

Analiza cene nekretnina

Seminarski rad u okviru kursa

Uvod u teoriju uzoraka

Matematički fakultet

Petar Zečević

mi16169@alas.matf.bg.ac.rs

26. jul 2020.

Sažetak

U radu je prikazana analiza srednje vrednosti cene nekretnina i cene metra kvadratnog nekretnina u okrugu King County pomoću raznih metoda uzorkovanja. Te metode su: prost slučajni uzorak, regresiono ocenjivanje i stratifikovani uzorak. Pokazalo se da su za relativno malu veličinu uzorka sva tri metoda dala pristojne rezultate.

Sadržaj

1	Uvod	2
1.1	Opis baze podataka	2
2	Analiza baze podataka	3
3	Teorijski uvod	4
3.1	Prost slučajni uzorak	6
3.2	Regresiono ocenjivanje	8
3.3	Stratifikovano uzorkovanje	8
4	Rezultati uzorkovanja	9
5	Zaključak	10
	Literatura	10
A	Dodatak	11

1 Uvod

Analiziranje cene nekretnina je traženo iz razloga što bi ljudima olakšala kupovinu stanova. Želimo da nad relativno malom bazom podataka, prodate kuće u periodu od godinu dana u okrugu King County, nađemo dobru metodu za uzorkovanje da bismo to mogli da primenimo nad nekom većom bazom (na primer nad celom državom).

Pokušaćemo da ocenimo kvalitet ocene srednje vrednosti cene prodatih nekretnina kao i cene kvadratnog metra prodatih nekretnina. Prvo ćemo da opišemo bazu podataka i analiziramo potencijalno korisna obeležja, zatim ćemo da ukratko teorijski objasnimo uzorkovanja koja smo koristili. Onda sledi analiza rezultata koje smo dobili. U dodatku se nalazi kod u programskom jeziku R koji je korišćen u radu.

1.1 Opis baze podataka

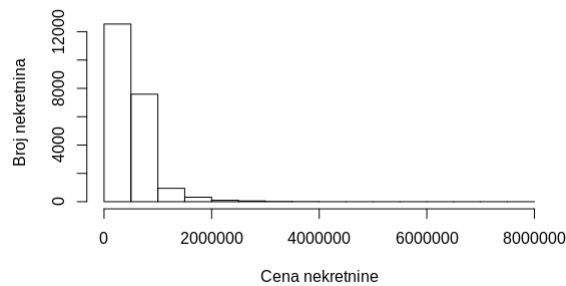
Baza predstavlja informacije o prodatim kućama u periodu od maja 2014. do maja 2015. godine u okrugu King County, SAD (uključujući Sietl) [1]. Sastoji se iz 21613 redova i 21 kolone. Kolone koje predstavljaju broj kvadratnih stopa su konvertovani u broj kvadratnih metara. Sledeći atributi se nalaze u bazi:

- **id** - jedinstveni identifikator reda. Nema nikakav značaj za istraživanje.
- **date** - datum prodaje. Izražen kao vremenski žig.
- **price** - cena po kojoj je nekretnina prodana. Ovo je ciljna promenljiva u našem istraživanju.
- **bedrooms** - broj soba. Postoje nekretnine bez soba što je verovatno greška
- **bathrooms** - broj kupatila. Ovaj atribut može uzeti vrednost na primer 2.5 jer treće kupatilo nema neki deo da bi bilo kompletno. Mi smo zaokruživali ove vrednosti jer je tako razumljivije.
- **sqft_living** - broj kvadratnih stopa u zatvorenom prostoru
- **sqft_lot** - broj kvadratnih stopa ukupno na placu
- **floors** - broj spratova. Može biti na primer 1.5 što znači da kuća ima jedan međusprat.
- **waterfronts** - indikator da li se nekretnina nalazi u blizini reke ili mora.
- **view** - ocena koliko dobar pogled ima
- **condition** - ocena u kakvom je stanju objekat
- **grade** - ocena koliko je kvalitetna izgradnja objekta
- **sqft_above** - broj kvadratnih stopa u zatvorenom prostoru ne računajući podrum
- **sqft_basement** - broj kvadratnih stopa podruma
- **yr_built** - godina izgradnje
- **yr_renovated** - godina poslednjeg renoviranja. Ako objekat nije renoviran onda je ova vrednost nula.
- **zipcode** - poštanski broj
- **lat** - geografska širina
- **long** - geografska dužina

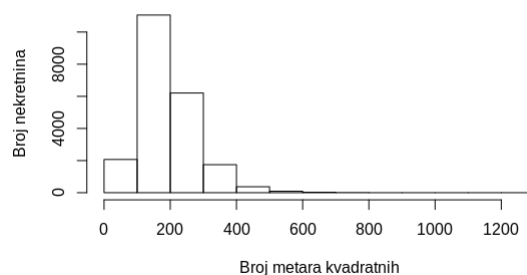
- **sqft_living15** - prosek broja kvadratnih stopa u zatvorenom prostoru za 15 najbližih kuća
- **sqft_lot15** - prosek broja kvadratnih stopa ukupno na placu za 15 najbližih kuća

2 Analiza baze podataka

U okrugu King County cene nekretnina se kreću od \$77,000 do \$7,700,000. 93% nekretnina koštaju ispod \$1,000,000 što možemo videti na slici 1. Na slici 2 može se videti histogram obeležja broja metara kvadratnih u zatvorenom prostoru nekretnine. Ostala obeležja za broj metara kvadratnih nisu dobro korelisana sa cenom nekretnine tako da na njih nije posvećeno mnogo pažnje.



