

# Statistička analiza podataka - projekt

Sapunanje

2023-01-04

## Statistika nogometaša engleske Premier lige

Studenti: Karlo Boroš, Petar Novak, Vlado Perković i Mislav Rendulić

### 1. Uvod

Ovaj projekt iz kolegija Statistička analiza podataka radili smo pod vodstvom asistenta *insert\_name*. Ovaj seminar ćemo podijeliti u par dijelova:

1. Uvod
2. Osnovna prilagodba podataka
3. Generalne informacije o podacima
4. XY-test
5. YZ-test
6. ZX-test
7. XZ-test
8. Rezultati

Imali smo sreće i dobili smo upravo zadatak koji smo i priželjkivali.

**Cilj** ovoga projekta je uzeti dane podatke i iz njih probati izvući zaključke i faktore koji mogu utjecati na rezultat, broj golova i sl. Naravno, nije potrebno naglasiti važnost korištenja ispravnih testova te dobivanje rezultata koji su validni.

### 2. Osnovna prilagodba podataka

Podatke je prvo potrebno učitati. Bitno je dobro ih proučiti kako ne bismo slučajno pogriješili u nekom zaključku. Nakon dobre analize možemo krenuti sa našim zadacima.

**izbrisati kasnije** *napomena ostatku ekipe: spremio sam podatke kao dataset.csv jer je ime dugo i ne očitava š pa je ovo najjednostavnije*

#poslije dodati include = FALSE, za sad nek ostane da vidimo sve sta treba

```
nogometasi <- read.csv('dataset.csv', encoding = "UTF-8")
nogometasi$Nation <- str_sub(nogometasi$Nation, -3)
head(nogometasi)
```

##		Player	Team	Nation	Pos	Age	MP	Starts	Min	X90s	Gls	Ast	G.PK	
## 1		Bukayo Saka	Arsenal	ENG	FW,MF	19	38	36	2,978	33.1	11	7	9	
## 2		Gabriel Dos Santos	Arsenal	BRA	DF	23	35	35	3,063	34.0	5	0	5	
## 3		Aaron Ramsdale	Arsenal	ENG	GK	23	34	34	3,060	34.0	0	0	0	
## 4		Ben White	Arsenal	ENG	DF	23	32	32	2,880	32.0	0	0	0	
## 5		Martin Ødegaard	Arsenal	NOR	MF	22	36	32	2,785	30.9	7	4	7	
## 6		Granit Khaka	Arsenal	SUI	MF,DF	28	27	27	2,327	25.9	1	2	1	
##		PK	PKatt	CrdrY	CrdrR	Gls.1	Ast.1	G.A	G.PK.1	G.A.PK	xG	npG	xG.1	
## 1	2	2	6	0	0.33	0.21	0.54	0.27	0.48	9.7	8.2	6.9	15.2	0.29

```
## 2 0 0 8 1 0.15 0.00 0.15 0.15 0.15 2.7 2.7 0.8 3.5 0.08
## 3 0 0 1 0 0.00 0.00 0.00 0.00 0.00 0.0 0.0 0.0 0.0 0.00
## 4 0 0 3 0 0.00 0.00 0.00 0.00 0.00 1.0 1.0 0.6 1.6 0.03
## 5 0 0 4 0 0.23 0.13 0.36 0.23 0.36 4.8 4.8 6.8 11.6 0.16
## 6 0 0 10 1 0.04 0.08 0.12 0.04 0.12 1.2 1.2 2.3 3.5 0.05
## xA.1 xG.xA npxG.1 npxG.xA.1
## 1 0.21 0.50 0.25 0.46
## 2 0.02 0.10 0.08 0.10
## 3 0.00 0.00 0.00 0.00
## 4 0.02 0.05 0.03 0.05
## 5 0.22 0.38 0.16 0.38
## 6 0.09 0.14 0.05 0.14
```

```
str(nogometasi)
```

```
## 'data.frame': 691 obs. of 30 variables:
## $ Player : chr "Bukayo Saka" "Gabriel Dos Santos" "Aaron Ramsdale" "Ben White" ...
## $ Team : chr "Arsenal" "Arsenal" "Arsenal" "Arsenal" ...
## $ Nation : chr "ENG" "BRA" "ENG" "ENG" ...
## $ Pos : chr "FW,MF" "DF" "GK" "DF" ...
## $ Age : int 19 23 23 22 28 28 24 21 20 ...
## $ MP : int 38 35 34 32 36 27 24 22 33 29 ...
## $ Starts : int 36 35 34 32 32 27 23 22 21 21 ...
## $ Min : chr "2,978" "3,063" "3,060" "2,880" ...
## $ X90s : num 33.1 34 34 32 30.9 25.9 22.5 21.3 21.3 20.7 ...
## $ GlS : int 11 5 0 0 7 1 2 1 10 6 ...
## $ Ast : int 7 0 0 0 4 2 1 3 2 6 ...
## $ G.PK : int 9 5 0 0 7 1 2 1 10 5 ...
## $ PK : int 2 0 0 0 0 0 0 0 0 1 ...
## $ PKatt : int 2 0 0 0 0 0 0 0 0 1 ...
## $ CrdY : int 6 8 1 3 4 10 6 0 1 3 ...
## $ CrdR : int 0 1 0 0 0 1 0 0 0 1 ...
## $ GlS.1 : num 0.33 0.15 0 0 0.23 0.04 0.09 0.05 0.47 0.29 ...
## $ Ast.1 : num 0.21 0 0 0 0.13 0.08 0.04 0.14 0.09 0.29 ...
## $ G.A : num 0.54 0.15 0 0 0.36 0.12 0.13 0.19 0.56 0.58 ...
## $ G.PK.1 : num 0.27 0.15 0 0 0.23 0.04 0.09 0.05 0.47 0.24 ...
## $ G.A.PK : num 0.48 0.15 0 0 0.36 0.12 0.13 0.19 0.56 0.53 ...
## $ xG : num 9.7 2.7 0 1 4.8 1.2 2.5 0.7 5.8 7.2 ...
## $ npxG : num 8.2 2.7 0 1 4.8 1.2 2.5 0.7 5.8 6.5 ...
## $ xA : num 6.9 0.8 0 0.6 6.8 2.3 1.3 1.9 2.2 3.3 ...
## $ npxG.xA : num 15.2 3.5 0 1.6 11.6 3.5 3.8 2.6 8 9.8 ...
## $ xG.1 : num 0.29 0.08 0 0.03 0.16 0.05 0.11 0.03 0.27 0.35 ...
## $ xA.1 : num 0.21 0.02 0 0.02 0.22 0.09 0.06 0.09 0.1 0.16 ...
## $ xG.xA : num 0.5 0.1 0 0.05 0.38 0.14 0.17 0.12 0.37 0.51 ...
## $ npxG.1 : num 0.25 0.08 0 0.03 0.16 0.05 0.11 0.03 0.27 0.31 ...
## $ npxG.xA.1 : num 0.46 0.1 0 0.05 0.38 0.14 0.17 0.12 0.37 0.47 ...
```

Nakon enkodiranja početnih podataka, postajala su odstupanja od stvarnih imena kod nekih igrača pa smo ta imena ručno ispravili.

Konverzija odigranih minuta iz char u numeric:

```
nogometasi$Min <- as.numeric(gsub(",", "", nogometasi$Min))
```

### 3. Pregled sezone

Ekipe koje su se natjecale u Premier Ligi u sezoni 2021/2022

```
#as_tibble(klubovi)
klubovi
```

```
## [1] "Arsenal"           "Aston Villa"
## [3] "Brentford"         "Brighton & Hove Albion"
## [5] "Burnley"           "Chelsea"
## [7] "Crystal Palace"    "Everton"
## [9] "Leeds United"      "Leicester City"
## [11] "Liverpool"         "Manchester City"
## [13] "Manchester United" "Newcastle United"
## [15] "Norwich City"      "Southampton"
## [17] "Tottenham Hotspur" "Watford"
## [19] "West Ham United"   "Wolverhampton Wanderers"
```

```
#cat('Preko rownames mozda nastimati da ne pise "value" + odluciti jel bi ispisali preko tibble-a ili o
```

```
najbolji_strijelci
```

#### Najbolji strijelci

##	Player	Team	Gls	Gls per 90 min
## 1	Mohamed Salah	Liverpool	23	0.75
## 2	Son Heung-min	Tottenham Hotspur	23	0.69
## 3	Cristiano Ronaldo	Manchester United	18	0.66
## 4	Harry Kane	Tottenham Hotspur	17	0.47
## 5	Sadio Mané	Liverpool	16	0.51

