$H_n = P_{n,-1}, \quad G_n = \lim_{k \rightarrow 0} P_{n,k} \quad A_n = P_{n,1} \quad \text{in} \quad K_n = P_{n,2}$$

ter

$$H_n \leq G_n \leq A_n \leq K_n \quad \text{ozziroma za } k \leq m \quad P_{n,k} \leq P_{n,m}$$

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 $\text{proj}_{\vec{u}} \vec{v}$ 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 kot 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.

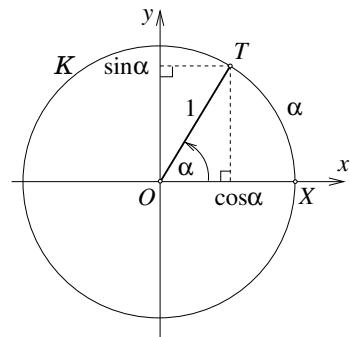
Pravimo ji tudi Cauchy-Schwarzova neenakost ali neenakost Bunjakovskega, glej http://en.wikipedia.org/wiki/Cauchy-Schwarz_inequality. To je ena izmed najpomembnejših neenakosti. Omogoča nam, da preverimo kolinearnost dveh vektorjev tako, da izračunamo vrednosti na levi in desni strani zgornje neenakosti 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. $|z| \cdot |w| = |zw|$.

Za $n = 3$ lahko dokažemo Cauchyjevo neenakost uporavbo 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



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 "prehoden" lok dolžine α (takrat je kot $\angle XOT$ enak α radijanov), ima točka T koordinati $(\cos \alpha, \sin \alpha)$. Funkcija sinus je pozitivna v prvem in drugem kvadrantu, funkcija cosinus pa v prvem in četrtem. Obe funkciji sta periodični s perido 2π (tj. 360°).

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 vektor-skega produkta enaka ploščini paralelograma, ki ga razpenjata vektorja \vec{u} in \vec{v} . **Lagrangeova identiteta** $|\vec{u}|^2 |\vec{v}|^2 = (\vec{u} \cdot \vec{v})^2 + |\vec{u} \times \vec{v}|^2$, ni nič drugega kot na drugačen način zapisana relacija $\sin^2 \varphi + \cos^2 \varphi = 1$ (ozioroma Pitagorjev izrek).

Za splošen $n \in \mathbb{N}$ lahko Cauchyjevo neenakost zapišemo tudi v naslednji obliki:

$$(a_1^2 + \cdots + a_n^2)(b_1^2 + \cdots + b_n^2) \geq (a_1 b_1 + \cdots + a_n b_n)^2,$$

kjer so $a_1, \dots, a_n, b_1, \dots, b_n$ poljubna realna števila. Lagrangeova identiteta za splošen n pa izgleda takole:

$$\left(\sum_{i=1}^n a_i^2 \right) \left(\sum_{i=1}^n b_i^2 \right) = \left(\sum_{i=1}^n a_i b_i \right)^2 + \sum_{i < j} (a_i b_j - a_j b_i)^2.$$

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 pričakovana vrednost njunega produkta skalarni produkt, tj.

$$\langle X, Y \rangle := \mathbb{E}(XY)$$

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

$$|\mathbb{E}(XY)|^2 \leq \mathbb{E}(X^2) \mathbb{E}(Y^2).$$

Naj bo $\mu = E(X)$ in $\nu = E(Y)$. Potem po Cauchyjevi neenakosti velja

$$\begin{aligned} |\text{Cov}(X, Y)|^2 &= |\mathbb{E}((X - \mu)(Y - \nu))|^2 = |\langle X - \mu, Y - \nu \rangle|^2 \\ &\leq \langle X - \mu, X - \mu \rangle \langle Y - \nu, Y - \nu \rangle = \mathbb{E}((X - \mu)^2) \mathbb{E}((Y - \nu)^2) = \text{D}(X) \text{D}(Y), \end{aligned}$$

kjer je D disperzija, Cov pa kovarianca.

Formalno je **vektorski prostor s skalarnim produkтом** (rečemo tudi *unitarni vektorski prostor*) vektorski prostor V nad poljubnim obsegom \mathbb{F} s skalarnim produkтом, 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*: $\langle x, y \rangle = \overline{\langle y, x \rangle}$.
- *Linearost na prvi koordinati*: $\langle ax, y \rangle = a\langle x, y \rangle$.
- *Pozitivna-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 produkтом.

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

$$0 \leq \langle u, u \rangle - |\langle u, v \rangle|^2 \cdot \langle v, v \rangle^{-1} \quad \text{ozziroma} \quad |\langle u, v \rangle|^2 \leq \langle u, u \rangle \cdot \langle v, v \rangle,$$

in končno po korenjenju neenakost, ki smo jo žeeli pokazati. \square

Trikotniško neenakost za vektorske prostore s skalarnim produkтом 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.

Dodatek B

Vadnica

Stari pregovor pravi: “[Brez muje se še čevelj ne obuje.](#)” Ni kaj, snov je najbolje utrjevati z reševanjem nalog. Sprva ni lahko, potem pa lahko postane celo zabavno.

B.1 Vaje za uvod

Za tiste, ki ste že nekoliko pozabili snov iz srednje šole¹ priporočam ponavljanje. Možnosti je veliko, npr.

- Nives Mihelič Erbežnik et al., Priprave na maturo, matematika, 1. izd., 1. natis. – Ljubljana: DZS, 2001. (Zbirka nalog za srednje šole) 399, str.

Za ogrevanje smo si izposodili nekaj nalog iz verjetnosti:

1. Kaj je poskus in kaj slučajni dogodek?
 - (a) Iz kompleta 32 kart izberemo 3 karte.
 - (b) Izvlečena karta iz komleta 32 kart je srčni kralj.
2. Zapisa $A \cap B \subset A$ in $A \cap B \subset B$ govorita o:
 - (a) vsoti dogodkov,
 - (b) produktu dogodkov.
3. Produkt nezdružljivih dogodkov A in B je:
 - (a) nemogoč dogodek,
 - (b) gotov dogodek,
 - (c) sestavljen dogodek.
4. Poljuben dogodek in njegova negacija sta:
 - (a) nezdružljiva in nasprotna dogodka,
 - (b) združljiva in nasprotna dogodka,
 - (c) združljiva in neodvisna dogodka.
5. Za dogodek A , da padejo pri metu igralne kocke več kot 3 pike, je nasprotni dogodek (negacija):
 - (a) da ne pade več pik kot 3,
 - (b) da ne padejo kvečjemu štiri pike.
6. Pošteno igralno kocko vržemo enkrat. Popolni sistem elementarnih dogodkov tega poskusa so dogodki, da:
 - (a) padejo šestice $\{e_6\}$,
 - (b) pade katerokoli število pik $\{e_1, e_2, e_3, e_4, e_5, e_6\}$.

¹ Pogosto se zgodi, da snov kot je verjetnost (ali celo statistika) učitelji v srednjih šolah sploh ne obdelajo, večinoma pa jo le na hitro preletijo.

7. Vržemo igralno kocko. Sestavljeni dogodek je:
- pade šestica,
 - pade sodo število pik, manjše od 5,
 - pade liho število pik, večje od 4.
8. Klasična in statistična verjetnost se razlikujeta po tem, da je ocena za verjetnost dogodkov:
- pri prvi dobljena računsko, pri drugi empirično,
 - pri prvi dobljena empirično, pri drugi računsko
9. Elementarni dogodki, ki jih obravnavamo s klasično definicijo verjetnosti:
- morajo sestavljati popoln simetričen sistem dogodkov,
 - ne smejo sestavljati popolnega simetričnega sistema dogodkov.
10. V posodi so 3 rdeče kroglice in 2 modri. Na slepo izvlečemo kroglico, jo vrnemo v posodo, jih premešamo in ponovno na slepo eno izvlečemo. Izid poskusa je elementarni dogodek, ki ga določata barvi dveh kroglic. Vseh elementarnih dogodkov v tem poskusu je:
- V_5^2 ,
 - C_2^5 .
11. V posodi sta 2 beli kroglici in 5 rdečih. Na slepo izvlečemo iz nje eno kroglico. Verjetnost, da je izvlečena kroglica rdeča, je:
- 2/7,
 - 5/7.
12. Iz kupa 32 igralnih kart izvlečemo 3 karte. Verjetnost, da je med njimi vsaj ena, je:
- 0·34,
 - 0·66.
13. Dva moška in pet žensk naključno razporedimo v vrsto. Verjetnost, da bosta moška sedela skupaj na začetku ali na

koncu vrste, je:

$$(a) 2/7, \quad (b) 2/21.$$

14. Hkrati vržemo 3 poštene igralne kocke. Verjetnost, da pade ena šestica, hkrati pa vsaka kocka pokaže različno število pik, je:
- $20/6^3$,
 - $5/18$.
15. V veslaškem klubu so 3 krmari in 10 veslačev, od tega 5 članov in 5 mladincev. Trener mora v naglici izbrati posadko za četverec s krmarjem, zato jo sestavi kar na slepo. Verjetnost, da bo izbral najboljšega krmarja, najboljšega člana in najboljšega mladinka, je:
- 2/45,
 - 3/45.
16. Problem rešujeta neodvisno drug od drugega dva učenca. Učenec A rešuje probleme tako, da je verjetnost za posamezno rešitev 0·9, učenec B pa z verjetnostjo 0·6. Verjetnost, da bo problem rešen, je:
- 0·04,
 - 0·96.
17. Probleme rešujeta neodvisno drug od drugega dva učenca. Učenec A rešuje probleme tako, da je verjetnost za posamezno rešitev 0·9, učenec B pa z verjetnostjo 0·6. Problem naj reši le en učenec. Verjetnost, da bo problem rešil samo učenec A , je:
- 0·857,
 - 0·143.
18. Pri poskusu sodelujejo štiri osebe. Vsaka si povsem naključno in neodvisno od drugih izbere neko naravno število, manjše od 100. Verjetnost, da bo vsaj eno od teh štirih števil deljivo s 3, je
- 1/3,
 - 2/3,
 - 65/81.
19. Študent obvlada 8 izpitnih vprašanj od 12-ih. Na izpitu mora na slepo izbrati

4 vprašanja. Pozitivno oceno doseže, če pravilno odgovori vsaj na 2 vprašanji. Verjetnost, da bo študent izpit opravil, je:

- (a) $2/3$, (b) $14/15$.

20. V prvi vrečki so 3 bele kroglice, v drugi pa 2 rdeči. Na slepo izberemo eno od obeh vrečk, v njej izberemo kroglico in jo prenesemo v sosednjo vrečko. Naposled še enkrat sežemo na slepo v eno od obeh vrečk in v njej vnovič na slepo izberemo kroglico. Verjetnost, da bo nazadnje izbrana kroglica bele barve, je:

- (a) $3/4$, (b) $25/48$.

21. V športnem oddelku gimnazije, ki ga obiskuje 24 dijakov, se jih 15 ukvarja z nogometom, 15 jih kolesari in 10 šahira. Nogomet in šah igra 6 dijakov. Z nogometom in kolesarjenjem se jih ukvarja 10, 7 pa jih kolesari in igra šah. Vse tri športe gojijo 3 dijaki. Izračunajte verjetnost dogodkov:

- (a) da se slučajno izbran dijak ne ukvarja z nobenim športom,
 (b) da slučajno izbran dijak kolesari in šahira, a ne igra nogometom,
 (c) da se med dvema slučajno izbranimi dijakoma eden ukvarja samo z nogometom, drugi pa samo kolesari.

22. Iz kompleta 52 igralnih kart na slepo izberemo tri.

- (a) Kolikšna je verjetnost, da so vse enake barve?
 (b) Kolikšna je verjetnost, da so vse tri različnih barv?

23. V n -kotniku na slepo izberemo dve oglisci. Kolikšno mora biti najmanj število n , da

bo verjetnost, da sta izbrani točki krajišči diagonale mnogokotnika, vsaj 0.95 ?

24. Črke besede **HIPERBOLA** napišemo na 9 listkov in jih premešamo. Nato na slepo izberemo 5 listkov in jih položimo na mizo v naključnem zaporedju. Izračunajte verjetnost naslednjih dogodkov:

- (a) sestavili smo besedo **LOPAR**,
 (b) sestavljeni besedi vsebuje 3 soglasnike in 2 samoglasnika,
 (c) sestavljeni besedi vsebuje črko **H**,
 (d) sestavljeni besedi vsebuje 4 soglasnike, se začne s črko **H** in konča s črko **A**.

25. Sedem različnih učbenikov na slepo pospravimo v dva predala: v večjega štiri, v manjšega tri učbenike. Kolikšna je verjetnost, da bosta določena dva učbenika znašla v istem predalu?

26. Na zabavi se je zbral 6 družin. Vsako družino prestavlja oče, mati in trije otroki. Za igro slepo izberemo dve odrasli osebi in enega otroka. Izračunajte verjetnost dogodkov:

- (a) da so osebe izbrane iz iste družine,
 (b) da so osebe izbrane iz treh različnih družin in sta odrasli osebi različnega spola.

27. Trener razdeli na slepo pet različnih štartnih številk petim atletom: Andreju, Borisu, Cenetu, Vidu in Žigi. Kolikšna je verjetnost, da dobi Andrej nižjo startno številko kot Vid in Žiga?

28. Štiri različna pisma za štiri različne napisne bomo na slepo zlepili v štiri ovojnice z njihovimi naslovi. Izračunaj verjetnost, da bosta natanko dve pismi prišli na pravi naslov.
29. Na treh kroglicah so številke 1, 3 in 5, na treh ploščicah pa 2, 4 in 6. Kroglice in ploščice postavimo v raven niz.
- (a) V koliko razporeditvah stojijo kroglice skupaj?
 - (b) V koliko razporeditvah se niz začne in konča s ploščico?
 - (c) Koliko štirimestnih števil lahko stavimo iz vseh števk, ki so na kroglicah in ploščicah?
 - (č) Kolikšna je verjetnost dogodka, da je štirimestno število večje od 5000?
30. V stolpnici stanuje 5 družin z enim otrokom, 3 družine s 3 otroki in 2 družini s 5 otroki. Zaradi anketiranja izberemo na slepo 3 družine. Izračunajte verjetnosti, da
- (a) imata dve izmed izbranih družin isto število otrok,
 - (b) imajo vse tri izbrane družine skupaj 7 otrok.
31. Kolikšna je verjetnost, da pri naključni permutaciji črk besede **ŠALA** enaki črki ne bosta stali druga zraven druge?
32. Naključen izbor treh črk v besedi **RAČUN** naj bo elementaren dogodek v poskusu.
- (a) Koliko je vseh različnih dogodkov v tem poskusu?
 - (b) Kolikšna je verjetnost, da sta med izbranimi tremi črkami vsaj dva soglasnika?
- (c) Kolikšna je verjetnost, da pri naključni permutaciji črk besede **RAČUN** sestavi nova beseda, v kateri oba samoglasnika nista sosednji črki?
33. Micka ima v vrečki pet kroglic, ki se nazven razlikujejo le po barvi: 3 so bele, 2 pa sta črni. Na slepo vlečemo kroglice eno za drugo iz vrečke. Kolikšna je verjetnost, da ji bo šele v tretjem poskusu uspelo prvič izvleči črno kroglico?
34. Na naši fakulteti so si študentje izbirne predmete izbrali takole: 55 študentov statistiko, 80 študentov kriptografijo, 75 študentov verjetnost, 25 študentov kriptografijo in verjetnost, 20 študenti statistiko in verjetnost, 5 študentov vse tri predmete.
- (a) Koliko študentov je v tem letniku?
 - (b) Koliko študentov je izbralo dva in koliko samo en predmet od vseh treh naštetih?
 - (c) Koliko študentov je izbralo kriptografijo in ne verjetnosti?
 - (č) Izračunaj verjetnost dogodka, da je slučajno izbran študent izbral statistiko.
 - (d) Izračunaj verjetnost dogodka, da je slučajno izbran študent izbral verjetnost, pri pogoju, da je izbral statistiko.
35. Pri poskusu sodelujejo štiri osebe. Vsaka si povsem naključno in neodvisno od drugih izbere neko naravno število, manjše od 100. Kolikšna je verjetnost, da bo vsaj eno od teh štirih števil deljivo s 3?