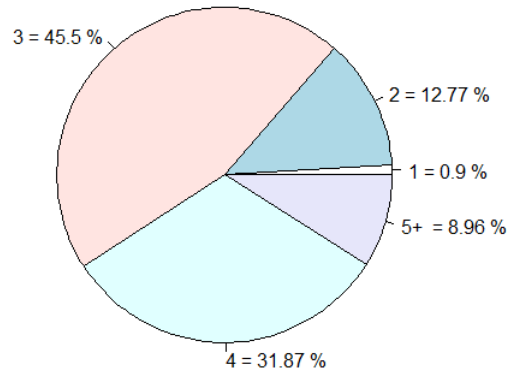
Slika 1: Histogram cena nekretnina



Slika 2: Histogram broja metara kvadratnih

Iz baze su izbačene nekretnine koje imaju nula soba ili kupatila jer verovatno predstavljaju neku grešku. Procenat nekretnina po broju soba može se videti na slici 3, a procenat nekretnina po broju kupatila na 4. Sve nekretnine sa više od 4 sobe su kategorisane kao 5+, a sa više od 3 kupatila u 4+ radi preglednijih slika. Vrednosti broja kupatila su zaokružena na

cele brojeve jer je nedefinisano šta bi na primer značilo da kuća ima 2.25 kupatila.



Slika 3: Procenat nekretnina po broju soba

Cena nekretnine ima najbolju korelaciju sa brojem metara kvadratnih u zatvorenom prostoru, 0.7. Graf korelacije je prikazan na slici 6. Procenat nekretnina po metrima kvadratnim u zatvorenom prostoru je prikazan na slici 5. Vrednosti su diskretizovane u intervale po $100m^2$.

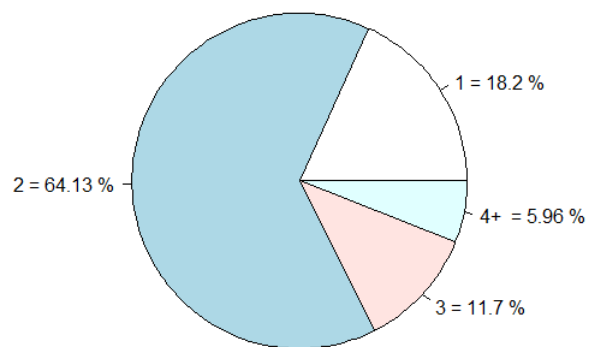
Prosečna cena metra kvadratnog u okrugu King County je 2843.03. Histogram je prikazan na slici 7.

Pošto imamo obeležja geografska širina i dužina, možemo aproksimativno da odredimo koja nekretnina se nalazi u gradu Sijetl, a koja van njega. Gledajući mapu King County okruga, uzimamo da se grad Sijetl nalazi na geografskoj širini u intervalu (47.501221, 47.735211) i geografskoj dužini (-122.443828, -122.238907). Procenat nekretnina u Sijetlu i van njega nalazi se na 8. Korelacija između ovog obeležja i cene nekretnina je samo 0.0955997 tako da ovo obeležje nećemo dalje koristiti. Korelacija između ovog obeležja i cene metra kvadratnog je bolja, 0.4337856, ali opet ne dovoljno dobra da bismo je iskoristili za regresiono ocenjivanje.

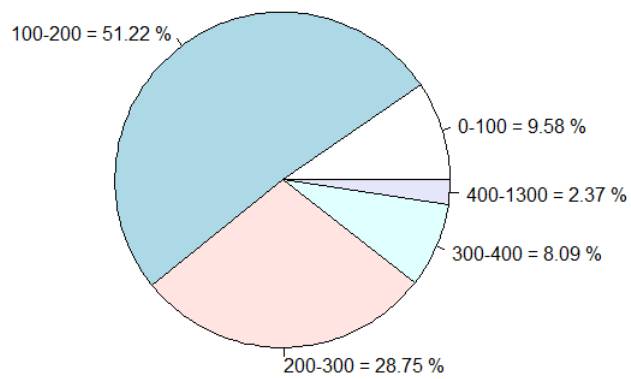
Po formuli iz [2] dobijamo da je dobra veličina uzorka 1411. To je 6.53% ukupne populacije.

