

Dokumentacija Zadaće 1 iz predmeta Mašinsko učenje

Beglerović Vedad

Beglerović Vildana

Buturović Lejla

Elektrotehnički fakultet Sarajevo

24. april 2022.

Sadržaj

1	Zadatak 1	4
1.1	Osnovne metode deskriptivne statistike	4
1.1.1	Upoznavanje sa setom podataka	4
1.2	Metode procjene lokacije i varijabilnosti podataka	5
1.2.1	Metode za procjenu lokacije	5
1.2.2	Metode za procjenu varijabilnosti	9
1.3	Metode za procjenu korelacije između varijabli	11
1.3.1	Korelacija između numeričkih varijabli	11
1.3.2	Korelacija između kategoričkih varijabli	14
1.4	Predprocesiranje podataka	18
1.4.1	Odbacivanje atributa sa visokim stepenom korelacije	18
1.4.2	Popunjavanje nedostajućih vrijednosti	20
1.4.3	Odbacivanje pronađenih outliera	24
2	Zadatak 2	25
2.1	Izgradnja modela klasifikacije	25
2.1.1	Drvo odlučivanja koje kao mjeru atributa selekcije koristi informacijsku dobit	26
2.1.2	Drvo odlučivanja koje kao mjeru atributa selekcije koristi gini indeks	27
2.1.3	C5.0 model klasifikacije	30
2.2	Predikcijski modeli sa metodom holdouta	31
2.2.1	Drvo odlučivanja koje kao mjeru atributa selekcije koristi informacijsku dobit sa metodom holdouta	31
2.2.2	Drvo odlučivanja koje kao mjeru atributa selekcije koristi gini indeks sa metodom holdouta	31
2.2.3	C5.0 model klasifikacije sa metodom holdouta	32
2.3	Predikcijski modeli sa metodom k-fold unakrsne validacije	33
2.3.1	Drvo odlučivanja koje kao mjeru atributa selekcije koristi informacijsku dobit sa metodom k-fold unakrsne validacije	33
2.3.2	Drvo odlučivanja koje kao mjeru atributa selekcije koristi gini indeks sa metodom k-fold unakrsne validacije	33
2.3.3	C5.0 model klasifikacije sa metodom k-fold unakrsne validacije	33
2.4	Predikcijski modeli sa metodom k-fold bootstrapping validacije	34
2.4.1	Drvo odlučivanja koje kao mjeru atributa selekcije koristi informacijsku dobit sa metodom k-fold bootstrapping validacije	34
2.4.2	Drvo odlučivanja koje kao mjeru atributa selekcije koristi gini indeks sa metodom k-fold bootstrapping validacije	34
2.4.3	C5.0 model klasifikacije sa metodom k-fold bootstrapping validacije	34
2.5	Balansiranje podataka	35

2.5.1	Drvo odlučivanja koje kao mjeru atributa selekcije koristi informacijsku dobit nakon balansiranja podataka	36
2.5.2	Drvo odlučivanja koje kao mjeru atributa selekcije koristi gini indeks nakon balansiranja podataka	36
2.5.3	C5.0 model nakon balansiranja	37
2.6	Ansambl tehnike za unaprjeđenje tačnosti klasifikacije	37
2.6.1	Bagging model	38
2.6.2	Boosting model uz korišćenje AdaBoost metode	38
2.6.3	Random forest model	39
3	Zadatak 3	40
4	Zadatak 4	41

1 Zadatak 1

1.1 Osnovne metode deskriptivne statistike

1.1.1 Upoznavanje sa setom podataka

Set podataka učitavamo naredbom:

```
podaci <- read.csv("customer_data_train.csv", fileEncoding = 'UTF-8')
```

Kako bismo se upoznali sa podacima možemo iskoristiti naredbu `head(podaci)` koja će nam prikazati prvih nekoliko instanci dataseta iz čega možemo zaključiti koje kolone imamo u setu.

```
head(podaci)
```

##	gender	Dependents	tenure	PhoneService	MultipleLines	InternetService
## 1	Male	No	8	Yes	No	<NA>
## 2	Female	<NA>	8	Yes	<NA>	Fiber optic
## 3	<NA>	No	21	Yes	No	Fiber optic
## 4	<NA>	<NA>	1	No	No phone service	DSL
## 5	<NA>	No	NA	<NA>	<NA>	Fiber optic
## 6	Male	No	69	Yes	No	No
##	StreamingTV	StreamingMovies	Contract	PaymentMethod		
## 1	Yes	Yes	<NA>	<NA>		
## 2	No	No	Month-to-month	Credit card (automatic)		
## 3	Yes	Yes	One year	Mailed check		
## 4	No	No	Month-to-month	Mailed check		
## 5	No	Yes	Month-to-month	Electronic check		
## 6	<NA>	<NA>	Two year	Bank transfer (automatic)		
##	MonthlyCharges	TotalCharges	DailyCharges	Churn		
## 1	NA	832.35	NA	Yes		
## 2	NA	548.90	NA	No		
## 3	104.55	2239.40	20.91	No		
## 4	35.90	35.90	7.18	No		
## 5	81.10	81.10	16.22	Yes		
## 6	19.30	1447.90	3.86	No		

Vidimo da imamo 14 atributa i već na prvi pogled uočavamo da imamo neke nedostajuće vrijednosti. Možemo zaključiti da imamo kategoričke i numeričke podatke kroz koje se detaljnije upoznajemo pomoću sljedećih naredbi:

```

1 cat("Kategorije spolova: ", unique(podaci$gender),
2 "\nKategorije izdržavanja (dependents): ", unique(podaci$Dependents),
3 "\nKategorije PhoneServices: ", unique(podaci$PhoneService),
4 "\nKategorije (brojevi) nekretnina mušterije koji koriste usluge: ", sort(unique(podaci$
   tenure)),
5 "\nKategorije MultipleLines: ", unique(podaci$MultipleLines),
6 "\nKategorije InternetService: ", unique(podaci$InternetService),
7 "\nKategorije StreamingTV: ", unique(podaci$StreamingTV),
8 "\nKategorije StreamingMovies: ", unique(podaci$StreamingMovies),
9 "\nKategorije ugovora: ", unique(podaci$Contract),
10 "\nKategorije metode plaćanja: ", unique(podaci$PaymentMethod),
11 "\nKategorije Churn: ", unique(podaci$Churn))

```

Rezultat izvršavanja koda je sljedeći:

Kategorije spolova: Male Female NA

Kategorije izdržavanja (dependents): No NA Yes Maybe

Kategorije PhoneServices: Yes No NA

Kategorije (brojevi) nekretnina mušterije koji koriste usluge: 0 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20
21 22 23 24 25 26 27 28 29 30 31 32 33 34 35 36 37 38 39 40 41 42 43 44 45 46 47 48 49 50 51 52 53 54 55 56
57 58 59 60 61 62 63 64 65 66 67 68 69 70 71 72

Kategorije MultipleLines: No NA No phone service Yes

Kategorije InternetService: NA Fiber optic DSL No

Kategorije StreamingTV: Yes No NA No internet service

Kategorije StreamingMovies: Yes No NA No internet service

Kategorije ugovora: NA Month-to-month One year Two year

Kategorije metode plaćanja: NA Credit card (automatic) Mailed check Electronic check Bank transfer (automatic)
abcd

Kategorije Churn: Yes No NA

1.2 Metode procjene lokacije i varijabilnosti podataka

1.2.1 Metode za procjenu lokacije

Da bismo dobili neke osnovne podatke o procjeni lokacije svih atributa iz našeg data seta možemo iskoristiti sljedeću naredbu:

```

1 summary <- lapply(podaci, summary)
2 summary

```

Nakon izvršavanja prethodnog isječka dobijamo sljedeći pregled:

<div> <div>\$gender</div> <div>Length Class Mode</div> <div>2000 character character</div> </div>							<div> <div>\$StreamingMovies</div> <div>Length Class Mode</div> <div>2000 character character</div> </div>						
<div> <div>\$Dependents</div> <div>Length Class Mode</div> <div>2000 character character</div> </div>							<div> <div>\$Contract</div> <div>Length Class Mode</div> <div>2000 character character</div> </div>						
<div> <div>\$tenure</div> <div>Min. 1st Qu. Median Mean 3rd Qu. Max.</div> <div>0.00 9.00 29.00 32.52 56.00 72.00</div> </div>							<div> <div>\$PaymentMethod</div> <div>Length Class Mode</div> <div>2000 character character</div> </div>						
<div> <div>\$PhoneService</div> <div>Length Class Mode</div> <div>2000 character character</div> </div>							<div> <div>\$MonthlyCharges</div> <div>Min. 1st Qu. Median Mean 3rd Qu. Max. NA's</div> <div>-1.22 35.00 70.30 64.41 89.40 118.65 256</div> </div>						
<div> <div>\$MultipleLines</div> <div>Length Class Mode</div> <div>2000 character character</div> </div>							<div> <div>\$TotalCharges</div> <div>Min. 1st Qu. Median Mean 3rd Qu. Max. NA's</div> <div>19.1 433.4 1415.4 2280.6 3751.7 10000.0 276</div> </div>						
<div> <div>\$InternetService</div> <div>Length Class Mode</div> <div>2000 character character</div> </div>							<div> <div>\$DailyCharges</div> <div>Min. 1st Qu. Median Mean 3rd Qu. Max. NA's</div> <div>-0.244 6.940 14.050 12.855 17.880 23.730 255</div> </div>						
<div> <div>\$StreamingTV</div> <div>Length Class Mode</div> <div>2000 character character</div> </div>							<div> <div>\$Churn</div> <div>Length Class Mode</div> <div>2000 character character</div> </div>						

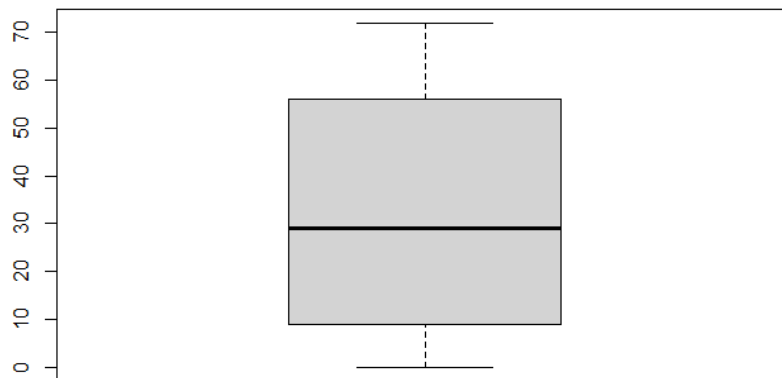
(a) (b)

Slika 1: Statistički podaci o atributima dataseta

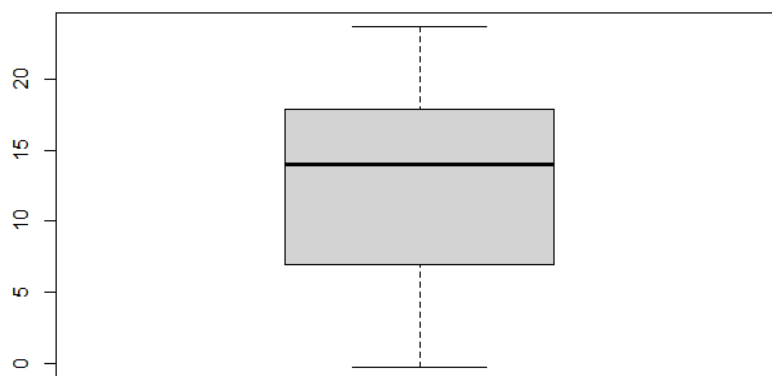
Za numeričke attribute, statističke podatke ćemo vizualizirati pomoću boxplot funkcije. Za tu svrhu možemo izdvojiti subset sa numeričkim vrijednostima:

```
1 num_cols <- subset(podaci, select = c(tenure, DailyCharges, MonthlyCharges, TotalCharges))
2 boxplots <- lapply(num_cols, boxplot)
```

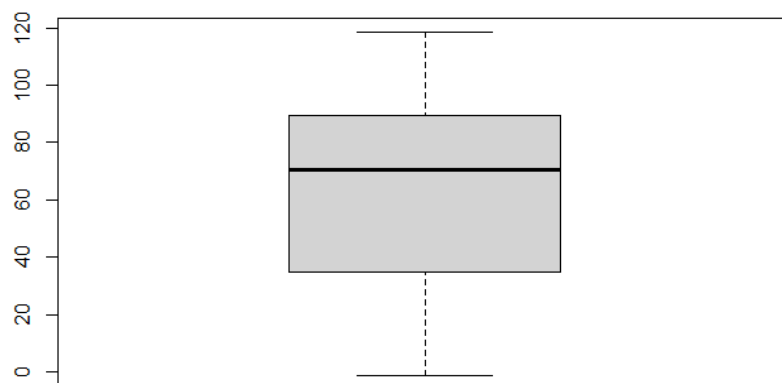
Na osnovu boxplot grafika atributa **'tenure'** (slika 2a) može se zaključiti da je distribucija podataka prilično ravnomjerna, s obzirom da se interkvartalni opseg nalazi nešto ispod sredine. Srednja vrijednost se ne nalazi na sredini opsega vrijednosti, već se nalazi nešto niže što znači da su vrijednosti broja zakupa teže ka nešto manjim vrijednostima. Srednja vrijednost se također ne nalazi na sredini između prvog i trećeg kvartala, već nešto malo ispod, što znači da više klijenata ima broj zakupa od 29 do 56. Za **dnevnu naplatu** (slika 2b) može se zaključiti da je distribucija podataka prilično ravnomjerna, s obzirom da se interkvartalni opseg nalazi skoro na sredini. Srednja vrijednost se ne nalazi na sredini opsega vrijednosti, već se nalazi nešto iznad što znači da vrijednosti dnevne naplate teže ka nešto većim vrijednostima. Srednja vrijednost se također ne nalazi na sredini između prvog i trećeg kvartala, već nešto malo iznad, što znači da više klijenata ima manju dnevnu naplatu (6.940 - 14.050). Za **mjesečnu naplatu** vrijedi ista analiza, ali sa drugačijim opsegom vrijednosti. Ukoliko zanemarimo opsege, možemo uočiti da su njihovi boxplot grafici identični (slike 2b i 2c) što nam ukazuje na potencijalnu jaku korelaciju. **Ukupna naplata** (slika 2d) ima neravnomjernu distribuciju podataka s obzirom da nema mnogo klijenata s ukupnom naplatom iznad interkvartalnog opsega (iznad 3752.175). Srednja vrijednost se ne nalazi na sredini opsega vrijednosti, već se nalazi blizu donje granice što znači da je totalna naplata uglavnom jako niska. Srednja vrijednost se također ne nalazi na sredini između prvog i trećeg kvartala, već blizu prvog kvartala, što znači da se mnogo više klijenata nalazi u donjem kraju interkvartalnog opsega (433.125 - 1415.425).