- 36.** Lokostrelec cilja v mirujočo tarčo. Na voljo ima 4 puščice, vrednost zadetka pri vsakem strelu je 0.6 (po zadetku streljanje prekine). Izračunaj verjetnost dogodkov:
- za zadetek porabi natanko 3 puščice,
 - za zadetek porabi največ 2 puščici,
 - tarča je zadeta,
 - lokostrelec porabi vse 4 puščice.
- 37.** Iz škatle, v kateri so 3 bele in 2 črni kroglici, vlečemo na slepo po eno kroglico brez vračanja, dokler ni število izvlečenih belih kroglic enako številu izvlečenih črnih kroglic, ali dokler v škatli ne zmanjka kroglic. Izračunajte verjetnost, da bo število izvlečenih belih kroglic enako številu izvlečenih črnih.
- 38.** V podjetju z velikim številom zaposlenih je 60% delavcev moških. 35% moških in 25% žensk ima visoko izobrazbo.
- Izračunajte verjetnost, da ima naključno izbrani delavec visoko izobrazbo.
 - Naključno izbrani delavec ima visoko izobrazbo. Kolikšna je verjetnost, da je to ženska?
 - Vsaj koliko delavcev moramo izbrati, da je med njimi z verjetnostjo večjo od 0.95 , vsaj eden z visoko izobrazbo.
- 39.** Lastnik stojnice na zabavišču je ugotovil, da je verjetnost, da nalkjučni gost s pikadom zadene v polno in s tem osvoji za nagrado plišastega medvedka, približno $1/10$. Najmanj kolikokrat mora tedaj oče Hinko vreči pikado, da bo verjetnost, da osvoji vsaj enega medvedka za svojo hčer,
- vsaj $1/2$? Računajte, da namerava oče vnaprej plačati število metov in ne glede na vmesne izide vreči vsa plačana pikada.
- 40.** Robotek stoji v spodnjem levem polju šahovnice velikosti 3×3 (ki jo narišemo na vrhu visoke stolpnice). Vsak njegov premik je naključen: z enako verjetnostjo se premakne zmeraj za eno polje bodisi v desno ali pa navzgor. Na ta način lahko robotek torej tudi zdrsne čez rob šahovnice, s čimer je njegove poti seveda konec. Kolikšna je verjetnost, da robotku uspe priti na zgornje desno polje?
- 41.** Vsak od šesterice prijateljev je pred kasaško dirko na slepo stavil na enega od treh konjev A, B ali C.
- Na koliko načinov je malhko v tem primeru stavilo teh šest prijateljev?
 - Kolikšna je verjetnost, da sta Albert in Bruno, da iz omenjene šesterice prijateljve, stavila na istega konja?
 - Kolikšna je verjetnost, da je vsaj eden od teh šestih prijateljev stavil na konja A?
 - Kolikšna je verjetnost, da jih je med njimi na konja A stavilo toliko, kot na oba ostala konja?
- 42.** Igralno kocko vržemo petkrat zapored.
- Kolikšna je verjetnost, da bo v prvem in zadnjem metu padlo enako število pik?
 - Kolikšna je verjetnost, da bo padlo pri tem sodo število pik vsaj enkrat?
 - Kolikšna je verjetnost, da bo šele

- v tretjem metu prvič padlo manj kakor tri pike?
- (č) Kolikšna je verjetnost, da bo padlo v vsakem naslednjem metu več pik kakor v prejšnjem?
43. Štirje igralci drug za drugim v krogu mečejo igralno kocko. Zmaga igralec, ki prvi vrže šestico. Izračunajte verjetnost, da zmaga igralec, ki igro začne?
44. Igralna kocka je prirejena tako, da so nekateri izidi bolj verjetno od drugih. Izidi enica, dvojka in trojka so enako verjetni, izida štirica in petica sta dvakrat bolj verjetna kot enica, šestica je trikrat verjetnejša od enice.
- (a) Izračunajte verjetnosti elementarnih dogodkov.
- (b) Izračunajte verjetnosti, da pade sodo število pik.
- (c) Izračunajte verjetnosti, da pade v treh poskusih vsaj enkrat šestica.
45. Pokaži, da iz $A \subset B$ sledi $AB = A$.
46. Pokaži, da iz $A \subset B$ sledi $A \cup B = B$.
47. Kocko vržemo sedemkrat. Kolikšna je verjetnost, da padejo več kot štiri šestice?
48. Katero število grbov je nabolj verjetno, če vržemo pošten kovanec
- (a) 70-krat, (b) 75-krat.
49. Tovarna izdeluje žarnice; med njimi je 4% takih, ki po kakovosti ne ustreza normi. Oceni verjetnost, da so med 50 kupljenimi žarnicami 4 neustrezne.
50. Trije lovci so hkrati ustrelili na divjega prašiča, ki je ubit z 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, 0·6?

B.2 Poskusi, dogodki in definicija verjetnosti

- Standardno kocko (s pikami od 1 do 6) želimo obtežiti tako, da ko jo bomo dvakrat (neodvisno) vrgli, bo vsota pik obeh metov zavzela vrednost od 2 do 12, vsako z enako verjetnostjo. Poišči tako obtežitev ali dokaži, da ne obstaja.
- Skupina sedmih moških in petih žensk se odpravlja na taborjenje. Imajo dva šotorja za tri osebe in tri šotore za dve osebi. Na koliko načinov se lahko razdelijo v šotore? Kaj pa, če naj bo v vsakem šotoru vsaj en moški in vsaj ena ženska? Kaj pa, če morajo biti v istem šotoru le osebki istega spola? (Ljudi in šotore med seboj ločimo.)
- Trije gusarji, Kuki, Luki in Muki, najdejo zaklad, štiri zlatnike. Razdelijo si jih na naslednji način: Kuki in Luki vržeta pošten kovanec. Če pade grb, dobi prvi zlatnik Kuki, sicer ga dobi Luki. Tako en gusar dobi prvi kovanec, ostala dva gusarja pa nato mečeta kovanec za drugi zlatnik. Za tretji zlatnik se na isti način potegujeta tista dva gusarja, ki nista dobila drugega zlatnika, ter za četrти zlatnik tista dva, ki nista dobila

tretjega zlatnika. Kaj je bolj verjetno, da ima Kuki isto število zlatnikov kot Luki ali kot Muki?

4. Janez gre z avtobusom v službo. Možno se je peljati z dvema progama, prva vozi na 5, druga pa na 7 minut in sta neodvisni. Če gre s prvo progo, potrebuje od trenutka, ko stopi na avtobus, pa do službe 15 minut, če gre z drugo, pa 13 minut. Recimo, da gre Janez na prvi avtobus, ki pride. Kolikšna je verjetnost, da bo do službe (skupaj s čakanjem) potreboval manj kot 18 minut?
5. 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

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

6. Na FRI nekoč niso hoteli vzeti v službo enega dobrega raziskovalca. Kot kaže, se ta s tem ni mogel sprijazniti in je po mnogih letih prišel na Fakulteto ter ugrabil 100 študentov - med njimi tudi Tebe. Zaprl vas je v veliko predavalnico in jo obdal z eksplozivom. Prosili ste ga, da vas izpusti, vendar se ga ni dalo omehčati saj tudi njega, kjudan njegovi veliki želji, da bi delal na FRI, daleč nazaj ali pa sedaj niso uslišali. Ker pa je imel tako zelo rad računalništvo (posebej verjetnost in statistiko, kjer se je naučil tudi kakšno novo strategijo), se je odločil, da vam vseeno ponudi naslednjo možnost:

- v drugi predavalnici je pripravil 100 omaric (oštrevilčenih od 1 do 100) in zbral vseh 100 vaših študentskih izkaznic;
- v vsako omarico je dal po eno izkaznico in vam - študentom povedal, da boste lahko šli en po en v to predavalnico z omaricami.
- vsak izmed vas bo lahko odprl največ 50 omaric (lahko tudi manj), jih nato zaprl in odšel ven;
- vsak študent, ki pride ven iz sobe z omaricami, se ne more pogovarjati (ali kako drugače komunicirati) z drugimi, saj dobi v usta veliko nogavico, čez glavo pa žakelj;
- če se bo zgodilo, da je vsak od vas videl med drugim tudi svojo izkaznico, potem vas bo izpustil, sicer pa bo šla cela stavba v zrak (pa čeprav ni dolgo tega kar je bila renovirana).

Študente je za trenutek zgrabila panika, saj so pomislili, da se jim ne piše dobro, potem pa so se spomnili, da si med njimi tudi Ti zvit(a) študent(ka), ki jim lahko

razložiš kakšno strategijo lahko ubrete, da bo verjetnost, da se rešite bistveno večja od $(1/2)^{100}$, recimo blizu $1/3$. Gotovo misliš, da gre samo za sanje, pa temu žal ni tako. Ali se še spomniš, kaj si jim povedal(a)?

7. Razprodano letalo premore 100 sedežev. Vkrcavanje poteka tako, da se najprej vkrca potnik, ki sedi na sedežu št. 1, nato potnik, ki sedi na sedežu št. 2, ... in na zadnje potnik, ki sedi na sedežu št. 100. Na žalost vkrcavanje ne poteka ravno po željah letalske družbe. Namreč 1. oseba ne spoštuje reda in se usede na poljubni sedež (vsakega z verjetnostjo $1/100$). Vsaka naslednja oseba, ki pride na letalo, se usede na svoj sedež, če je le-ta prost. Če pa je ta zaseden, si tudi ona med preostalimi prostimi sedeži izbere enega naključno (vsakega z enako verjetnostjo).

Kolikšna je verjetnost, da 100-ti potnik sedi na sedežu št. 100 (torej na svojem sedežu)?

- (a) $1/100$,
- (b) $1/99$,
- (c) $1/e$,
- (d) $1/2$,
- (e) $1 - 1/e$,
- (f) $98/99$,
- (g) $99/100$,
- (h) nobena izmed naštetih možnosti.

B.3 Pogojna verjetnost

1. Pepe deli karte pri taroku. Ko deli talon (6 kart), pogleda, ali je v njem škis (ena izmed 54 kart, kolikor jih je vseh skupaj). Če je škis v talonu, mu ga s pogojno verjetnostjo 50% uspe neopaženo vtihotapiti med svojih 12 kart, z verjetnostjo 30% mu to ne uspe (a tudi nihče nič ne opazi), z verjetnostjo 20% pa ga razkrinkajo. Recimo, da soigralci niso opazili nič sumljivega. Kolikšna je pogojna verjetnost, da ima Pepe škisa?
2. Osebi A in B mečeta kovanec, pri katerem grb pade z verjetnostjo p . Zmaga igralec, ki prej vrže grb. Najprej kovanec vrže oseba A in če pade grb zmaga, sicer poda kovanec osebi B . Če oseba B vrže grb zmaga, sicer poda kovanec nazaj osebi A in igra se nadaljuje, dokler nekdo ne zmaga.
 - (a) Denimo, da je kovanec pravičen ($p = 1/2$). Kolikšna je verjetnost, da zmaga oseba A ? Kolikšna, da zmaga oseba B ?
 - (b) Posplošimo igro na k igralcev. Torej najprej vrže kovanec oseba A_1 , nato $A_2, \dots, A_k, A_1, \dots$ dokler ena oseba ne vrže grba. Kolikšna je verjetnost, da zmaga oseba A_i ?

- (c) Denimo, da je zmagala oseba A_i . Kolikšna je pogojna verjetnost, da je zmagala v drugem krogu (oseba A_i je v drugem poskusu vrga grb)?
3. Med dvajsetimi kovanci sta dva kovanca z dvema grboma ter dva kovanca z dvema ciframi. Ostali kovanci imajo en grb in eno cifro. Naključno si izberemo en kovanec in ga vržemo 5-krat.
- (a) Kolikšna je verjetnost, da vržemo natanko 5 grbov?
- (b) Recimo, da smo vrgli 5 grbov. Kolikšna je pogojna verjetnost, da smo povlekli kovanec z dvema grboma?
4. Vsak izmed 10 študentov si izbere bodisi belo bodisi črno kroglico, vsako enako verjetno. Nato 8-krat seže v to posodo in vsakič izvleče kroglico, ki je ne vrne. Izkaže se, da je izzrebal 5 belih in 3 črne.
- (a) Kolikšna je verjetnost, da sta v posodi ostali še dve črni kroglici?
- (b) Študent Anže ve, da je v posodo prispeval belo kroglico. Koliko je zanj pogojna verjetnost, da sta v posodi ostali še dve črni kroglici?
- (c) Študent Blaž ve, da je v posodo vrgel črno kroglico (ne ve pa, kaj ve Anže). Koliko je zanj pogojna verjetnost, da sta v posodi še dve črni kroglici?
5. V prvi posodi sta dve beli in ena črna kroglica, v drugi posodi pa dve črni in ena bela kroglica. Najprej Marjetica vrže goljufivi kovanec, na katerem grb pade z verjetnostjo 40%. Če pade grb, z izvleče eno kroglico iz prve posode, sicer izvleče eno kroglico in druge posode. Izvlečene kroglice ne vrne nazaj v posodo. Za njo pride Trdoglav, ki izvleče eno kroglico iz iste posode iz katere je c kroglico izvlekla Marjetica.
- (a) Kolikšna je verjetnost, da je Trdoglav izvlekel črno kroglico?
- (b) Kolikšna je pogojna verjetnost, da je Trdoglav izvlekel črno kroglico, če je Marjetica izvlekla črno kroglico?
- (c) Kolikšna je pogojna verjetnost, da je Marjetica izvlekla črno kroglico, če je Trdoglav izvlekel črno kroglico?
6. Gусар želi poiskati zaklad, ki se nahaja na enem izmed n otokov. Otoke začne raziskovati enega za drugim, tako da za naslednji otok izbere enega od preostalih, vse z enako verjetnostjo. Ker želi gусар v čim krajšem času poiskati zaklad je malce površen in zato najde zaklad na posameznem otoku le z verjetnostjo p , če se zaklad na tem otoku tudi nahaja. Naj bo k celo število, $0 \leq k \leq n$. Kolikšna je verjetnost, (a) da gусар ne najde zaklada na enem izmed prvih k otokov, in (b) da se s zaklad ne nahaja na nobenem izmed prvih k otokov, če vemo, da gусар ni našel zaklada na prvih k otokih.