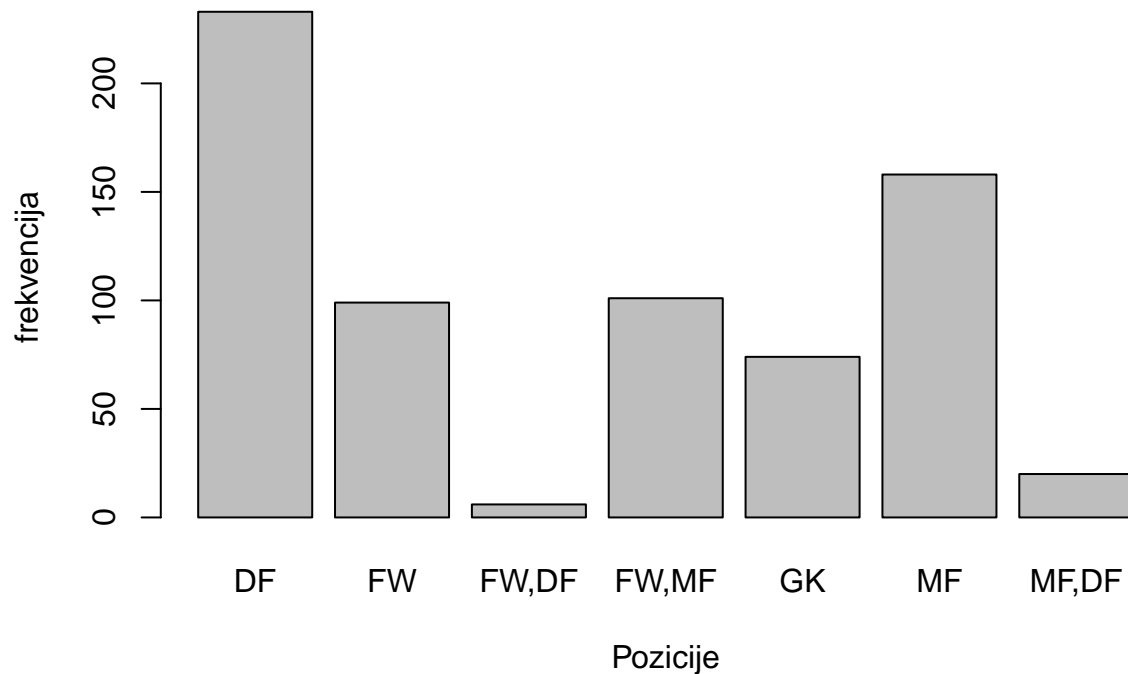
```
najbolji_asistenti
```

#### Najbolji asistenti

##	Player	Team	Ast	Ast per 90 min
## 1	Mohamed Salah	Liverpool	13	0.42
## 2	Trent Alexander-Arnold	Liverpool	12	0.38
## 3	Mason Mount	Chelsea	10	0.38
## 4	Harvey Barnes	Leicester City	10	0.43
## 5	Andrew Robertson	Liverpool	10	0.35
## 6	Jarrod Bowen	West Ham United	10	0.30

**Pozicije igrača** Vizualizacija razdiobe igrača po pozicijama:

```
nogometasi %>% select(Pos) %>% summarise(uniPos = ifelse(Pos == "DF,FW", "FW,DF", ifelse(Pos == "MF,FW"
#c("GK", "DF", "MF,DF", "MF", "FW,MF", "FW,DF", "FW"))
barplot(table(popravak), xlab = "Pozicije", ylab = "frekvencija")
```



Primijetimo veliki broj obrambenih igrača što i ima smisla kada pogledamo da ekipe najčešće igraju s 4 igrača u obrani. Neki igrači su igrali pozicije beka i napadačkog krila pa spadaju u skupinu “FW,DF” koja je na prvi pogled dosta neuobičajena.

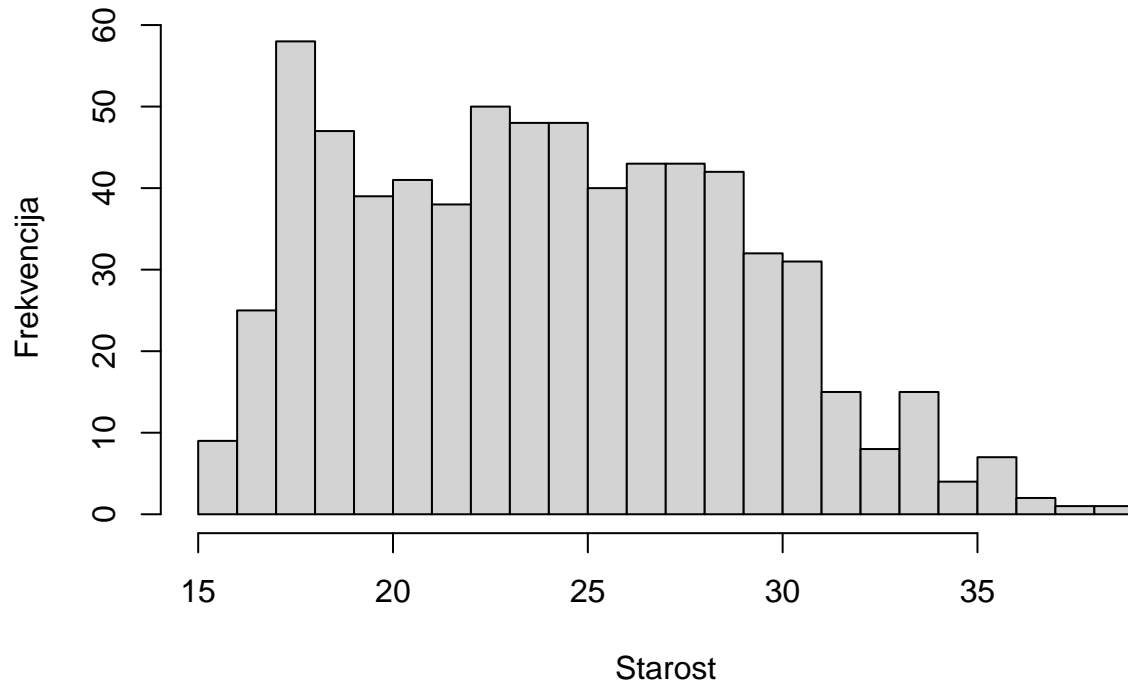
```
distrStarosti <- hist(nogometasi$Age,
  breaks = 20,
  main="Razdioba starosti igrača",
  xlab="Starost",
  ylab='Frekvencija'
)
```

#### Godine igrača

```
## Warning in title(main = main, sub = sub, xlab = xlab, ylab = ylab, ...):
## conversion failure on 'Razdioba starosti igrača' in 'mbcsToSbcs': dot
## substituted for <c4>

## Warning in title(main = main, sub = sub, xlab = xlab, ylab = ylab, ...):
## conversion failure on 'Razdioba starosti igrača' in 'mbcsToSbcs': dot
## substituted for <8d>
```

## Razdioba starosti igra..a

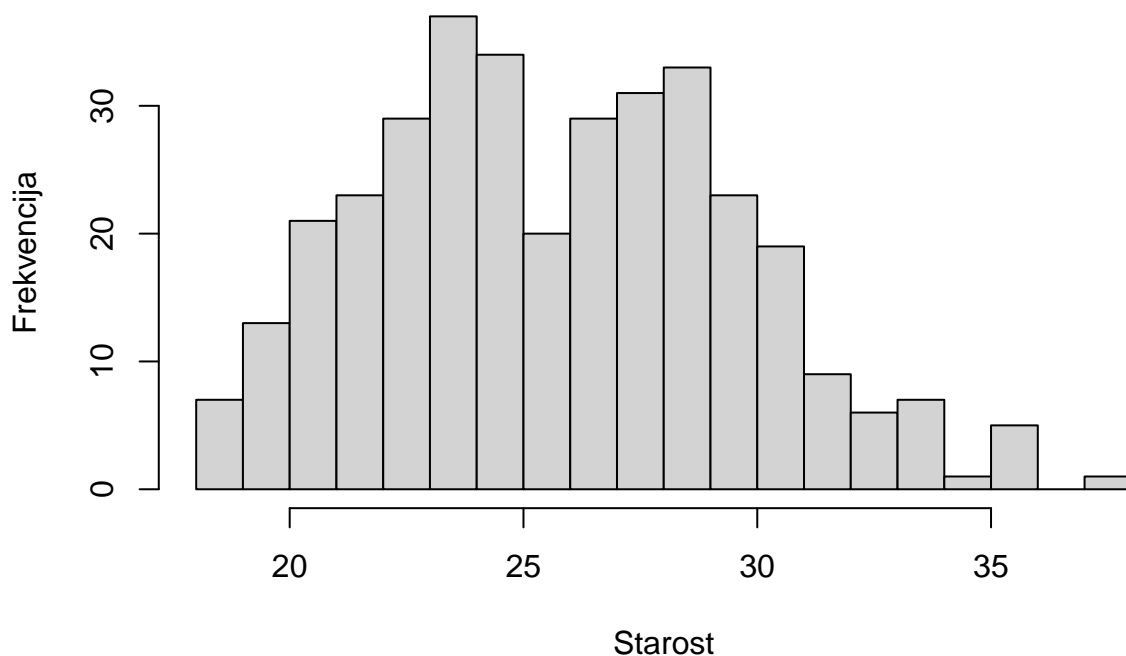


```
x <- nogometasi %>% filter(!is.na(X90s) & X90s >= 9.5)
distrStarosti <- hist(x$Age,
  breaks = 20,
  main="Starost igrača sa 25%+ minutaze",
  xlab="Starost",
  ylab='Frekvencija'
)
```

```
## Warning in title(main = main, sub = sub, xlab = xlab, ylab = ylab, ...):
## conversion failure on 'Starost igrača sa 25%+ minutaze' in 'mbcsToSbcs': dot
## substituted for <c4>

## Warning in title(main = main, sub = sub, xlab = xlab, ylab = ylab, ...):
## conversion failure on 'Starost igrača sa 25%+ minutaze' in 'mbcsToSbcs': dot
## substituted for <8d>
```

## Starost igra..a sa 25%+ minutaze



4.