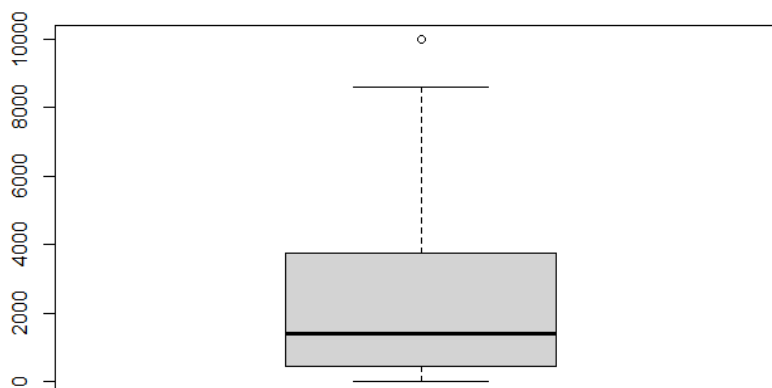
(a) Tenure



(b) Dnevna naplata



(c) Mjesečna naplata



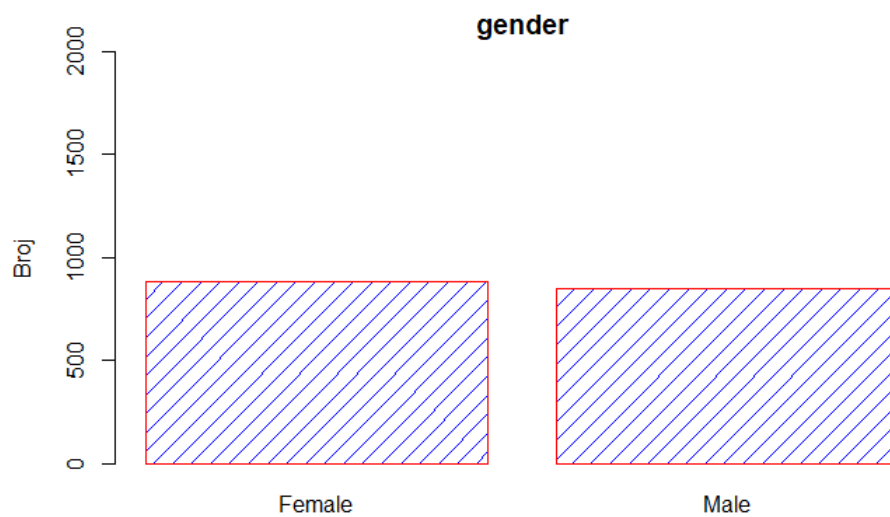
(d) Ukupna naplata

Slika 2: Boxplot grafici numeričkih atributa

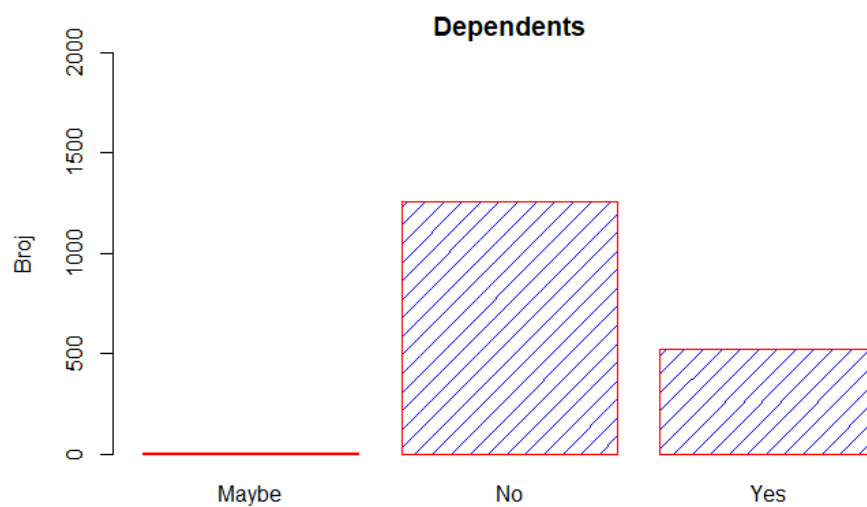
Za statistički prikaz kategoričkih atributa koristit ćemo bar plot:

```
1 for (i in colnames(podaci)){
2   if (class(podaci[[i]]) == "character")
3     barplot(table(podaci[[i]]),
4             main=colnames(podaci[i]),
5             ylim=c(0, length(podaci[[i]])),
6             ylab="Broj",
7             border="red",
8             col="blue",
9             density=10
10          )
11 }
```

Navedeni kod rezultira sa 10 grafika:

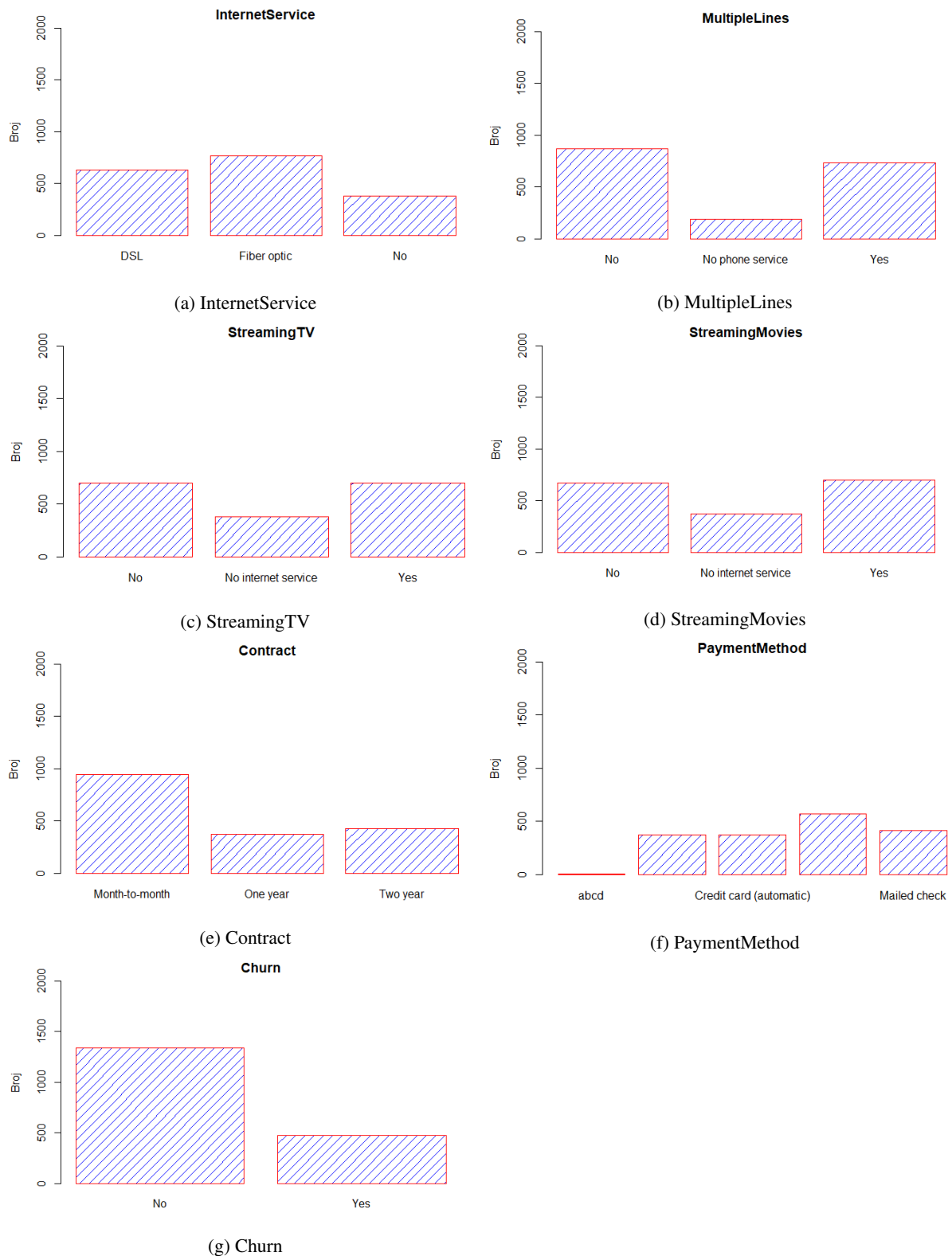


(a) Gender



(b) Dependents

Slika 3: Spol i



Slika 4: Barplot grafici kategoričkih atributa

1.2.2 Metode za procjenu varijabilnosti

```

1 trimmed_mean <- function(vektor, p)
2 {
3   vektor <- sort(vektor)

```

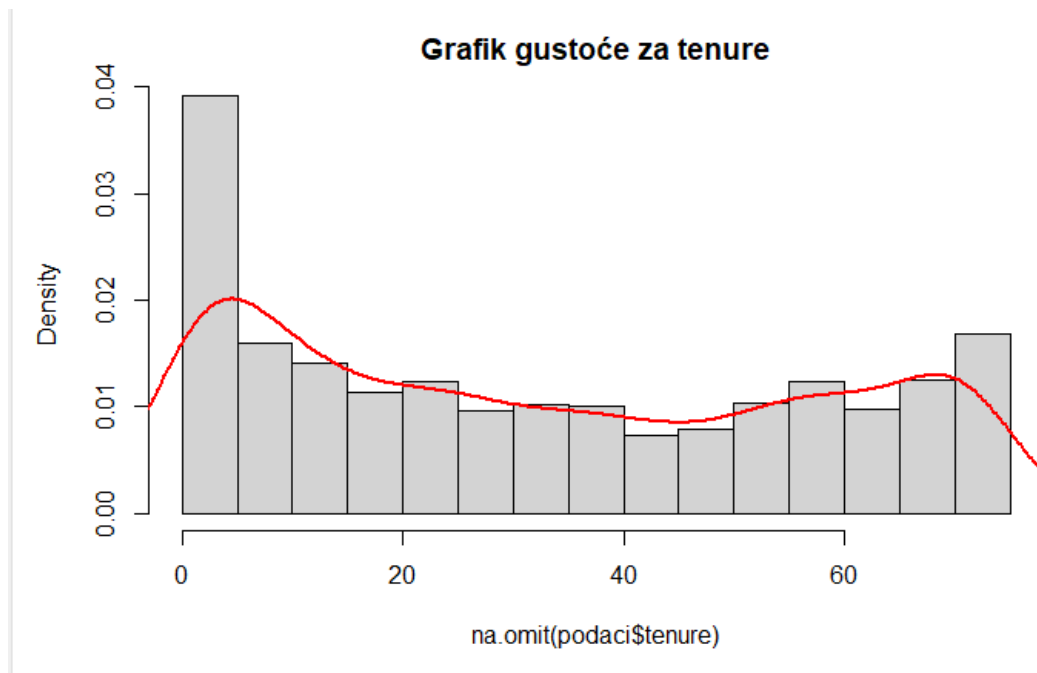
```

4 trimmed_vektor <- vektor[(1 + p) : (length(vektor) - p)]
5 return (sum(trimmed_vektor) / length(trimmed_vektor))
6 }
7
8 cat("Mean za broj nekretnina koje mu terija posjeduje:", mean(podaci_brojcani_bez_na$tenure),
  "\n")
9 cat("Trimmed_mean za broj nekretnina koje mu terija posjeduje (p=1):", trimmed_mean(podaci_
  brojcani_bez_na$tenure, 1), "\n")
10 cat("Trimmed_mean za broj nekretnina koje mu terija posjeduje (p=2):", trimmed_mean(podaci_
  brojcani_bez_na$tenure, 2), "\n")
11 cat("Trimmed_mean za broj nekretnina koje mu terija posjeduje (p=4):", trimmed_mean(podaci_
  brojcani_bez_na$tenure, 4), "\n")
12 cat("Median za broj nekretnina koje mu terija posjeduje:", median(podaci$tenure), "\n\n")
13
14 cat("Mean za dnevnu naplatu:", mean(podaci_brojcani_bez_na$DailyCharges), "\n")
15 cat("Trimmed_mean za dnevnu naplatu (p=1):", trimmed_mean(podaci_brojcani_bez_na$DailyCharges,
  1), "\n")
16 cat("Trimmed_mean za dnevnu naplatu (p=2):", trimmed_mean(podaci_brojcani_bez_na$DailyCharges,
  2), "\n")
17 cat("Trimmed_mean za dnevnu naplatu (p=4):", trimmed_mean(podaci_brojcani_bez_na$DailyCharges,
  4), "\n")
18 cat("Median za dnevnu naplatu:", median(podaci_brojcani_bez_na$DailyCharges), "\n\n")
19
20 cat("Mean za mjese nu naplatu:", mean(podaci$MonthlyCharges), "\n")
21 cat("Trimmed_mean za mjese nu naplatu (p=1):", trimmed_mean(podaci$MonthlyCharges, 1), "\n")
22 cat("Trimmed_mean za mjese nu naplatu (p=2):", trimmed_mean(podaci$MonthlyCharges, 2), "\n")
23 cat("Trimmed_mean za mjese nu naplatu (p=4):", trimmed_mean(podaci$MonthlyCharges, 4), "\n")
24 cat("Median za mjese nu naplatu:", median(podaci$MonthlyCharges), "\n\n")

```

U navedenom isječku prikazano je izračunavanje mean, median i trimmed_mean vrijednosti brojčanih atributa, ali bez instanci sa na vrijednostima brojčanih atributa, tako da se ovi parametri trebaju uzeti s rezervom, jer ukoliko bi se pojedinačno računalo za svaki brojčani atribut vrlo je moguće da će se podaci razlikovati.

Grafik gustoće za tenure prikazan je na sljedećoj slici:



Slika 5: Grafik gustoće ta tenure

1.3 Metode za procjenu korelacije između varijabli

1.3.1 Korelacija između numeričkih varijabli

Za određivanje korelacije između numeričkih varijabli, koristili smo Pearsonov koeficijent korelacije. Odredili smo korelaciju svakog numeričkog atributa sa svakim, odnosno korelacije između atributa: 'tenure', 'DailyCharges', 'MonthlyCharges' i 'TotalCharges'.

```

1 library(reshape2)
2 library(ggplot2)
3 # funkcija za odsijecanje gornjeg dijela matrice
4 get_upper_tri <- function(cormat){
5   cormat[lower.tri(cormat)]<- NA
6   return(cormat)
7 }
8 # korelacijski koeficijent koriste i NA podatke ima rezultnu vrijednost NA
9 podaci_brojcani <- subset(podaci, select = c(tenure, DailyCharges, MonthlyCharges, TotalCharges))
10 podaci_brojcani_bez_na <- na.omit(podaci_brojcani)
11 # kreiranje korelacijske matrice za sve podatke
12 cormat <- round(cor(podaci_brojcani_bez_na), 2)
13 # odsijecanje gornjeg dijela matrice
14 upper_tri <- get_upper_tri(cormat)
15 melted_cormat <- melt(upper_tri, na.rm = TRUE)
16
17 ggheatmap <- ggplot(data = melted_cormat, aes(Var2, Var1, fill = value))+
18   geom_tile(color = "white")+
19   scale_fill_gradient2(low = "blue", high = "red", mid = "white",
20     midpoint = 0.5, limit = c(0,1), space = "Lab",
21     name="Pearson\nCorrelation") +

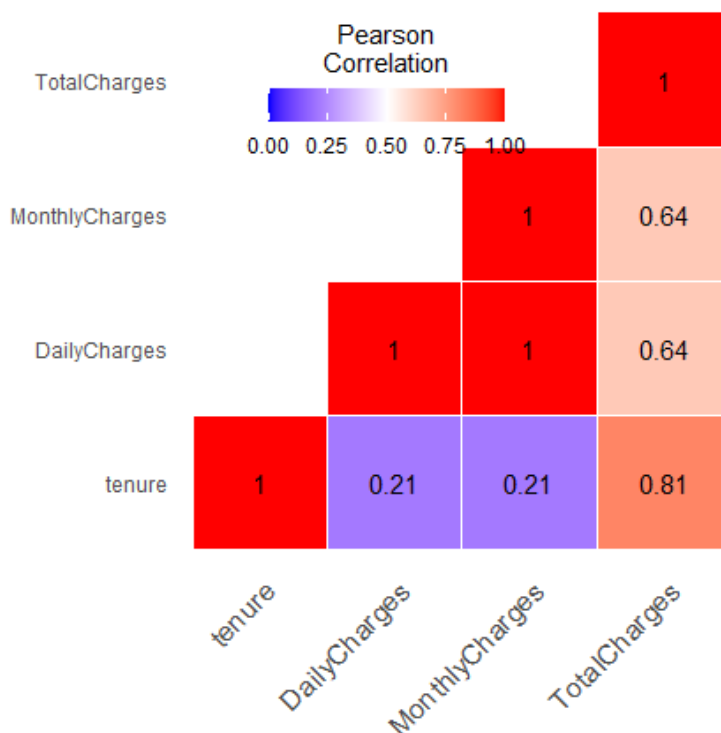
```

```

22 theme_minimal()+
23 theme(axis.text.x = element_text(angle = 45, vjust = 1,
24   size = 12, hjust = 1))+
25 coord_fixed()
26 ggheatmap +
27 geom_text(aes(Var2, Var1, label = value), color = "black", size = 4) +
28 theme(
29   axis.title.x = element_blank(),
30   axis.title.y = element_blank(),
31   panel.grid.major = element_blank(),
32   panel.border = element_blank(),
33   panel.background = element_blank(),
34   axis.ticks = element_blank(),
35   legend.justification = c(1, 0),
36   legend.position = c(0.6, 0.7),
37   legend.direction = "horizontal")+
38 guides(fill = guide_colorbar(barwidth = 7, barheight = 1,
39   title.position = "top", title.hjust = 0.5))
40
41 ggheatmap

```

U liniji 9 smo u varijablu `podaci_brojčani` izdvojili numeričke atribute, zatim smo u liniji 10 pomoću funkcije `na.omit` odstranili instance (klijente) koje nemaju definisane vrijednosti jedne ili više navedenih brojčanih atributa i preostale instance pohranili u varijablu `podaci_brojčani_bez_na`. U liniji koje slijede kreirali smo heat mapu koja pokazuje koeficijente korelacije između svakih od atributa. Heat mapa je prikazana na slici 6.