3 Teorijski uvod

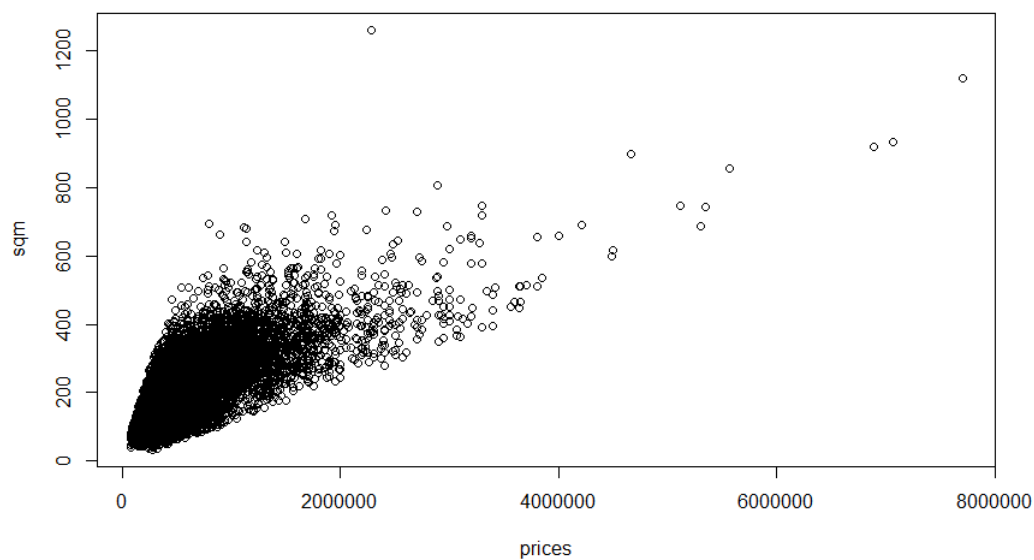
U radu ćemo koristiti prost slučajan i stratifikovan uzorak u kome se iz stratuma izvlače uzorci prostim slučajnim uzorkovanjem. Takođe ćemo koristiti regresiono ocenjivanje prostog slučajnog uzorkovanja. U nastavku sledi teorijska podloga ovih uzorkovanja.



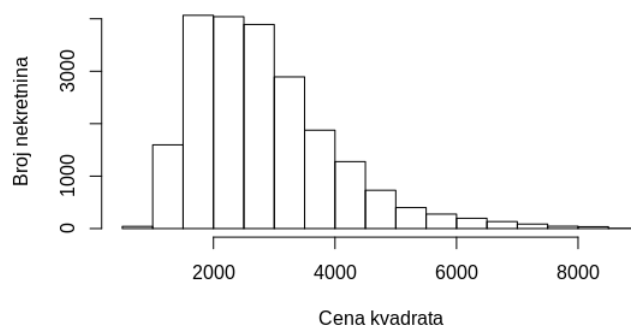
Slika 4: Procenat nekretnina po broju kupatila



Slika 5: Procenat nekretnina po kvadratnim metrima



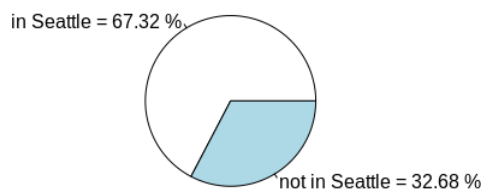
Slika 6: Korelacija između cene i broja kvadratnih metara



Slika 7: Histogram cena metra kvadratnog

3.1 Prost slučajni uzorak

Prost slučajni uzorak bez ponavljanja je plan izbora uzorka koji iz populacije od N jedinica, nasumično bira n različitih jedinica tako da svaka kombinacija ima jednaku verovatnoću [3]. Verovatnoća da određena



Slika 8: Procenat nekretnina u Sijetlu i van njega

kombinacija bude izabrana iznosi:

$$p = \begin{cases} \left(\frac{N}{n} \right)^{-1} & \text{ako je obim uzorka jednak } n \\ 0 & \text{inace} \end{cases}$$

Ocena populacijske srednje vrednosti iznosi:

$$\hat{m}_Y = \frac{1}{n} \sum_{k \in S} y_k$$

Ova ocena je nepristrasna. Disperzija ove ocene i ocena te disperzije su:

$$D(\hat{m}_Y) = \frac{\sigma^2}{n} \left(1 - \frac{n}{N} \right)$$

$$D(\hat{m}_Y) = \frac{S^2}{n} \left(1 - \frac{n}{N} \right)$$

Kada su u pitanju intervalne ocene, aproksimativni $100(1 - \alpha)\%$ interval poeverenja izgleda:

$$I_{\hat{m}_Y} = [\hat{m}_Y - z_{1-\frac{\alpha}{2}} \sqrt{D(\hat{m}_Y)}, \hat{m}_Y + z_{1-\frac{\alpha}{2}} \sqrt{D(\hat{m}_Y)}]$$

- m_Y - populacijska srednja vrednost
- S - prost slučajan uzorak
- σ^2 - populacijska disperzija
- S^2 - uzoračka disperzija
- $z_{1-\frac{\alpha}{2}}$ - vrednost $(1 - \frac{\alpha}{2})$ kvantila standardne normalne raspodele. Ukoliko je veličina uzorka manja od 30 onda se koristi studentova raspodela.