B.4 Bernoullijevo zaporedje neodvisnih poskusov

B.5 Slučajne spremenljivke in porazdelitve

B.6 Slučajni vektorji

B.7 Funkcije slučajnih spremenljivke in vektorjev

1. Hkrati vržemo tri kovance in za tem še enkrat vržemo tiste kovance, na katerih je padel grb. Meti so med seboj neodvisni, verjetnosti, da pade grb je na vseh kovancih enaka 50%. Naj bo X število kovancev, na katerih je padla cifra v prvem metu in naj bo Y število kovancev, na katerih je padla s cifra v drugem metu.
 - (a) Zapišite porazdelitev slučajnega vektorja (X, Y) .
 - (b) Poišcite robni porazdelitvi in porazdelitev produkta XY .
 - (c) Sta slučajni spremenljivki X in Y neodvisni?
 - (d) Izračunajte pričakovano vrednost $E(X - 2Y)$.

B.8 Momenti in kovarianca

1. Fakulteta za računalništvo in informatiko ima nad pritličjem 8 različnih nadstropij ($M_1, 1, M_2, \dots, M_4, 4$), ki so bila nekoč označena s številkami od 1 do 8. Denimo, da so tudi sedaj označena tako. V dvigalo vstopi 6 profesorjev, med katerimi si vsak izbere nadstropje naključno, vsakega z isto verjetnostjo in neodvisno od ostalih profesorjev.
 - (a) Kolikšno je pričakovano število postankov dvigala?
 - (b) Naj bo H najvišje nadstropje, v katerem se dvigalo ustavi. Poišcite tako število (mediano) m , za katerega velja $P(H \leq m) \leq 1/2 \leq P(H \leq m)$.
2. Dragić, Nachbar in Lakovič tekmujejo v metanju trojk. Vsak izmed njih vrže na koš enkrat. Prvi na koš vrže Dragić, za njim Nachbar in na koncu se Lakovič. Preden začnejo metati na koš, da vsak izmed njih v kapo 100EUR. Kdor koš zgreši, mora znesek v kapi podvojiti. Kdor koš zadane, dobi ves denar iz kape, nato pa vsi prispevajo začetnih 100EUR. Kar je v kapi po koncu metov, si razdelijo na enake dele. Recimo, da Dragić zadane trojko z verjetnostjo 25%, Nachbar z verjetnostjo 50% in Lakovič z verjetnostjo 75%. Koliko ima Dragić v povprečju dobička oz. izgube?
3. Kako sta povezani mediana zaporedja in mediana slučajne spremenljivke?
4. Predpostavi, da je graf gostote verjetnosti $p(x)$ zvezne slučajne spremenljivke simetričen. Kaj lahko v tem primeru poveš o mediani in pričakovani vrednosti (mat. upanje) te slučajne spremenljivke?

B.9 Karakteristične funkcije in limitni izreki

B.10 Opisna statistika

1. Množico, ki jo statistično opazujemo, imenujemo:
 - (a) populacija,
 - (b) vzorec,
 - (c) statistična spremenljivka (znak).
 2. Značilnosti populacije kot celote imenujemo:
 - (a) statistične enote,
 - (b) statistične spremenljivke,
 - (c) statistične parametri.
 3. Frekvenca je število, ki pove:
 - (a) koliko je v opazovani populaciji vseh možnih vrednosti statistične spremenljivke;
 - (b) kolikokrat je v opazovani populaciji nastopila določena vrednost statistične spremenljivke.
 4. Frekvenčna distribucija je prikaz:
 - (a) zbranih vrednosti statistične spremenljivke v ustrezni tabeli;
 - (b) zbranih vrednosti statistične spremenljivke v ustremnem grafikonu.
 5. Frekvenca posameznega razreda je enaka:
 - (a) številu enot opazovane populacije;
 - (b) številu vrednosti številske spremenljivke, ki spada v tisti razred.
 - (c) količniku med številom vrednosti številske spremenljivke, ki spada v tisti razred, in številom vseh enot populacije.
 6. Širina razreda je enaka:
- (a) aritmetični sredini obeh mej razreda,
 - (b) razliki med zgornjo in spodnjo mejo razreda.
7. Porazdelitev absolutnih in porazdelitev relativnih frekvenc lahko prikažemo s:
 - (a) frekvenčnimi poligoni,
 - (b) histogrami,
 - (c) frekvenčnimi kolači.
 8. Frekvenčni kolač (struktturni krog) se stavljajo krožni izseki. Ti kažejo
 - (a) deleže enot, ki sodijo v posamezne razrede;
 - (b) število enot, ki sodijo v posamezne razrede.
 9. Uteženo aritmetično sredino koristno uporabljamo, če:
 - (a) statistična spremenljivka zavzame različne vrednosti na več enotah populacije;
 - (b) statistična spremenljivka zavzame isto vrednost na več enotah populacije;
 - (c) so vrednosti statistične spremenljivke razdeljene v razrede.
 10. Varianca, mera za razpršenost posameznih vrednosti, je:
 - (a) število, ki je kvadratni koren iz povprečja kvadratov odklonov posameznih vrednosti statistične spremenljivke od aritmetične sredine;
 - (b) število, ki je enako povprečju kvadratov odklonov posameznih vre-

dnosti statistične spremenljivke od aritmetične sredine;

11. V razredu s 25 učenci je 8% odličnih, 28% pravdobrih, trije so nezadostni, sedem je zadostnih. Ostali so dobri. Izračunajte povprečno oceno razreda in narišite histogram frekvenc ocen.
12. V razredu je 25 učencev. Učenka Andreja je računala povprečno število točk pri šolski nalogi. Pri prvem računanju se je zmotila: ni upošteva svojega dosežka in dobila povprečje 74,5 točk. Ko je napako popravila, je dobila povprečje 75 točk. Koliko točk je dosegla Andreja pri šolski nalogi?
13. Za izhodišče vzemimo stavek:

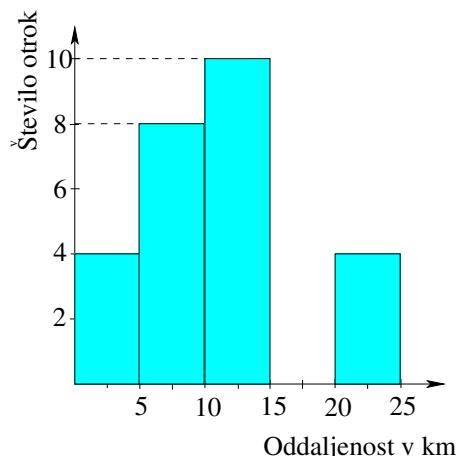
Je zvito kakor kozji rog.

- (a) V zgornjem stavku opazujte števila črk v posameznih besedah. Izračunajte povprečno število črk v besedah tega stavka in standarno deviacijo števila črk v besedah.
- (b) Opazujte besede v zgornjem stavku še enkrat: tokrat štejte, koliko samoglasnikov vsebuje ta in ona beseda. Ali je razpršenost teh podatkov večja ali manjša od razpršenosti podatkov o dolžini

posameznih besed iz prejšnjega vprašanja?

14. V nekem razredu je 30 otrok. Njihovi domovi so od šole oddaljeni največ 25 km. Oddaljenosti so prikazane s frekvenčnim histogramom, v katerem manjka podatek za oddaljenost od 15 km do 20 km. Mejne oddaljenosti 5 km, 10 km, ... stejemo v histogramu k manjšim vrednostim, na primer $15 \text{ km} \in (10 \text{ km}, 15 \text{ km}]$.

 - (a) Koliko domov otrok je od šole oddaljenih od 5 km do 10 km?
 - (b) Koliko domov otrok je od šole oddaljenih od 15 km do 20 km?
 - (c) Koliko so domovi otrok povprečno oddaljeni od šole? Izračunajte povprečno vrednost, ki jo lahko dobite iz podatkov v histogramu.



B.11 Vzorčenje

B.12 Cenilke

B.13 Intervali zaupanja

B.14 Preverjanje statističnih domnev

1. Proizvajalec meritcev električne moči, ki se uporablja za uravnavanje pragov energije pri podatkovno-komunikacijskih sistemih, trdi, da ob normalno delujoči proizvodni liniji proizvede največ 10% nedelujočih meritcev.

Trgovec je pravkar dobil pošiljko 25 omenjenih meritcev. Recimo, da želi preveriti domnevo $H_0 : p = 0.10$ napram $H_1 : p > 0.10$, kjer je p pravi delež pokvarjenih meritcev. Pri testu uporabi $y \geq 6$ za prag zavrnitve (y je seveda pravo število pokvarjenih meritcev v pošiljki).

- (a) Določite vrednost α za ta test.
- (b) Določite β , če je $p = 0.2$. Koliko je moč testa za to vrednost p ?
- (c) Določite β , če je $p = 0.4$. Koliko je moč testa za to vrednost p ?
2. S pomočjo računalniške simulacije so raziskovali učinek napak na strojih na performanse proizvodnega sistema (Industrial Engineering, Aug. 1990). Študija se je osredotočila na sistem z enim strojem. Povprečni čas med začetkom dveh procesiranj je 1.25 minute pri konstantnem času procesiranja 1 minuta. Naprava je pokvarjena 10% časa. Po $n = 5$ neodvisnih simulacijah dolžine 160 ur je povprečna produktivnost na teden (40-urni delavnik) $\bar{X} = 1908.8$ izdelkov. Za sistem brez napak je povprečna produktivnost 1920 izdelkov. Če predviedevamo, da je standardna deviacija 5 simulacij $s = 18$ izdelkov, testiraj domnevo, da je resnična povprečna produktivnost manjša od 1920 izdelkov. Testirajte pri značilnosti $\alpha = 0.05$.
3. Rezultati druge raziskave o nacionalnem zdravju in prehranjevanju v ZDA so pokazali, da imajo ljudje v starosti med pol leta in 74 let v krvi povprečno koncentracijo svinca $14\mu\text{g}/\text{dl}$ (Analytical Chemistry, feb. 1986). Poleg tega so ugotovili, da imajo črnski otroci v starosti do pet let znatno višje koncentracije od ostalih.

V naključnem vzorcu 200 črnskih otrok v starosti pod pet let je bila ugotovljena povprečna koncentracija svinca v krvi $21\mu\text{g}/\text{dl}$ s standardnim odklonom $10\mu\text{g}/\text{dl}$. Je to dovolj, da lahko trdimo, da je pravo povprečje pri črnski populaciji pod pet let res večje od $14\mu\text{g}/\text{dl}$? Testirajte z vrednostjo $\alpha = 0.01$.

4. Dovoljena koncentracija PCB-ja (nevarna substanca) v vodi po standardu, ki ga je postavila EPA, je 5 delcev na milijon. Večji proizvajalec PCB-ja, ki se uporablja za električno izolacijo, spušča manjše količine PCB-ja skupaj z odpadno vodo. Uprava tovarne je izdala navodila, da je treba ustaviti proizvodnjo, če povprečna koncentracija PCB-ja v odplakah preseže 3 delce/milijon. Analiza 50 naključnih vzorcev odpadne vode je pokazala naslednje rezultate:

$$\bar{y} = 3.1 \text{ delcev na milijon} \quad \text{in} \quad s = 0.5 \text{ delcev na milijon}$$

- (a) Ali so zgornji statistični rezultati analize dovoljeni dokaz, da je proizvodnjo potrebno ustaviti? Uporabite $\alpha = 0.01$.
- (b) Če bi bili vi menedžer tovarne, ali bi uporabili večjo ali manjšo vrednost za α v testu pod točko (a)? Pojasnite!
5. Vrtanje "globokih lukenj" je družina procesov, ki omogočajo vrtanje lukenj, ki so vsaj desetkrat globlje kot je premer svedra. Eden bistvenih problemov pri globokem vrtanju je zastajanje izvrtanih okruškov v utorih svedra. Izveden je bil eksperiment, s katerim so preučili uspešnost globokega vrtanja v primeru, ko se okruški sprijemajo (Journal of Engineering for Industry, maj 1993). Globina izvrtanih lukenj pri 50 vrtanjih je v povprečju znašala $y = 81.2$ mm s standardno deviacijo vzorca $s = 50.2$ mm. Testirajte domnevo, da se pravo povprečje izvrtanih lukenj μ razlikuje od 75 mm. Uporabite stopnjo značilnosti $\alpha = 0.01$.

6. Inštitut za okolje in tehnologijo je izdal študijo o onesnaženosti zemlje na Nizozemskem. Skupno so zbrali, posušili in analizirali za prisotnost cianida 72.400 gramov vzorcev zemlje. Z infrardečo mikroskopsko metodo so merili koncentracijo cianida v miligramih na kilogram zemlje v vzorcih. Povprečna koncentracija cianida v vzorcih je bila $y = 84$ mg/kg, standardna deviacija pa $s = 80$ mg/kg.

Uporabite to informacijo za testiranje domneve, da je resnično povprečje koncentracije cianida v zemlji na Nizozemskem manjše od 100 mg/kg pri stopnji značilnosti $\alpha = 0.10$.

7. Ali tekmovanje med različnimi odseki za raziskovanje in razvoj (R&R) v sklopu ameriškega ministrstva za obrambo, ki delajo na istem projektu, izboljša produktivnost? Za odgovor na to vprašanje so raziskali produktivnost pri 58 projektih, dodeljenih večim oddelkom, ter pri 63 projektih, dodeljenih enemu samemu oddelku (IEEE Transaction on Engineering Management, feb. 1990). Glede na rezultate so dobili povprečen faktor uspešnosti pri prvih (tekmovalnih) projektih enak 7.62, pri drugih (netekmovalnih) pa 6.95.

- (a) Postavite ničelno in alternativno domnevo za odločanje, ali povprečje uspešnosti pri tekmovalnih projektov presega povprečje uspešnosti pri netekmovalnih projektih.