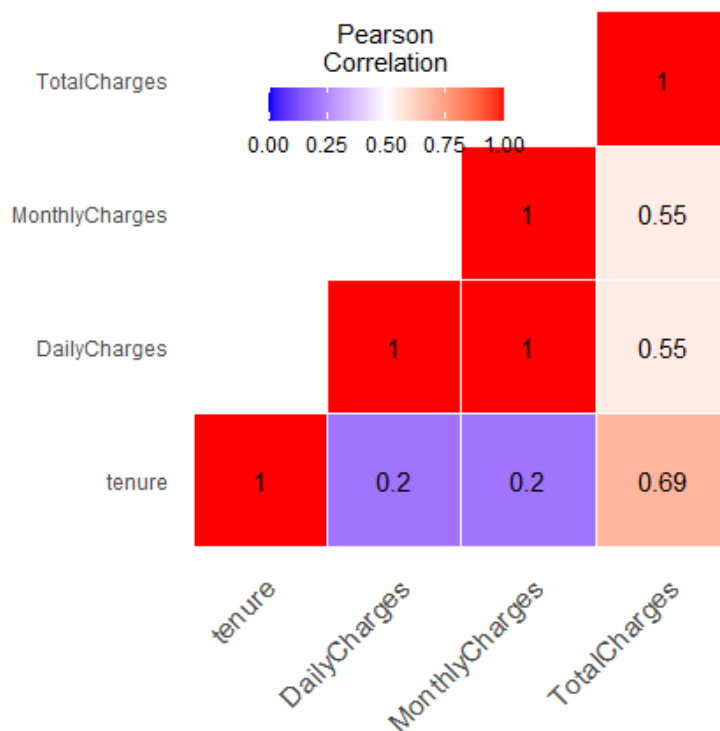
Slika 6: Pearsonov koeficijent korelacije između numeričkih atributa

Svi Pearsonovi koeficijenti korelacija su pozitivni što znači da se varijable pomjeraju u istom smjeru, odnosno da im vrijednosti uzajamno rastu ili uzajamno opadaju. Dnevna naplata i mjesečni troškovi klijenata su u potpunoj korelaciji (koeficijent 1), pa je potrebno izbaci jedan od navedenih atributa. Smatra se da su u jakoj korelaciji svi atributi čiji Pearsonov koeficijent ima vrijednost veću od 0,7 ukoliko se kreću u istom smjeru, odnosno manju od -0,7 ukoliko se kreću u suprotnim smjerovima. Prema tome, koeficijent koji iznosi 0.81 indicira da su u jakoj korelaciji broj zakupa (tenure) i ukupna naplata (TotalCharges), pa je potrebno izbaci jedan od ovih atributa.

Možemo provesti i analizu korelacije u slučaju da ne izbacujemo NA vrijednosti. Da bismo to učinili, prvo ćemo popuniti nedostajuće vrijednosti sa medijanima odgovarajućeg atributa. To ćemo učiniti sljedećim kodom:

```
1 podaci_brojcani <- podaci_brojcani %>%
2 mutate_if(is.numeric, function(x) ifelse(is.na(x), median(x, na.rm = T), x))
```

Ukoliko sada nad varijablom podaci_brojcani izvršimo analizu korelacije dobit ćemo sljedeći grafik:



Slika 7: Pearsonov koeficijent korelacije između numeričkih atributa (NA popunjene medijanom)

Vidimo da se koeficijent korelacije između atributa TotalCharges i tenure malo smanjio, no i dalje je jak i iznosi približno 0.7. S obzirom da je koeficijent na granici za izbacivanje atributa, možemo još i napraviti subset koji će se sastojati iz atributa 'tenure' i 'TotalCharges' bez instanci koje posjeduju bar jednu NA vrijednost i izračunati njihovu korelaciju. Navedeno je prikazano u sljedećem isječku:

```
1 total_tenure <- subset(podaci, select = c(tenure, TotalCharges))
2 total_tenure_bez_na <- na.omit(total_tenure)
3 cor <- cor.test(total_tenure_bez_na$tenure, total_tenure_bez_na$TotalCharges, method = "
  pearson")
4 cat("Pearsonov koeficijent korelacije tenure-TotalCharges: ", cor$estimate)
```

Rezultat izvršavanja isječka je:

Pearsonov koeficijent korelacije tenure-TotalCharges: 0.8143779

Dakle, definitivno ćemo izbaciti ili atribut TotalCharges ili atribut tenure.

1.3.2 Korelacija između kategoričkih varijabli

Za određivanje korelacije između kategoričkih varijabli koristili smo chi-square koeficijent korelacije. Odredili smo korelaciju svakog kategoričkog atributa sa svakim i koeficijente smjestili u korelacijsku matricu. Prije određivanja chi-square izbačeni su outlieri, s obzirom da je to jedan od uslova da bi analiza bila uspješna.

```

1 chi_matrica <- matrix(0,ncol(podaci), ncol(podaci))
2 rownames(chi_matrica) <- colnames(podaci)
3 colnames(chi_matrica) <- colnames(podaci)
4 critical_matrica <- matrix(0,ncol(podaci), ncol(podaci))
5 rownames(critical_matrica) <- colnames(podaci)
6 colnames(critical_matrica) <- colnames(podaci)
7 razlika_matrica <- matrix(0,ncol(podaci), ncol(podaci))
8 rownames(razlika_matrica) <- colnames(podaci)
9 colnames(razlika_matrica) <- colnames(podaci)
10 for (k in colnames(podaci)){
11   atribut1 <- sort(unique(podaci[[k]]))
12   for (l in colnames(podaci)){
13     atribut2 <- sort(unique(podaci[[l]]))
14     if(class(podaci[[k]]) == "character" & class(podaci[[l]]) == "character"){
15       data2 <- matrix(0, length(atribut1), length(atribut2))
16       rownames(data2) <- atribut1
17       colnames(data2) <- atribut2
18       cat("\nMatrica korelacije atributa ", colnames(podaci[k]), " i ", colnames(podaci[l]),
19         ":\n\n")
20       for (i in 1 : length(atribut1)) {
21         for (j in 1 : length(atribut2)) {
22
23           redovi <- subset(podaci, (podaci[[k]] == atribut1[i] & podaci[[l]] ==
24             atribut2[j]))
25           data2[i, j] = length(redovi[[k]])
26         }
27       }
28       print(data2)
29       chi <- chisq.test(data2)
30       chi_matrica[k,l] <- chi$statistic
31       critical <- qchisq(p = chi$p.value, df = chi$parameter)
32       critical_matrica[k,l] <- critical
33       if(chi$statistic - critical >0){
34         razlika_matrica[k,l] <- chi_matrica[k,l] - critical
35       }
36     }
37   }
38   cat("Chi matrica:\n\n")
39   print(chi_matrica)
40   cat("Matrica krticnih vrijednosti:\n\n")
41   print(critical_matrica)
42   cat("Matrica razlike chi matrice i kriticnih vrijednosti:\n\n")
43   #print(razlika_matrica)
44   razlika_matrica <- razlika_matrica[,!colnames(razlika_matrica) %in% c("DailyCharges", "tenure",
45     "MonthlyCharges", "TotalCharges")]
46   razlika_matrica <- razlika_matrica[!rownames(razlika_matrica) %in% c("DailyCharges", "tenure",

```

```

    "MonthlyCharges", "TotalCharges"),,]
46 print(razlika_matrica)
47 razlika_df <- as.data.frame(t(razlika_matrica))
48 print(razlika_df)
49 normalize <- function(x) {
50   return (round((x - min(x)) / (max(x) - min(x)),5))
51 }
52 for (i in colnames(razlika_df)){
53   razlika_df[[i]]<-normalize(razlika_df[[i]])
54 }
55 print(razlika_df)
56 razlika_matrica <- as.matrix(razlika_df)
57 print(razlika_matrica)

```

U linijama 1 - 9 deklariramo matrice: `chi_matrica`, `critical_matrica` i `razlika_matrica` i postavimo njihove vrijednosti na 0. Zatim u for petlji u liniji 10 prolazimo kroz sve atribute dataframea i u svakoj iteraciji petlje po jedan pohranjujemo u varijablu 'atribut1'. Zatim u liniji 12 otvaramo novu petlju kojom također prolazimo kroz sve atribute i pohranjujemo ih u varijablu atribut 2. Ukoliko su u pitanju kategorički atributi (uslov u liniji 14) formira se korelacijska matrica. Zatim se matrica šalje u funkciju `chisq.test` u liniji 28 i u chi matricu se dodaje vrijednost `chi$statistic`. Kritične vrijednosti se dodaju u odgovarajuća polja u critical matrici. Konačno, u liniji 33 se u matricu razlike dodaje razlika između chi-square vrijednosti i kritične vrijednosti ukoliko je razlika pozitivna, odnosno ukoliko su odgovarajući atributi u korelaciji. Na slici 8 su prikazane neke od korelacijskih matrica. Na slici 10 je prikazana heat mapa korelacije kategoričkih atributa.

Matrica korelacije atributa gender i Dependents

	No	Yes
Female	543	224
Male	542	229

Matrica korelacije atributa gender i PhoneService

	No	Yes
Female	79	685
Male	66	665

Matrica korelacije atributa gender i MultipleLines

	No	No phone service	Yes
Female	376	90	331
Male	365	71	312

Matrica korelacije atributa gender i InternetService

	DSL	Fiber optic	No
Female	285	345	147
Male	253	317	175

Matrica korelacije atributa gender i StreamingTV

	No	No internet service	Yes
Female	318	153	314
Male	287	171	290

Matrica korelacije atributa gender i StreamingMovies

	No	No internet service	Yes
Female	290	156	320
Male	286	170	286

Matrica korelacije atributa gender i Contract

	Month-to-month	One year	Two year
Female	426	167	180
Male	393	154	185

(a)

Matrica korelacije atributa Dependents i StreamingMovies

	No	No internet service	Yes
No	452	203	443
Yes	137	136	174

Matrica korelacije atributa Dependents i Contract

	Month-to-month	One year	Two year
No	673	222	199
Yes	162	116	174

Matrica korelacije atributa Dependents i PaymentMethod

	Bank transfer (automatic)	Credit card (automatic)	Electronic check	Mailed check
No	205	208	417	247
Yes	119	117	91	122

Matrica korelacije atributa Dependents i Churn

	No	Yes
No	774	372
Yes	411	57

Matrica korelacije atributa PhoneService i gender

	Female	Male
No	79	66
Yes	685	665

Matrica korelacije atributa PhoneService i Dependents

	No	Yes
No	99	52
Yes	994	396

Matrica korelacije atributa PhoneService i PhoneService

	No	Yes
No	169	0

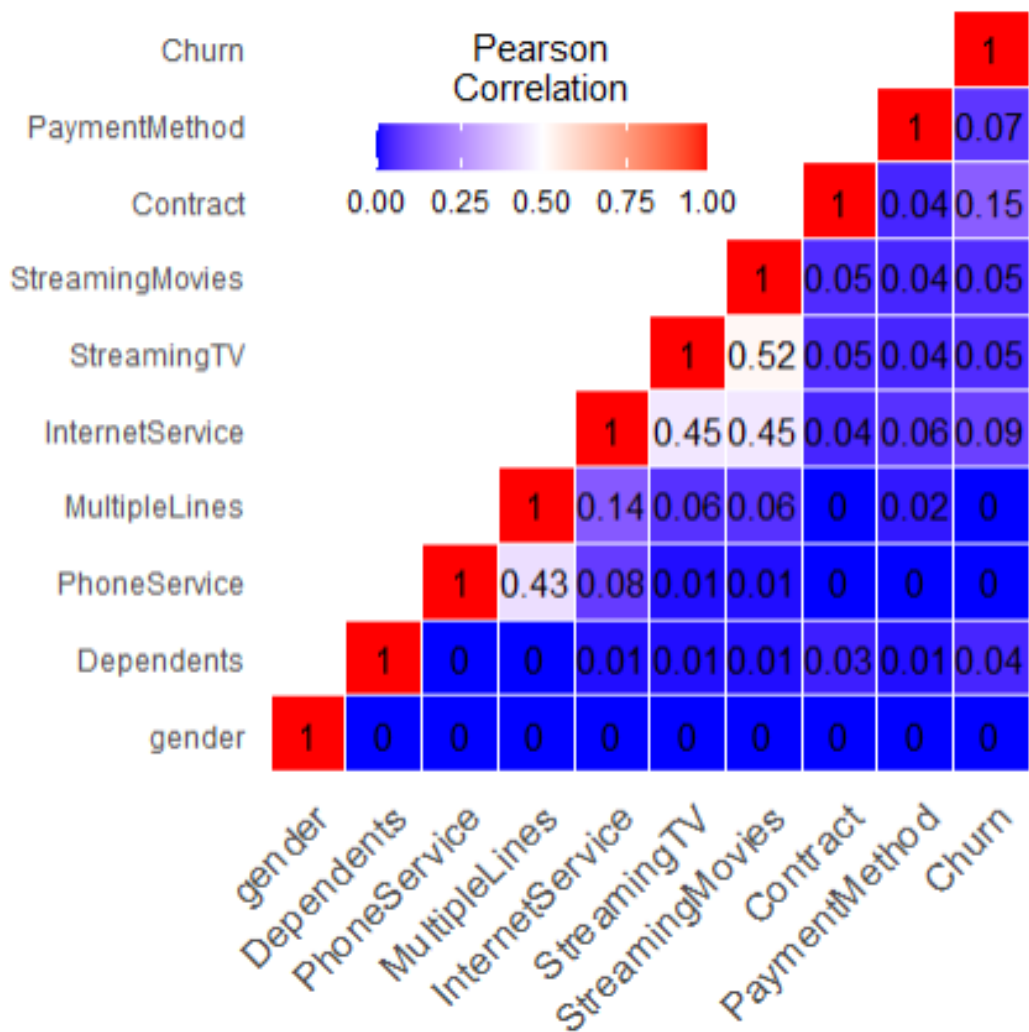
(b)