Postoji li razlika u broju odigranih minuta mladih igrača (do 25 godina) među premierligaškim ekipama? Podijelimo igrače...

```
mladi <- nogometasi %>% filter(Age <= 25)
cat("Broj mladih igrača do 25 godina iznosi: ", nrow(mladi), "\n")
```

```
## Broj mladih igrača do 25 godina iznosi: 403
```

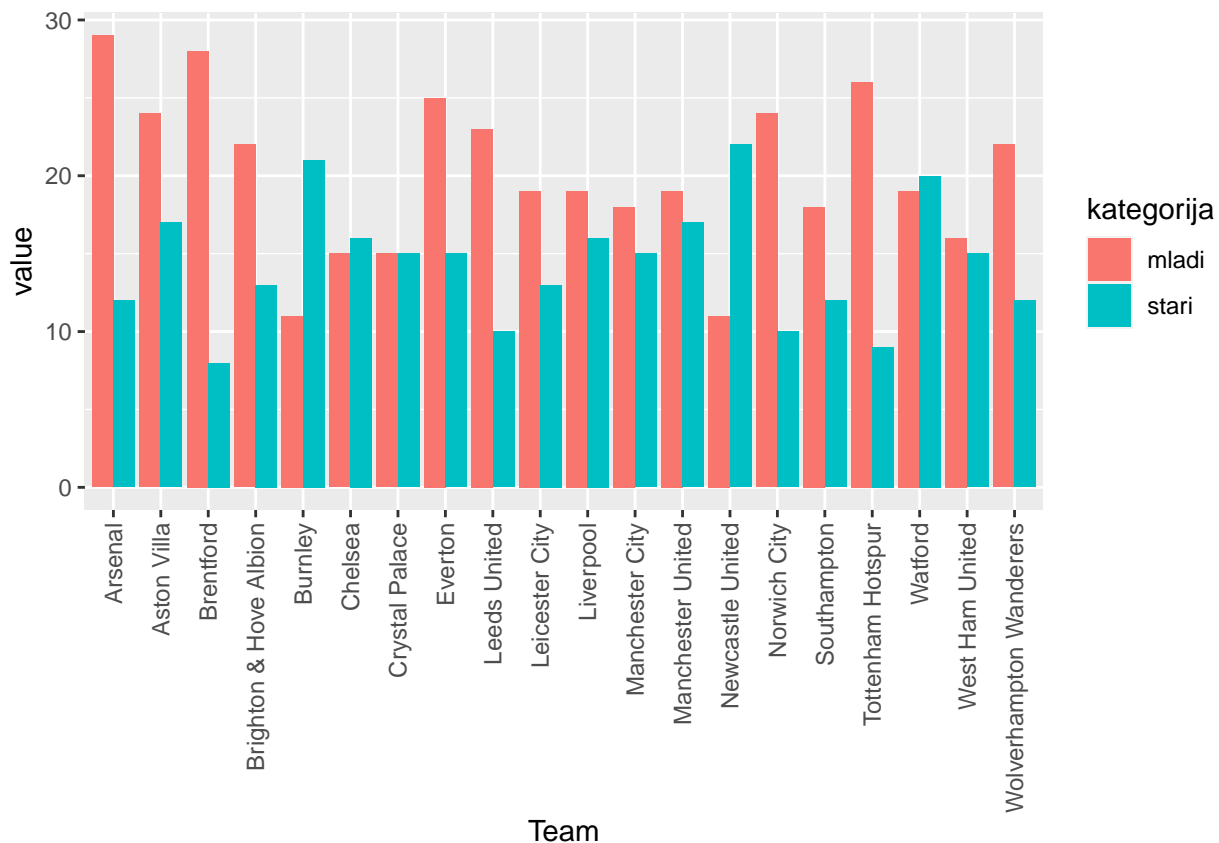
```
stari <- nogometasi %>% filter(Age > 25)
cat("Broj igrača iznad 25 godina iznosi: ", nrow(stari))
```

```
## Broj igrača iznad 25 godina iznosi: 284
```

Vizualizirajmo podjelu igrača u samim klubovima:

```
nogometasi_god <- nogometasi %>% summarise(mladi = ifelse(Age <= 25, 1, 0), Team) %>% group_by(Team) %>%
  summarise(mladi = sum(mladi))

ggplot(nogometasi_god, aes(x=Team, y=value, fill=kategorija)) +
  geom_bar(stat="identity", position="dodge") +
  theme(axis.text.x = element_text(angle = 90, hjust = 1, vjust = 0.5))
```



```
#scale_fill_brewer(palette = "Set1")
```

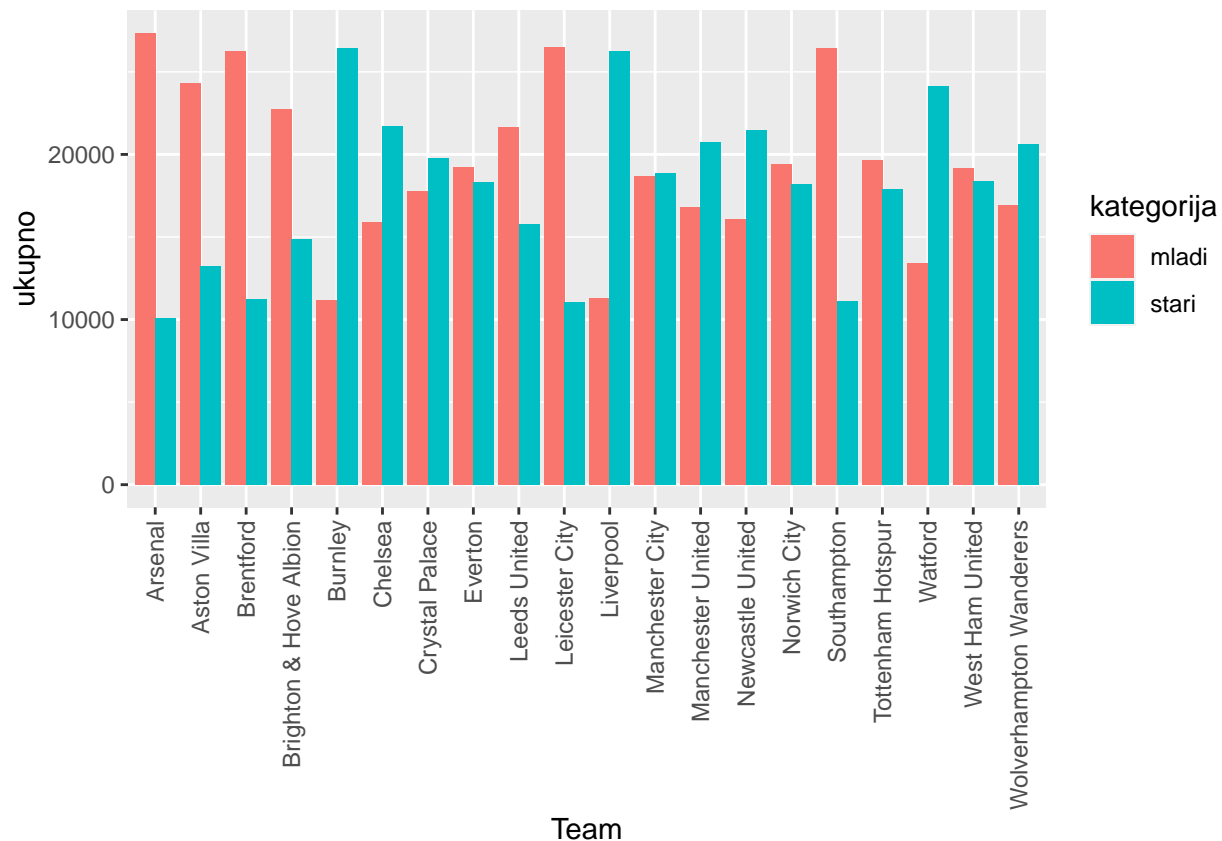
Vidimo da klubovi pretežno nastoje priključivati mlađe igrače u ekipu uz iznimke timova Burnley i Newcastle United.