3.2 Regresiono ocenjivanje

Regresiono ocenjivanje je tehnika ocenjivanja, koja koristi neke dodatne informacije radi postizanja veće preciznosti ocena [3]. Koristi se u slučajevima kada se veza između glavnog i pomoćnog obeležja najbolje može opisati regresionom pravom koja ne prolazi kroz koordinatni početak, odnosno u situacijama kada se smatra da između obeležja postoji približno linearna veza, koja ne prolazi kroz koordinatni početak. Da bismo mogli da koristimo regresiono ocenjivanje, potrebno je da nam bude poznata suma ili srednja vrednost pomoćnog obeležja. Prvo definišemo koeficijent korelacije ciljnog i pomoćnog obeležja kao

$$\rho = \frac{\sum_{i=1}^N (x_i - \bar{x})(y_i - \bar{y})}{(N-1)s_x s_y}$$

i uzorački koeficijent korelacije kao

$$\hat{\rho} = \frac{\sum_{i \in S} (x_i - \bar{x}_n)(y_i - \bar{y}_n)}{(n-1)s_n(x)s_n(y)}$$

Ocena srednje vrednosti obeležja x je data sa

$$\hat{x}_{LR} = \hat{x}_n + \hat{b}(\bar{y} - \bar{y}_n)$$

gde je

$$\hat{b} = \hat{\rho} \frac{s_n(x)}{s_n(y)}$$

Disperzija ove ocene je

$$D(\hat{x}_{LR}) = \frac{\sigma_d^2}{n} \left(1 - \frac{n}{N}\right)$$

a ocena te disprezije:

$$D(\hat{\hat{x}}_{LR}) = \frac{\bar{s}_d^2}{n} \left(1 - \frac{n}{N}\right)$$

gde su

$$\sigma_d^2 = \frac{1}{N-1} \sum_{k=1}^N (x_k - (\bar{x} + b(y_k - \bar{y})))^2$$

$$\bar{s}_d^2 = \frac{1}{n-1} \sum_{k \in S} (x_k - (\bar{x}_n + b(y_k - \bar{y}_n)))^2$$

3.3 Straifikovano uzorkovanje

Stratifikacija podrazumeva podelu populacije na delove koji se nazivaju stratumi, pri čemu treba formirati relativno homogene, među sobom različite stratume, što znači da vrednosti obeležja, koje je predmet istraživanja, treba da budu približne na jedinkama unutar jednog stratumu, ali da se bitno razlikuju među stratumima. Dobijeni stratumi su disjunktni i treba da zadovoljavaju uslov pokrivenosti [3]. Odnosno treba da važi:

$$N_1 + N_2 + \dots + N_L = N$$

gde je $N_h, \forall h = 1, \bar{L}$ veličina stratumu, a N ukupna veličina populacije. Uzorak veličine n dobijen iz stratifikovane populacije sadrži iz svakog stratumu uzorak veličine $n_h, \forall h = 1, \bar{L}$, gde važi da je:

$$n_1 + n_2 + \dots + n_L = n$$

Ocena za uzoračku sredinu je

$$\bar{x}_{str} = \frac{1}{N} \sum_{h=1}^L \frac{N_h^2 s_h^2}{n_h} \left(1 - \frac{n_h}{N_h}\right)$$

Disperzija te ocene i ocena te disperzije su:

$$D(\bar{x}_{str}) = \frac{1}{N^2} \sum_{h=1}^L \frac{N_h^2 s_h^2}{n_h} \left(1 - \frac{n_h}{N_h}\right)$$

$$D(\hat{\bar{x}}_{str}) = \frac{1}{N^2} \sum_{h=1}^L \frac{N_h^2 s_{n_h}^2}{n_h} \left(1 - \frac{n_h}{N_h}\right)$$

Nakon stratifikacije, iz svakog stratuma biraju se uzorci, unapred određenih obima. Uzorci se biraju međusobno nezavisno iz različitih stratuma i nije neophodno koristiti isti plan uzorkovanja za sve stratumе.

Veličine uzoraka iz stratuma mogu da se uzimaju proporcionalno broju jedinki u tom stratumu, odnosno:

$$\frac{n_1}{N_1} = \dots = \frac{n_L}{N_L} = \frac{n}{N} \Rightarrow n_h = \frac{n}{N} N_h, h = 1, \dots, L$$

Drugi način za odabir veličine uzoraka je Nejmanov raspored. On uzima u obzir disperzije u stratumima:

$$n_h = \frac{N_h s_h}{\sum_{h=1}^L N_h s_h} n, h = 1, \dots, L.$$

4 Rezultati uzorkovanja

Za ocenu srednje vrednosti nad ciljnim obeležjem cena nekretnine vršimo prosto slučajno uzorkovanje, regresiono ocenjivanje gde je pomoćna promenljiva broj metara kvadratnih i stratifikovano uzorkovanje gde su stratumi podeljeni po brojevima kupatila u objektu (ovde se takođe broj kupatila zaokružuje na ceo broj i svi objekti sa preko 3 kupatila se svrstavaju u isti stratum). Za ocenu srednje vrednosti nad ciljnim obeležjem cena metra kvadratnog nekretnine vršimo prosto slučajno uzorkovanje i stratifikovano uzorkovanje gde su stratumi podeljeni po brojevima kupatila u objektu.