Slika 8: Neke od korelacijskih matrica

Formirana normalizirana matrica razlike chi-square vrijednosti i kritičnih vrijednosti:

	gender <dbl>	Dependents <dbl>	PhoneService <dbl>	MultipleLines <dbl>	InternetService <dbl>	StreamingTV <dbl>	StreamingMovies <dbl>	Contract <dbl>	PaymentMethod <dbl>
gender	1.00000	0.00000	0.00014	0.00002	0.00000	0.00000	0.00000	0.00000	0.00089
Dependents	0.00000	1.00000	0.00118	0.00000	0.01227	0.00815	0.00823	0.02817	0.00954
PhoneService	0.00014	0.00114	1.00000	0.43350	0.08275	0.01237	0.01291	0.00000	0.00000
MultipleLines	0.00002	0.00000	0.89999	1.00000	0.14012	0.05683	0.05523	0.00404	0.01915
InternetService	0.00271	0.02712	0.17313	0.14031	1.00000	0.45491	0.45180	0.04104	0.06185
StreamingTV	0.00118	0.01744	0.02671	0.05705	0.45501	1.00000	0.52498	0.04746	0.04221
StreamingMovies	0.00077	0.01694	0.02698	0.05440	0.44464	0.51659	1.00000	0.05427	0.03910
Contract	0.00000	0.05529	0.00000	0.00394	0.03903	0.04600	0.05372	1.00000	0.03833
PaymentMethod	0.00301	0.02815	0.00036	0.02789	0.08911	0.06109	0.05765	0.05707	1.00000
Churn	0.00000	0.03896	0.00000	0.00196	0.04206	0.02302	0.02444	0.08022	0.02454

Slika 9: Matrica razlike kategoričkih atributa



Slika 10: Heat mapa korelacije kategoričkih atributa

1.4 Predprocesiranje podataka

1.4.1 Odbacivanje atributa sa visokim stepenom korelacije

Prva dva atributa koja su u korelaciji su 'MonthlyCharges' i 'DailyCharges'. Ova dva atributa imaju iste **koefficiente korelacije sa ostalim atributima** što se vidi na slikama 6 i 7. Da bismo odlučili koji atribut izbaciti, provjerit ćemo koji atribut ima više NA vrijednosti jer smatramo da je bolje izbaciti atribut koji ima više nedostajućih vrijednosti.

```

1 cat("Broj NA vrijednosti u koloni MonthlyCharges: ",
2 length(subset(podaci, {is.na(podaci$MonthlyCharges)})$MonthlyCharges),
3 "\nBroj NA vrijednosti u koloni DailyCharges: ",
4 length(subset(podaci, {is.na(podaci$DailyCharges)})$DailyCharges))

```

Nakon izvršavanja isječka dobijamo sljedeći rezultat:

Broj NA vrijednosti u koloni MonthlyCharges: 256

Broj NA vrijednosti u koloni DailyCharges: 255

S obzirom da imaju približan broj NA vrijednosti, izbacit ćemo atribut dnevne naplate jer mjesečna naplata ima približnije vrijednosti ukupnoj naplati nego dnevna naplata. Ovo ćemo uraditi pomoću sljedećeg koda:

```
1 podaci <- subset(podaci, select = -c(DailyCharges))
```

Analognu analizu radimo i za attribute 'tenure' i 'TotalCharges':

```
1 cat("Broj NA vrijednosti u koloni TotalCharges: ",  
2 length(subset(podaci, {is.na(podaci$TotalCharges)}))$TotalCharges),  
3 "\nBroj NA vrijednosti u koloni tenure: ",  
4 length(subset(podaci, {is.na(podaci$tenre)}))$tenure))
```

Nakon izvršavanja isječka dobijamo sljedeći rezultat:

Broj NA vrijednosti u koloni TotalCharges: 276

Broj NA vrijednosti u koloni tenure: 0

Najbitniju ulogu u izbacivanju atributa igra **stepen korelacije sa ostalim atributima**, pa tako vidimo na slici 6 da atribut 'TotalCharges' ima visok stepen korelacije sa 'MonthlyCharges' i 'DailyCharges' (0,64), dok atribut 'tenure' ima nešto niže koeficijente (0.21). Također, atribut 'TotalCharges' ima više nedostajućih vrijednosti. Iz navedenih razloga izbacujemo atribut 'TotalCharges'.

```
1 podaci <- subset(podaci, select = -c(TotalCharges))
```

Svi parovi kategoričkih atributa čija je vrijednost nula nisu u korelaciji. Možemo primijetiti da atributi **phoneService** i **MultipleLines** imaju nešto viši stepen korelacije (0,43) što ima smisla jer više linija ne može imati neko ko ne koristi telefonske usluge. Također, parovi **InternetService** i **StreamingTV** (0,45), **InternetService** i **StreamingMovies** (0,45), **StreamingTV** i **StreamingMovies** (0,52). Možemo izbaciti kategoriju **MultipleLines** s obzirom da ima veću korelaciju i sa ostalim atributima u odnosu na phoneService. Sasvim je svejedno da li ćemo odbaciti **StreamingMovies** ili **StreamingTV**, s obzirom da su im korelacije sa ostalim atributima jednake. Izbacit ćemo kategoriju koja ima više NA vrijednosti.

```
cat("Broj NA vrijednosti u koloni StreamingMovies: ",length(subset(podaci, {is.na(podaci$
  StreamingMovies)}))$StreamingMovies), "\nBroj NA vrijednosti u StreamingTV: ",length(subset
  (podaci, {is.na(podaci$StreamingTV)}))$StreamingTV))
```

Kod iznad daje rezultat:

Broj NA vrijednosti u koloni StreamingMovies: 252

Broj NA vrijednosti u StreamingTV: 224

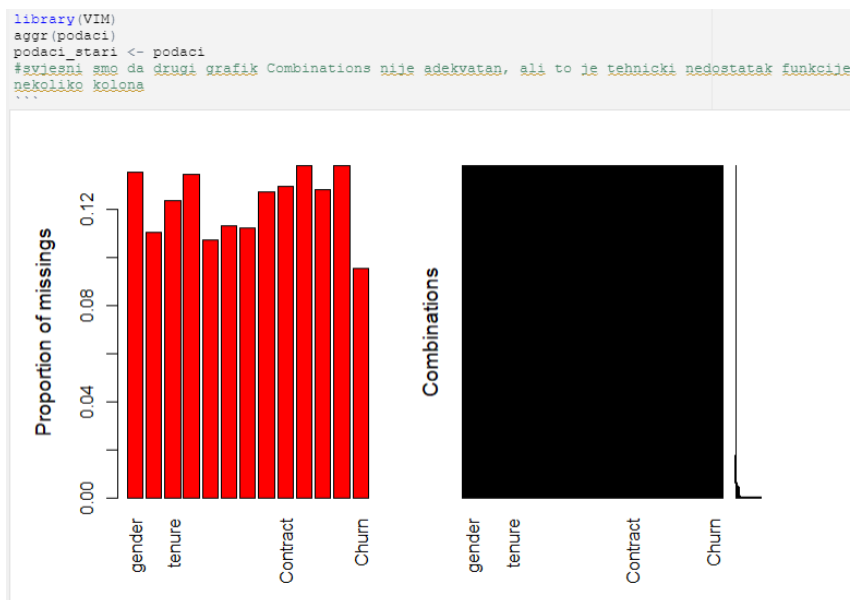
Dakle, izbacit ćemo **StreamingMovies**. Možemo još izbaciti ili **StreamingTV** ili **InternetService**. Internet-Service ima nešto veću korelaciju sa ostalim varijablama, pa ćemo nju izbaciti.

Izbacivanje dovršavamo sljedećim kodom:

```
podaci <- subset(podaci, select = -c(MultipleLines, StreamingMovies, InternetService))
```

1.4.2 Popunjavanje nedostajućih vrijednosti

Analiza na vrijednosti pokazuje sljedeću situaciju:



Slika 11: Analiza NA vrijednosti

Vidimo da u svim atributima imamo nedostajućih vrijednosti, stoga ih je potrebno popuniti. Oko tri četvrtine instanci ima bar jednu nedostajuću vrijednost što može biti problem za naše modele:

```

Vidimo da se u navedenim kategorijama nalaze NA vrijednosti, što znači da imamo nedostajućih vrijednosti. U nastavku ćemo izvršiti analizu distribucije nedostajućih
vrijednosti: <!-- -->

```{r}
#length(subset(podaci, (is.na(podaci$gender)))$MultipleLines)
#length(subset(podaci, is.na(podaci$MultipleLines)))
cat("Broj instanci dataseta: ", length(podaci$gender) , "\n",
#ovdje fali MonthlyCharges u subsetu jer je u chunku iznad izbacena ta kolona
"Broj instanci sa NA vrijednostima u bilo kojoj koloni: ", length(subset(podaci, (is.na(gender) | is.na(Dependents) | is.na(tenure) | is.na(PhoneService) |
is.na(MultipleLines) | is.na(InternetService) | is.na(StreamingTV) | is.na(StreamingMovies) | is.na(Contract) | is.na(PaymentMethod) | is.na(TotalCharges) |
is.na(MonthlyCharges) | is.na(Churn)))$MultipleLines))
...

```

Broj instanci dataseta: 2000  
 Broj instanci sa NA vrijednostima u bilo kojoj koloni: 1527

**Slika 12:** Broj instanci sa bar jednim nedostajućim podatkom

Brojčane attribute koje nećemo odbaciti probali smo popuniti na primitivniji način (samo medijanom). To smo uradili pomoću sljedećeg isječka koda:

```

1 library(tidyverse)
2 podaci <- podaci %>%
3 mutate_if(is.numeric, function(x) ifelse(is.na(x), median(x, na.rm = T), x))

```

Drugi način koji smo pokušali je da atribut tenure popunimo medijanom, a da atribut MonthlyCharges popunimo prosječnom vrijednosti po jednoj kategoriji atributa tenure (s obzirom da ima diskretne vrijednosti). Naši modeli su oba slučaja davali slične rezultate.

```

1 podaci_stari <- podaci
2 not_na <- subset(podaci, is.na(tenure) == FALSE)$tenure
3 median <- median(not_na)
4 for (i in 1 : length(podaci$tenure))
5 {
6 if (is.na(podaci$tenure[i]) == TRUE)
7 {
8 podaci$tenure[i] <- median
9 }
10 }
11
12 podaci_kategorija <- podaci
13 podaci_kategorija$tenure <- factor(podaci_kategorija$tenure)
14 not_na <- subset(podaci_kategorija, is.na(MonthlyCharges) == FALSE)
15 srednje_vrijednosti <- c()
16 svi_ratinzi <- levels(podaci_kategorija$tenure)
17 for (i in 1 : length(svi_ratinzi))
18 {
19 srednja_vrijednost <- median(subset(not_na, tenure ==
20 svi_ratinzi[i])$MonthlyCharges)
21 srednje_vrijednosti <- append(srednje_vrijednosti, srednja_vrijednost)
22 }
23 for (i in 1 : length(podaci_kategorija$MonthlyCharges))
24 {
25 if (is.na(podaci_kategorija$MonthlyCharges[i]) == TRUE)
26 {
27 index <- podaci_kategorija$tenure[i]
28 podaci_kategorija$MonthlyCharges[i] <- srednje_vrijednosti[index]

```

```

29 }
30 }
31 podaci_kategorija$tenure <- as.numeric(as.character(podaci_kategorija$tenure))
32 podaci <- podaci_kategorija

```

Vrijednosti kategoričkih varijabli smo popunili na sljedeći način:

```

1 male<-0
2 dependent<-0
3 for (i in 1 : length(podaci$gender))
4 {
5
6 #popunjavanje StreamingTV na No internet service ako je InternetService No
7 if (is.na(podaci$StreamingTV[i]) == TRUE & is.na(podaci$InternetService[i]) == FALSE &
8 podaci$InternetService[i]=="No")
9 {
10 podaci$StreamingTV[i] <- "No internet service"
11 }
12 #popunjavanje PhoneService na Yes ako je MultipleLines yes
13 if (is.na(podaci$PhoneService[i]) == TRUE & is.na(podaci$MultipleLines[i]) == FALSE & podaci
14 $MultipleLines[i]=="Yes")
15 {
16 podaci$PhoneService[i] <- "Yes"
17 }
18
19 if (is.na(podaci$gender[i]) == TRUE)
20 {
21 if (male == 1){
22 podaci$gender[i] <- "Female"
23 male<-0
24 } else {
25 podaci$gender[i] <- "Male"
26 male<-1
27 }
28 }
29
30 if (is.na(podaci$Dependents[i]) == TRUE)
31 {
32 if (dependent == 1){
33 podaci$Dependents[i] <- "No"
34 dependent<-0
35 } else {
36 podaci$Dependents[i] <- "Yes"
37 dependent<-1
38 }
39 }
40
41 #ako su PhoneService i MultipleLines oba Na
42 if (is.na(podaci$PhoneService[i]) == TRUE & is.na(podaci$MultipleLines[i])==TRUE){
43 podaci$PhoneService[i]<- 'No'
44 }
45 }

```

```

40 podaci$MultipleLines[i]<- 'No phone service '
41 }
42 #ako su InternetService i StreamingTV i StreamingMovies svi Na
43 if (is.na(podaci$InternetService[i]) == TRUE & is.na(podaci$StreamingTV[i])==TRUE & is.na(
44 podaci$StreamingMovies[i])==TRUE){
45 podaci$InternetService[i]<- 'No'
46 podaci$StreamingTV[i]<- 'No internet service '
47 podaci$StreamingMovies[i]<- 'No internet service '
48 }
49 #ako su InternetService i StreamingTV Na i StreamingMovies No
50 if (is.na(podaci$InternetService[i]) == TRUE & is.na(podaci$StreamingTV[i])==TRUE & podaci$
51 StreamingMovies[i]=='No'){
52 podaci$InternetService[i]<- 'No'
53 podaci$StreamingTV[i]<- 'No internet service '
54 }
55 #ako su InternetService i StreamingMovies Na i StreamingTV No
56 if (is.na(podaci$InternetService[i]) == TRUE & is.na(podaci$StreamingMovies[i])==TRUE &
57 podaci$StreamingTV[i]=='No'){
58 podaci$InternetService[i]<- 'No'
59 podaci$StreamingMovies[i]<- 'No internet service '
60 }
61 if (is.na(podaci$InternetService[i]) == TRUE){
62 podaci$InternetService[i]<- 'DSL '
63 }
64 if (is.na(podaci$PhoneService[i]) == TRUE){
65 podaci$PhoneService[i]<- 'Yes '
66 }
67 if (is.na(podaci$MultipleLines[i]) == TRUE){
68 podaci$MultipleLines[i]<- 'Yes '
69 }
70 if (is.na(podaci$StreamingTV[i]) == TRUE){
71 podaci$StreamingTV[i]<- 'Yes '
72 }
73 if (is.na(podaci$StreamingMovies[i]) == TRUE){
74 podaci$StreamingMovies[i]<- 'Yes '
75 }
76 if (is.na(podaci$Contract[i]) == TRUE){
77 podaci$Contract[i]<- 'Month-to-month '
78 }
79 if (is.na(podaci$PaymentMethod[i]) == TRUE){
80 podaci$PaymentMethod[i]<- 'Electronic check '
81 }
82 }

```

### 1.4.3 Odbacivanje pronađenih outliera

Postoji mnogo tehnika pomoću kojih možemo uočiti outliere. Jedna od tehnika za kategoričke attribute ili attribute sa diskretnim vrijednostima jeste pomoću analize barplot grafika. Pa tako na grafiku 3b možemo uočiti kako imamo zanemariv broj klijenata čija je vrijednost atributa 'Dependents' 'Maybe'. U sljedećem isječku pronalazimo te instance:

```
print(subset(podaci , Dependents == "Maybe"))
```

	gender <chr>	Dependents <chr>	tenure <int>	PhoneService <chr>	MultipleLines <chr>	InternetService <chr>
550	Female	Maybe	61	Yes	No	DSL

Slika 13: Outlier atributa 'Dependents'

Na slici 4f uočavamo da postoji outlier kategorije 'PaymentMethod' koji ima vrijednost 'abcd':

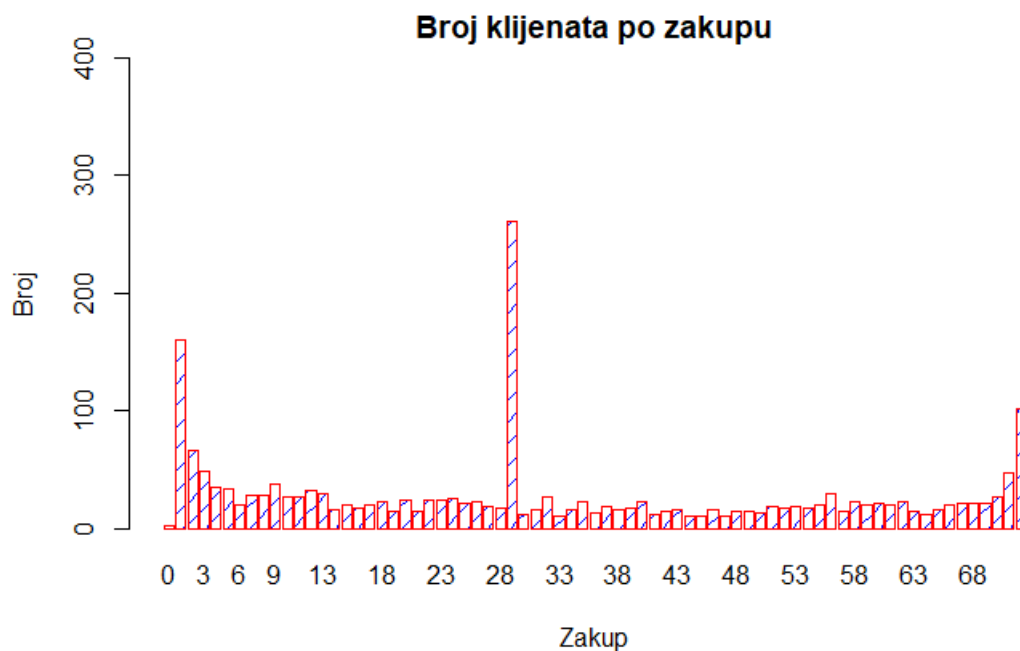
```
print(subset(podaci , PaymentMethod == "abcd"))
```

	gender <chr>	Dependents <chr>	tenure <int>	PhoneService <chr>	MultipleLines <chr>	InternetService <chr>	StreamingService <chr>
70	NA	No	1	Yes	No	No	No intern

Slika 14: Outlier atributa 'PaymentMethod'

Preostaje nam da ispitamo da li ima outliera u atributu 'tenure' i za tu svrhu pravimo barplot dijagram.





Slika 15: Barplot atributa 'tenure'

Možemo primijetiti da outlier ima vrijednost atributa tenure 0. U pitanju su dvije instance:

	gender <chr>	Dependents <chr>	tenure <int>	PhoneService <chr>	MultipleLines <chr>	InternetService <chr>
76	Female	Yes	0	Yes	Yes	DSL
1632	Female	NA	0	NA	No	DSL

Slika 16: Outlieri atributa 'tenure'

Pronađene outlieri ćemo ukloniti iz dataframea sljedećom naredbom:

```
podaci <- podaci[!rownames(podaci) %in% c("550", "70", "76", "1632"),]
```

Isto je bilo moguće i koristeći subset funkciju sa kriterijima, ali s obzirom da nema puno instanci koje treba ukloniti navedeni način je prihvatljiv.

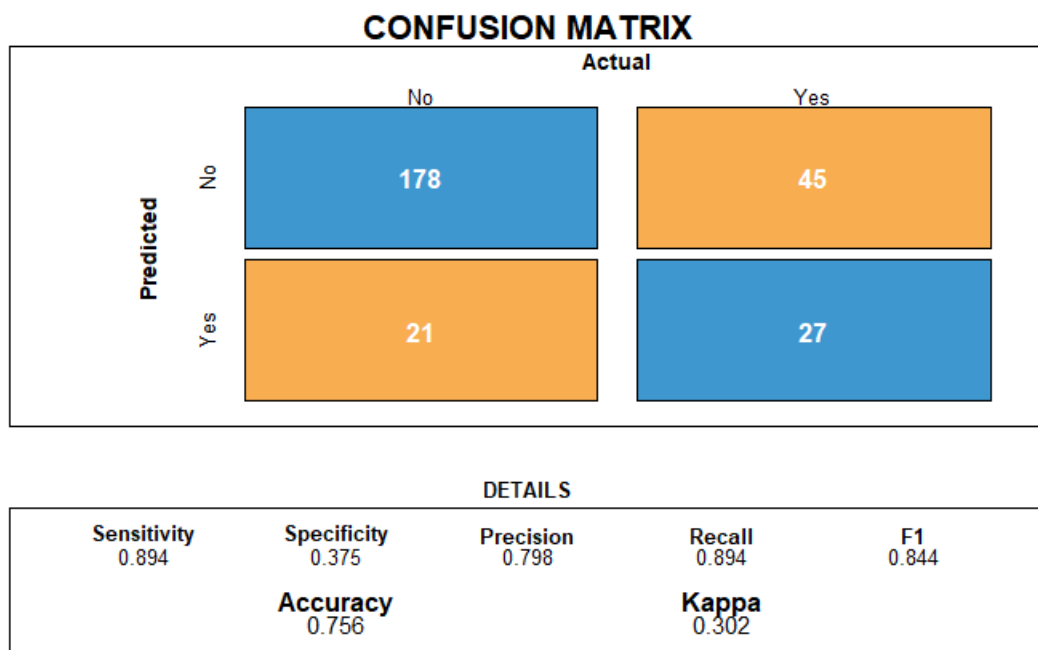
## 2 Zadatak 2

Prije kreiranja modela uradili smo faktorizaciju kategoričkih varijabli i podjeli skup podataka na trening i testni podskup. Zbog broja instanci u skupu, podjelu smo uradili tako da se 85% podataka nalazi u trening podskupu, a 15% u testnom podskupu.

### 2.1 Izgradnja modela klasifikacije

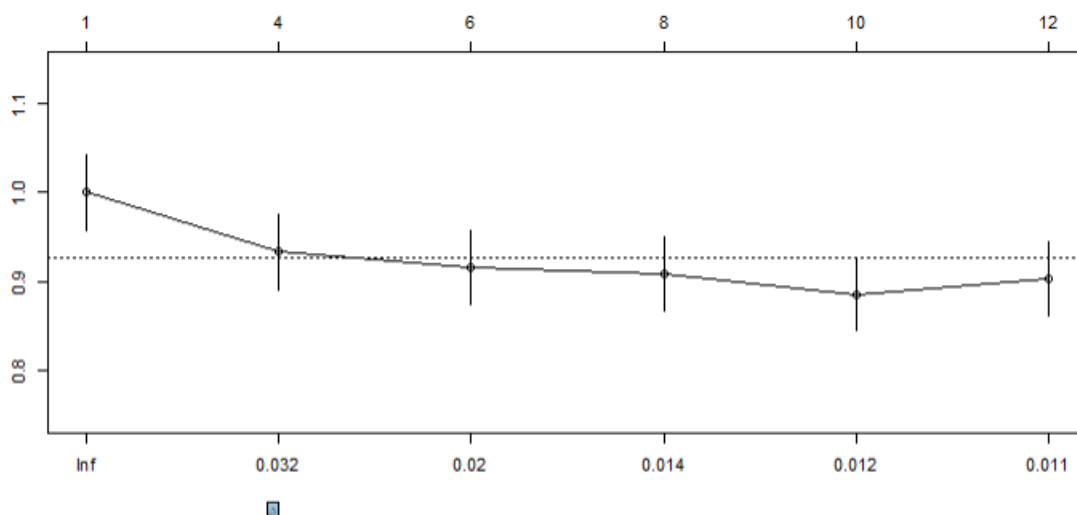
\*Napomena: Vrijednosti Sensivity i Specifity su zamijenjene.

### 2.1.1 Drvo odlučivanja koje kao mjeru atributa selekcije koristi informacijsku dobit



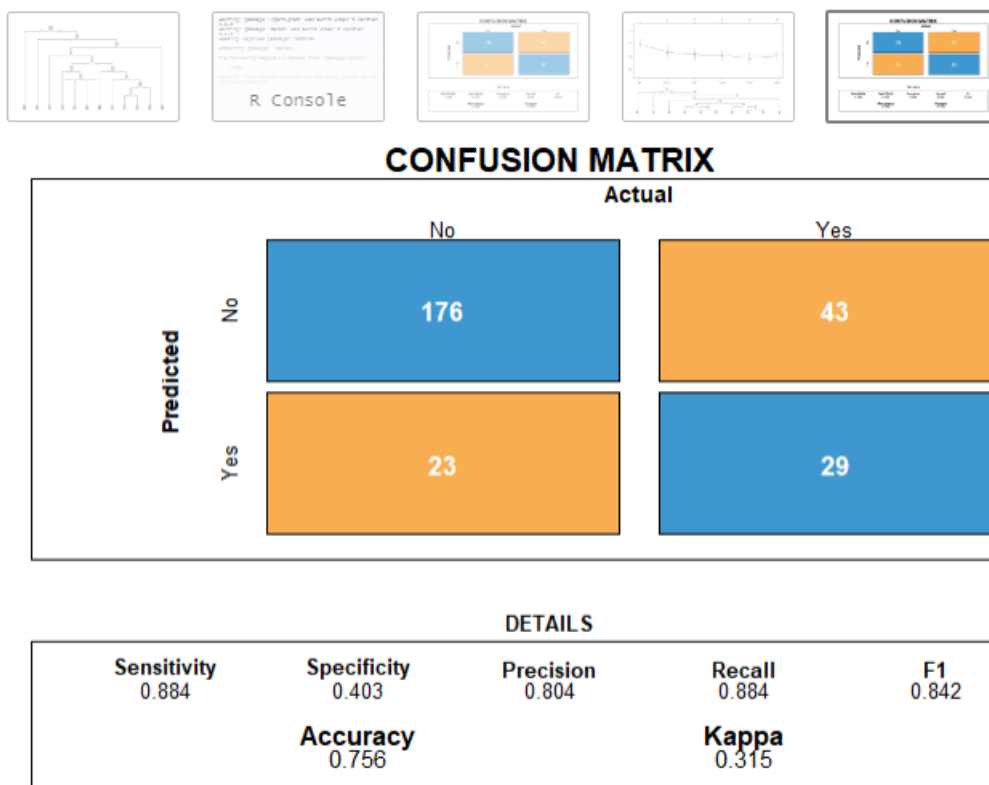
**Slika 17:** Konfuzijska matrica za drvo odlučivanja koje kao mjeru atributa selekcije koristi informacijsku dobit

Kao što možemo vidjeti, Sensitivity je niska što znači da klasifikator ne klasifikuje ispravno stvarno pozitivne instance.



**Slika 18:** Roc kriva za drvo odlučivanja koje kao mjeru atributa selekcije koristi informacijsku dobit

Na ovaj model smo primijeli čišćenje stabla i dobili smo bolje rezultate za Sensivity.



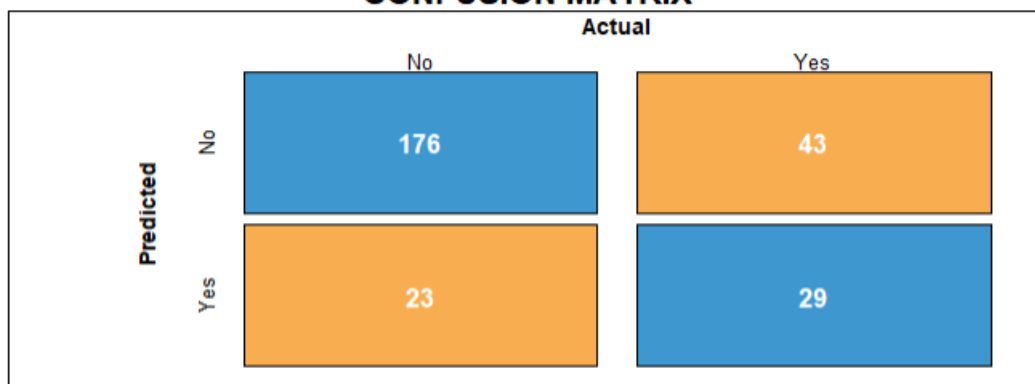
Slika 19: Konfuzijska matrica za drvo odlučivanja koje kao mjeru atributa selekcije koristi informacijsku dobit

### 2.1.2 Drvo odlučivanja koje kao mjeru atributa selekcije koristi gini indeks

Za model koji koristi gini indeks smo dobili slične rezultate kao kod modela koji koristi informacijsku dobit.



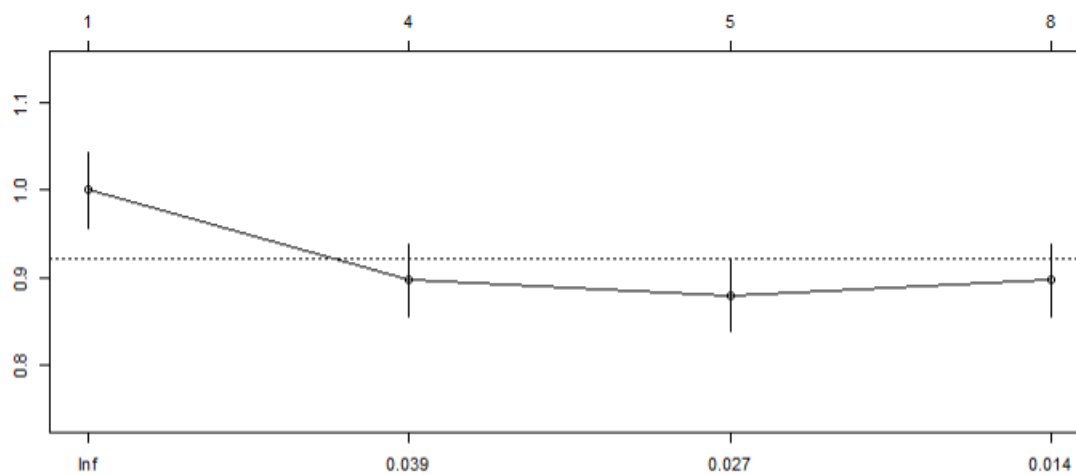
## CONFUSION MATRIX



## DETAILS

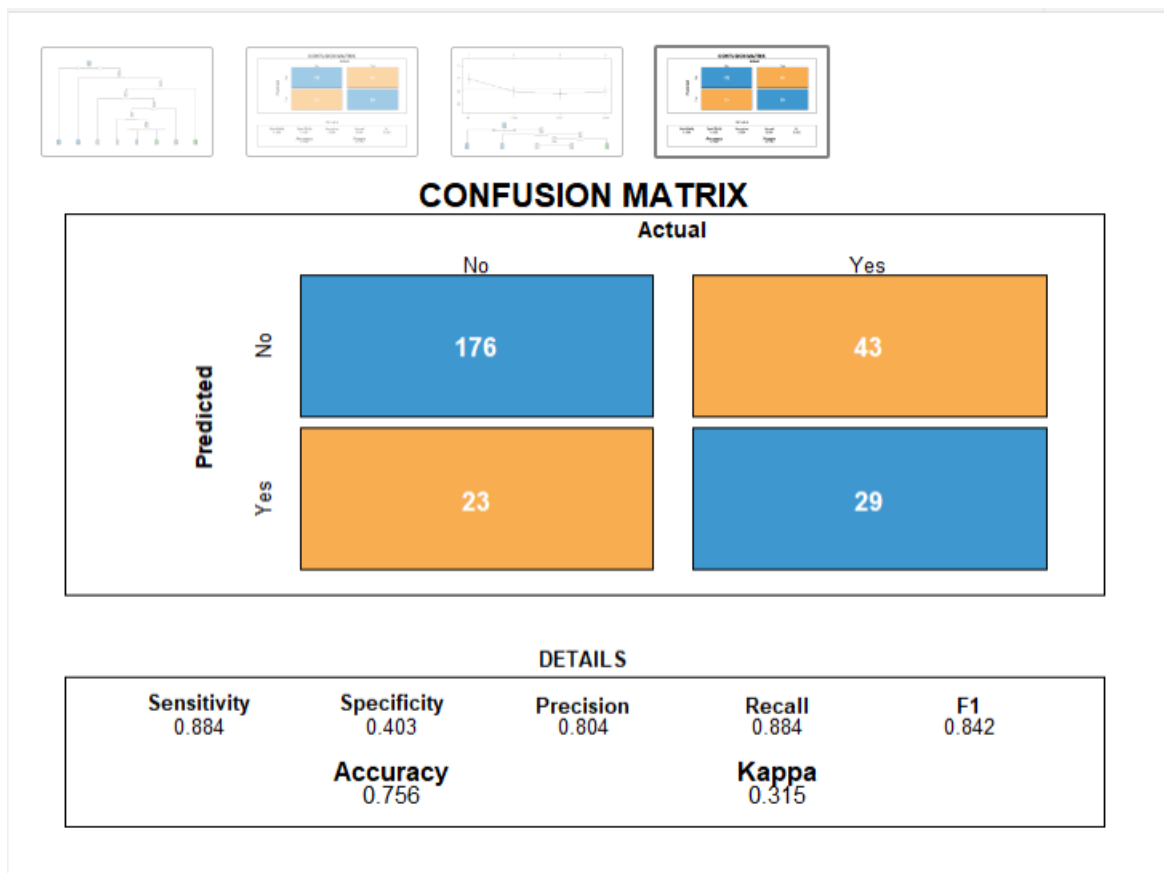
<b>Sensitivity</b> 0.884	<b>Specificity</b> 0.403	<b>Precision</b> 0.804	<b>Recall</b> 0.884	<b>F1</b> 0.842
<b>Accuracy</b> 0.756			<b>Kappa</b> 0.315	

**Slika 20:** Drvo odlučivanja koje kao mjeru atributa selekcije koristi gini indeks



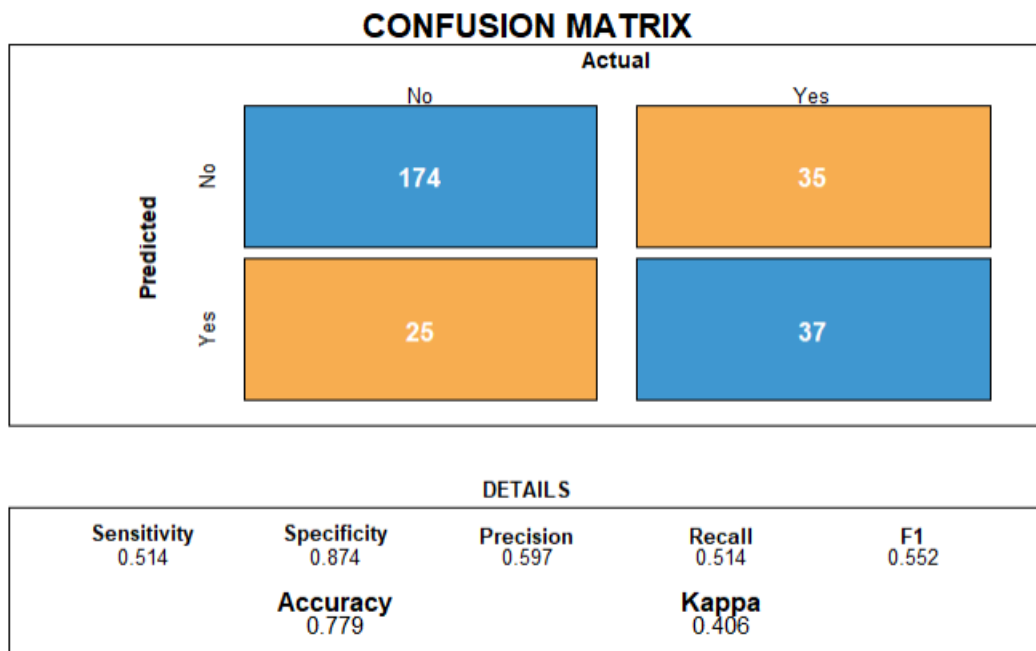
**Slika 21:** Roc kriva za drvo odlučivanja koje kao mjeru atributa selekcije koristi gini indeks

Čišćenje stabla kod modela koji koristi gini indeks nije poboljšalo rezultate.



**Slika 22:** Konfuzijska matrica za drvo odlučivanja koje kao mjeru atributa selekcije koristi gini indeks

### 2.1.3 C5.0 model klasifikacije

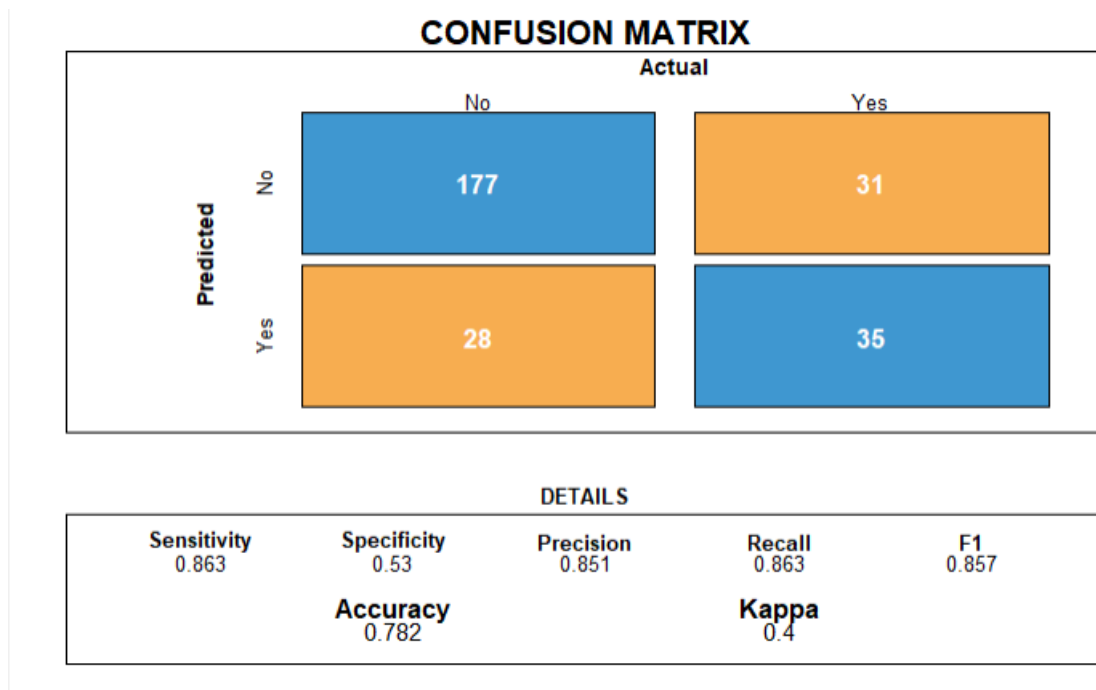


**Slika 23:** C5.0 model klasifikacije

## 2.2 Predikcijski modeli sa metodom holdouta

### 2.2.1 Drvo odlučivanja koje kao mjeru atributa selekcije koristi informacijsku dobit sa metodom holdouta

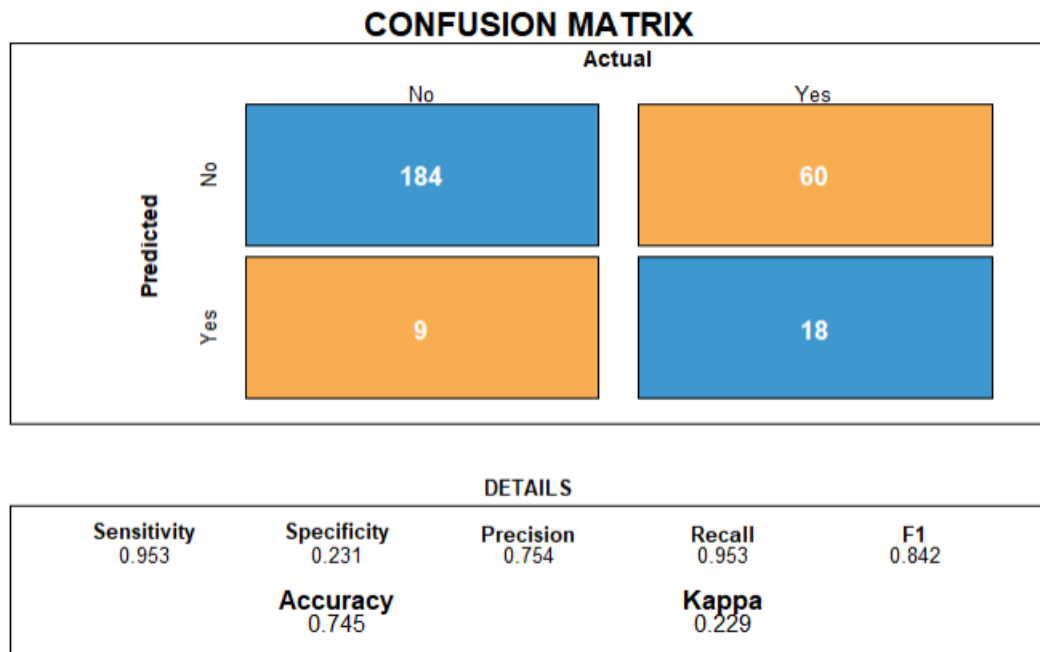
Metoda holdouta pokazuje bolju vrijednosti svih parametara osim parametra Sensivity.



**Slika 24:** Konfuzijska matrica za drvo odlučivanja koje kao mjeru atributa selekcije koristi informacijsku dobit sa metodom holdouta

### 2.2.2 Drvo odlučivanja koje kao mjeru atributa selekcije koristi gini indeks sa metodom holdouta

Metodom holdouta smo dobili lošiji Sensivity za model koji koristi gini indeks.



**Slika 25:** Drvo odlučivanja koje kao mjeru atributa selekcije koristi gini indeks sa metodom holdouta

### 2.2.3 C5.0 model klasifikacije sa metodom holdouta

Confusion Matrix and Statistics	
	Reference
Prediction	No Yes
No	484 126
Yes	50 62
Accuracy : 0.7562	
95% CI : (0.7232, 0.7871)	
No Information Rate : 0.7396	
P-Value [Acc > NIR] : 0.1648	
Kappa : 0.2717	
McNemar's Test P-Value : 1.574e-08	
Sensitivity : 0.32979	
Specificity : 0.90637	
Pos Pred Value : 0.55357	
Neg Pred Value : 0.79344	
Prevalence : 0.26039	
Detection Rate : 0.08587	
Detection Prevalence : 0.15512	
Balanced Accuracy : 0.61808	
'Positive' Class : Yes	

**Slika 26:** C5.0 model klasifikacije sa metodom holdouta



## 2.3 Predikcijski modeli sa metodom k-fold unakrsne validacije

### 2.3.1 Drvo odlučivanja koje kao mjeru atributa selekcije koristi informacijsku dobit sa metodom k-fold unakrsne validacije

Za k-fold unakrsnu validaciju smo napravili funkciju koju smo pozivali sa vrijednostima parametra  $k = 5$  i  $k = 10$ . Parametar  $k$  pokazuje koliko puta će se ponoviti klasifikacija. Za obje vrijednosti parametara dobili smo slične rezultate.