Pogledajmo sada koliko te iste mlade igrače timovi zapravo i koriste...

```
nogometasi_min <- nogometasi %>% filter(!is.na(Age)) %>% summarise(Team, kategorija = ifelse(Age <= 26,
```

```
## `summarise()` has grouped output by 'Team', 'kategorija'. You can override
## using the `.groups` argument.
```

```
ggplot(nogometasi_min, aes(x=Team, y=ukupno, fill=kategorija)) +
  geom_bar(stat="identity", position="dodge") +
  theme(axis.text.x = element_text(angle = 90, hjust = 1, vjust = 0.5))
```



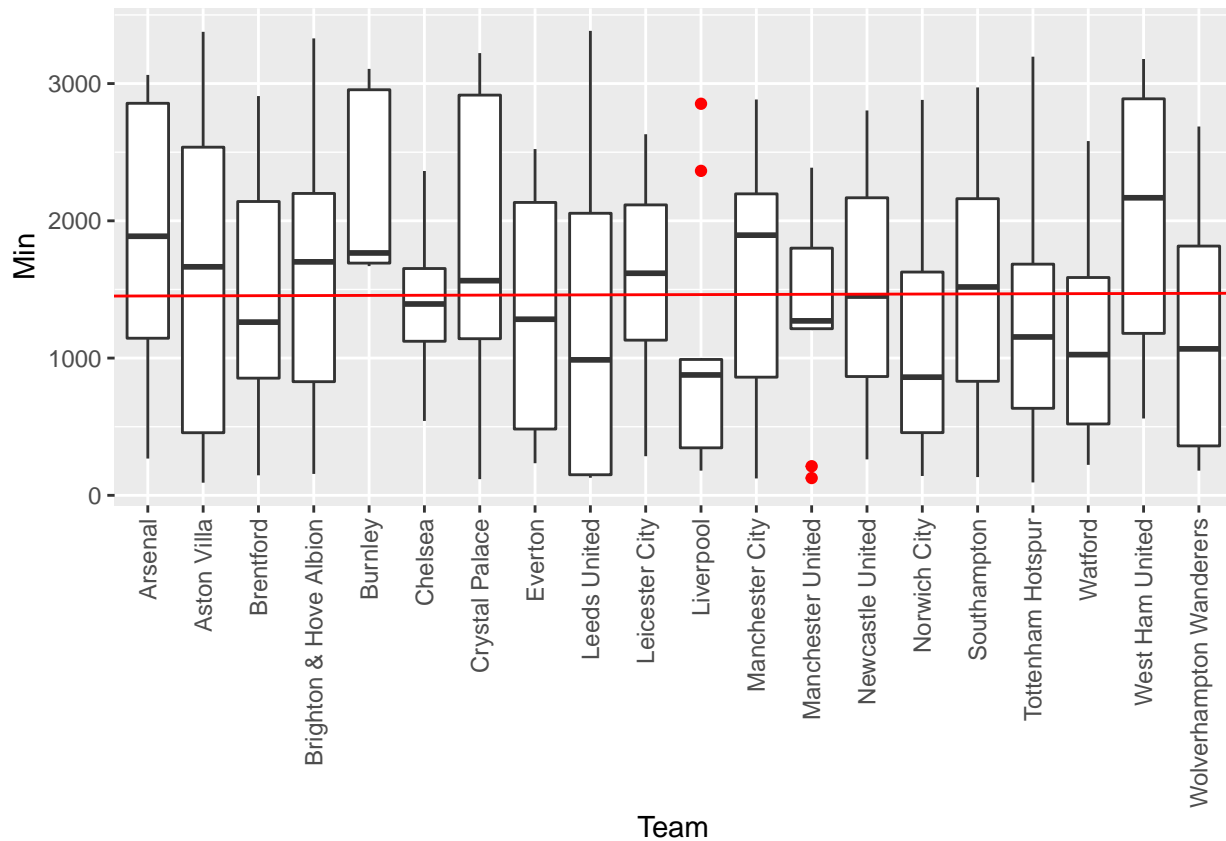
Kod analize u obzir ćemo uzeti mlade igrače koji su upisali barem 90 minuta.

```
mladi90 <- nogometasi %>% filter(Age <= 25 & Min >= 90)
```

Pogledajmo koliko su u prosjeku klubovi davali minuta svojim mladim igračima

```
ggplot(mladi90, aes(x = Team, y = Min)) +
  geom_boxplot(outlier.color = "red") +
  geom_abline(intercept = mean(mladi90$Min), col = "Red") +
  theme(axis.text.x = element_text(angle = 90, hjust = 1, vjust = 0.5))
```





Testirat ćemo homogenost varijance raspodjele minuta mladih igrača po klubovima:

```
bartlett.test(mladi90$Min ~ mladi90$Team)
```

```
##
## Bartlett test of homogeneity of variances
##
## data: mladi90$Min by mladi90$Team
## Bartlett's K-squared = 12.618, df = 19, p-value = 0.8575
```

Sada je potrebno testirati normalnost distribucije odigranih minuta za igrače do 25 godina ukupno i po klubovima:

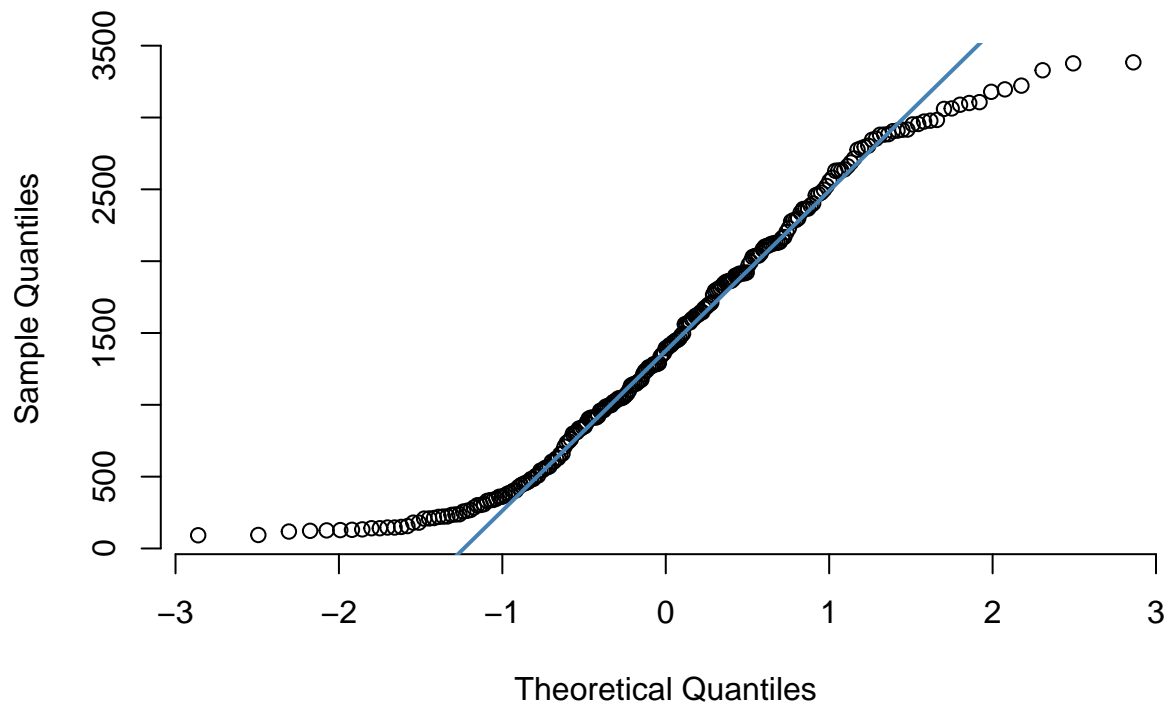
```
qqnorm(mladi90$Min, pch = 1, frame = FALSE, main='Odigrane minute za igrače do 25 godina')
```

```
## Warning in title(...): conversion failure on 'Odigrane minute za igrače do 25
## godina' in 'mbcsToSbcs': dot substituted for <c4>

## Warning in title(...): conversion failure on 'Odigrane minute za igrače do 25
## godina' in 'mbcsToSbcs': dot substituted for <8d>
```

```
qqline(mladi90$Min, col = "steelblue", lwd = 2)
```