Prava srednja vrednost ciljnog obeležja je 540259. Prema tome, dobili smo da je stratifikovano uzorkovanje dalo najpribližniju vrednost. Međutim, što se tiče ocene disperzije ocena, ovo uzorkovanje ima najgori rezultat. Najbolji rezultat daje ocena disperzije regresionog ocenjivanja. Što se tiče intervala poverenja, i ovde se regresiono ocenjivanje pokazuje kao najbolje mada ni stratifikovano uzorkovanje ne daje loš rezultat. Rezultati se mogu videti u tabeli 1.

Što se tiče prosečne cene metra kvadratnog, gde je prava srednja vrednost 2843.03, tu se bolje pokazalo prosto slučajno uzorkovanje mada stratifikovano uzorkovanje nije mnogo lošije. Rezultati se mogu videti u tabeli 2.

Tabela 1: Ocene srednje vrednosti, ocene disperzija tih ocena i intervali tih ocena za tri ocenjivanja.

plan	\hat{m}_Y	$\hat{D}(\hat{m}_Y)$	$I_{\hat{m}_Y}$
PSU	531920.8	69334651	[515600.7, 548240.9]
REG	544836.3	41271541	[532244.9, 557427.7]
STRAT	538549.5	86723407	[520297.2, 556801.7]

Tabela 2: Ocene srednje vrednosti, ocene disperzija tih ocena i intervali tih ocena za obeležje cena metra kvadratnog.

plan	\hat{m}_Y	$\hat{D}(\hat{m}_Y)$	$I_{\hat{m}_Y}$
PSU	2848.116	905.7124	[2789.130, 2907.101]
STRAT	2822.609	1116.006	[2757.133, 2888.085]

5 Zaključak

Možemo primetiti da svi načini uzorkovanja daju prihvatljive rezultate. Što znači da lako možemo oceniti koliko prosečno košta stan u nekom gradu, okrugu ili zemlji sa relativno malim uzorkom u odnosu na populaciju i možemo primeniti bilo koji od metoda primenjenih u radu.

U ovom konkretnom primeru (King County okrug) vidimo da nije velika razlika u cenama između velikog grada i okruga, što ne znači da u nekom drugom okrugu to ne bi bio slučaj. Da bi se ova hipoteza proverila potrebno je raditi na dosta više baza podataka.

Literatura

- [1] harlfoxem. House Sales in King County, USA. on-line at: <https://www.kaggle.com/harlfoxem/housesalesprediction>.
- [2] Robert V. Krejcie and Daryle V. Morgan. Determining sample size for research activities. on-line at: https://home.kku.ac.th/sompong/guest_speaker/KrejcieandMorgan_article.pdf.
- [3] Astrea Camstra Reinder Banning and Paul Kottnerus. *Sampling Theory, Sampling design and estimation models*. Statistics Netherlands, 2012.

