

```

1  ##### OPŠTE INFORMACIJE O BAZI #####
2
3
4  # Baza podataka je preuzeta sa sledećeg linka:
5  # https://github.com/fivethirtyeight/data/tree/master/alcohol-consumption
6
7  # U datoteci baze možemo pronaći podatke o potrošnji alkohola u 193 zemlje tokom 2010.
  godine.
8  # Potrošnja alkohola je iskazana kroz četiri varijable:
9
10 # 1. prosečna potrošnja piva po glavi stanovnika,
11 # koja je izražena kroz broj konzumiranih limenki piva (cans of beer),
12 # u ovoj anlizi ova varijabla će nositi naziv potrošnja limenki piva;
13
14 # 2. prosečna potrošnja vina po glavi stanovnika,
15 # koja je izražena kroz broj konzumiranih čaša vina (glasses of wine),
16 # u ovoj anlizi ova varijabla će nositi naziv potrošnja čaša vina;
17
18 # 3. prosečna potrošnja žestokih pića po glavi stanovnika
19 # koja je izražena kroz broj konzumiranih čašica žestokih pića (shots of spirits),
20 # u ovoj anlizi ova varijabla će nositi naziv potrošnja čašica žestokih pića:
21
22 # 4. prosečan unos čistog alkohola po glavi stanovnika, koji je izražen u litrima,
23 # i u ovoj analizi ova varijabla će nositi naziv unos čistog alkohola.
24
25 # Važno je napomenuti da su čaše, čašice i limenke samo standardizovane mere
26 # za uobičajeni način konzumiranja datih pića a ne pravi podaci o načinu njihovog
  konzumiranja.
27 # Primera radi ako je neko popio dve litre piva to je iskazano preko 4 limenke,
28 # bez obzira na to da li je to pivo u realnosti konzumirano putem krigele, flaše, ili
  limenke.
29
30
31 ##### UČITAVANJE I UPOZNAVANJE SA KARATERISTIKAMA BAZE #####
32
33
34 # originalna baza
35 orig_baza <- read.csv("data/drinks.csv", stringsAsFactors = FALSE)
36
37 # baza za sređivanje
38 baza_n <- orig_baza
39
40 # Upoznavanje sa osnovnim karakteristikama baze
41 ncol(baza_n)
42 nrow(baza_n)
43
44 tail(baza_n, 10)
45 head(baza_n, 10)
46
47 summary(baza_n)
48 str(baza_n)
49
50 ##### SREĐIVANJE BAZE #####
51
52 # 1. Provera prisutnosti nedostajućih vrednosti.
53
54 is.na(baza_n) # klasičan način
55 sum(is.na(baza_n)) # pregledniji način, ako je vrednost nula onda ih nema
56
57 # Očigledno da nema nedostajućih vrednosti u NA formatu.
58 # Ali to ne znači da one stvarno ne postoje.
59
60 # Sve opservacije koje imaju vrednost nula,
61 # su potencijalno nedostajuće vrednosti ili rezultat lošeg merenja.
62 # Primera radi teško je zamisliti zemlju u kojoj stvarno nema nikakve potrošnje alkohola,
63 # i gde nijedan njen stanovnik ne konzumira ni kap piva, vina ili žestokih pića.
64 # Čak i kad bi postojala zabrana konzumiranja pića,
65 # to ne znači da bi se ona nužno poštovala u realnim okolnostima.
66 # Zbog toga smatram je potrebno eliminisati sve opservacije tog tipa.

```

```

67
68
69 # 2. Proces eliminacije
70
71 bazal <- replace(baza_n, baza_n == 0, NA)
72 baza <- na.omit(bazal)
73 summary(baza)
74
75 # Prvo su zamenjene sve 0 vrednosti sa NA vrednostima.
76 # Nakon toga je primenjena funkcija na.omit,
77 # koja eliminiše sve observacije koje imaju vrednos NA.
78 # Na ovaj način je eliminisatno 38 zemalja,
79 # sa potencijalno "problematičnim vrednostima".
80 # Broj elimnisanih zemalja može delovati kao preobiman,
81 # ali ne treba zaboraviti da da je,
82 # ova baza obuhvatila skoro sve zemlje sveta i da samim tim
83 # imamo idalje prilično dobar uzorak,
84 # u kome zasigurno nema nedostajućih vrednosti.