## Odigrane minute za igra..e do 25 godina



```
require(nortest)
```

```
## Loading required package: nortest
```

```
lillie.test(mladi90$Min[mladi90$Team == "Arsenal"])
```

```
##
```

```
## Lilliefors (Kolmogorov-Smirnov) normality test
```

```
##
```

```
## data: mladi90$Min[mladi90$Team == "Arsenal"]
```

```
## D = 0.18585, p-value = 0.2112
```

```
lillie.test(mladi90$Min[mladi90$Team == "Aston Villa"])
```

```
##
```

```
## Lilliefors (Kolmogorov-Smirnov) normality test
```

```
##
```

```
## data: mladi90$Min[mladi90$Team == "Aston Villa"]
```

```
## D = 0.18316, p-value = 0.3232
```

```
lillie.test(mladi90$Min[mladi90$Team == "Brentford"])
```

```
##
```

```
## Lilliefors (Kolmogorov-Smirnov) normality test
```

```
##
```

```
## data: mladi90$Min[mladi90$Team == "Brentford"]
```

```
## D = 0.12846, p-value = 0.6017
```

```
lillie.test(mladi90$Min[mladi90$Team == "Brighton & Hove Albion"])
```

```
##
```

```
## Lilliefors (Kolmogorov-Smirnov) normality test
```

```
##
```

```
## data:  mladi90$Min[mladi90$Team == "Brighton & Hove Albion"]
## D = 0.11693, p-value = 0.9435
```

```
lillie.test(mladi90$Min[mladi90$Team == "Burnley"])
```

```
##
##  Lilliefors (Kolmogorov-Smirnov) normality test
##
## data:  mladi90$Min[mladi90$Team == "Burnley"]
## D = 0.34198, p-value = 0.05652
```

```
lillie.test(mladi90$Min[mladi90$Team == "Chelsea"])
```

```
##
##  Lilliefors (Kolmogorov-Smirnov) normality test
##
## data:  mladi90$Min[mladi90$Team == "Chelsea"]
## D = 0.16879, p-value = 0.5134
```

```
lillie.test(mladi90$Min[mladi90$Team == "Crystal Palace"])
```

```
##
##  Lilliefors (Kolmogorov-Smirnov) normality test
##
## data:  mladi90$Min[mladi90$Team == "Crystal Palace"]
## D = 0.24401, p-value = 0.127
```

```
lillie.test(mladi90$Min[mladi90$Team == "Leeds United"])
```

```
##
##  Lilliefors (Kolmogorov-Smirnov) normality test
##
## data:  mladi90$Min[mladi90$Team == "Leeds United"]
## D = 0.17923, p-value = 0.1555
```

```
lillie.test(mladi90$Min[mladi90$Team == "Leicester City"])
```

```
##
##  Lilliefors (Kolmogorov-Smirnov) normality test
##
## data:  mladi90$Min[mladi90$Team == "Leicester City"]
## D = 0.14925, p-value = 0.4919
```

```
lillie.test(mladi90$Min[mladi90$Team == "Liverpool"])
```

```
##
##  Lilliefors (Kolmogorov-Smirnov) normality test
##
## data:  mladi90$Min[mladi90$Team == "Liverpool"]
## D = 0.31804, p-value = 0.009129
```

```
lillie.test(mladi90$Min[mladi90$Team == "Manchester City"])
```

```
##
##  Lilliefors (Kolmogorov-Smirnov) normality test
##
## data:  mladi90$Min[mladi90$Team == "Manchester City"]
## D = 0.24065, p-value = 0.188
```

```
lillie.test(mladi90$Min[mladi90$Team == "Manchester United"])
```

```
##  
## Lilliefors (Kolmogorov-Smirnov) normality test  
##  
## data: mladi90$Min[mladi90$Team == "Manchester United"]  
## D = 0.2188, p-value = 0.2451
```

```
lillie.test(mladi90$Min[mladi90$Team == "Newcastle United"])
```

```
##  
## Lilliefors (Kolmogorov-Smirnov) normality test  
##  
## data: mladi90$Min[mladi90$Team == "Newcastle United"]  
## D = 0.13818, p-value = 0.9242
```

```
lillie.test(mladi90$Min[mladi90$Team == "Norwich City"])
```

```
##  
## Lilliefors (Kolmogorov-Smirnov) normality test  
##  
## data: mladi90$Min[mladi90$Team == "Norwich City"]  
## D = 0.1879, p-value = 0.09276
```

```
lillie.test(mladi90$Min[mladi90$Team == "Southampton"])
```

```
##  
## Lilliefors (Kolmogorov-Smirnov) normality test  
##  
## data: mladi90$Min[mladi90$Team == "Southampton"]  
## D = 0.13111, p-value = 0.7402
```

```
lillie.test(mladi90$Min[mladi90$Team == "Tottenham Hotspur"])
```

```
##  
## Lilliefors (Kolmogorov-Smirnov) normality test  
##  
## data: mladi90$Min[mladi90$Team == "Tottenham Hotspur"]  
## D = 0.12067, p-value = 0.7712
```

```
lillie.test(mladi90$Min[mladi90$Team == "Watford"])
```

```
##  
## Lilliefors (Kolmogorov-Smirnov) normality test  
##  
## data: mladi90$Min[mladi90$Team == "Watford"]  
## D = 0.13074, p-value = 0.8984
```

```
lillie.test(mladi90$Min[mladi90$Team == "West Ham United"])
```

```
##  
## Lilliefors (Kolmogorov-Smirnov) normality test  
##  
## data: mladi90$Min[mladi90$Team == "West Ham United"]  
## D = 0.20373, p-value = 0.5104
```

```
lillie.test(mladi90$Min[mladi90$Team == "Wolverhampton Wanderers"])
```

```
##
```

```
## Lilliefors (Kolmogorov-Smirnov) normality test
##
## data:  mladi90$Min[mladi90$Team == "Wolverhampton Wanderers"]
## D = 0.2024, p-value = 0.1538
```

Na razini značajnosti od 5% jedino Liverpool pravi probleme kod normalnosti. Iako varijanca i sredina ne odudaraju, uzorak ima izražene stršeće vrijednosti (Trent i Jota).

```
mladi90 %>% filter(Team == "Liverpool") %>% select(Player, Min)
```

```
##           Player  Min
## 1 Trent Alexander-Arnold 2853
## 2         Diogo Jota 2364
## 3      Ibrahima Konaté  990
## 4         Luis Díaz  958
## 5      Curtis Jones  851
## 6      Kostas Tsimikas  877
## 7      Harvey Elliott  346
## 8         Joe Gomez  331
## 9      Caoimhín Kelleher  180
```

```
mladi90bezL <- mladi90 %>% filter(Team != "Liverpool")
```

Sada kada smo pretpostavili homogenost varijance, normalnost i nezavisnost provest ćemo ANOVA test: Nulta hipoteza: Raspodijela minuta igračima do 25 godina se ne razlikuje po klubovima Alternativna hipoteza: Raspodijela minuta igračima do 25 godina razlikuje se u barem jednom klubu  $\alpha = 0.05$ .

```
anova(lm(Min ~ Team, data = mladi90bezL))
```

```
## Analysis of Variance Table
##
## Response: Min
##           Df      Sum Sq Mean Sq F value Pr(>F)
## Team       18  17052162   947342  1.1298  0.3251
## Residuals 209 175244441   838490
```

**Zaključci:** Ne odbacujemo nultu hipotezu da se raspodijela minuta razlikuje po klubovima.

Liverpool nismo uvrstili u test jer nismo mogli pretpostaviti normalnost, ali ni za tu ekipu ne možemo reći da značajno odstupa od prosjeka.

```
mean(mladi90$Min[mladi90$Team == "Liverpool"])
```

```
## [1] 1083.333
```

```
mean(mladi90$Min)
```

```
## [1] 1451.515
```

```
...
```

## 5.

**Dobivaju li u prosjeku više žutih kartona napadači ili igrači veznog reda?** Uzmimo za početak prosječne vrijednosti dobijenih žutih kartona kao motivaciju za statističko ispitivanje. Budući da se neki igrači smatraju i veznjacima i napadačima, njih nećemo ubrajati u ispitivanje kako bi mogli pretpostaviti nezavisnost.