- (b) Določite območje zavrnitve za test pri vrednosti $\alpha = 0.05$.
- (c) Izkazalo se je, da pri zgornjem testu p -vrednost leži v intervalu med 0.02 in 0.03. Kakšen je pravilen zaključek?
8. Inštitut za okolje in prostor je naredil študijo o insekticidih, ki se uporabljajo v nasadu orhidej v dolini San Joaquin v Kaliforniji. Zbrali so vzorce zraka v nasadu in jih testirali vsak dan v obdobju najbolj intenzivnega škropljenja. V spodnji tabeli so prikazane količine oksonov in thionov v zraku (v ng/m^3) in razmerje med količino thionov in oksonov. Primerjaj povprečje razmerij med oksoni in thioni v megleh in jasnih/oblačnih pogojih v nasadu orhidej z uporabo domneve. Uporabi stopnjo značilnosti $\alpha = 0.05$.
- | Datum | Vremenske razmere | Thioni | Oksoni | Razmerje oksoni/thioni |
|---------|-------------------|--------|--------|------------------------|
| 15. jan | megla | 38.2 | 10.3 | 0.270 |
| 17. | megla | 28.6 | 6.9 | 0.241 |
| 18. | megla | 30.2 | 6.2 | 0.205 |
| 19. | megla | 23.7 | 12.4 | 0.523 |
| 20. | megla | 62.3 | - | vzorec izgubljen |
| 20. | jasno | 74.1 | 45.8 | 0.618 |
| 21. | megla | 88.2 | 9.9 | 0.112 |
| 21. | jasno | 46.4 | 27.4 | 0.591 |
| 22. | megla | 135.9 | 44.8 | 0.330 |
| 23. | megla | 102.9 | 27.8 | 0.270 |
| 23. | oblačno | 28.9 | 6.5 | 0.225 |
| 25. | megla | 46.9 | 11.2 | 0.239 |
| 25. | jasno | 44.3 | 16.6 | 0.375 |
9. Tovarna bi rada ugotovila kateri od naslednjih dveh virov energije – plin in elektrika – bo dal cenejšo energijo. Ena mera za ekonomično proizvodnjo energije (v angl. se imenuje plant investment perdelivered quad) je izračunana tako, da vzamemo vsoto vloženega denarja (dolarji) v določen vir s strani neke tovarne in jo delimo s količino energije (v kvadrilionih britanskih termalnih enot). Manjši kot je kvocient manj plača tovarna za dobavljeni elektriko. Izbran je bil naključni vzorec 11ih tovarn, ki uporabljata električno energijo

1	2	3	4	5	6	7	8	9	10	11
204.15	0.57	62.76	89.72	0.35	85.46	0.78	0.65	44.38	9.28	78.60

in 16ih tovarn, ki uporabljajo plin:

1	2	3	4	5	6	7	8
0.78	16.66	74.94	0.01	0.54	23.59	88.79	0.64
9	10	11	12	13	14	15	16
0.82	91.84	7.20	66.64	0.74	64.67	165.60	0.36

Nato je bil izračunan zgoraj omenjeni kvocient za vsako tovarno in podatki vnešeni v tabele, sledi pa še izpis iz MINITABA za analizo podatkov.

```
TWOSAMPLE T FOR electric VS gas
      N    MEAN   STDEV    SE MEAN
electric 11    52,4   62,4      19
gas       16    37,7   49,0      12
```

```
95 PCT CI FOR MU electric - MU gas: (-30, 59)
```

```
TTEST MU electric = MU gas (VS NE) : T=0,68      P=0,50      DF= 25
POOLED STDEV = 54,8 (standardna e)
```

- (a) Ali ti podatki predstavljajo dovolj dober pokazatelj za stopnjo značilnosti $\alpha = 0,05$, ki kaže na razliko med povprečjem kvocientov, ki uporabljajo plin in tistimi, ki uporabljajo elektriko?
- (b) Kakšne predpostavke so potrebne, da bo postopek veljaven?
- (c) Preveri, če so predpostavke iz točke (b) smiselno izpolnjene. Kako to vpliva na upravičenost rezultata v točki (a)?
10. Raziskovalci so naredili eksperiment za določitev vpliva puščavskih granivorjev (živali, ki jedo semena) na gostoto in porazdelitev semen v prsti (Ecology, Dec. 1979). Določene vrste puščavskih glodalcev delajo zaloge semen na površju zemlje, zato je bil eksperiment zasnovan tako, da so raziskovali ali zaradi takšnih zalog semen v povprečju na tistem območju iz semen zraste več mladih rastlinic kot na sosednjih kontrolnih območjih. Locirali so 40 majhnih območij, kjer so glodalci kopičili semena in jih pokrili z mrežo, da so glodalci preprečili ponoven dostop. Zamrežili so tudi sosednja območja za kontrolo. Potem so opazovali število semen, ki je vzklilo na zamreženih območjih. V spodnji tabeli je povzetek zbranih podatkov. Ali so rezultati eksperimenta zadosten dokaz (pri stopnji značilnosti $\alpha = 0,05$), da je povprečno število vzkaljenih semen na območjih, kjer so jih kopičili glodalci bistveno večje kot na kontrolnih območjih?

Območja z zalogami	Kontrolna območja
$n_1 = 40$	$n_2 = 40$
$y_1 = 5,3$	$s_2 = 2,7$
$s_1 = 1,3$	$s_2 = 0,7$

11. Percepcija govora pretežno gluhih oseb se zanaša predvsem na branje z ust, tj., zaznavanje pogovornega jezika z opazovanjem artikuliranih gibov, izrazov na obrazu ter gest sogovornika. Ali se da percepcijo govora izboljšati tako, da bralcu glasu vizualno

predstavimo podatek o poudarku zlogov. Da bi raziskali ta fenomen je 10 oseb z normalnim sluhom sodelovalo v eksperimentu pri katerem so morali ustno ponoviti zvočno predvajane stavke, katerih informacijo niso videli na video-monitorju. (Journal of the Acoustical Society of America, Feb. 1986). Stavki so bili predstavljeni osebam pod naslednjima pogojema:

- (1) branje govora z informacijo o frekvenci in amplitudi govornega signala (oznaka $S + F + A$) ter
- (2) samo branje govora (oznaka S).

Za vsakega od 10-ih oseb, je bila izračunana razlika med procentom pravilno reproducirane vsebine pod pogojem $S + F + A$ in pod pogojem S . Povprečje in standardni odklon razlik sta naslednja:

$$\bar{d} = 20\cdot 4 \quad \text{in} \quad s_d = 17\cdot 44.$$

Testiraj domnevo, da je povprečje procentov pravilnih vsebin pri pogoju $S + F + A$ presega ustrezno povprečje pod pogoju S . Privzami $\alpha = 0\cdot 05$.

12. Tetraklorodibenzo-p-dioksin (TCDD) je visoko-toksična substanca, ki jo najdemo v industrijskih odpadkih. Znanstveniki so naredili študijo s katero so določili količino TCDD-ja v tkivih volovskih žab, ki živijo na območju Rocky Branch Creek v Arkanzasu, za katerega se ve, da je kontaminirano z TCDD (Chemosphere, Feb. 1986). Merili so količino TCDD-ja (v delcih na trilijon) v različnih tkivih štirih samic volovskih žab. Za vsako žabo so izmerili razmerje med količino TCDD-ja v tkivu in v nožni misici. Relativno razmerje med koncentracijo TCDD-ja v jetrih in jajčnikih s je podana v spremljajoči preglednici. Raziskovalci po zbranih podatkih sklepajo, da je povprečna koncentracija TCDD-ja v jajčnikih samic volovskih žab večja kot povprečna koncentracija v jetrih. Testirajte to trditev z uporabo $\alpha = 0\cdot 05$.

Žaba	Jetra	Jajčniki
A	11·0	34·2
B	14·6	41·2
C	14·3	32·5
D	12·2	26·2

13. V Merckovih raziskovalnih laboratorijih so izvedli poskus ocene učinka novega zdravila, pri čemer so uporabili prirejen plavalni labirint. Devetnajstim oplojenim jezovnim podganam so dali 12·5 mg zdravila. Iz vsakega legla so za plavanje v labirintu naključno izbrali enega moškega ter enega ženskega mladiča. Vsak podganji mladič je postavljen v vodo na enem koncu labirinta, plava pa, dokler uspešno ne najde izhoda na nasprotnem koncu. Če mladiču ne uspe najti izhoda v določenem časovnem intervalu,

ga spet postavijo na začetek in mu dajo še eno priložnost za pobeg. Poskus ponavljamo vse dokler vsak mladič uspešno ne pobegne trikrat. Število plavanj, potrebnih za tri uspešne pobege posameznega mladiča, je predstavljeno v spodnji razpredelnici. Ali lahko sklepamo, da je povprečno število potrebnih pobegov različno pri moških oz. ženskih mladičih?

Leglo	Moški	Ženski	Leglo	Moški	Ženski
1	8	5	11	6	5
2	8	4	12	6	3
3	6	7	13	12	5
4	6	3	14	3	8
5	6	5	15	3	4
6	6	3	16	8	12
7	3	8	17	3	6
8	5	10	18	6	4
9	4	4	19	9	5
10	4	4			

TEST MU = 0.000 proti MU \$not = \$ 0.000

	N	E	STDEV	SE MEAN	T	P VALUE
Swimoiff	19	0.368	3.515	0.806	0.46	0.65

15. Raziskovalci z Univerze v Rochesteru so študirali trenje, ki nastane v procesu zajemanja papirja iz kasete v fotokopirnem stroju (Journal of Engineering for Industry, May 1993). Poskus je vseboval opazovanje zamika posameznih listov papirja v kupu, s katerega je fotokopirni stroj zajemal. Posamezno zajemanje je označeno kot uspešno, če se noben list z izjemo vrhnjega ni zamaknil za več kot 25% celotne predvidene dolžine premika. V kupu stotih listov je bilo zajemanje uspešno 94-krat. Predviden delež uspešnosti postopka je 0·9. Pri stopnji značilnosti $\alpha = 0\cdot1$ oceni, ali zanesljivost presega 0·9.
16. Ameriška agencija za znanost, NSF, je v študiji 2237-ih doktorandov s področij tehniških ved na ameriških univerzah ugotovila, da je bilo med njimi 607 ameriških državljanov (revija Science, 24. september 1993). S stopnjo značilnosti $\alpha = 0\cdot01$ testiraj domnevo, da dejanski odstotek doktorskih nazivov, podeljenih tujim (neameriškim) državljanom na ameriških univerzah, presega 0·5.
17. Zaradi dvomov o zadostni varnosti v letalskem prometu je ameriška zvezna agencija za letalstvo (FAA) uvedla sankcije zoper letalske družbe, ki na preizkusih varnosti dobijo nezadostno oceno. Ena izmed serij poskusov, izpeljanih na mednarodnem letališču v Los Angelesu (LAX), je pokazala, da so varnostniki odkrili zgolj 72 izmed 100 lažnih orožij, ki so jih inspektorji FAA bodisi nosili pri sebi bodisi spravili v osebno prtljago. Kot so zatrdirili z FAA, je bila stopnja odkritih primerov občutno pod državnim povprečjem, ki znaša 0·8. Testiraj to domnevo, pri čemer za stopnjo značilnosti vzemi $\alpha = 0\cdot1$.

18. Oddelek za mehaniko organizacije ASEE vsakih 10 let opravi nacionalno raziskavo o poddiplomskem poučevanju mehanike na kolidžih in univerzah. Leta 1985 so na 66 izmed 100 kolidžev poučevali statiko tekočin v poddiplomskem programu strojništva. Leta 1975 pa je ta predmet poučevalo 43% kolidžev (Engineering Education, april 1986).

Izvedite test, s katerim boste določili ali se je odstotek kolidžev, ki poučujejo statiko tekočin, povečal med leti 1975 in 1985, če sklepamo da je bilo leta 1975 v raziskavo pravtako vključenih 100 kolidžev. Uporabite stopnjo značilnosti $\alpha = 0.01$.

19. V študiji, katere namen je bilo ugotoviti vpliv t.i. multifunkcijske delovne postaje (MFDP) na način dela zaposlenih (Datamation, 15 februar 1986), sta sodelovali dve skupini svetovalcev iz varnostne agencije s sedežem v St. Louisu, in sicer skupina 12-ih svetovalcev, ki trenutno uporablajo programsko opremo MFDP, ter kontrolna skupina 25-ih, ki je ne uporablja. Ena izmed vprašanj, zastavljenih udeležencem, se je dotikalo virov informacij. V testni skupini uporabnikov MFDP so štirje izjavili, da je računalnik njihov glavni vir informacij, medtem ko sta v kontrolni skupini isto izjavila dva njena člana.

- (a) Ali lahko sklepamo o razliki med deležem uporabnikov MFDP, ki se med vsemi viri informacij najbolj zanašajo na računalnik, ter istim deležem med neuporabniki MFDP? Test izvedi z $\alpha = 0.1$.
- (b) Ali sta vzorca dovolj velika, da lahko uporabimo aproksimacijski postopek iz točke (a)?

20. Sisteme za ogrevanje hiš, ki izkoriščanje sončno energijo lahko razdelimo na dve skupini - pasivne in aktivne. Pri pasivnih sistemih hiša deluje kot kolektor sončne energije, pri aktivnih sistemih pa se uporablja napredna mehanska oprema, ki pretvarja sončne žarke v toploto. V raziskavi je bilo vključenih 50 domov z pasivnimi sistemi in 50 domov z aktivnimi sistemi. Raziskovalci so zabeležili število domov, ki so porabili manj kot 200 galon kuričnega olja za ogrevanje v zadnjem letu. Ali rezultati raziskave pri stopnji značilnosti $\alpha = 0.02$ nakazujejo na to, da se pasivni in aktivni sistemi razlikujejo po učinkovitosti?

	št. domov	št. domov, ki so porabili manj kot 200 galon kuričnega olja za ogrevanje
pasivni sončni sistemi	50	37
aktivni sončni sistemi	50	46

21. Najpogostejsa metoda za dezinfekcijo vode za pitje je prosta rezidualna klorinacija. Pred kratkim pa je kot alternativa precej pozornosti pritegnil postopek preamoniacije (tj. dodajanja amoniaka vodi pred dodajanjem prostega klora). V eni študiji je bila preamoniacija izvedena na 44-ih vzorcih vode in se je izkazalo, da je imel indeks odtočne

umazanosti povprečje 1·8 ter standardni odklon 0·16 (American Water Works Journal, januar 1986).

Ali lahko sklepamo, da varianca indeksa odtočne umazanosti vodnih vzorcev, dezinficiranih s preamoniacijo, presega 0·0016? (Vrednost 0·0016 predstavlja doslej znano varianco omenjenega indeksa.) Pri testu vzemi $\alpha = 0·01$.

22. Poliklorinirani bifenili (PCB-ji), ki se uporabljajo v proizvodnji velikih električnih transformatorjev in kondenzatorjev so izjemno nevarni onesnaževalci okolja. Agencija za zaščito okolja (EPA) preizkuša novo napravo za merjenje koncentracije PCB-jev v ribah. Za testiranje natančnosti nove naprave so naredili 7 merjenj koncentracije PCB-ja na isti ribi. Podatki so zabeleženi v spodnji tabeli (v delcih na milijon):

1	2	3	4	5	6	7
6·2	5·8	5·7	6·3	5·9	5·8	6·0

Ali nova naprava ustreza specifikacijam Agencije za zaščito okolja, če le-ta zahteva, da naprave za merjenje koncentracije PCB-ja merijo z varianco manjšo od 0·1. Testirajte pri stopnji značilnosti $\alpha = 0·05$.