```
10-fold validacija
Najveća tačnost: 0.8166667 , fold: 2, najveća kappa: 0.441906 , fold: 6
Najmanja tačnost: 0.7333333 , fold: 5, najmanja kappa: 0.2270992 , fold: 9
Srednja tačnost: 0.7815838, srednja kappa: 0.3214565

5-fold validacija
Najveća tačnost: 0.8199446 , fold: 1, najveća kappa: 0.4414558 , fold: 1
Najmanja tačnost: 0.7479224 , fold: 2, najmanja kappa: 0.2521456 , fold: 4
Srednja tačnost: 0.781602, srednja kappa: 0.3324738
```

**Slika 27:** Konfuzijska matrica za drvo odlučivanja koje kao mjeru atributa selekcije koristi informacijsku dobit sa metodom k-fold unakrsne validacije

### 2.3.2 Drvo odlučivanja koje kao mjeru atributa selekcije koristi gini indeks sa metodom k-fold unakrsne validacije

I za model koji koristi gini indeks smo dobili slične vrijednosti srednje tačnosti za vrijednosti parametara  $k=5$  i  $k=10$ .

```
10-fold validacija
Najveća tačnost: 0.8287293 , fold: 4, najveća kappa: 0.4925387 , fold: 4
Najmanja tačnost: 0.7222222 , fold: 6, najmanja kappa: 0.1680532 , fold: 6
Srednja tačnost: 0.7815316, srednja kappa: 0.3342524

5-fold validacija
Najveća tačnost: 0.800554 , fold: 5, najveća kappa: 0.3701047 , fold: 5
Najmanja tačnost: 0.7534626 , fold: 2, najmanja kappa: 0.2905983 , fold: 3
Srednja tačnost: 0.7815897, srednja kappa: 0.333075
```

**Slika 28:** Drvo odlučivanja koje kao mjeru atributa selekcije koristi gini indeks sa metodom k-fold unakrsne validacije

### 2.3.3 C5.0 model klasifikacije sa metodom k-fold unakrsne validacije

Kod C5.0 modela smo dobili bolju srednju tačnost kappe za vrijednost parametra  $k = 10$ .

```
10-fold validacija
Najveća tačnost: 0.8277778 , fold: 3, najveća kappa: 0.5735249 , fold: 3
Najmanja tačnost: 0.7333333 , fold: 2, najmanja kappa: 0.2785571 , fold: 2
Srednja tačnost: 0.7788091, srednja kappa: 0.4160015

5-fold validacija
Najveća tačnost: 0.7950139 , fold: 1, najveća kappa: 0.4490596 , fold: 1
Najmanja tačnost: 0.7229917 , fold: 5, najmanja kappa: 0.185556 , fold: 2
Srednja tačnost: 0.7633102, srednja kappa: 0.3356188
```

**Slika 29:** C5.0 model klasifikacije sa metodom k-fold unakrsne validacije

## 2.4 Predikcijski modeli sa metodom k-fold bootstrapping validacije

### 2.4.1 Drvo odlučivanja koje kao mjeru atributa selekcije koristi informacijsku dobit sa metodom k-fold bootstrapping validacije

Kao i kod k-fold unakrsne validacije i kod k-fold bootstrapping validacije smo pozivali funkciju za različite vrijednosti parametra k.