Moramo pripaziti i na činjenicu da postoji podosta igrača s vrlo malo minuta odigrano, stoga ima smisla gledati igrače koji su u cijeloj sezoni sveukupno barem 50% minuta odigrali.

```
//samo napadaci i samo veznjaci
// napadaci su "FW", "FW,MF", "FW, DF"
// veznjaci su "MF", "MF,FW", "MF,DF"
#pitanje: ima li smisla ukljucivat MF,DF pod veznjake kao i FW,DF pod napadače
veznjaci <- nogometasi %>% filter(Pos == "MF" | Pos == "MF,FW" | Pos == "MF,DF") %>% filter(!is.na(X90s))
napadaci <- nogometasi %>% filter(Pos == "FW" | Pos == "FW,MF" | Pos == "FW,DF") %>% filter(!is.na(X90s))
cat("Prosječan broj žutih kartona igrača veznog reda iznosi: ", mean(veznjaci$CrdY, na.rm = T), "\n")

## Prosječan broj žutih kartona igrača veznog reda iznosi: 4.723077
cat("Prosječan broj žutih kartona napadača iznosi: ", mean(napadaci$CrdY, na.rm = T))

## Prosječan broj žutih kartona napadača iznosi: 3.765957
Vizualiziramo li podatke pomoću box plot:
boxplot(veznjaci$CrdY, napadaci$CrdY,
        names = c('broj žutih kartona veznih igrača', 'broj žutih kartona napadača'),
        main='Box plot raspodjele žutih kartona među veznjacima i napadačima')

## Warning in axis(side = base::quote(1), at = base::quote(1:2), labels =
## base::quote(c("broj žutih kartona veznih igrača", : conversion failure on 'broj
## žutih kartona veznih igrača' in 'mbcsToSbcs': dot substituted for <c5>

## Warning in axis(side = base::quote(1), at = base::quote(1:2), labels =
## base::quote(c("broj žutih kartona veznih igrača", : conversion failure on 'broj
## žutih kartona veznih igrača' in 'mbcsToSbcs': dot substituted for <be>

## Warning in axis(side = base::quote(1), at = base::quote(1:2), labels =
## base::quote(c("broj žutih kartona veznih igrača", : conversion failure on 'broj
## žutih kartona veznih igrača' in 'mbcsToSbcs': dot substituted for <c4>

## Warning in axis(side = base::quote(1), at = base::quote(1:2), labels =
## base::quote(c("broj žutih kartona veznih igrača", : conversion failure on 'broj
## žutih kartona veznih igrača' in 'mbcsToSbcs': dot substituted for <8d>

## Warning in axis(side = base::quote(1), at = base::quote(1:2), labels =
## base::quote(c("broj žutih kartona veznih igrača", : conversion failure on 'broj
## žutih kartona veznih igrača' in 'mbcsToSbcs': dot substituted for <c5>

## Warning in axis(side = base::quote(1), at = base::quote(1:2), labels =
## base::quote(c("broj žutih kartona veznih igrača", : conversion failure on 'broj
## žutih kartona veznih igrača' in 'mbcsToSbcs': dot substituted for <be>

## Warning in axis(side = base::quote(1), at = base::quote(1:2), labels =
## base::quote(c("broj žutih kartona veznih igrača", : conversion failure on 'broj
## žutih kartona veznih igrača' in 'mbcsToSbcs': dot substituted for <c4>

## Warning in axis(side = base::quote(1), at = base::quote(1:2), labels =
## base::quote(c("broj žutih kartona veznih igrača", : conversion failure on 'broj
## žutih kartona veznih igrača' in 'mbcsToSbcs': dot substituted for <8d>

## Warning in axis(side = base::quote(1), at = base::quote(1:2), labels =
## base::quote(c("broj žutih kartona veznih igrača", : conversion failure on 'broj
## žutih kartona napadača' in 'mbcsToSbcs': dot substituted for <c5>

## Warning in axis(side = base::quote(1), at = base::quote(1:2), labels =
## base::quote(c("broj žutih kartona veznih igrača", : conversion failure on 'broj
## žutih kartona napadača' in 'mbcsToSbcs': dot substituted for <be>
```

```

## Warning in axis(side = base::quote(1), at = base::quote(1:2), labels =
## base::quote(c("broj žutih kartona veznih igrača", : conversion failure on 'broj
## žutih kartona napadača' in 'mbcsToSbcs': dot substituted for <c4>

## Warning in axis(side = base::quote(1), at = base::quote(1:2), labels =
## base::quote(c("broj žutih kartona veznih igrača", : conversion failure on 'broj
## žutih kartona napadača' in 'mbcsToSbcs': dot substituted for <8d>

## Warning in axis(side = base::quote(1), at = base::quote(1:2), labels =
## base::quote(c("broj žutih kartona veznih igrača", : conversion failure on 'broj
## žutih kartona napadača' in 'mbcsToSbcs': dot substituted for <c5>

## Warning in axis(side = base::quote(1), at = base::quote(1:2), labels =
## base::quote(c("broj žutih kartona veznih igrača", : conversion failure on 'broj
## žutih kartona napadača' in 'mbcsToSbcs': dot substituted for <be>

## Warning in axis(side = base::quote(1), at = base::quote(1:2), labels =
## base::quote(c("broj žutih kartona veznih igrača", : conversion failure on 'broj
## žutih kartona napadača' in 'mbcsToSbcs': dot substituted for <c4>

## Warning in axis(side = base::quote(1), at = base::quote(1:2), labels =
## base::quote(c("broj žutih kartona veznih igrača", : conversion failure on 'broj
## žutih kartona napadača' in 'mbcsToSbcs': dot substituted for <8d>

## Warning in (function (main = NULL, sub = NULL, xlab = NULL, ylab = NULL, :
## conversion failure on 'Box plot raspodjele žutih kartona među veznjacima i
## napadačima' in 'mbcsToSbcs': dot substituted for <c5>

## Warning in (function (main = NULL, sub = NULL, xlab = NULL, ylab = NULL, :
## conversion failure on 'Box plot raspodjele žutih kartona među veznjacima i
## napadačima' in 'mbcsToSbcs': dot substituted for <be>

## Warning in (function (main = NULL, sub = NULL, xlab = NULL, ylab = NULL, :
## conversion failure on 'Box plot raspodjele žutih kartona među veznjacima i
## napadačima' in 'mbcsToSbcs': dot substituted for <c4>

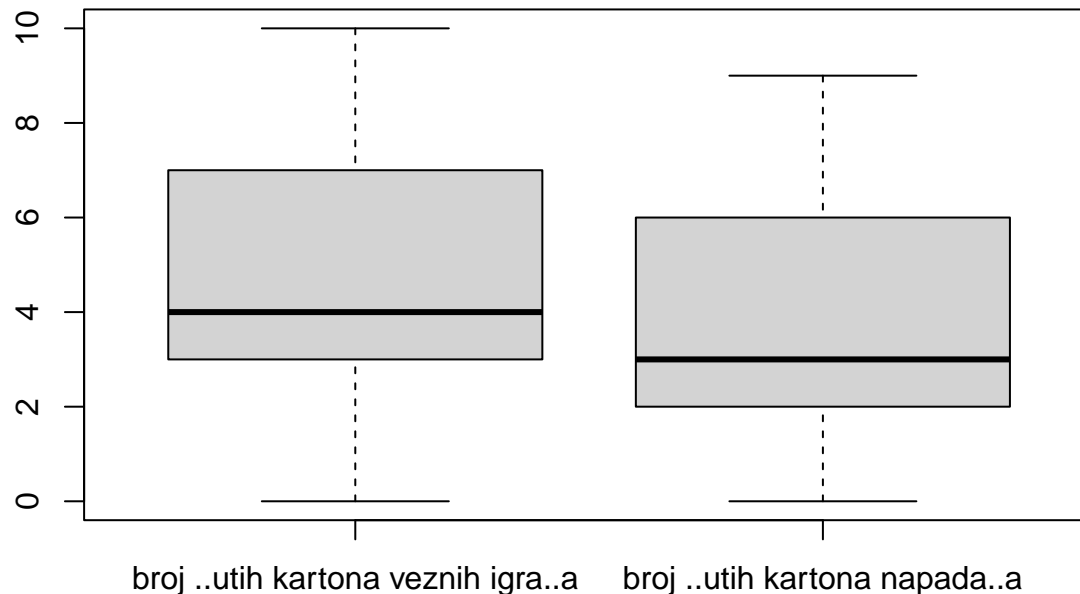
## Warning in (function (main = NULL, sub = NULL, xlab = NULL, ylab = NULL, :
## conversion failure on 'Box plot raspodjele žutih kartona među veznjacima i
## napadačima' in 'mbcsToSbcs': dot substituted for <91>

## Warning in (function (main = NULL, sub = NULL, xlab = NULL, ylab = NULL, :
## conversion failure on 'Box plot raspodjele žutih kartona među veznjacima i
## napadačima' in 'mbcsToSbcs': dot substituted for <c4>

## Warning in (function (main = NULL, sub = NULL, xlab = NULL, ylab = NULL, :
## conversion failure on 'Box plot raspodjele žutih kartona među veznjacima i
## napadačima' in 'mbcsToSbcs': dot substituted for <8d>

```

## Box plot raspodjele ..utih kartona me..u veznjacima i napada..ima



dobijemo bolju sliku stvarne raspodjele žutih kartona u kojoj vidimo neke indikacije da bi mogla postojati razlika u broju žutih kartona. Ovakvo ispitivanje bismo mogli provesti klasičnim t-testom, no prvo se moramo uvjeriti da raspodjele kartona dolaze iz normalne razdiobe.