A Dodatak

```
1000 install.packages("ggplot2")
1001 library(ggplot2)
1002 library(scales)

1004 #otklanjam o notifikaciju sa realnih brojeva
options(scipen=999)

1006 #konvertujemo stope kvadratne u metre kvadratne
1008 f2_to_m2 = (0.3048)^2

1010 N = length(kc_house_data$id)

1012 #izbacujemo redove koji imaju 0 soba ili 0 kupatila jer su to
      verovatno greske
indexes_wo_ano = c()
1014 for (i in 1:N){
      if (kc_house_data$bedrooms[i] > 0 && kc_house_data$bathrooms[i]
        > 0.5)
1016     {
1018         indexes_wo_ano = c(indexes_wo_ano, i)
      }
    }

1020 house_data_prep = kc_house_data[indexes_wo_ano,]

1022 N = length(house_data_prep$id)

1024 #pravimo kolone za metre kvadratne umesto kolona za stope
      kvadratne
house_data_prep$sqm_living = house_data_prep$sqft_living*f2_to_m2
1026 house_data_prep$sqm_lot = house_data_prep$sqft_lot*f2_to_m2
house_data_prep$sqm_above = house_data_prep$sqft_above*f2_to_m2
1028 house_data_prep$sqm_basement = house_data_prep$sqft_basement*
      f2_to_m2
house_data_prep$sqm_living15 = house_data_prep$sqft_living15*
      f2_to_m2
1030 house_data_prep$sqm_lot15 = house_data_prep$sqft_lot15*f2_to_m2
house_data_prep = subset(house_data_prep, select=-c(sqft_living,
      sqft_lot, sqft_above, sqft_basement, sqft_living15,
      sqft_lot15))
1032 #dimnames(house_data_prep)

1034 #histogram cena
hist(house_data_prep$price, xlab="Cena nekretnine", ylab = "Broj
      nekretnina", main = NULL)
1036 #ggplot(house_data_prep, aes(x=price)) + geom_histogram()

1038 hist(house_data_prep$sqm_living, xlab="Broj metara kvadratnih",
      ylab = "Broj nekretnina", main = NULL)

1040 #procenat nekretnina koje kostaju manje od 1,000,000 dolara
n_lt_1m = length(house_data_prep[house_data_prep$price<1000000,1])
1042 perc_n_lt_1m = n_lt_1m/N*100
perc_n_lt_1m

1044 #minimalna i maksimalna cena
1046 min(house_data_prep$price)
max(house_data_prep$price)

1048 #racunamo cene jednog kvadrata svake nekretnine kao i njihovu
      prosechnu ocenu
house_data_prep$price_of_m2 = house_data_prep$price/
      house_data_prep$sqm_living
house_data_prep$price_of_m2
1052 hist(house_data_prep$price_of_m2, xlab = "Cena kvadrata", ylab = "
      Broj nekretnina", main = NULL)

1054 #sve kuće sa 4 ili više soba kategorisemo zajedno
bedrooms = house_data_prep$bedrooms
1056 for (i in 1:length(bedrooms)){
      if (bedrooms[i] > 4){
1058         bedrooms[i] = 5
      }
```

```

    }
1060 }

1062 #pravimo pie chart za broj soba
mytable = table(bedrooms)
1064 lbls = paste(names(mytable), "=", round(mytable/sum(mytable)*100,
      digits = 2), "%", sep=" ")
      lbls[5] = paste("5+ ", "=", round(mytable[5]/sum(mytable)*100,
      digits = 2), "%", sep=" ")
1066 pie(mytable, labels = lbls,
      main=NULL)
1068

1070 #sve kuće sa 3 ili više kupatila kategorisemo zajedno
bathrooms = house_data_prep$bathrooms
for (i in 1:length(bathrooms)){
1072   bathrooms[i] = round(bathrooms[i])
      if (bathrooms[i] > 3){
1074     bathrooms[i] = 4
      }
1076 }

1078 #pravimo pie chart za broj kupatila
mytable = table(bathrooms)
1080 lbls = paste(names(mytable), "=", round(mytable/sum(mytable)*100,
      digits = 2), "%", sep=" ")
      lbls[4] = paste("4+ ", "=", round(mytable[4]/sum(mytable)*100,
      digits = 2), "%", sep=" ")
1082 pie(mytable, labels = lbls,
      main=NULL)
1084

1086 #diskretizujemo kolonu sa metrima kvadratnim
house_data_prep$sqm_living_bin = cut(house_data_prep$sqm_living,
      breaks = c(0,100,200,300,400,1300), labels = c("0-100", "
      100-200", "200-300", "300-400", "400-1300"))

1088 #pravimo pie chart sa diskretizovanim podacima
bins = house_data_prep$sqm_living_bin
mytable = table(bins)
1090 lbls = paste(names(mytable), "=", round(mytable/sum(mytable)*100,
      digits = 2), "%", sep=" ")
1092 pie(mytable, labels = lbls,
      main=NULL)
1094

1096 #pravimo grafik korelacije izmedju cene i metara kvadratnih
      nekretnine
plot(house_data_prep$price, house_data_prep$sqm_living, xlab = "
      prices", ylab = "sqm")
cor(house_data_prep$price, house_data_prep$sqm_living)
1098

1100 #racunamo idealnu velicinu uzorka
n = ceiling(qchisq(0.95, 7)*N*0.5*0.5/(0.05^2*(N-1)) + qchisq
      (0.95, 7)*0.5*0.5)
n
1102 procenat_velicine_uzorka = n/N*100

1104 #postavljamo seed na fiksnu vrednost zbog vise pokretanja
set.seed(6)
1106

1108 #Vadjenje prostog slucajnog uzorka
sample_psu = sample(1:N, n)

1110 #Ocena srednje vrednosti prostog slucajnog uzorka
xn_ocena_psu = mean(house_data_prep$price[sample_psu])
1112

1114 #Ocena disprezije ocene srednje vrednosti prostog slucajnog uzorka
s2_psu = var(house_data_prep$price[sample_psu])
d_ocena_psu = s2_psu/n*(1-n/N)
1116
xn_ocena_psu
1118 d_ocena_psu

1120 #Disperzija srednje vrednosti prostog slucajnog uzorka
var(house_data_prep$price)/n*(1-n/N)
1122