23. (Nadaljevanje 8. naloge.) Spomnite se, da je bila izvedena študija za primerjavo povprečnih razmerij med oksoni in thioni v Kalifornijskem nasadu orhidej v dveh vremenskih pogojih - megla in jasno/oblačno vreme. Testirajte predpostavko, da so variance v različnih vremenskih pogojih enake, ker je to pogoj da je primerjava povprečij veljavna. Testirajte pri stopnji značilnosti $\alpha = 0·05$.
24. (Nadaljevanje 11. naloge - o branju govora z ust.) Za osebe z normalnim sluhom je bil izveden še en eksperiment za primerjavo razpršenosti v zaznavanju pogovornega jezika v skupini tistih brez izkušenj v branju z ust in skupino tistih z izkušnjami v branju z ust. Vzorec je sestavljal 24 neizkušenih in 12 izkušenih oseb. Vsaka oseba je morala ustno ponoviti stavke pod različnimi pogoji, med katerim je bil tudi dodatek za branje govora s podatkom o poudarku zlogov. Povzetek rezultatov (procent pravilnih zlogov) za obe skupini je podan v tabeli. Testirajte ali obstaja razlika med variancama v procentu pravilno reproduciranih zlogov med obema skupinama. Pri testu porabite $\alpha = 0·10$.

neizkušeni	izkušeni
$n_1 = 24$	$n_2 = 12$
$\bar{y}_1 = 87·1$	$\bar{y}_2 = 86·1$
$s_1 = 8·7$	$s_2 = 12·4$

B.15 Bivariatna analiza in regresija

B.16 Časovne vrste in trendi

Dodatek C

Zgodovina in biografije

Kombinatoriko je definiral Jacob Bernoulli v svoji knjigi *Ars Conjectandi*, ki je izšla po njegovi smrti v Baslu leta 1713. Osrednja ideja kombinatorike je *preštevanje*:

1. koliko bo treba prebrati za izpit?¹
2. diagonale n -kotnika
3. število 7. mestnih dvojiških števil: $2^7 = 128$
4. L. Euler in N. Bernoulli sta si dopisovala o *problemu zamenjave pisem* in ga tudi rešila.
5. Več faz, kombinatorično drevo ...

Zgodovina: Rimski patricij Boetius v 5. stoletju, indijski matematik Bhaskara v 12. stoletju, Židovski učenjak iz Avignona Levi ben Gerson v 14. stoletju (ki je znal izračunati permutacije, kombinacije in variacije brez ponavljanja). Leta 1494 Luca Pacioli v knjigi *Suma de aritmeticae* na koliko načinov se lahko 10 oseb vsede v prvo vrsto z 10imi sedeži. 16. stoletje: Cardan, Tartaglio, Leibniz, Buteom; 17. stoletje: Pascal Fermat, Jacob Bernoulli ...

ime	dela
1. Pierre de Fermat	(1601-1665) pisma (Pascalu), <i>Varia opera mathematica</i>
2. Blaise Pascal	(1623-1662) pisma (Fermatu)
3. Jacob Bernoulli	(1654-1705) Umetnost ugibanja 1713 (<i>Ars Conjectandi</i>)
4. Abraham de Moivre	(1667-1754) <i>The Doctrine of Chances</i> 1718, 1733
5. Thomas Bayes	(1702-1761)
6. Leonhard Euler	(1707-1783) veliko knjig o algebri, analizi,... letno 800 str.
7. Joseph Louis Lagrange	(1736-1813)
8. Pierre-Simon Laplace	(1749-1827) <i>Nebesna mehanika, Analitična teorija verjetnosti</i>
9. Frederic Gauss	(1777-1855) <i>Disquisitiones Arithmeticae</i> , 1801
10. Simeon Poisson	(1781-1840)
11. Augustin Louis Cauchy	(1789-1857) 789 člankov
12. Arthur Cayley	(1821-1895) okoli 1000 objavljениh del
13. Emile Borel	(1871-1956)
14. William Sealy Gosset	(1876-1937)
15. Sir Ronald A. Fisher	(1890-1962)
16. Egon Sharpe Pearson	(1895-1980)
17. Frank Plumpton Ramsey	(1903-1930)
18. Andrei Kolmogorov	(1903-1987)
19. Paul Erdős	(1913-1996)

¹Pa poskusimo na hitro oceniti, kaj imamo pred seboj: 36 vrstic \times 86 črk je $\doteq 3K$, 100 strani je torej 300K, 200 strani 600K in čez 300 strani skoraj 1M.

C.1 Pierre de Fermat (1601-1665)

Iz članka M. Omladič, Zagonetni Fermat, Presek ???

1. Mirno življenje. V malem mestecu Beaumont-de-Lomagne na skrajnjem jugu Francije, v vojvodini Gascogni, se je 4. avgusta 1601 rodil eden največjih matematikov vseh časov Pierre de Fermat. Umrl je dobrih 63 let kasneje, 12. januarja 1665 v bližnjem kraju Cartresu, kjer je tudi pokopan. O življenju človeka, katerega genij še danes občudujemo, vemo le malo. Njegov oče Dominique je bil trgovec z usnjem in beaumonški uradnik, mati Claire de Long pa je izvirala iz pravnische družine. Šolal se je v rojstnem kraju in v bližnjem Toulousu, kjer je končal študij prava ter postal sodni uradnik. Kasneje se je tudi za stalno naselil v Toulousu, se poročil in imel pet otrok. Njegova pravnische kariera je bila počasna in mirna ter je dosegla svoj vrh takrat, ko je bil imenovan za kraljevega svetnika v lokalnem parlamentu v Toulousu. Nekateri viri pričajo o tem, da se je odlikoval po svojem poštenju, obzirnosti in zglednem vedenju. Svoj prosti čas je Fermat posvečal matematiki in dosegal izredne rezultate.

Danes velja za začetnika diferencialnega računa in skupaj z drugim znamen francoskim matematikom Blaiseom Pascalom za tvorca *teorije verjetnosti*. Njegova verjetno najpomembnejša odkritja pa sodijo v področje matematike, ki mu danes pravimo teorija števil. Na žalost nimamo pravega pregleda nad celotnim njegovim matematičnim delom, saj je svoje rezultate redko objavljal. Nekaj o njegovem delu lahko izvemo iz pisem, ki jih je pisal drugim matematikom in fizikom svojega časa. Dopisoval si je s Pascalom, Renejem Descartesom in mnogimi drugimi, po nekaterih virih celo z Isaacom Newtonom. Navadno je Fermat svoje matematične domisleke pripisal kar na robove knjig, ki jih je prebiral. Njegov sin Samuel je imel težko delo, ko je poskušal po očetovi smrti zbrati njegovo delo in ga izdati v knjižni obliki. Kljub temu je knjiga leta 1679 izšla pod naslovom *Varia opera mathematica*, ta knjiga pa je hkrati tudi edino pričevanje, ki nam je ostalo o Fermatovem delu. Precejšnje število Fermatovih trditev je ostalo nedokazanih, bodisi da so se dokazi izgubili, bodisi jih zagonetni, a genialni mož ni nikdar zapisal. Cela desetletja so minila, preden so najboljši matematiki Evrope uspeli razvozlati nekatere Fermatove uganke. Ena med njimi je še dobrih 300 let po njegovi smrti ostala nerazrešena, to je sloviti "poslednji izrek".

2. Tangenta na krivuljo. Tangenta na krivuljo je eden najpomembnejših pojmov, ki jih je vpeljal Fermat in prav to ga postavlja na čelo tistih, ki jih danes imamo za začetnike diferencialnega računa. Zamislimo si poljubno krivuljo v ravnini in izberimo točko T na krivulji. Tisti premici, ki bo šla skozi to točko in se bo najbolje prilegalna krivulji v bližini te točke, bomo rekli tangenta na krivuljo v dani točki. Vprašanje je, kako priti do te premice? Izberimo si na krivulji še dve nadaljnji točki P in Q , eno npr. desno, drugo pa levo od točke T . Premica skozi ti dve točki je *sekanta krivulje*. Če bosta točki začeli drseti po krivulji vsaka s svoje strani proti točki T , pa bo sekanta najbrž vse bližje tangenti! In ko bosta končno točki zdrsnili v T , bo sekanta skočila v tangentu! Že prav, boste rekli, toda kako to izračunati? ...

3. Fermatov princip. Fermat pa se ni ukvarjal samo z matematiko, ampak je posvetil svoj čas včasih tudi nekaterim fizikalnim problemom. Med njegova najpomembnejša odkritja s tega področja štejemo danes njegovo trditev o širjenju svetlobe v optiki. To je znameniti Fermatov princip, ki pravi, da se svetloba širi po taki poti, za katero porabi najmanj časa. Fermatov sodobnik Descartes je ta princip v svojih pismih žolčno napadal in poskušal pri tem mirnega Fermata celo žaliti. Descartes je namreč v zvezi s širjenjem svetlobe zagovarjal svoje lomne in odbojne zakone, ki so sicer pravilni, ni pa s previdel, da Fermat vidi dlje, saj so Descartesovi zakoni samo ena izmed posledic Fermatovega principa.

Prejšnje razmišljjanje o tangentah nas je dovolj podkovalo, da bomo znali to razumeti! ...

4. Teorija števil in poslednji izrek. Za konec naj navedem še dva Fermatova problema iz teorije števil. Najprej moram omeniti znameniti Fermatov izrek, ki ga je kakih 20 let po Fermatovi smrti prvi dokazal nemški matematik in filozof Gottfried Wilhelm Leibniz. Njegova vsebina je kaj preprosta: *Pri poljubnem naravnem številu n in praštevilu p je število*

$$M(n, p) = n^p - n$$

deljivo s p! Da bi se prepričali o pravilnosti te trditve, si oglejmo nekatere primere. Tako je npr. $M(n, 2) = n^2 - n = n(n - 1)$. Če je n deljiv z 2, tedaj je očitno tudi $M(n, 2)$ deljiv z 2, če pa n ni deljiv z 2, tedaj je $n - 1$ prav gotovo deljiv z 2 in spet je $M(n, 2)$ deljiv z 2. S podobnim premislekom se prepričamo, da velja Fermatov izrek v naslednjih primerih:

$$\begin{aligned} M(n, 3) &= n^3 - n = n(n - 1)(n + 1), \\ M(n, 5) &= n^5 - n = n(n - 1)(n + 1)(n - 2)(n + 1) + 5(n^3 - n), \\ M(n, 7) &= n^7 - n = n(n - 1)(n + 1)(n - 2)(n + 2)(n - 3)(n + 3) + 7(n^3 - n)(2n^2 - 5). \end{aligned}$$

Ta izrek sodi, kot sem že omenil, v posebno vejo matematike, imenovano teorija števil, ki se ukvarja predvsem z naravnimi in celimi števili. Začetki te matematične teorije segajo še v stari vek, saj je že staro grški matematik Diofant iz Aleksandrije zastavljal in reševal probleme, ki sodijo v to vejo. Fermat je imel v svoji knjižnici tudi neko izdajo Diofantove matematike iz leta 1621, ki jo je vneto prebiral in včasih tudi zapisoval vanjo na robove svoje ideje. In prav na robu te knjige so našli trditev, ki se je glasila nekako tako: *Če je n poljubno naravno število večje kot 2, potem ne obstajajo nobena taka cela števila x, y in z, da bi bila izpolnjena enačba $x^n + y^n = z^n$.* Ob tej trditvi je Fermat zapisal v isto knjigo, da pozna čudovit dokaz za to, da pa je na robu knjige premalo prostora, da bi lahko dokaz zapisal. S tem je Fermat zastavil matematičnemu svetu uganko, ki je doslej ni še nihče razrešil. Največji matematiki sveta so se ukvarjali s tem “poslednjim izrekom”, toda nihče ga doslej ni znal niti dokazati, niti ovreči! Jasno je, da pri $n = 2$ celo številske rešitve zgornje enačbe obstajajo! Rešitve so tedaj npr. $x = 3, y = 4, z = 5; x = 5, y = 12, z = 13; x = 7, y = 24, z = 25; x = 8, y = 15, z = 17$; itd., nove rešitve pa dobimo iz teh, če vsa tri števila množimo z istim faktorjem. Tako so rešitve tudi $x = 6, y = 8, z = 10; x = 9, y = 12, z = 15; x = 12, y = 16, z = 20$; itd. Vidimo, da je rešitev kar neskončno. Toda že za $n = 3, 4, 5$ in še za nekatera druga naravna števila n so matematiki dokazali, da zgornja enačba sploh nima rešitev.

To pa še ni vse. Fermat je trdil še mnogo več!

C.2 Blaise Pascal (1623-1662)

Po http://sl.wikipedia.org/wiki/Blaise_Pascal velja Blaise Pascal za enega največjih neuresničenih talentov v zgodovini matematike. S svojo izjemno matematično nadarjenostjo in s prav neverjetno intuicijo v geometriji bi ob nekoliko drugačnih življenjskih okoliščinah v znanosti zapustil najbrž daleč globlje sledi. Žal mu je dokaj ustvarjalne energije pobral njegov nenehen boj s krhkim zdravjem, razen tega pa se je v zrelih ustvarjalnih letih večinoma ukvarjal bolj z religioznimi vprašanji kot z znanostjo. Ostal je eden najkontroverznejših genijev v vsej zgodovini zahodne misli.

Že kot dete se je moral Pascal boriti z boleznimi, menda so se dolgo bali celo za njegovo življenje. V šolo kasneje prav zaradi svojega krhkega zdravja ni hodil, pač pa je imel domačega

učitelja. Njegov oče se je namreč bal, da bi se njegov zagnanec zaradi učenja v šoli preveč izčrpal. Domačemu učitelju je celo naročil, naj da dečku predvsem dobro humanistično izobrazbo s poudarkom na latinskem in grškem jeziku, izogiba pa se naj zahtevne matematike. Kljub temu se je Pascal odločil žrtvovati ves svoj prosti čas za študij geometrije. V vsega nekaj tednih je samostojno in na povsem samosvoj odkril več izrekov evklidske geometrije. Oče po vsem tem seveda ni mogel več skrivati navdušenja nad sinovim talentom. Podaril mu je *Evklidove Elemente* in mu priskrbel najboljše pogoje za matematično izobrazbo. Leta 1640 je Pascal kot komaj šestnajstletni mladenič napisal izvrstno razpravo iz projektivne geometrije z naslovom *Razprava o stožnicah* (*Essai pour les coniques*), dolgo vsega eno stran. Na njej se je predstavil z enim od svojih sploh največjih odkritij v matematiki – z znamenitim izrekom o šestkotniku, včrtanemu stožnici, kasneje znamen tudi pod priljubljenim imenom mistični heksagram. Ko je razprava prišla v roke Renéu Descartesu, se mu je zdela naravnost briljantna in menda dolgo ni mogel verjeti, da jo je resnično napisal tak mladenič. Podobno je bil kasneje leta 1676 ob njenem prebiranju začuden tudi nemški matematik Gottfried Wilhelm Leibniz. Žal razprava ni bila nikdar objavljena in danes velja za izgubljeno.