```
10-fold bootstrap
Najveća tačnost: 0.8277778 , fold: 3, najveća kappa: 0.4149196 , fold: 8
Najmanja tačnost: 0.7444444 , fold: 4, najmanja kappa: 0.200237 , fold: 9
Srednja tačnost: 0.7872222, srednja kappa: 0.3165408

5-fold bootstrap
Najveća tačnost: 0.7833333 , fold: 2, najveća kappa: 0.3129771 , fold: 4
Najmanja tačnost: 0.7611111 , fold: 5, najmanja kappa: 0.2271978 , fold: 5
Srednja tačnost: 0.7705556, srednja kappa: 0.282632
```

**Slika 30:** Konfuzijska matrica za drvo odlučivanja koje kao mjeru atributa selekcije koristi informacijsku dobit sa metodom k-fold bootstrapping validacije

### 2.4.2 Drvo odlučivanja koje kao mjeru atributa selekcije koristi gini indeks sa metodom k-fold bootstrapping validacije

```
10-fold bootstrap
Najveća tačnost: 0.8055556 , fold: 1, najveća kappa: 0.3697479 , fold: 4
Najmanja tačnost: 0.7388889 , fold: 3, najmanja kappa: 0.2634762 , fold: 5
Srednja tačnost: 0.775, srednja kappa: 0.3176853

5-fold bootstrap
Najveća tačnost: 0.8055556 , fold: 3, najveća kappa: 0.3435791 , fold: 3
Najmanja tačnost: 0.7611111 , fold: 5, najmanja kappa: 0.2461832 , fold: 2
Srednja tačnost: 0.7805556, srednja kappa: 0.3059102
```

**Slika 31:** Drvo odlučivanja koje kao mjeru atributa selekcije koristi gini indeks sa metodom k-fold bootstrapping validacije

### 2.4.3 C5.0 model klasifikacije sa metodom k-fold bootstrapping validacije

```
10-fold bootstrap
Najveća tačnost: 0.8444444 , fold: 6, najveća kappa: 0.542068 , fold: 6
Najmanja tačnost: 0.7777778 , fold: 1, najmanja kappa: 0.3939394 , fold: 3
Srednja tačnost: 0.8005556, srednja kappa: 0.471508

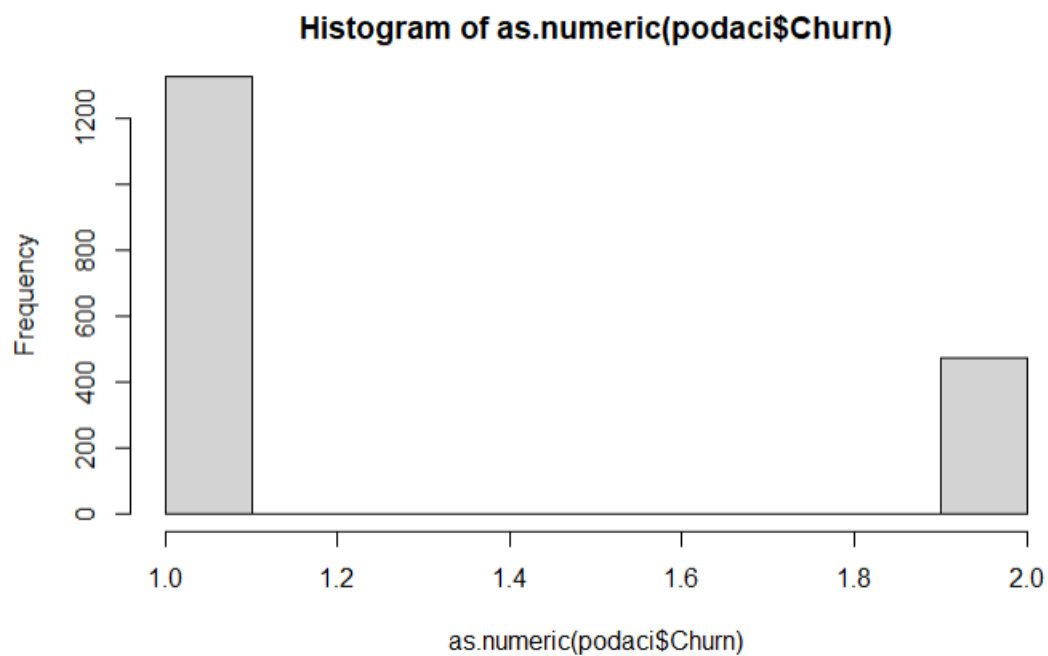
5-fold bootstrap
Najveća tačnost: 0.8416667 , fold: 4, najveća kappa: 0.5823836 , fold: 4
Najmanja tačnost: 0.7861111 , fold: 3, najmanja kappa: 0.4340547 , fold: 3
Srednja tačnost: 0.8105556, srednja kappa: 0.497575
```

---

**Slika 32:** C5.0 model klasifikacije sa metodom k-fold bootstrapping validacije

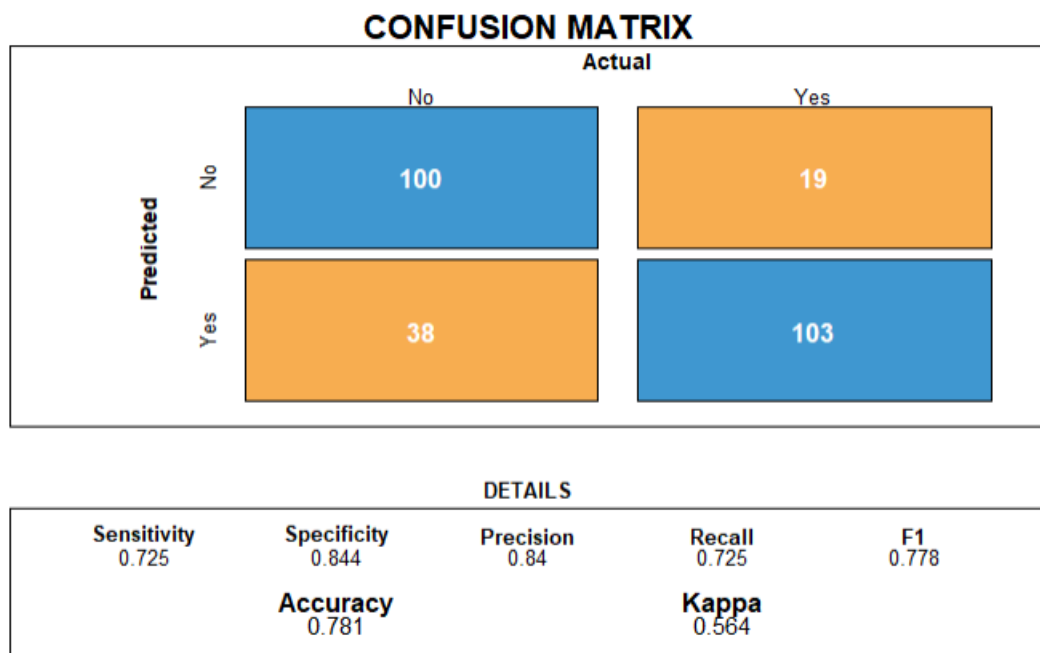
## 2.5 Balansiranje podataka

Na histogramu možemo uočiti da u datasetu postoji više instanci koje imaju No vrijednost za varijablu Churn. Stoga je potrebno izvršiti balansiranje dataseta. Za vrijednost parametra N smo uzeli 2600 kako bi se izbalansirao set.



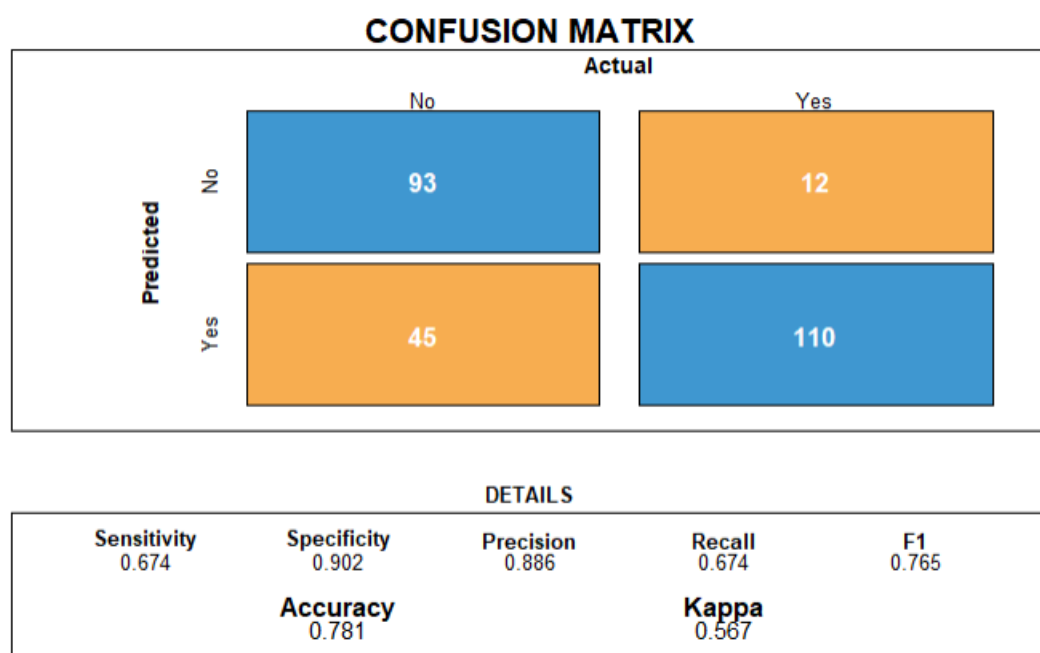
**Slika 33:** Histogram prije balansiranja

### 2.5.1 Drvo odlučivanja koje kao mjeru atributa selekcije koristi informacijsku dobit nakon balansiranja podataka



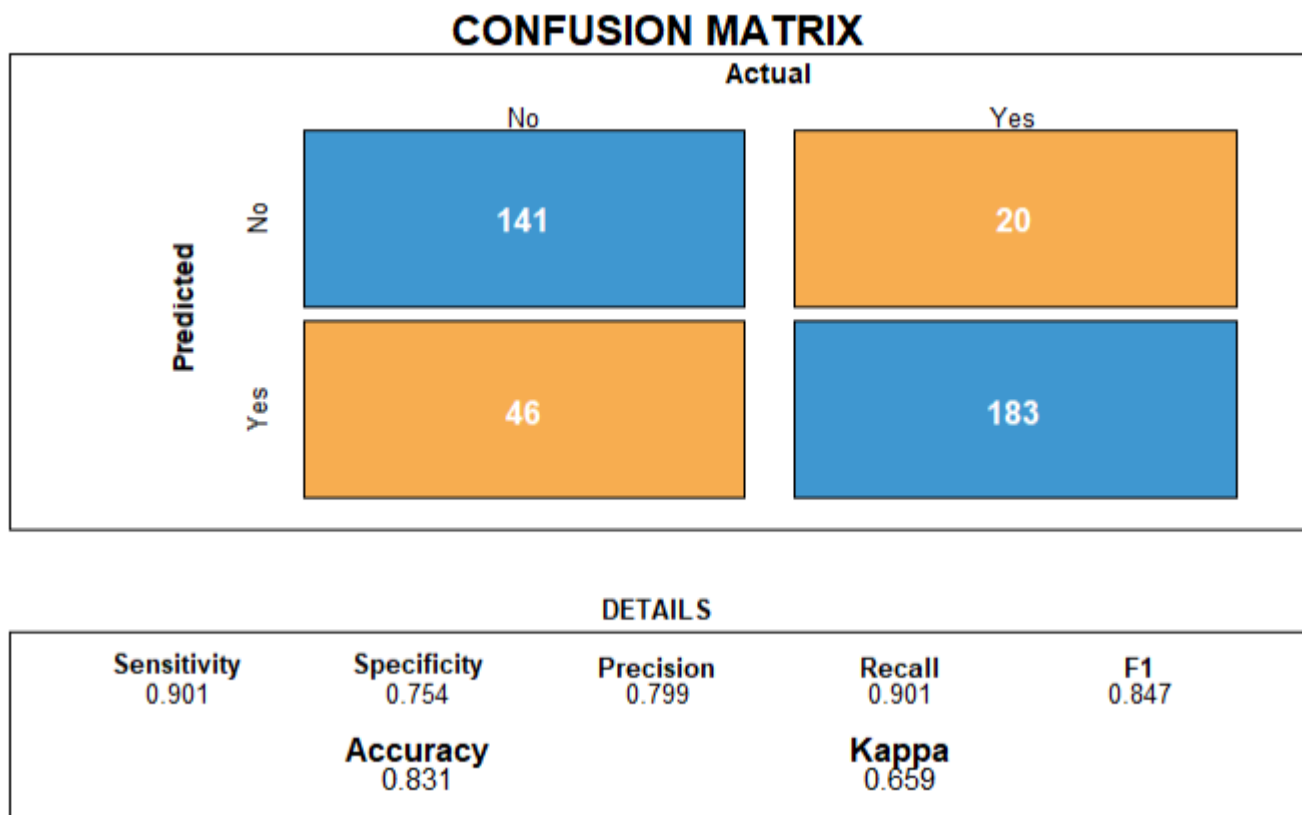
**Slika 34:** Drvo odlučivanja koje kao mjeru atributa selekcije koristi informacijsku dobit nakon balansiranja podataka

### 2.5.2 Drvo odlučivanja koje kao mjeru atributa selekcije koristi gini indeks nakon balansiranja podataka



**Slika 35:** Drvo odlučivanja koje kao mjeru atributa selekcije koristi gini indeks nakon balansiranja podataka

### 2.5.3 C5.0 model nakon balansiranja



Slika 36: C5.0 model nakon balansiranja

Za sve modele smo dobili bolje vrijednosti parametara nakon treniranja nad balansiranim podacima.

## 2.6 Ansambl tehnike za unaprjeđenje tačnosti klasifikacije

Za parametar nbagg smo uzeli vrijednost 50 što znači da će se kreirati 50 modela, a za parametar minsplite smo uzeli vrijednost 4 što označava minimalni broj instanci za granjanje.

### 2.6.1 Bagging model

```
Confusion Matrix and Statistics

 Reference
Prediction No Yes
No 127 3
Yes 11 119

 Accuracy : 0.9462
 95% CI : (0.9113, 0.9703)
No Information Rate : 0.5308
P-Value [Acc > NIR] : < 2e-16

 Kappa : 0.8923

McNemar's Test P-value : 0.06137

 Sensitivity : 0.9203
 Specificity : 0.9754
 Pos Pred Value : 0.9769
 Neg Pred Value : 0.9154
 Prevalence : 0.5308
 Detection Rate : 0.4885
 Detection Prevalence : 0.5000
 Balanced Accuracy : 0.9478

 'Positive' Class : No
```

Slika 37: Bagging model

### 2.6.2 Boosting model uz korišćenje AdaBoost metode

Na slici su prikazani parametri za boosting model koji koristi AdaBoost metodu. Vrijednost parametra mfinal je 50.