```

```

1124 #Racunjanje intervala poverenja za psu
aplha = 0.05

1126 #Posto je n>30 onda koristimo normalnu raspodelu
z = qnorm(1 - aplha/2)
1128 interval = c(xn_ocena_psu - z*sqrt(d_ocena_psu), xn_ocena_psu + z*
sqrt(d_ocena_psu))
interval

1130 #Regresiono ocenjivanje gde je pomocna promenljiva broj metara
kvadratnih zatvorenog prostora
1132 uzorak_zareg = sample(1:N, n)

1134 #Racunanje pomocnih promenljivih ro i b
ro = cor(house_data_prep$price[uzorak_zareg],
house_data_prep$sqm_living[uzorak_zareg])
1136 b = ro*sqrt(var(house_data_prep$price[uzorak_zareg])/sqrt(var(
house_data_prep$sqm_living[uzorak_zareg])))

1138 #Racunanje ocene srednje vrednosti regresionim ocenjivanjem
xn_lr = mean(house_data_prep$price[uzorak_zareg])+b*(mean(
house_data_prep$sqm_living[uzorak_zareg] - mean(
house_data_prep$sqm_living[uzorak_zareg])))

1140 #Racunanje ocene disperzije srednje vrednosti regresionim
ocenjivanjem
1142 se = 1/(n-1)*sum((house_data_prep$price[uzorak_zareg] - (mean(
house_data_prep$price[uzorak_zareg]) + b*(
house_data_prep$sqm_living[uzorak_zareg] - mean(
house_data_prep$sqm_living[uzorak_zareg]))))^2)
d_ocena_reg = se/n*(1-n/N)

1144 #Racunanje intervala poverenja za regresiono ocenjivanje
1146 interval_reg = c(xn_lr - z*sqrt(d_ocena_reg), xn_lr + z*sqrt(
d_ocena_reg))
interval_reg

1148 xn_lr
1150 d_ocena_reg

1152 #Prava srednja vrednost
mean(house_data_prep$price)

1154 #Pravimo stratumе od obelezja bathrooms pa hocemo da
diskretizujemo to obelezje
1156 #Takodje sve nekretnine sa 4 ili vise kupatila svrstavamo u jedan
stratum
house_data_prep$bathrooms_rounded = round(
house_data_prep$bathrooms)
1158 for (i in 1:N){
if (house_data_prep$bathrooms_rounded[i] > 4){
1160 house_data_prep$bathrooms_rounded[i] = 4
}
}
1162 unique(house_data_prep$bathrooms_rounded)

1164 #Pravimo listu stratuma
stratumi = list()
stratumi[[1]] = house_data_prep[house_data_prep$bathrooms_rounded
== 1,]
1168 stratumi[[2]] = house_data_prep[house_data_prep$bathrooms_rounded
== 2,]
stratumi[[3]] = house_data_prep[house_data_prep$bathrooms_rounded
== 3,]
1170 stratumi[[4]] = house_data_prep[house_data_prep$bathrooms_rounded
== 4,]
stratumi[[1]]

1172 #Cuvamo disperzije svakog stratuma i velicine
1174 disprezije_stratuma = c(sqrt(var(stratumi[[1]]$price)), sqrt(var(
stratumi[[2]]$price)), sqrt(var(stratumi[[3]]$price)), sqrt(
var(stratumi[[4]]$price)))
N_stratuma = c(length(stratumi[[1]]$price), length(stratumi[[2]]
$price), length(stratumi[[3]]$price), length(stratumi[[4]]

```

```

$price))
1176 #Nejmanovom rapodelom odredjujemo velicine uzorka iz svakog
      stratum
1178 n_stratuma = round(n*(N_stratuma*disprezije_stratuma)/sum(
      N_stratuma*disprezije_stratuma))

1180 #Posto smo odmah dobili zeljeni ukupni broj elemenata u uzoraku,
      ne moramo da dodajemo ili oduzimamo elemente
      n_stratuma
1182 sum(n_stratuma) == n

1184 #Uzimamo prost slucajan uzorak iz svakog stratum
indeksi_strat = c()
1186 indeksi_strat[[1]] = sample(1:N_stratuma[1], n_stratuma[1],
      replace = F)
      indeksi_strat[[2]] = sample(1:N_stratuma[2], n_stratuma[2],
      replace = F)
1188 indeksi_strat[[3]] = sample(1:N_stratuma[3], n_stratuma[3],
      replace = F)
      indeksi_strat[[4]] = sample(1:N_stratuma[4], n_stratuma[4],
      replace = F)

1190 #Racunamo ocene srednje vrednosti svakog stratum
xn_strat = c()
1192 xn_strat[1] = mean(house_data_prep$price[indeksi_strat[[1]])
1194 xn_strat[2] = mean(house_data_prep$price[indeksi_strat[[2]])
      xn_strat[3] = mean(house_data_prep$price[indeksi_strat[[3]])
1196 xn_strat[4] = mean(house_data_prep$price[indeksi_strat[[4]])

1198 #Racunamo ocene uzoracke disperzije svakog stratum
uzoracka_disperzija_stratuma = c()
1200 uzoracka_disperzija_stratuma[1] = var(house_data_prep$price[
      indeksi_strat[[1]])
      uzoracka_disperzija_stratuma[2] = var(house_data_prep$price[
      indeksi_strat[[2]])
1202 uzoracka_disperzija_stratuma[3] = var(house_data_prep$price[
      indeksi_strat[[3]])
      uzoracka_disperzija_stratuma[4] = var(house_data_prep$price[
      indeksi_strat[[4]])

1204 #Racunamo ocenu srednje vrednosti dobijenu stratifikovanim
      uzorkovanjem
1206 xn_st = 1/N*sum(N_stratuma*xn_strat)
      xn_st

1208 #Racunamo ocenu disperzije te ocene
1210 d_ocena_st = 1/N^2*sum((N_stratuma^2*uzoracka_disperzija_stratuma)
      /n_stratuma*(1-n_stratuma/N_stratuma))
      d_ocena_st

1212 #Racunamo interval poverenja te ocene
1214 interval_strat = c(xn_st - z*sqrt(d_ocena_st), xn_st + z*sqrt(
      d_ocena_st))
      interval_strat

1216 #Aproksimativno delimo nekretnine na one koje su u Sijetlu i one
      koje su van njega
1218 lat_min = 47.501221
      long_max = -122.238907
1220 lat_max = 47.735211
      long_min = -122.443828
1222 house_data_prep$in_seattle = ifelse(house_data_prep$lat > lat_min
      & house_data_prep$lat < lat_max & house_data_prep$long >
      long_min & house_data_prep$long < long_max, 1, 0)
      unique(house_data_prep$in_seattle)

1224 mytable = table(house_data_prep$in_seattle)
1226 lbls = paste(c("in Seattle", "not in Seattle"), "=", round(mytable
      /sum(mytable)*100, digits = 2), "%", sep=" ")
      pie(mytable, labels = lbls,
1228 main=NULL)