Pri osemnajstih je Pascal izumil in sestavil računski stroj Pascaline, prvi takšne vrste v zgodovini, ki je zmogel seštevati in odštevati. Z njim je pomagal očetu pri njegovi obsežni reorganizaciji davčnega sistema v mestu Rouen, kjer je dobil mesto glavnega davkarja. Leta 1649 je ustanovil podjetje za njihovo proizvodnjo, a so bile naprave predrage in premalo zanesljive. Le 7 se jih je ohranilo do današnjih dni. Nesrečni dogodek v novembru leta 1654, ko se je Pascal s konjsko vprego prevrnil in komaj ostal živ, pa je močno spremenil nadaljnji potek njegovega življenja. V njem je namreč zaznal skrivnostno opozorilo, da so njegove matematične dejavnosti vse prej kot po Božji volji. Odločil se je, da se poslej posveti izključno premišljevanju o veri in Bogu. Kmalu zatem se je pridružil skupini zavzetih kristjanov, ki se je zbirala okrog samostana Port Royal, in ji ostal zvest do smrti. Lotil se je pisanja zagovora krščanstva, a je njegovo delo ostalo nedokončano. Prijatelji so po smrti zbrali gradivo na lističih in ga izdali pod naslovom *Misli*.

C.3 Matematični Bernoulli

Redke družine so prispevale k matematiki toliko, kot Bernoullijevi iz švicarskega Basla. Kar sedem Bernoullijev iz treh generacij med leti 1680 in 1800 je bilo odličnih matematikov. Pet izmed njih je pomagalo graditi novo matematično teorijo verjetnosti.

Jacob (1654–1705) in **Johann** (1667–1748) sta bila otroka uspešnega švicarskega trgovca, vendar sta matematiko študirala proti volji svojega praktičnega očeta. Oba sta bila med najboljšimi matematiki svojega časa, vendar je bil Jacob tisti, ki se je osredotočil na verjetnost. Bil je prvi, ki je jasno videl idejo z dolgoročnimi povprečji kot način za merjenje slučajnosti (“Zakon velikih števil (neodvisnih poskusov)” je en od stebrov moderne teorije verjetnosti) in avtor “druge” knjige o verjetnosti z naslovom *Ars Conjectandi*, ki je bila objavljena šele po njegovi smrti l. 1713.

Johannov sin **Daniel** (1700–1782) ter Jacobov in Johannov nečak **Nicholas** (1687–1759) sta se prav tako ukvarjala z verjetnostjo. Nicholas je opazil, da lahko z verjetnostjo pojasnimo vzorec v rojstvu dečkov in deklic. Kljub temu da se je tudi sam upiral očetovim željam, je Johann želel, da bi njegov sin Daniel postal trgovec ali zdravnik, a Daniela to ni odvrnilo in je vseeno postal še en matematik iz družine Bernoullijev. Na področju verjetnosti se je ukvarjal s pravičnim določanjem cen v igrah naključij, poleg tega pa je dokazal učinkovitost

cepiva proti kozam.

Matematična družina Bernoullijev je, tako kot njihovi glasbeni sodobniki iz družine Bach, nenevaden primer nadarjenosti za določeno področje, ki se pokaže v zaporednih generacijah. Delo Bernoullijev je pomagalo verjetnosti, da je od svojega rojstva v svetu hazarda zrasla do spoštovanega orodja svetovne uporabnosti.

C.4 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 Jacob (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>

C.5 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.

C.6 Leonhard Euler (1707-1783)

Euler je bil verjetno en izmed najbolj plodovitih matematikov. Govorilo se je: “Euler je pisal matematiko tako neumorno kot človek diha.” Rodil se je leta 1707 v Baslu v Švici kot sin protestantskega duhovnika, ki je tudi sam študiral matematiko. Njegov genij se je razvil že zelo zgodaj. Obiskoval je univerzo v Baslu, kjer je kot 16 leten hkrati diplomiral iz umetnosti in magistriral iz filozofije. V času svojega bivanja v Baslu je imel srečo, da ga je en dan na teden inštruiral matematiko izjemni Johann Bernoulli. Sprva je na očetovo željo študiral teologijo, potem pa ga je z 18-imi prevzela matematika, ki jo je začel raziskovati. Vseeno pa je očetov vpliv ostal prisoten, tako da je bil Euler globoko pobožen. V različnih časovnih obdobjih je poučeval na Akademiji znanosti Sv. Petra (v Rusiji), Univerzi v Baslu in berlinski akademije znanosti. Eulerjeva energija za delo je bila praktično neomejena. Njegova zbrana dela obsegajo 60-80 volumnov/knjig velikega formata, hkrati pa se verjame, da je bilo veliko njegovih del izgubljenih. Posebej presenetljivo je, da je bilo kljub slepoti v zadnjih 17ih letih njegovega življenja, eno njegovih najbolj produktivnih obdobij. Njegov nezmotljiv spomin je bil fenomenalen. Njegova sposobnost reševanja problemov v glavi je bila takorekoč neverjetna. V glavi je izpeljal gibanje planetov s katerim se je ukvarjal Isaac Newton, enkrat pa je v glavi izpeljal zapletene izračune, da bi razrešil spor med dvema študentoma, katerih izračuni so se razlikovali na 50em decimalnem mestu. Eulerjev največji prispevek je bila sistematizacija matematike. Eulerjev genij je dal skladnost matematični pokrajni. Bil je prvi matematik, ki je v fizikalne probleme vpeljal vso moč analize. Prispeval je takorekoč na vsako področje matematike, kakor tudi teorije optike, gibanja planetov, elektrike, magnetizma in splošne mehanike.

Po članku M. Juvana in P. Potočnika, Najpomembnejši matematiki, Presek,....: Leonhard Euler se je rodil leta 1707 v Baslu v Švici. Tam je tudi študiral, sprva na očetovo željo teologijo, ker pa ga ta ni pritegnila, se je prepisal na matematiko. Študij je dokončal leta 1726, leta kasneje pa se je od pravil v Sankt Petersburg, glavno mesto carske Rusije, kjer je dobil službo na matematično-fizikalnem oddelku Akademije znanosti. Tam je ostal do leta

1741, ko se je na povabilo pruskega kralja Friderika II (Velikega) preselil v Berlin. V Berlinu je ostal 25 let, ko pa so se njegovi odnosi s kraljem poslabšali, se je leta 1766 vrnil v Sankt Petersburg. Tam je ostal do svoje smrti leta 1783. Euler je bil izredno plodovit znanstvenik. Napisal je preko 800 knjig, člankov in drugih matematičnih, pa tudi astronomskih in fizikalnih del. Čeprav je imel že pri tridesetih letih težave z vidom, po vrnitvi v Rusijo pa je popolnoma oslepel, je ostal izredno delaven vse do smrti. Imel je izvrsten spomin, kar mu je omogočalo, da je ob pomoči sodelavcev in sinov kljub slepoti nadaljeval z delom. Pomembno je prispeval k vsem dejavnostim tedanje matematike, zelo vplivni pa so bili tudi njegovi učbeniki, v katerih je sistematično zbral in uredil tedanje znanje z različnih področij matematike. Tako je npr. vpeljal pisavo $f(x)$ za vrednost funkcije v točki x , e za (Eulerjevo) število, ki je osnova naravnih logaritmov, π za razmerje med obsegom in premerom kroga, znak \sum za zapisovanje vsot, itd. Po njem se imenuje tudi Eulerjeva funkcija φ iz teorije števil (vrednost $\varphi(n)$ pove, koliko števil od 1 do n je tujih s številom n), Eulerjeva formula, ki v osnovni različici podaja zvezo med številom oglišč, robov in lic pri poliedrih, Eulerjevi sprehodi v teoriji grafov (to so sprehodi, ki vsebujejo vsako povezavo grafa natanko enkrat) idr.

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 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).

C.7 Joseph Louis Lagrange (1736-1813)

Lagrange je bil sin javnega uslužbenca in se je rodil v Torinu, Italiji. S 16imi leti je začel študirati matematiko kot samouk in pri 19ih je bil izbran za profesorja na kraljevski artilerijski šoli v Torinu. Naslednje leto je poslal Eulerju rešitve nekaterih slavnih problemov, pri tem pa je uporabil nove metode, ki so se kasneje razcvetale v vejo analize, ki se imenuje analiza varijacij. Te metode in Lagrangova uporaba pri problemih nebesne mehanike so bili tako izjemni, da so ga pri 25ih nekateri smatrali za največjega matematike tistih časov.

C.8 Pierre-Simon Laplace (1749–1827)

Po članku: Marija Vencelj, Pierre-Simon Laplace (1749-1827) — ob 250-letnici rojstva, *Presek* 26/4 (1998/99), 213–217.

Na predek in izpopolnjevanje matematike sta tesno povezana z blagostanjem države.

Napoleon 1

Pierre-Simon Laplace se je rodil 23. marca 1749 v kraju Beaumonten-Auge v francoski Normandiji. O njegovem otroštvu je malo znanega, saj se je zelo sramoval svojih preprostih staršev in napravil vse, da bi skril svoje kmečko poreklo. Bistremu dečku iz vaške šole so omogočili, da je od sedmega do šestnajstega leta kot zunanjji učenec obiskoval benediktinski kolegij v Beaumont-en-Auge. Nato je za dve leti odšel na univerzo v Caenu, kjer je odkril svoj matematični talent. Namesto, da bi ustregel staršem in nadaljeval šolanje na teološki fakulteti, je 1768. leta za vedno stresel s svojih čevljev beaumontski prah in od šel v Pariz osvojil matematični svet. Z d'A lembertovo pomočjo je postal profesor matematike na vojaški šoli v Parizu. Povprečno nadarjene kadete iz uglednih družin je poučeval geometrijo, trigonometrijo, elementarno analizo in statistiko, kar je bilo daleč pod njegovimi ambicijami in sposobnostmi, vendar mu je plača omogočila, da je lahko ostal v Parizu.

Kot je bilo pričakovati, je Laplace vso svojo energijo in sposobnosti usmeril v cilj, doseči matematični sloves. Po petih letih Pariza je, komaj 24-leten, postal dopisni član Akademije znanosti. Condorcet, ki je nekaj pred tem postal stalni tajnik akademije, je zapisal, da še nikoli do tedaj akademija ni prejela od tako mladega kandidata in v tako kratkem času toliko pomembnih člankov s tako različnih in zahtevnih področij.

V poznih sedemdesetih letih 18. stoletja se je Laplaceov sloves razširil tudi izven majhnega kroga matematikov, ki so lahko razumeli njegovo delo. Tako je leta 1785, star 36 let, postal redni član Akademije znanosti in bil od poznih osemdesetih let dalje ena vodilnih osebnosti akademije. V letu svoje izvolitve za rednega člana akademije je Laplace, kot član izpitne komisije za kraljevo artilerijsko šolo, srečal tudi osebo, ki je kasneje odločilno vplivala na njegovo javno delovanje. To je bil tedaj šestnajstletni kadet z imenom Napoleon Bonaparte.

Leta 1788 se je Laplace poročil z dvajset let mlajšo Marie-Charlotte de Courty de Romanges. Imela sta dva otroka. Skozi revolucijo je Laplace takorekoč pojezdil na konju in videl marsikaj, vendar relativno brez skrbi za življenje. Imel je nekaj pomembnih položajev. Bil je član Komiteja za uteži in mere. Sodeloval je pri organizaciji Ecole normal in Ecole polytechnique ter bil na obeh visokih šolah tudi profesor. Po revoluciji se je strastno posvetil politiki. Na tem področju se je izkazal kot pravi genij, saj se je dobro znašel v nemirnih političnih vodah in vedno zнал pluti s tokom. Zadnja leta življenja je prebil Laplace na svojem lepem posestvu v Arcueilu blizu Pariza. Umrl je po kratki bolezni, star 78 let.

Opisati Laplaceovo znanstveno delo bi bil zaradi njegove obsežnosti in zahtevnosti za Presek prehud zalogaj. Zato si ga oglejmo le v grobih črtah. Laplace je od vseh matematikov osemnajstega stoletja prišel najbliže temu, kar imenujemo uporabna matematika. Vendar moramo celo v njegovem primeru vzeti pojem uporabna v zelo omejenem pomenu. Ukvajal se je predvsem s teorijo verjetnosti in nebesno mehaniko in vprašanje je, za kako praktični lahko štejemo v njegovi dobi ti dve področji. Lahko pa z gotovostjo trdimo, da je Laplace v prvi vrsti videl v matematiki orodje, ki ga je genialno priredil za vsak posebni problem, ki se je pojavil. Bil je velik filozof, ki je želel spoznati naravo in je v ta namen izkoristil višjo matematiko. Vendar je tudi s strogo teoretičnega vidika njegov prispevek matematiki velik, predvsem v potencialni teoriji in teoriji verjetnosti. Ni pretirano reči, da teorija verjetnosti dolguje Laplaceu več kot kateremukoli drugemu matematiku.

Laplaceovo znanstveno življenje bi lahko razdelili na štiri etape, od katerih sta se prvi dve odvijali pod starim monarhističnim režimom in zadnji dve v čas u francoske revolucije - Napoleonovega režima in restavracije. V prvem obdobju od leta 1768 do 1778 je Laplace vzhajal na znanstvenem obzorju, pisal članke o problemih integralnega računa, matematični astronomiji, o vesolju, teoriji iger na srečo in aposteriorni verjetnosti. V tem obdobju rasti je ustvaril svoj stil, sloves in filozofski položaj. Oblikoval je določene matematične tehnike in zastavil program raziskav na dveh področjih, verjetnosti in nebesni mehaniki, na katerih je matematično deloval do konca življenja. V drugem obdobju (1778-1789) je ti dve področji obogatil s pomembnimi rezultati (po katerih je slaven), ki jih je kasneje vključil v svoji veliki deli **Nebesna mehanika** (*Mécanique céleste*) in **Analitična teorija verjetnosti** (*Theorie analytique des probabilités*). Uporabljene matematične postopke je večinoma sam uvedel in razvil. Med najpomembnejše sodijo rodovne funkcije, transformacija, ki po njem nosi ime Laplaceova, formula za razvoj determinante (ki se tudi imenuje po Laplaceu), variacija konstant za iskanje aproksimativnih rešitev diferencialnih enačb v astronomiji in posplošena gravitacijska funkcija. V tem času je Laplaceovo pozornost pritegnila tudi fizika. Deloma mu je prav sodelovanje z Lavoiserjem v topotni teoriji na strežaj odprlo vrata v vplivno znanstveno srenjo. V tretjem obdobju (1789-1805), času revolucije in vlade direktorija, je Laplace dosegel svoj vrh. V zgodnjih devetdesetih letih 18. stoletja je napisal obsežno serijo člankov o planetarni astronomiji in sodeloval pri pripravi metričnega merskega sistema. V drugi polovici devetdesetih let je bil najvplivnejša osebnost oddelka za eksaktne znanosti na Institut de France (ki je nadomestil bivšo Akademijo znanosti); imel je močan položaj v svetu Ecole polytechnique, od koder je izšla prva generacija matematičnih fizikov. V letih 1799 do 1805 so izšli prvi širje deli Nebesne mehanike, v katerih je posplošil zakone mehanike za njihovo uporabo pri obravnavi gibanja in števila nebesnih teles. V Laplaceovem delu četrtega obdobja (1805-1827) opazimo elemente vzpona in pojemanja intelektualne moči. Zreli, morda že tudi starajoči se Laplace, je skupaj z Bertholletom ustanovil neformalno šolo Societe d'Arcueil. Toda ta se ni ukvarjala z astronomijo. V središču njenega zanimanja je bila fizika: kapilarnost, teorija topote, optika delcev in hitrost zvoka. Čeprav je Societe d'Arcueil imela morda nekoliko preveč 'šolski' značaj, ni dvoma o njenem velikem prispevku k matematizaciji znanosti.