Normalnost ćemo provjeriti histogramom i qq plotom.

```
hist(veznjaci$CrdY,  
     breaks=seq(min(veznjaci$CrdY, na.rm = T),max(veznjaci$CrdY, na.rm = T)+1,0.25),  
     main='Histogram količine žutih kartona igrača veznog reda',  
     xlab='broj žutih kartona')
```

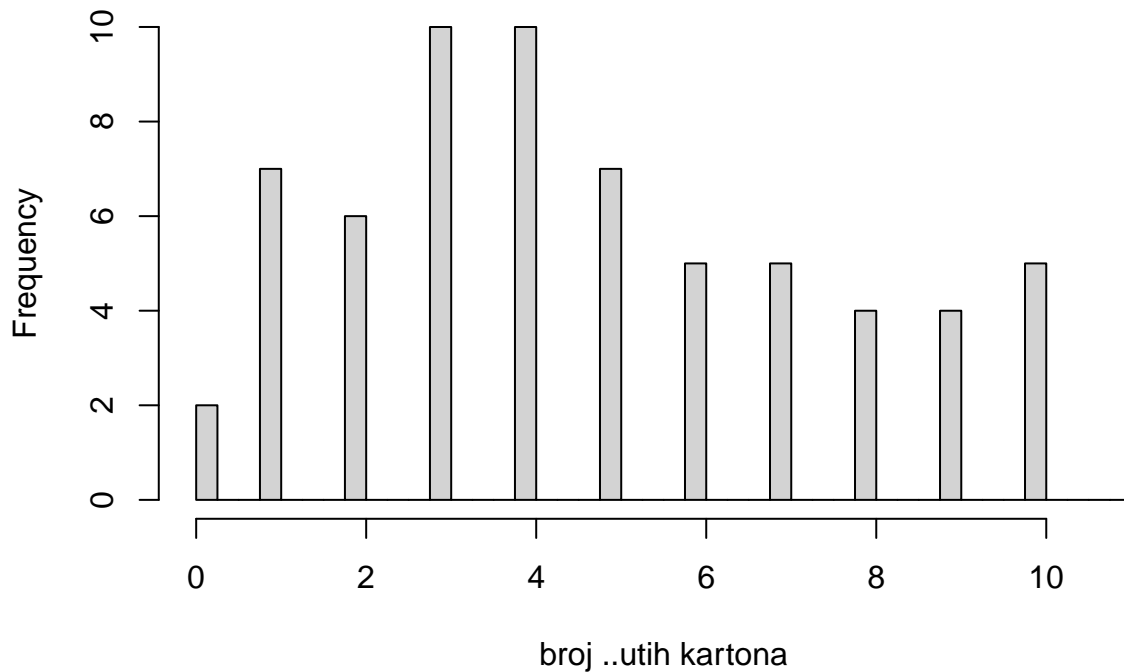
```
## Warning in title(main = main, sub = sub, xlab = xlab, ylab = ylab, ...):  
## conversion failure on 'Histogram količine žutih kartona igrača veznog reda' in  
## 'mbcsToSbcs': dot substituted for <c4>  
  
## Warning in title(main = main, sub = sub, xlab = xlab, ylab = ylab, ...):  
## conversion failure on 'Histogram količine žutih kartona igrača veznog reda' in  
## 'mbcsToSbcs': dot substituted for <8d>  
  
## Warning in title(main = main, sub = sub, xlab = xlab, ylab = ylab, ...):  
## conversion failure on 'Histogram količine žutih kartona igrača veznog reda' in  
## 'mbcsToSbcs': dot substituted for <c5>  
  
## Warning in title(main = main, sub = sub, xlab = xlab, ylab = ylab, ...):  
## conversion failure on 'Histogram količine žutih kartona igrača veznog reda' in  
## 'mbcsToSbcs': dot substituted for <be>  
  
## Warning in title(main = main, sub = sub, xlab = xlab, ylab = ylab, ...):  
## conversion failure on 'Histogram količine žutih kartona igrača veznog reda' in  
## 'mbcsToSbcs': dot substituted for <c4>  
  
## Warning in title(main = main, sub = sub, xlab = xlab, ylab = ylab, ...):  
## conversion failure on 'Histogram količine žutih kartona igrača veznog reda' in  
## 'mbcsToSbcs': dot substituted for <8d>  
  
## Warning in title(main = main, sub = sub, xlab = xlab, ylab = ylab, ...):
```



```
## conversion failure on 'broj žutih kartona' in 'mbcsToSbcs': dot substituted for
## <c5>

## Warning in title(main = main, sub = sub, xlab = xlab, ylab = ylab, ...):
## conversion failure on 'broj žutih kartona' in 'mbcsToSbcs': dot substituted for
## <be>
```

## Histogram količine žutih kartona igrača vrnog reda



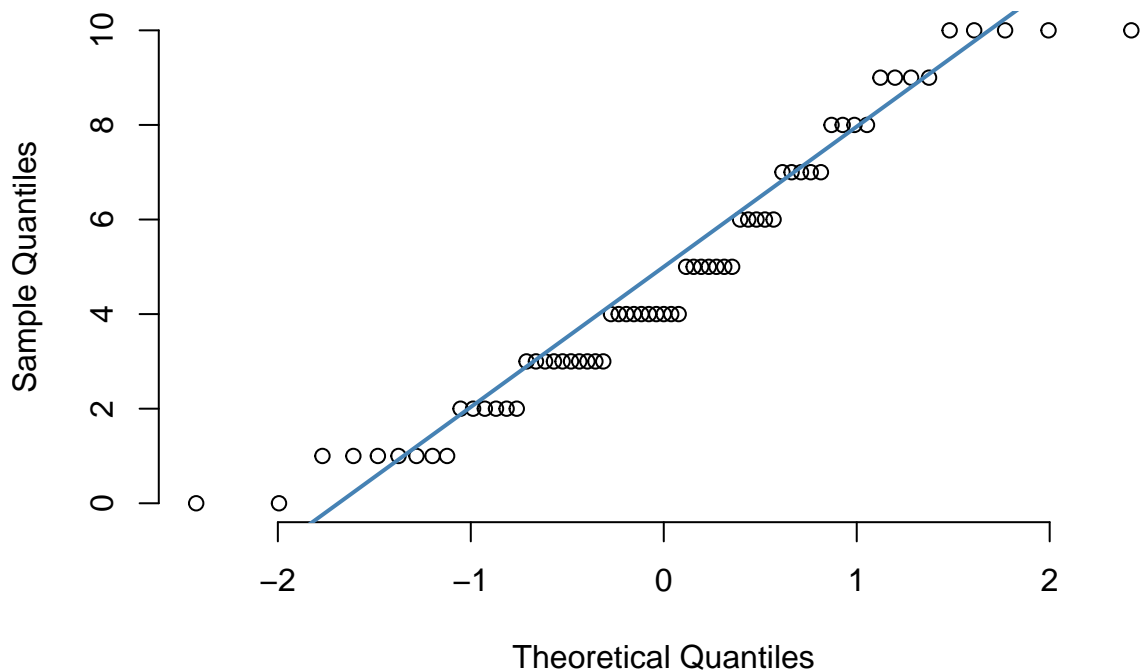
```
qqnorm(veznjaci$CrdY, pch = 1, frame = FALSE, main='igrači vrnog reda')
```

```
## Warning in title(...): conversion failure on 'igrači vrnog reda' in
## 'mbcsToSbcs': dot substituted for <c4>
```

```
## Warning in title(...): conversion failure on 'igrači vrnog reda' in
## 'mbcsToSbcs': dot substituted for <8d>
```

```
qqline(veznjaci$CrdY, col = "steelblue", lwd = 2)
```