```

```

1230 #Zelimo da vidimo da li postoji korelacije izmedju cene i toga da
      li se nekretnina nalazi u Sijetlu
      #Ispostavlja se da nema korelacije
1232 cor(house_data_prep$price, house_data_prep$seattle)

1234 #Isto to samo za obelezje cena metra kvadratnog
      cor(house_data_prep$price_of_m2, house_data_prep$seattle)
1236
      #U ovom delu koda ocenjijemo srednju vrednost za obeleyje cena
      metra kvadratnog
1238 #Koristimo PSU i strafifikovano uzorkovanje po broju kupatila
      #Veoma slicno kao do sad tako da nece biti komentara
1240 sample_psu_2 = sample(1:N, n)

1242 xn_ocena_psu_2 = mean(house_data_prep$price_of_m2[sample_psu_2])

1244 s2_psu_2 = var(house_data_prep$price_of_m2[sample_psu_2])
      d_ocena_psu_2 = s2_psu_2/n*(1-n/N)
1246
      xn_ocena_psu_2
1248 d_ocena_psu_2
      mean(house_data_prep$price_of_m2)
1250
      interval_psu_2 = c(xn_ocena_psu_2 - z*sqrt(d_ocena_psu_2),
      xn_ocena_psu_2 + z*sqrt(d_ocena_psu_2))
1252 interval_psu_2

1254 indeks_i_strat_2 = c()
      indeks_i_strat_2[[1]] = sample(1:N_stratuma[1], n_stratuma[1],
      replace = F)
1256 indeks_i_strat_2[[2]] = sample(1:N_stratuma[2], n_stratuma[2],
      replace = F)
      indeks_i_strat_2[[3]] = sample(1:N_stratuma[3], n_stratuma[3],
      replace = F)
1258 indeks_i_strat_2[[4]] = sample(1:N_stratuma[4], n_stratuma[4],
      replace = F)

1260 xn_strat_2 = c()
      xn_strat_2[1] = mean(house_data_prep$price_of_m2[indeksi_strat_2
      [[1]])
1262 xn_strat_2[2] = mean(house_data_prep$price_of_m2[indeksi_strat_2
      [[2]])
      xn_strat_2[3] = mean(house_data_prep$price_of_m2[indeksi_strat_2
      [[3]])
1264 xn_strat_2[4] = mean(house_data_prep$price_of_m2[indeksi_strat_2
      [[4]])

1266 uzoracka_disperzija_stratuma_2 = c()
      uzoracka_disperzija_stratuma_2[1] = var(
      house_data_prep$price_of_m2[indeksi_strat_2[[1]])
1268 uzoracka_disperzija_stratuma_2[2] = var(
      house_data_prep$price_of_m2[indeksi_strat_2[[2]])
      uzoracka_disperzija_stratuma_2[3] = var(
      house_data_prep$price_of_m2[indeksi_strat_2[[3]])
1270 uzoracka_disperzija_stratuma_2[4] = var(
      house_data_prep$price_of_m2[indeksi_strat_2[[4]])

1272 xn_st_2 = 1/N*sum(N_stratuma*xn_strat_2)
      xn_st_2
1274
      d_ocena_st_2 = 1/N^2*sum((N_stratuma^2*
      uzoracka_disperzija_stratuma_2)/n_stratuma*(1-n_stratuma/
      N_stratuma))
1276 d_ocena_st_2

1278 interval_strat = c(xn_st_2 - z*sqrt(d_ocena_st_2), xn_st_2 + z*
      sqrt(d_ocena_st_2))
      interval_strat

```

Listing 1: Kod korišćen u radu