```
Confusion Matrix and Statistics

 Reference
Prediction No Yes
No 110 20
Yes 28 102

 Accuracy : 0.8154
 95% CI : (0.7628, 0.8606)
No Information Rate : 0.5308
P-value [Acc > NIR] : <2e-16

 Kappa : 0.6308

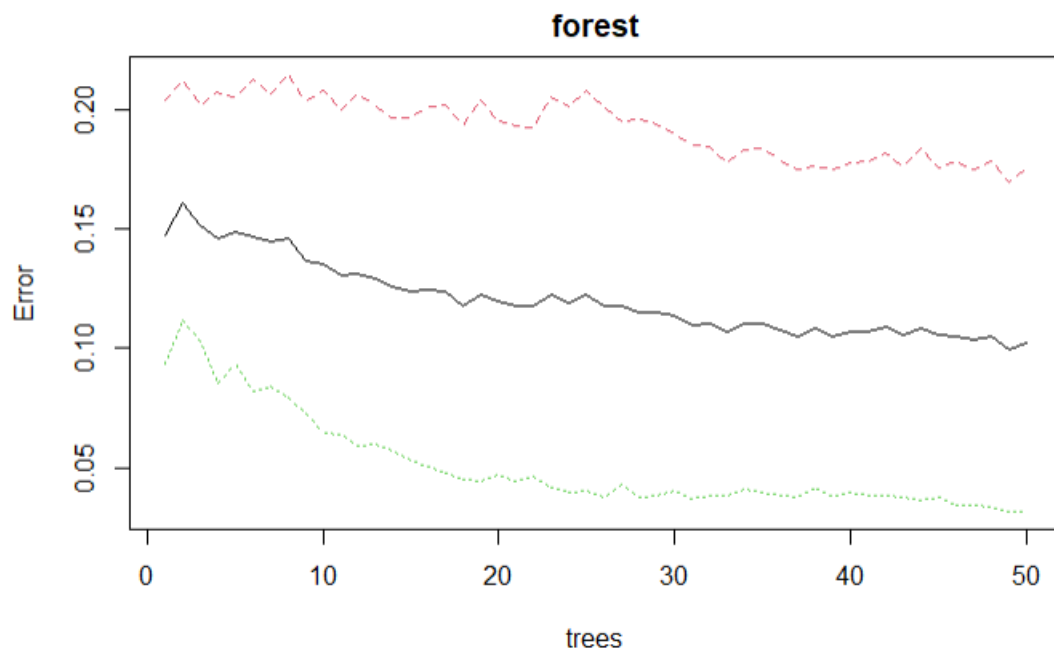
McNemar's Test P-value : 0.3123

 Sensitivity : 0.7971
 Specificity : 0.8361
 Pos Pred Value : 0.8462
 Neg Pred Value : 0.7846
 Prevalence : 0.5308
 Detection Rate : 0.4231
 Detection Prevalence : 0.5000
 Balanced Accuracy : 0.8166
```

Slika 38: Boosting model uz korišćenje AdaBoost metode

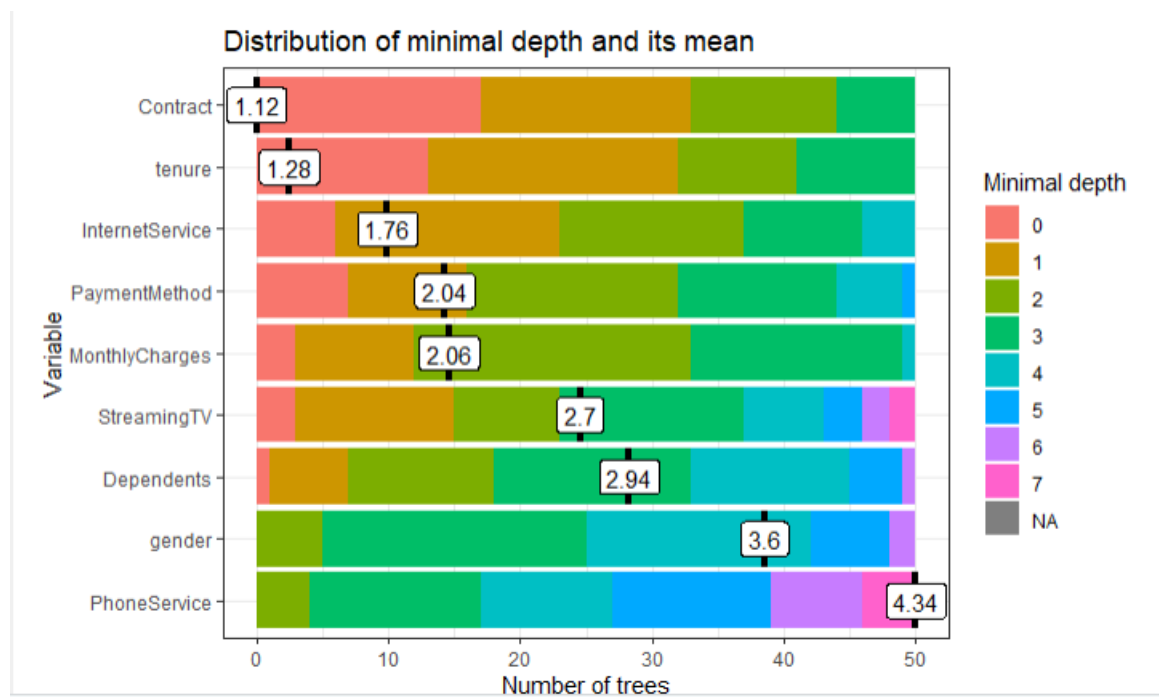
### 2.6.3 Random forest model

Za ovaj model smo kkoristi vrijednost parametra mtree = 50.



Slika 39: Prikaz ovisnosti greške klasifikacije o broju modela

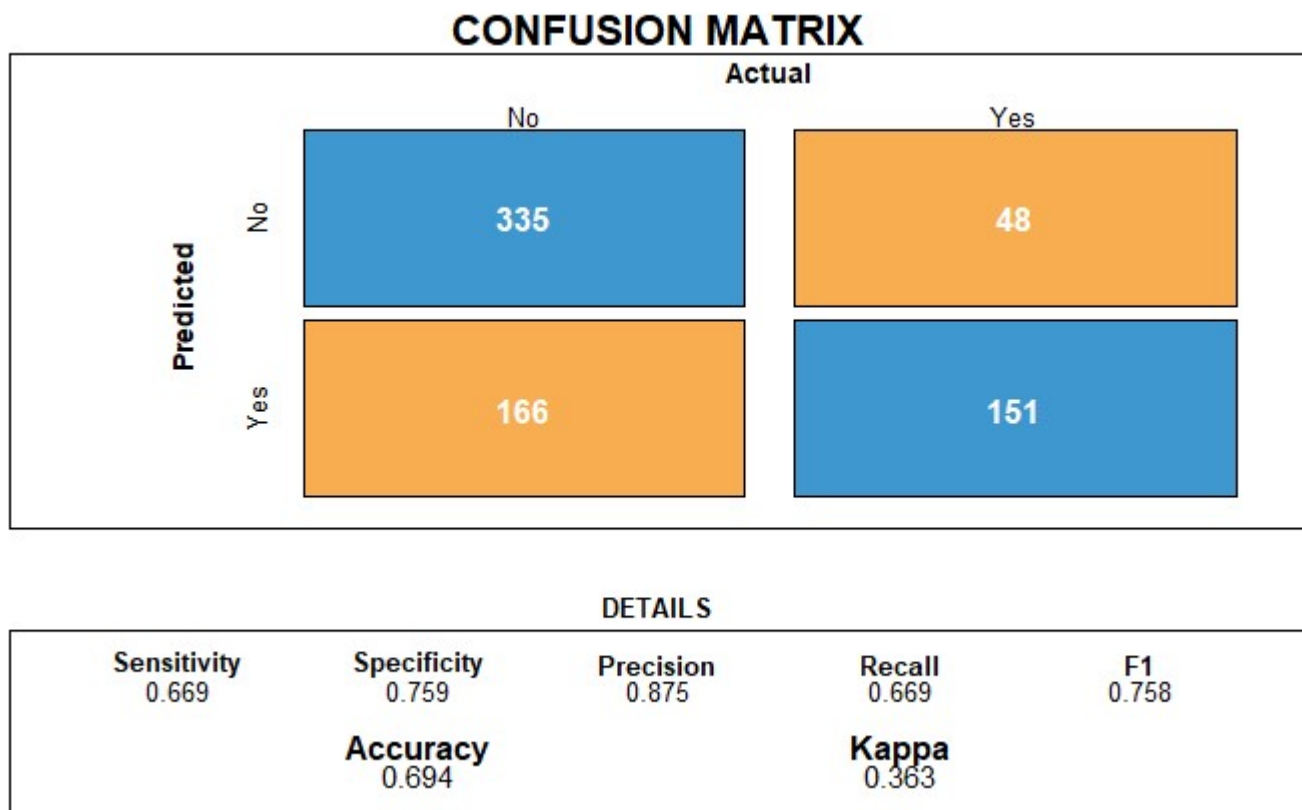
Na osnovu stepena značaja svih atributa možemo zaključiti da su Contract i tenure najvažniji atributi.



Slika 40: Prikaz stepena značaja svih atributa u random forest modelu

### 3 Zadatak 3

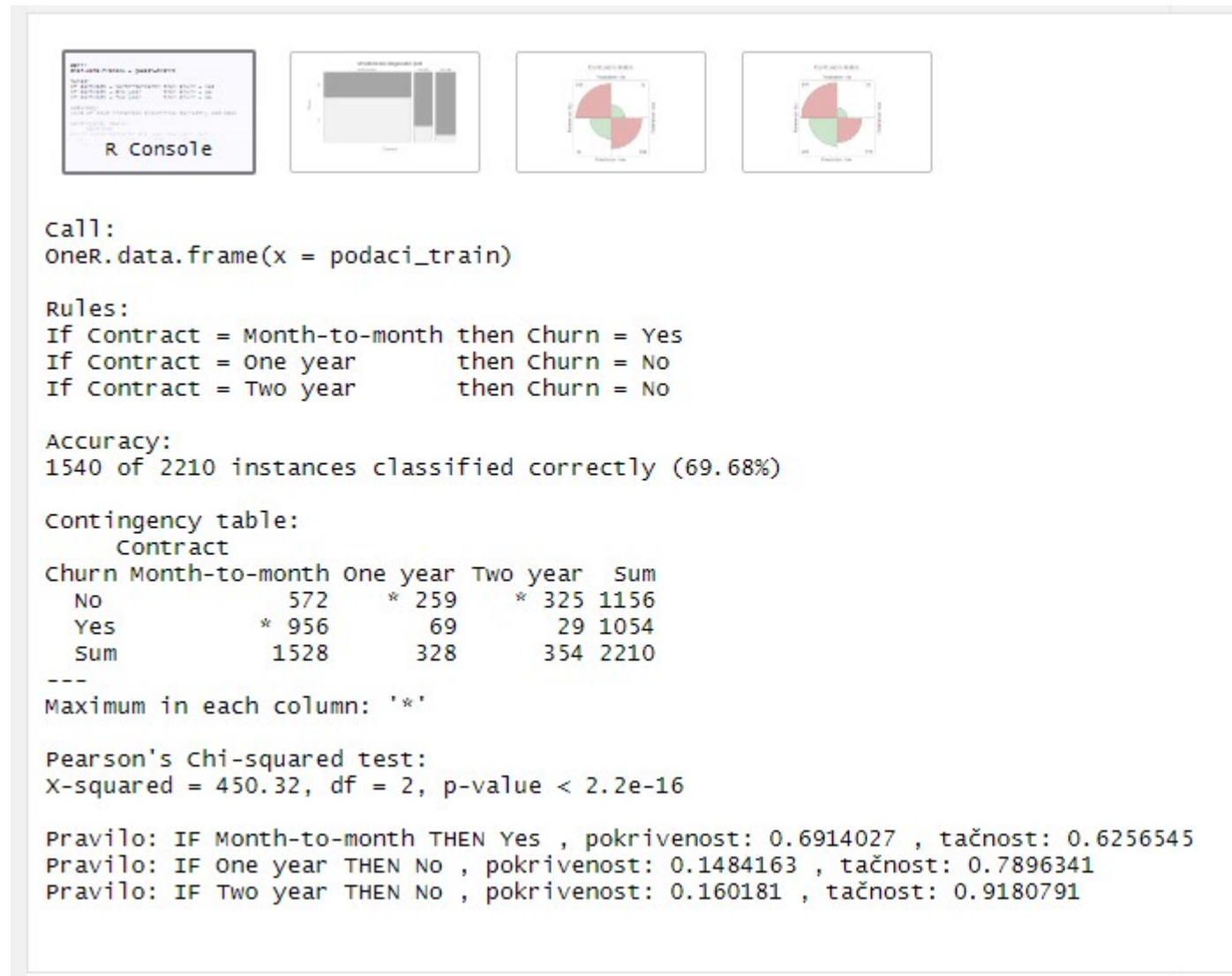
Najbolje rezultate u drugom zadatku nam je dao bagging model. Međutim nad testnim skupom bolje rezultate nam je dao C5.0 pa smo prikazali njegove rezultate.



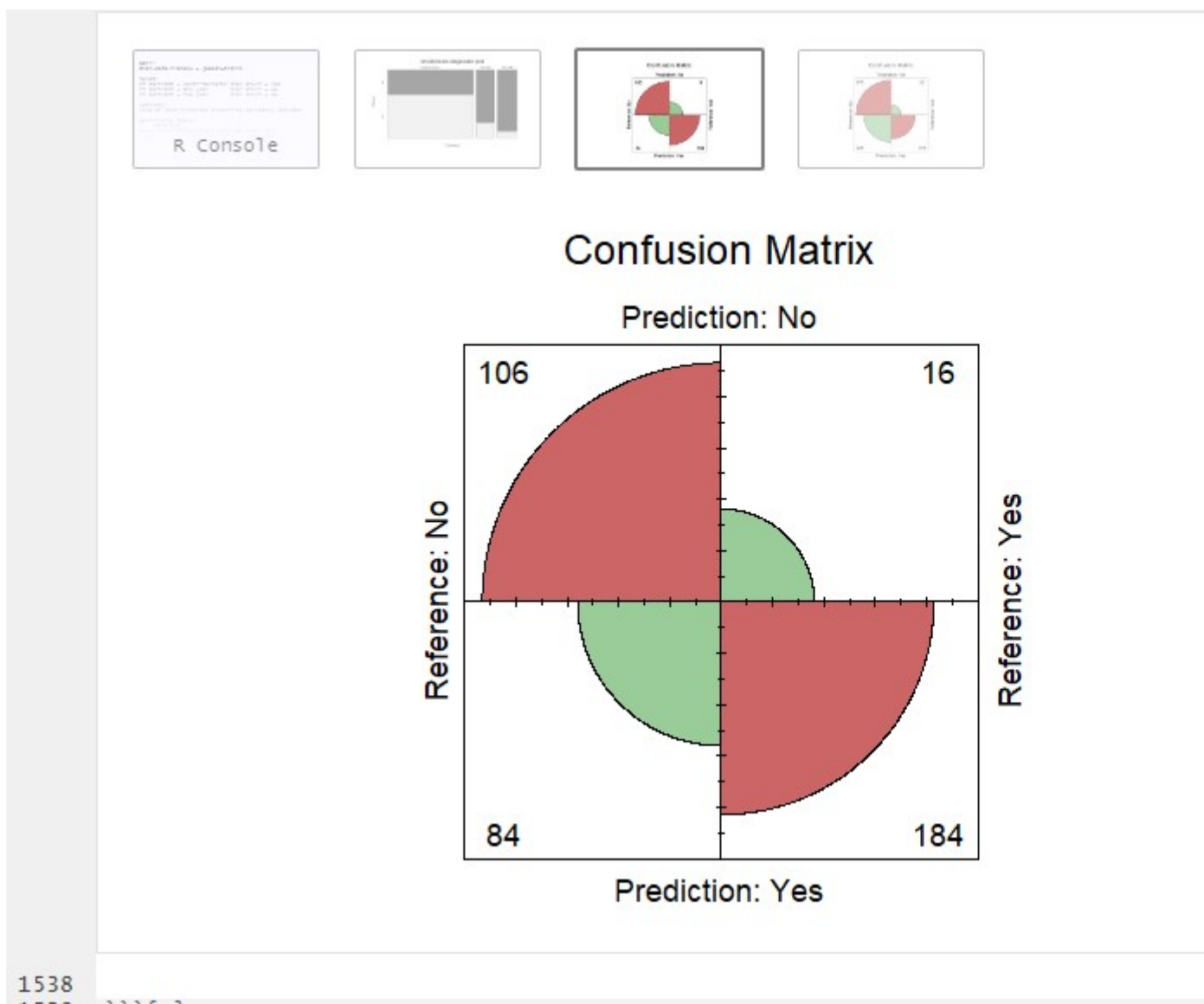
**Slika 41:** Primjena C5.0 modela nad testnim skupom



## 4 Zadatak 4



Slika 42: Model OneR klasifikatora



**Slika 43:** Prikaz konfuzijske matrice sa proporcijama tačnih i netačnih klasifikacija