85
86
87 ##### ZANIMLJIVOSTI #####
88
89
90 # U ovom delu su istaknuti određeni zanimljiv podaci,
91 # koji nisu rezultat ozbiljne statističke analize,
92 # već prostog izvlačenja podataka iz baze.
93 # Pa samim tim nemaju status istraživačkog nalaza,
94 # već zanimljiv, a možda i korisne informacije.
95
96 # Za sve podatke koji će u ovom delu biti istaknuti,
97 # se podrazumeva da je potrošnja računata po glavi stanovnika.
98
99 # 1. Gde je najviše konzumiran alkohol tokom 2010. godine,
100 # U Srbiji, Bugarskoj ili Rusiji?
101
102 baza[baza$country == "Serbia", ]
103 baza[baza$country == "Bulgaria", ]
104 baza[baza$country == "Russian Federation", ]
105
106 # Očigledno je najviše konzumiran u Rusiji (11.5 litara čistog alkohola),
107 # u kojoj dominira potrošnja čašica žestokih pića (326).
108 # Mada je zanimljiv podatak da je u Srbiji
109 # konzumirano znatno više limenki piva (283) i čaša vina (127),
110 # u odnosu na Rusiju (247/73) i Bugarsku (231/94).
111
112
113 # 2. U kojoj zemlji je najviše konzumirano pivo tokom 2010. godine?
114 baza[(which.max(baza$beer_servings)), c(1,2)]
115 # Odgovor je pomalo začuđujući, reč je o Namibiji.
116
117 # 3. U kojoj zemlji je najviše konzumirano vino tokom 2010. godine?
118 baza[(which.max(baza$wine_servings)), c(1,4)]
119 # Odgovor ne bi trebalo da nas čudi, reč je o Francuskoj koja je poznata po vinima.
120
121 # 4. U kojoj zemlji su najviše konzumirana žestoka pića tokom 2010. godine?
122 baza[(which.max(baza$spirit_servings)), c(1,3)]
123 # Odgovor isto može da bude začuđujući, pošto je reč o Grenadi.
124
125 # 5. U kojoj zemlji je najviše konzumiran alkohol tokom 2010. godine.
126 baza[(which.max(baza$total_litres_of_pure_alcohol)), c(1,5)]
127 # Reč je o Belorusiji. Odgovor je verovatno očekivan.
128
129 # 6. Lista zemalja sa najmanje konzumiranim pivom tokom 2010. godine.
130 baza[baza$beer_servings == 1, c(1,2)]
131
132 # 7. Lista zemalja sa najmanje konzumiranim vinom tokom 2010. godine.
133 baza[baza$wine_servings == 1, c(1,4)]
134
135 # 8. Lista zemalja u kojoj su najmanje konzumirana žestoka pića tokom 2010. godine.

```

```

136 baza[baza$spirit_servings == 1, c(1,3)]
137
138 # 9. Za kraj možemo videti u kojoj zemlji je najmanje konzumiran alkohol tokom 2010.
godine?
139 baza[(which.min(baza$total_litres_of_pure_alcohol)), c(1,5)]
140 # Najmanje je konzumiran na Komorima.
141
142
143 ##### ISTRAŽIVAČKA PITANJA #####
144
145 # 1. Da li postoji statistički značajna veza između:
146 # - potrošnje: limenki piva i čaša vina
147 # - potrošnje: čašica žestokih pića i limenki piva
148 # - potrošnje: čaša vina i čašica žestokih pića
149
150 # 2. Koliko dobro mere potrošnje: limenki piva, čaša vina i čašica žestokih pića,
151 # predviđaju ukupan unos čistog alkohola.
152
153 # Za prvo istraživačko pitanje koristiće se korelacioni testovi,
154 # a za drugo metod višestruke regresije.
155
156
157 ##### PROVERA NORMALNOSTI DISTRIBUCIJE #####
158
159
160 # Ali pre primene pomenutih statističkih tehnika,
161 # treba proveriti normalnost distribucija datih varijabli,
162 # radi odabira adekvatnih statističkih testova.
163
164 # To radimo pomoću Shapiro-Wilk testa.
165
166 shapiro.test(baza$beer_servings)
167 shapiro.test(baza$spirit_servings)
168 shapiro.test(baza$wine_servings)
169 shapiro.test(baza$total_litres_of_pure_alcohol)
170
171 # p vrednost za sve varijable je znatno manja od 0.01
172 # usled čega odbacujemo nultu hipotezu u korist alterativne,
173 # i dolazimo do zaključka da date varijable nemaju normalnu distribuciju,
174 # pa je poželjno koristiti neparametraskе tehnike tamo gde je to moguće.