Po letu 1810 je Laplaceovo zanimanje spet pritegnila verjetnost. V Analitični teoriji verjetnosti, ki je izšla 1812. leta, je zbral in posplošil rezultate svojih zgodnjih raziskovalnih let in dodal pomembne novosti za njihovo uporabo, npr. metodo najmanjših kvadratov. Kasnejšim izdajam je dodal analizo verodostojnosti prič, izbora nepristrane sodne porote in volilnih teles ter napak pri statistični obdelavi geodetskih in meteoroloških podatkov.

V letih 1823-1825 je postopoma izšel peti del Nebesne mehanike, ki je dejansko samostojno

delo. Medtem ko vsebujejo prvi štirje deli Laplaceove rezultate, ki jih je objavil že pri stari akademiji, je vsebina petega dela nov pomemben fizikalni material, ki je nastal v tem obdobju. Laplace je obe svoji veliki deli pospremil s preprosto razlago, namenjeno inteligentni francoski (ne ozko strokovni) publiku. Razlaga vesolja (*Exposition du systém du monde*) je izšla leta 1796 pred izidom Nebesne mehanike, Filozofska razprava o verjetnosti (*Essai philosophique sur les probabilités*) pa leta 1814 kot uvod v drugo izdajo Analitične teorije verjetnosti.

Prav je, da opišemo Laplacea tudi po drugi strani. Čeprav sodi med matematike francoske revolucije, dejansko ni sodeloval v revolucionarnih aktivnostih. Kakor je imel zelo stroga merila glede znanstvenih resnic, kaže, da je bil v politiki brez pomislekov. To ne pomeni, da se je plašno umaknil v ozadje. Brez strahu se je odkrito družil z znanstvenimi kolegi, ki so bili v kriznem obdobju politično sumljivi. Pravijo celo, da sta se skupaj z Lagrangeom izognila giljotini samo zato, ker ju je Napoleon potreboval za izračunavanje poti topovskih izstrelkov in pripravo zalog za pohod v Egipt. Laplace pa je bil tudi grob koristolovec. Po vsakem padcu vlade je dobil boljši položaj. Prav nič težko mu ni bilo čez noč preleviti se iz divjega republikanca v stastnega monarhista, ko se je Napoleon proglašil za cesarja. Napoleonova odlikovanja vseh vrst so krasila nestanovitna Laplaceova prsa, celo Veliki križ častne legije in Orden prijateljstva, Napoleon ga je imenoval tudi za grofa cesarstva. Zasedal je številne ugledne položaje, za kratek čas ga je Napoleon, ki je bil velik občudovalec znanstvenikov, celo postavil za notranjega ministra.² Ko pa je Napoleon padel, je Laplace brez oklevanja podpisal listino, s katero je obsodil svojega dobrotnika in nemudoma svojo vdanost prenesel na Ludvika XVIII. Poplačan je bil s sedežem v Zgornjem domu, k čemur je sodil tudi plemiški naziv markiz de Laplace, in bil imenovan za predsednika odbora za reorganizacijo École polytechnique.

C.9 Johann Carl Friedrich Gauss (1777–1855)

Še kot najstnik je rešil 2000 let odprt problem konstrukcije 17-kotnika samo s šestilom in ravnilom. Njegovo delo “*Disquisitiones Arithmeticae*” (1801) je verjetno najpomembnejša knjiga na področju teorije števil, ki je bila kdajkoli napisana. Gojil je strast do popolnosti, ki ga je gnala k poliranju in ponovni obdelavi namesto, da bi objavil nedokončane izdelke – njegov moto je bil “Malo, a zrele” tako da je veliko njegovih odkritij ostalo skritih v dnevnikih, ki so ostali neobjavljeni za časa njegovega življenja. Med njegovimi dosežki so Gaussova zvonasta krivulja napake (osnova v teoriji verjetnostnega računa), geometrijska interpretacija kompleksnih števil in njihova osrednja vloga v matematiki, razvoj metod za karakterizacijo ploskev glede na to katere krivulje vsebujejo ter odkritje ne-Evklidskih geometrij 30 let pred objavo s strani drugih. Odkril je heliotrope, bifilar magentometer in elektrotelegraf.

Po članku M. Juvana in P. Potočnika, Najpomembnejši matematiki, Presek, ...: Johann Carl Friederich Gauss se je rodil leta 1777 v Braunschweigu v Nemčiji. Že kot sedem letni deček je opozoril na svoje matematične sposobnosti. Znana je anekdota o mladem Gaussu, ki je v šoli v nekaj sekundah izračunal vsoto vseh naravnih števil med 1 in 100. Seštevanja se seveda ni lotil neposredno (kot je pričakoval učitelj), temveč je opazil, da lahko teh 100 števil združimo v 50 parov, tako da je vsota vsakega para enaka 101. Gaussova kariera je povezana z univerzo v Gottingenu. Tam je namreč leta 1795 pričel svoj matematični študij. Gottingen je zapustil leta 1798 in kmalu zatem objavil svoje prvo pomembno odkritje: konstrukcijo

²Laplace ni pokazal nobene nadarjenosti za uradovanje in Napoleon se je bojda celo rogal na njegov račun, “da je prinesel duh neskončno majhnega v upravljanje državnih zadev.”

pravilnega 17-kotnika s šestilom in ravnih ter kriterij, ki pove, za katera naravna števila n je konstrukcija pravilnega n -kotnika s šestilom in ravnih možna. Svojo doktorsko disertacijo s področja algebре je leta 1799 predložil univerzi v Helmstedtu. V Cottingen se je za vedno vrnil leta 1807, ko je bil imenovan za direktorja tamkajšnjega observatorija. To imenovanje ga je rešilo finančnih skrbi ob vzdrževanju številne družine in mu hkrati dalo dobre možnosti za stalno znanstveno delo. Najpomembnejša Gaussova dela sodijo na področje algebре in teorije števil, ukvarjal pa se je tudi z drugimi matematičnimi področji. Veliko svojega časa je posvečal tudi povsem praktičnim problemom. Med drugim je računal orbite nebesnih teles ter se ukvarjal z Zemljinim magnetizmom in geodezijo. Umrl je leta 1855 v Cottingenu v 78. letu starosti.

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 ‘till 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 superficies 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.

C.10 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,

C.11 Augustin Louis Cauchy (1789-1857)

Cauchyjeva zgodnja izobrazba je bila pridobljena od njegovega očeta, odvetnika in mojstra klasike. L. 1805 se je vpisal na L’Ecole Polytechnique, da bi študiral inžinirstvo, zaradi slabega zdravja pa so mu svetovali, da se skoncentriira na matematiko. Njegova velika dela so se pričela 1811 z serijo briljantnih rešitev nekaterih težkih odprtih problemov. L. 1814 je napisal obravnavo integralov, ki je postala osnova za moderno kompleksno analizo; l. 1816 je sledil članek o širjenju valov v tekočinah, za katerega je dobil nagrado francoske akademije;

l. 1822 je napisal članek, ki predstavlja osnovo moderne teorije elastičnosti. Cauchyjevi matematični prispevki v naslednjih 35 letih so bili briljantni ter so obsegali preko 700 člankov (26 vol.). S Cauchyjem se je začela doba moderne analize. V matematiko je vpeljal standarde preciznosti in strogosti, o katerih se Leibnizu in Newtonu niti sanjati ni moglo.

<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.

C.12 Arthur Cayley (1821-1895)

Je en najpomembnejših matematikov 19. stoletja. Objavil je okoli 1000 del in ga uvrščamo med najplodovitejše matematike vseh časov.

C.13 Emile Borel (1871–1956)

French mathematician educated at the École Normale Supérieure (1889-1892) who received is Docteur 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 than thirty books and nearly 300 papers, including a number of popular works of high scientific quality, and devoted more than fifty papers to history and philosophy of sciences.

<http://scienceworld.wolfram.com/biography/Borel.html>

C.14 William Sealy Gosset (1876-1937)

Leta 1908 je pod psevdonimom 'Student' objavil članek, v katerem je vpeljal Studentovo t -porazdelitev. He is famous as a statistician, born in Canterbury, England. attended Winchester College. After graduating chemistry and mathematics at New College, Oxford in 1899, he joined the Dublin brewery of Arthur Guinness & Son. Guinness was a progressive agro-chemical business and Gosset would apply his statistical knowledge both in the brewery and on the farm—to the selection of the best yielding varieties of barley. Gosset acquired that knowledge by study, trial and error and by spending two terms in 1906–7 in the biometric laboratory of Karl Pearson. Gosset and Pearson had a good relationship and Pearson helped Gosset with the mathematics of his papers. Pearson helped with the 1908 papers but he had little appreciation of their importance. The papers addressed the brewer's concern with small samples, while the biometrician typically had hundreds of observations and saw no urgency in developing small-sample methods. Another researcher at Guinness had previously published a paper containing trade secrets of the Guinness brewery. To prevent further disclosure of confidential information, Guinness prohibited its employees from publishing any papers regardless of the contained information. This meant that Gosset was unable to publish his works under his own name. He therefore used the pseudonym Student for his publications to avoid their detection by his employer. Thus his most famous achievement is now referred to as Student's t -distribution, which might otherwise have been Gosset's t -distribution.

Gosset had almost all of his papers including The probable error of a mean published in Pearson's journal Biometrika using the pseudonym Student. However, it was R. A. Fisher who appreciated the importance of Gosset's small-sample work, after Gosset had written to him to say I am sending you a copy of Student's Tables as you are the only man that's ever likely to use them! Fisher believed that Gosset had effected a "logical revolution". Ironically the t -statistic for which Gosset is famous was actually Fisher's creation. Gosset's statistic was $z = t/\sqrt{n-1}$. Fisher introduced the t -form because it fit in with his theory of degrees of freedom. Fisher was also responsible for the applications of the t -distribution to regression.

Gosset's interest in barley cultivation led him to speculate that design of experiments should aim, not only at improving the average yield, but also at breeding varieties whose yield was insensitive (robust) to variation in soil and climate. This principle only occurs in the later thought of Fisher and then in the work of Genichi Taguchi in the 1950s. In 1935, he left Dublin to take up the position of Head Brewer, in charge of the scientific side of production, at a new Guinness brewery at Park Royal in North West London. He died in Beaconsfield, England of a heart attack.

Gosset was a friend of both Pearson and Fisher, an achievement, for each had a massive ego and a loathing for the other. Gosset was a modest man who cut short an admirer with the comment that "Fisher would have discovered it all anyway."

C.15 Sir Ronald A. Fisher (1890-1962)

Ideje in metode, ki jih danes poznamo pod imenom *statistika*, so v devetnajstem in dvajsetem stoletju izumili ljudje, ki so se ukvarjali s problemi, pri katerih je bilo potrebno analizirati velike količine podatkov. Astronomija, biologija, družboslovne vede in celo geodezija lahko trdijo, da so igrale pomembno vlogo pri rojstvu statistike. Če pa si kdo zasluži naziv "oče statistike", je to Sir Ronald A. Fisher. Fisherjevi zapiski so opredelili statistiko kot posebno področje proučevanja, katerega metode lahko uporabimo v številnih panogah. Sistematisiral je matematično teorijo statistike in izumil številne nove metode. Slučajeni primerjalni eksperiment je najbrž njegov največji dosežek.

Kot druge statistične pionirje so tudi Fisherja gnale zahteve praktičnih problemov. Od leta 1919 naprej je delal na področju kmetijstva na terenu v Rothamstedu v Angliji. Kako naj razporedimo sajenje različnih vrst pridelka ali pa uporabo različnih gnojil, da jih lahko primerjamo? Ker se rodovitnost in druge lastnosti spreminjajo, ko se premikamo po polju, so eksperimentatorji uporabljali zamotane vzorce, ki so posnemali šahovnico. Fisher je imel boljšo idejo: namenoma uredimo parcele naključno.

C.16 Egon Sharpe Pearson (1895-1980)

Egon Sharpe Pearson (Hampstead, 11 August 1895 – Midhurst, 12 June 1980) was the only son of Karl Pearson, and like his father, a leading British statistician. He went to Winchester School and Trinity College, Cambridge, and succeeded his father as professor of statistics at University College London and as editor of the journal Biometrika. Pearson is best known for development of the Neyman-Pearson lemma of statistical hypothesis testing. He was President of the Royal Statistical Society in 1955–56 and was awarded its Guy Medal in Gold in 1955. He was awarded a CBE in 1946. He was elected a Fellow of the Royal Society in Mar 1966. His candidacy citation read: "Known throughout the world as co-author of the Neyman-Pearson theory of testing statistical hypotheses, and responsible for many important contributions to problems of statistical inference and methodology, especially in the development and use of the likelihood ratio criterion. Has played a leading role in furthering the applications of statistical methods - for example, in industry, and also during and since the war, in the assessment and testing of weapons."

C.17 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.

C.18 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 Poincaré.

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.

C.19 Paul Erdős (1913-1996)

<http://www.maa.org/mathland/mathland.html>

<http://www.britannica.com/EBchecked/topic/191138/Paul-Erdos:PaulHoffman>

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."

Dodatek D

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")`

D.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 izvajjanju 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()`.

D.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.

D.3 Najosnovnejše o spremenljivkah

Prireditev: `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))` izbriše 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`.

D.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.

D.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`.

D.6 Podatkovne strukture

D.6.1 Vektorji

Primer: Konstrukcija vektorja: `c(7, 8, 9)` da isto kot `7:9` ali `seq(7, 9)`. \diamond

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")`. \diamond

Če ni določeno drugače, ukaz `c` izpusti vrednosti `NULL`. Z ukazom `[...]` dobimo vse 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.

D.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.
- Nasprost ukaz `array` deluje na poljubnorazsežnih tabelah. V drugem argumentu določimo številke razsežnosti, ki ostanejo.
- Ukaz `mapply(f, c(x11, ..., x1n), ..., c(xm1, ..., xmnn))` vrne `c(f(x11, ..., xm1), ..., f(x1n, ..., xmnn))`.
- 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)).`

D.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.

D.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 nastavitev `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", "cetrta"), c("prvi", "drugi", "tretji")))`

Vrsticam in stolpcem lahko imena pripomemo ali spremojemo tudi naknadno, recimo: `dimnames(x) <- list(c("prvi", "drugi", "tretji"))`. ◇

POZOR! Vektor ni eno-razsežna tabela, zato mu ne moremo pripometi nastavitev `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 `[...]`. Sveda pa lahko še vedno uporabimo tudi številke.