175
176 ## Histogrami ##
177
178 # Pomoću histograma možemo grafički
179 # predstaviti distribuciju datih varijabli,
180 # kao bismo imali stekli što bolji uvid o njima.
181
182 ## Instalacija i pozivanje ggplot paketa ##
183
184 # Instalacija nije neophodna ako je paket već instaliran,
185 # zato je i data u vidu komentara a ne komande, u narednom redu:
186 # install.packages("ggplot2")
187
188 # Isti princip će važiti za svako instaliranje paketa u ovoj analizi.
189
190 # pozivanje paketa je obavezan korak
191 library(ggplot2)
192
193 # 1. Histogram potrošnje lmenki piva.
194
195 ggplot(baza, aes( x = baza$beer_servings))+
196   geom_histogram(aes(y = ..density..), bins = 12, colour = "white", fill = "grey75") +
197   geom_density(aes(y = ..density..), colour = "blueviolet") +
198   ggtitle("Potrošnja piva u 2010. godini") +
199   xlab("broj konzumiranih limenki piva")+
200   ylab("gustina")
201
202 # Ovde možemo zapaziti da je distribucija zakrivljena ulevo,
203 # iz čega možemo izvesti zaključak,

```

```

204 # da je u velikom broju zemalja potrošnja limenki piva niska.
205
206 # 2. Histogram potrošnje čaša vina.
207
208 ggplot(baza, aes(x = baza$wine_servings))+
209   geom_histogram(aes(y = ..density..), bins = 12, colour = "white", fill = "grey75") +
210   geom_density(aes(y = ..density..), colour = "darkred") +
211   ggtitle("Potrošnja vina u 2010. godini") +
212   xlab("broj konzumiranih čaša vina") +
213   ylab("gustina")
214
215 # Sličan je slučaj kao i sa pivom,
216 # samo što je u ovom slučaju zakrivljenje još intezivnije.
217 # Iz čega možemo zaključiti da je potrošnja čaša vina
218 # u još većm broju zemalja niska nego što je to slučaj sa pivom.
219
220 # 3. Histogram potrošnje čašica žestokih pića.
221
222 ggplot(baza, aes(x = baza$spirit_servings)) +
223   geom_histogram(aes(y = ..density..), bins=12, colour = "white", fill = "grey75") +
224   geom_density(aes(y = ..density..), colour = "dodgerblue1") +
225   ggtitle("Potrošnja žestokih pića u 2010. godini") +
226   xlab("broj konzumiranih čašica žestokih pića") +
227   ylab("gustina")
228
229 # Isti slučaj kao i sa vinom i pivom,
230 # distribucija je zakrivljena ulevo
231 # što znači da u većini zemalja imamo,
232 # nisku potrošnu čašica žestokog pića.
233
234 # 4. Histogram unosa čistog alkohola.
235
236 ggplot(baza, aes(x = baza$total_litres_of_pure_alcohol)) +
237   geom_histogram(aes(y = ..density..), bins = 12, colour = "white", fill = "grey75") +
238   geom_density(aes(y = ..density..), colour = "slateblue") +
239   ggtitle("Unos čistog alkohola u 2010. godini") +
240   xlab("čist alkohol izražen u litrima") +
241   ylab("gustina")
242
243 # Zakrivljenje je znatno manje izraženo nego u predhodnim varijablama.
244 # U ovom slučaju imamo velki broj,
245 # kako zemalja sa niskim vrednostima unosa čistog alkohola,
246 # tako i zemalja sa srednjim vrednostima (opseg 5-10),
247 # dok je najmanji broj zemalja sa izrazito visokim vrednostima ove varijalbe.
248
249
250 ##### KORELACIJA #####
251
252 pvs <- data.frame(baza$beer_servings, baza$spirit_servings, baza$wine_servings)
253
254 # Matrica za korelaciju koja je sastavljena od
255 # varijabli nad kojima želimo da primenimo korelacione testove.