Oznake komponent upoštevata tudi ukaza `plot` in `boxplot`.

D.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 nastavljamo z ukazom `names`.

D.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 nastavitevijo `levels`:

`table(factor(vektor, levels=vrednosti))`.

Pri vektorju s predpisanimi vrednosti 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.

D.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 nastavljamo z ukazom `names`.

Primer: `nase_zene <- nasi_mozje[,c(žena", "starost", "ime", "st_otrok")]`
`names(nase_zene) <- c("ime", "starost", "moz", "st_otrok")`



D.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.

D.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.

D.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 nastavitevijo `sep`. V nasprotju z ukazom print ukaz `sep` ne zaključi 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 zavitih oklepajev in vrne vrednost zadnjega.

D.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 priteja 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.

D.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.

D.8 Knjižnice z dodatnimi funkcijami

Ukaz `library()` vrne seznam knjižnic, ki so na voljo. Ukaz `library(knjižnica)` naloži ustrezno knjižnico.

D.9 Vhod in izhod

D.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 nastavitevami, 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 nastavitevijo `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 nastavitev **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 posevnici 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.

D.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).
```

◇

D.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 narekova je. 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`.

D.9.4 Izvoz grafike

Datoteko najprej odpremo s klicem ustreznega gonilnika, npr. `postscript(datoteka)` ali `pdf(datoteka)`. Nato kličemo ustreerne ukaze, npr. `plot`. Nazadnje datoteko zapremo z ukazom `dev.off()`.

D.10 Verjetnostne porazdelitve

Porazdelitve v R-u računamo z ukazom: `predpona+porazdelitev(vrednost, parametri)`. Parametri so opcije ustreznih funkcij, 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 binoomsko `Bi(10, 0.4)`. \diamond

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.

D.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 nastavitev: `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)).
```

D.12 Statistično sklepanje

D.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 (ozioroma 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 = β).
```

ali tudi:

```
binom.test(c(k, n - k), conf.level = β).
```

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 domneva:

izbira `"two.sided"` ali `"t"` ali izpustitev določila pomeni, da alternativna domneva trdi $p \neq p_0$.

Izbira `"less"` ali `"l"` pomeni, da alternativna domneva trdi $p < p_0$. Izbira `"greater"` ali `"g"` pa pomeni, da alternativna domneva trdi $p > p_0$.

D.12.2 Primerjava verjetnosti dveh poskusov/deležev v dveh populacijah

Testiramo ničelno domnevo, 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 domneve $p_1 = p_2$ izvedemo z ukazom:

```
prop.test(c(k1, k2), c(n1,n2), alternative = "two.sided" ali "less" ali "greater").
```

Določilo `alternative` pove, kaj je alternativna domneva:

- Izbira `"two.sided"` ali `"t"` ali izpustitev določila pomeni, da alternativna domneva trdi $p_1 \neq p_2$.
- Izbira `"less"` ali `"l"` pomeni, da alternativna domneva trdi $p_1 < p_2$.
- Izbira `"greater"` ali `"g"` pa pomeni, da alternativna domneva trdi $p_1 > p_2$.

D.12.3 Primerjava verjetnosti več poskusov/deležev več populacijah

Testiramo ničelno domnevo, 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 uspelih poskusov posamezne vrste oziroma števila enot iz posamezne populacije z določeno lastnostjo pa vektor k . Test izvedemo z ukazom: `prop.test(k, n)`.

D.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 = β).
```

Hipotezo, da je $\mu = \mu_0$, testiramo z ukazom:

```
t.test(v, alternative = "two.sided" ali "less" ali "greater", μ = μ₀).
```

Privzeta vrednost določila μ je 0.

Določilo `alternative` pove, kaj je alternativna domneva:

- Izbira `"two.sided"` ali `"t"` ali izpustitev določila pomeni, da alternativna domneva trdi $\mu \neq \mu_0$.
- Izbira `"less"` ali `"l"` pomeni, da alternativna domneva trdi $\mu < \mu_0$.
- Izbira `"greater"` ali `"g"` pa pomeni, da alternativna domneva trdi $\mu > \mu_0$.

D.12.5 Test mediane

Testiramo domnevo, 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 > μ₀ })), sum(sapply(v, function(x) { x != μ₀ })), 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").
```

D.12.6 Primerjava porazdelitev dveh spremenljivk

1. Dve (približno) normalni spremenljivki 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 domnevo, 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 domneva 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 domneva trdi $\mu_1 \neq \mu_2$.
- Izbira `"less"` ali `"l"` pomeni, da alternativna domneva trdi $\mu_1 < \mu_2$.
- Izbira `"greater"` ali `"g"` pa pomeni, da alternativna domneva 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 domneva:

- izbira `"two.sided"` ali `"t"` ali izpustitev določila pomeni, da alternativna domneva trdi, da sta porazdelitvi različni.
- Izbira `"less"` ali `"l"` pomeni, da alternativna domneva trdi, da je prva spremenljivka stohastično večja od druge.
- Izbira `"greater"` ali `"g"` pa pomeni, da alternativna domneva 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 domnevo, 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)`.

D.12.7 Koreliranost

Zanima nas koreliranost dveh spremenljivk na isti populaciji. Vzamemo vzorec, vrednosti prve spremenljivke naj tvorijo vektor v , vrednosti druge pa vektor w .

Interval zaupanja za koreacijski 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 koreacijski koeficient nas zanima.

- Izbira `"pearson"` ali `"p"` ali opustitev določila pomeni običajni Pearsonov koreacijski koeficient, ki je primeren za številske spremenljivke, porazdeljene (približno) normalno.
- Izbira `"spearman"` ali `s` pomeni Spearmanov,
- izbira `"kendall"` ali `k` pa Kendallov koreacijski koeficient;

Slednja testa sta primerna za urejenostne spremenljivke. Spearmanov koreacijski koeficient (ρ) je lažji za računanje, o Kendallovem koeficientu (τ) pa dosti statistikov meni, da je verodostojnejši.

Hipotezo, da sta spremenljivki nekorelirani (ozioroma 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 domneva:

- Izbira `"two.sided"` ali `t` ali izpustitev določila pomeni, da alternativna domneva trdi, da sta spremenljivki odvisni (ozioroma korelirani).
- Izbira `"less"` ali `l` pomeni, da alternativna domneva trdi, da sta spremenljivki pozitivno asociirani (ozioroma korelirani).
- Izbira `"greater"` ali `g` pa pomeni, da alternativna domneva trdi, da sta spremenljivki negativno asociirani (ozioroma korelirani).

Računanje p -vrednosti za Kendallov koreacijski 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, 159
časovna vrsta, 193
škatla z brki, 118
štetje, 263
- algebra dogodkov, 24
analiza variance, 174
asimetrija, 98
- Bayes, T. (1702-1761), 267
Bayesov obrazec, 35
Bernoullijev obrazec, 38
Bernoullijeve zaporedje, 38
binomski simbol, 230
bivariantna analiza, 179
- Borel, E., 275
Borelove množice, 84
Buffon, G. (1707-1788), 19
- Cauchy, A.L. (1789-1857), 274
Cayley, A. (1821-1895), 275
cenilka
 asimptotično nepristranska, 137
 dosledna, 136
 nepristranska, 137
 točkovna, 136
centil, 118
Centralni limitni izrek, 102, 128
centralni limitni zakon, 102
- De Moivre, A. (1667-1754), 267
De Moivrov točkovni obrazec, 44
disperzija, 93
dogodek, 10
- elementaren, 12
enak, 11
gotov, 11
način, 11
nasproten, 12
nemogoč, 11
neodvisna, 29
nezdružljiva, 12
osnoven, 12
produkt, 12
sestavljen, 12
slučajen, 10
vsota, 11
- domneva
 enostavna, 156
 formalen postopek za preverjanje, 165
 neparametrična, 156
 sestavljena, 156
domneva (statistična), 155
- enakomerna zvezna, 57
Erdős, P. (1913-1996), 278
Erdőseva probabilistična metoda, 212
Euler, L. (1707-1783), 268
- faktorijel, 224
Fermat, P. (1601-1665), 264
Fisher, Sir Ronald, A. (1890-1962), 277
formula za popolno verjetnost, 34
frekvenca razreda, 115
 relativna, 115
funkcija, 224
 bijektivna, 224

- Gama, 62
injektivna, 224
karakteristična, 100
napake, 59
pogojna verjetnostna, 89
slučajnega vektorja, 87
surjektivna, 224
- Gauss, J.C.F. (1777–1855), 272
Gosset, W.S. (1876-1937), 276
gostota verjetnosti, 56
- histogram, 115
- interval zaupanja, 147
delež, velik vzorec, 153
delež, znani odklon, 153
kvocient varianc, majhen vzorec, 154
pričakovana vrednost z znanim odklonom, 150
pričakovana vrednost, majhen vzorec, 150
pričakovana vrednost, velik vzorec, 150
razlika deležev, velik vzorec, 153
razlika deležev, znan odklon, 153
razlika pričakovanih vrednosti, majhen vzorec, 151
razlika pričakovanih vrednosti, majhen vzorec, z enakima znanima odklonoma, 151
teoretična interpretacija, 149
varjanca, majhen vzorec, 154
za razliko pričakovanih vrednosti, ujemajoči
se pari, majhen vzorec, 152
za razliko pričakovanih vrednosti, ujemajočimera
se pari, velik vzorec, 152
za razliko pričakovanih vrednosti, veliki vzorci, 151
za razliko pričakovanih vrednosti, znani odkloni, 151
- Jacobijeva determinanta, 88
- koeficient
asimetrije, 121
Cramerjev, 181
Jaccardov, 182
koreacijski, 96
parcialne korelacije, 186
Pearsonov, 181, 182
Sokal Michenerjev, 182
sploščenosti, 121
Yulov, 182
- Kolmogorov, A. (1903-1987), 278
- kombinacije, 230
kombinatorika, 263
kontingenčna tabela, 27
konvolucija, 87
koreliranost, 95
kovariančna matrika, 97
kovarianca, 95
kumulativa, 114
kvantil, 98
kvartil, 98, 118
kvartilni razmik, 98
- Lagrange, J.L. (1736-1813), 270
Lagrangeova indentiteta, 239
Laplace, S.P. (1749–1827), 270
Laplaceov intervalski obrazec, 58
Latinski kvadrat, 216
list, 114
- mediana, 98, 117, 118
- lokacija, 117
oblika, 98, 120
razpršenost, 119
- metoda
drsečih sredin, 197

- momentov, 141
- najmanjših kvadratov, 189
- največjega verjetja, 142
- moč statističnega testa, 159
- modus, 117
- moment, 97
 - centralni, 97
 - vzorčni začetni, 136
 - začetni, 97
- napaka
 - 1. vrste, 159
 - 2. vrste, 159
- neenakost
 - Čebiševa, 101
 - Cauchyjeva, 238
- ogiva, 115
- osnovni izrek statistike, 124
- parameter, 110
- Pascal, B. (1623-1662), 265
- Pearson, E. S. (1895-1980), 277
- permutacija, 224
 - ciklična (cikel), 225
 - liha, 226
 - s ponavljanjem, 226
 - soda, 226
 - transpozicija, 225
- poddogodek, 11
- pogojna gostota, 90
- pogojna porazdelitvena funkcija, 89
- Poisson, S. (1781–1840), 274
- pojasnjena varianca (ANOVA), 191
- popoln sistem dogodkov, 13
- populacija, 3, 109
- porazdelitev
 - binomska, 48
 - eksponentna, 61
- Fisherjeva, 135
- frekvenčna, 113
- Gama, 62
- geometrijska, 52
- hi-kvadrat, 63
- hipergeometrijska, 54
- negativna binomska, 52
- normalna, 57
- Pascalova, 52
- Poissonova, 51
- polinomska, 73
- razlike vzorčnih deležev, 146
- razlike vzorčnih povprečij, 145
- standardizirana normalna, 58
- Studentova, 134
- vzorčne disperzije, 131
- vzorčnega povprečja, 131, 143
- vzorčnih deležev, 144
- porazdelitvena funkcija, 46
 - robna, 71
- porazdelitveni zakon, 45
- poskus, 10
 - dvo-stopenjski, 34
- povprečje, 126
 - populacije, 117
 - uteženo, 47
 - vzorčno, 117
- pravilo vsote, 17
- preverjanje domneve
 - dlež, 172
 - porazdelitev spremenljivke, 176
 - povprečje razlik, 171
 - pričakovana vrednost, 166
 - razlika deležev, 172
 - razlika pričakovanih vrednosti, 169
 - variance, 174
- pričakovana vrednost, 91
- slučajnega vektorja, 96

- prikaz
krožni, 113
Q-Q, 118
stolpčni, 113
pristranost, 137
- Ramsey, F. P. (1903-1930), 277
Ramseyjeva teorija, 209
ranžirana vrsta, 113
rang, 113
razbitje, 34
razpon, 115, 117
razred, 113
regresija, 186
regresijska premica, 189
 druga, 190
 prva, 190
regresijski koeficient, 187
relativna frekvenca, 13
- skalarni produkt, 238
skoraj gotova konvergenca, 100
slučajna odstopanja, 174
slučajna spremenljivka, 45
 diskretna, 46
 zvezna, 46, 56
slučajni vektor, 71
slučajni vzorec
 enostavnji, 124
sploščenost, 98
spremenljivka, 109
 številska, 111
 absolutna, 112
 imenska, 111
 opisna, 111
 razmernostna, 112
 razmična, 112
 urejenostna, 112
sredina
 aritmetična, 237
 geometrična, 237
 harmonična, 237
 kvadratna, 237
 potenčna, 237
sredinska mera, 126
 geometrijska, 126
 povprečna vrednost, 126
standardizacija, 94, 122
standardna deviacija, 94
standardni odklon, 94
statistična enota, 109
statistika, 109, 110
 analitična, 3
 inferenčna, 109
 opisna, 3, 109, 111
steblo, 114
Stirlingov obrazec, 234
stopnja tveganja, 159
stopnja zaupanja, 159
stopnja značilnosti, 159
- teorija vzorčenja, 124
trend, 198
- urejeno zaporedje, 113
- varianca, 93
vektorski produkt, 238
verjetnost, 24
 definicija, 9
 klasična definicija, 15
 pogojna, 27
 statistična definicija, 14
verjetnostna funkcija, 72
verjetnostna konvergenca, 100
verjetnostna tabela, 72
verjetnostni model, 15
verjetnostni prostor, 24

vzorčna disperzija, 126

popravljena, 126

vzorčna statistika, 129

vzorčne ocene, 125

vzorčni odklon, 126

vzorčni prostor, 15

vzorčni razmah, 126

vzorčno povprečje, 126

vzorec, 3, 109, 123

mediana, 126

modus, 126

zakon velikih števil

šibki, 100

Bernoullijev, 59

krepki, 100

zamenjalna šifra, 205