256
257
258 # Instalacija i pozivanje paketa za korelaciju
259 # install.packages("Hmisc")
260 library("Hmisc")
261
262 ## Korelacioni testovi nad matricom ##
263
264 rcorr(as.matrix(pvs), type = c("spearman"))
265
266 n <- rcorr(as.matrix(pvs), type = c("spearman"))
267 print(n$P, digits = 15) # tačniji prikaz p vrednosti
268
269 ## Nalazi korelacionih testova ##
270
271 # 1. Veza između potrošnje limenki piva i potrošnje čaša vina,
272 # istražena je pomoću Spirmanovog ro koeficijenta korelacije,

```

```

273 # izračunata je pozitivna korelacija srednje jačine između dve promenjive,
274 # r = 0.62, n = 155, p < 0,01
275 # a to znači da sa porastom potrošnje limenki piva,
276 # raste i potrošnja čaša vina.
277
278 # 2. Veza između potrošnje čašica žestokog pića i potrošnje limenki piva
279 # istražena je pomoću Spirmanovog ro koeficijenta korelacije,
280 # izračunata je pozitivna korelacija srednje jačine između dve promenjive,
281 # r = 0.50, n = 155, p < 0,01
282 # a to znači da sa porastom potrošnje limenki piva,
283 # raste i potrošnja čašica žestokih pića.
284
285 # 3. Veza između čaša vina i čašica žestokog pića
286 # istražena je pomoću Spirmanovog ro koeficijenta korelacije,
287 # izračunata je pozitivna korelacija srednje jačine između dve promenjive,
288 # r = 0.39, n = 155, p < 0,01
289 # a to znači da sa porastom potrošnje čaša vina,
290 # raste i potrošnja čašica žestokih pića.
291
292
293 ## Vizualizacija korelacije u ggplot-u (dijagrami raspršenosti) ##
294
295 # Instalacija i pozivanje ggplot-a
296 # install.packages("ggplot2")
297 library(ggplot2)
298
299 # Prikaz dijagrama
300
301 # 1. Dijagram korelacije potrošnje čaša vina i limenki piva
302
303 ggplot(baza, aes(x = baza$beer_servings, y = baza$wine_servings)) +
304   geom_point(size = 3, shape = 2, colour = "blue") +
305   ggtitle("Korelacija potrošnje limenki piva i čaša vina") +
306   xlab("potrošnja limenki piva") +
307   ylab("potrošnja čaša vina")
308
309 # 2. Dijagram korelacije potrošnje čašica žestokih pića i limenki piva
310
311 ggplot(baza, aes(x = baza$beer_servings, y = baza$spirit_servings)) +
312   geom_point(size = 3, shape = 1, colour = "red3")+
313   ggtitle("Korelacija potrošnje limenki piva i čašica žestokih pića") +
314   xlab("potrošnja limenki piva") +
315   ylab("potrošnja čašica žestokih pića")
316
317 # 3. Dijagram korelacije potrošnje čaša vina i čašica žestokih pića
318
319 ggplot(baza, aes(x = baza$wine_servings, y = baza$spirit_servings)) +
320   geom_point(size = 3, shape = 0, colour = "purple") +
321   ggtitle("Korelacija potrošnje čaša vina i čašica žestokih pića") +
322   xlab("potrošnja čaša vina") +
323   ylab("potrošnja čašica žestokih pića")
324
325 ## Vizualizacija u ggcorrplot-u (kvadrati i krugovi) ##
326
327 # Instaliranje i pozivanje datog paketa paketa
328 # install.packages("ggcorrplot")
329 library(ggcorrplot)
330
331 #Sređivanje imena kolona i matrica za korelaciju
332 names(pvs) = c("potrošnja krigli piva", "potrošnja čašica žestokih pića", "potrošnja
čaša vina")
333 viz <- cor(pvs, method = "spearman")
334
335 #Vizalizacija korelacije
336 ggcorrplot(viz, lab = TRUE) # Prvi način, kvadrat
337 ggcorrplot(viz, method = "circle") # drugi način krug
338
339
340 ##### REGRESIJA #####

```

```

341
342
343 # Korelacioni testovi predstavljaju
344 # prvi korak za izradu, regresionih modela
345
346
347 # Zato prvo treba napraviti korelacionu matricu
348 # a nakon toga sprovesti korelacione testove.
349
350 reg <- data.frame(baza$total_litres_of_pure_alcohol, baza$beer_servings,
baza$wine_servings, baza$spirit_servings)
351
352 # Instalacija i pozivanje paketa za korelaciju
353 # install.packages("Hmisc")
354 library("Hmisc")
355
356 ## Korelacioni testovi ##
357
358 rcorr(as.matrix(reg), type = c("spearman"))
359 m <- rcorr(as.matrix(reg), type = c("spearman"))
360 print(m$P, digits = 5)
361
362
363 ## Nalazi ##
364
365 # Spirmanov ro koeficijent korelacije,
366 # je pokazao da postoji pozitivna korelacija između varijable
367 # ukupnog unosa čistog alkohola i svih drugih varijabli potrošnje alkohola.
368 # Najjača je korelacija sa potrošnjom limenki piva (r = 0.85).
369 # Dok je sa potrošnjom čaša vina (r = 0.70)
370 # i čašica žestokog pića (r = 0.66) ona umerene jačine.
371 # Iz svega pomenutog možemo izvesti zaključak
372 # da sa porastom potrošnje limenki piva, čaša vina i čašica žestokih pića,
373 # raste i ukupan unos čistog alkohola.
374 # Ovakav nalaz je logičan i očekivan.
375 # Sve pomenute veze su statistički značajne p < 0.01
376
377 # Iz ovakvih nalaza možemo zaključiti da je najsvrsishodnije
378 # napraviti tri regresiona modela: jedan sa svim varijablama,
379 # drugi sa pivom i vinom i treći samo sa pivom.
380
381 # Izrada regresionih modela:
382
383 # 1. Prvi model, sve tri varijable:
384
385 Model1 <- lm(baza$total_litres_of_pure_alcohol ~ baza$wine_servings +
baza$spirit_servings + baza$beer_servings )
386 summary(Model1)
387
388 # Model je odličan, pošto objašnjava 88% varijabiliteta ukupnog unosa čistog alkohola.
389 # r^2 = 0.88, F = 385, p < 0.01 (za ceo model i sve koeficijente)
390
391 # 2. Drugi model, bez žestokih pića:
392
393 Model2 <- lm (baza$total_litres_of_pure_alcohol ~ baza$wine_servings +
baza$beer_servings)
394 summary(Model2)
395
396 # I ovaj model je prilično dobar pošto objašnjava 75% varijabiliteta
397 # ukupnog unosa čistog alkohola, sa dve varijable
398 # r^2 = 0.75, F = 234.5, p < 0.01 (za ceo model i sve koeficijente)
399
400 # 2.1. Vizualizacija drugog modela
401
402 ggplot (Model2, aes (y=baza$total_litres_of_pure_alcohol, x=baza$beer_servings,
color=baza$wine_servings)) +
403   geom_point (size = 4) +
404   stat_smooth(method = "lm", se = FALSE, colour = "red", size = 1.3 )+
405   ggtitle("Model broj 2") +

```

```

406     xlab("potrošnja limenki piva") +
407     ylab("ukupan unos čistog alkohola") +
408     labs(color = "potrošnja čaša vina")
409
410 # 3. Treći model, samo limenke piva:
411
412 Model3 <- lm (baza$total_litres_of_pure_alcohol ~ baza$beer_servings)
413 summary(Model3)
414
415 # Ovaj model je takođe veoma dobar pošto objašnjava 66 % varjabiliteta
416 # ukupnog unosa čistog alkohola sa samo jednom varijablom
417 #  $r^2 = 0.66$ ,  $F = 304.9$ ,  $p < 0.01$  (za ceo model i koeficijente)
418
419 # 3.1. Vizualizacija trećeg modela
420
421 ggplot(Model3, aes(y = baza$total_litres_of_pure_alcohol, x = baza$beer_servings)) +
422   geom_point(size = 3, color = "azure4" )+
423   stat_smooth(method = "lm", se = FALSE) +
424   ggtitle("Model broj 3") +
425   xlab("potrošnja limenki piva") +
426   ylab("ukupna unos čistog alkohola")
427
428 # Zaključak
429
430 # Sva tri modela su odlična ali je prvi ipak najbolji,
431 # jer objašnjava najveći procenat varijabiliteta ukupnog unosa čistog alkohola.
432 # Ovaj model pokazuje da su
433 # varijable potrošnje čaša vina, limenki piva i čašica žestokog pića,
434 # odlčni prediktori vrednosti varijable ukupnog unosa čistog alkohola.
435 # Takav nalaz je očekivan pošto je reč o najpopularnijim vrstama pića.
436
437 ##### ZAVRŠETAK ANALIZE #####
438
439 # Analizu sproveo: Novak Tešić
440 # GitHub profil: https://github.com/NovakTestic

```