## igra..i veznog reda



```
hist(napadaci$CrdY,
     breaks=seq(min(napadaci$CrdY, na.rm = T),max(napadaci$CrdY, na.rm = T)+1,0.25),
     main='Histogram količine žutih kartona napadača',
     xlab='broj žutih kartona')
```

```
## Warning in title(main = main, sub = sub, xlab = xlab, ylab = ylab, ...):
## conversion failure on 'Histogram količine žutih kartona napadača' in
## 'mbcsToSbcs': dot substituted for <c4>

## Warning in title(main = main, sub = sub, xlab = xlab, ylab = ylab, ...):
## conversion failure on 'Histogram količine žutih kartona napadača' in
## 'mbcsToSbcs': dot substituted for <8d>

## Warning in title(main = main, sub = sub, xlab = xlab, ylab = ylab, ...):
## conversion failure on 'Histogram količine žutih kartona napadača' in
## 'mbcsToSbcs': dot substituted for <c5>

## Warning in title(main = main, sub = sub, xlab = xlab, ylab = ylab, ...):
## conversion failure on 'Histogram količine žutih kartona napadača' in
## 'mbcsToSbcs': dot substituted for <be>

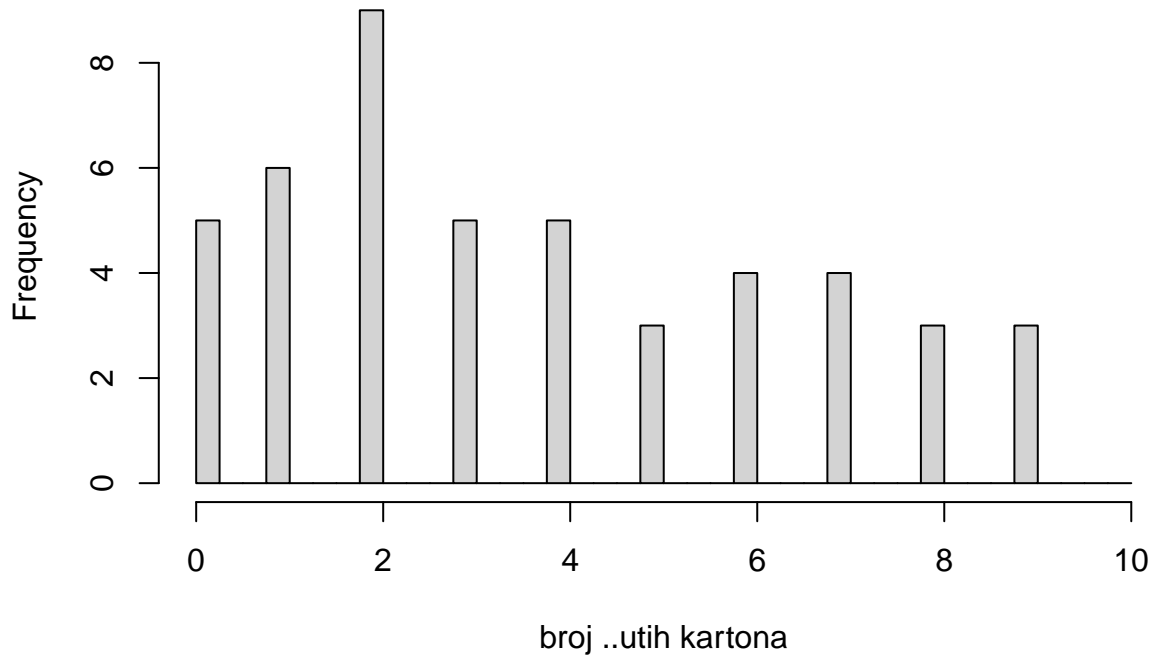
## Warning in title(main = main, sub = sub, xlab = xlab, ylab = ylab, ...):
## conversion failure on 'Histogram količine žutih kartona napadača' in
## 'mbcsToSbcs': dot substituted for <c4>

## Warning in title(main = main, sub = sub, xlab = xlab, ylab = ylab, ...):
## conversion failure on 'Histogram količine žutih kartona napadača' in
## 'mbcsToSbcs': dot substituted for <8d>

## Warning in title(main = main, sub = sub, xlab = xlab, ylab = ylab, ...):
## conversion failure on 'broj žutih kartona' in 'mbcsToSbcs': dot substituted for
## <c5>
```

```
## Warning in title(main = main, sub = sub, xlab = xlab, ylab = ylab, ...):
## conversion failure on 'broj žutih kartona' in 'mbcsToSbcs': dot substituted for
## <be>
```

## Histogram količine žutih kartona napadača

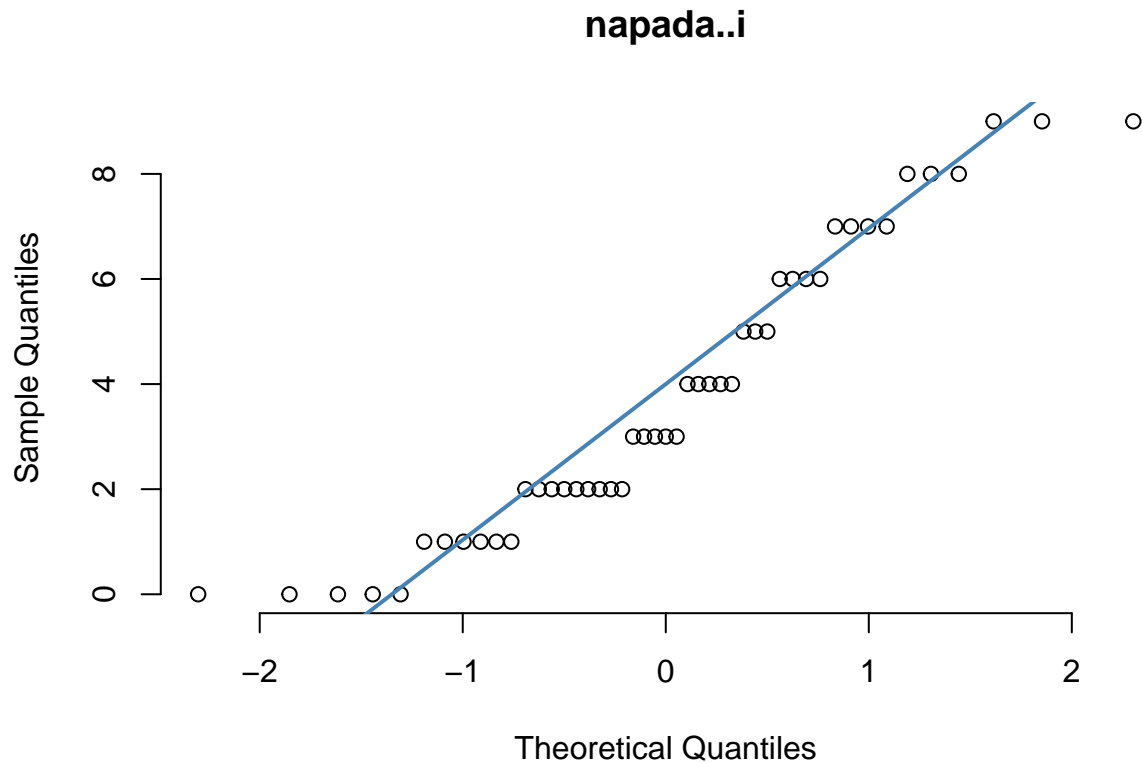


```
qqnorm(napadaci$CrđY, pch = 1, frame = FALSE, main='napadači')
```

```
## Warning in title(...): conversion failure on 'napadači' in 'mbcsToSbcs': dot
## substituted for <c4>
```

```
## Warning in title(...): conversion failure on 'napadači' in 'mbcsToSbcs': dot
## substituted for <8d>
```

```
qqline(napadaci$CrđY, col = "steelblue", lwd = 2)
```



Budući da znamo da je t-test poprilično robustan, dajemo si za pravo koristiti ga iako gore prikazane razdiobe nisu distribuirane normalnom razdiobom, no nisu ni predaleko od iste.

Provjeravamo jesu li varijance uzoraka značajno različite:

```
cat("Varijanca broja žutih kartona kod veznjaka iznosi: ", var(veznjaci$CrdY), "\n")
```

```
## Varijanca broja žutih kartona kod veznjaka iznosi: 7.984615
```

```
cat("Varijanca broja žutih kartona kod napadača iznosi: ", var(napadaci$CrdY))
```

```
## Varijanca broja žutih kartona kod napadača iznosi: 7.617946
```

ispitajmo...

```
var.test(veznjaci$CrdY, napadaci$CrdY)
```

```
##
```

```
## F test to compare two variances
```

```
##
```

```
## data: veznjaci$CrdY and napadaci$CrdY
```

```
## F = 1.0481, num df = 64, denom df = 46, p-value = 0.876
```

```
## alternative hypothesis: true ratio of variances is not equal to 1
```

```
## 95 percent confidence interval:
```

```
## 0.6024446 1.7798549
```

```
## sample estimates:
```

```
## ratio of variances
```

```
## 1.048132
```

Ne odbacujemo  $H_0$  koja kaže da su varijance jednake. Dakle koristit ćemo **t-test za dva uzorka s pretpostavkom jednakih varijanci**.

$H_0$ : broj žutih kartona između veznjaka i napadača je jednak.

H1: broj žutih kartona kod veznjaka veći je od onog kod napadača.

Odabir H1 motiviran je saznanjem da očekujemo da veznjaci imaju više žutih kartona.

```
t.test(veznjaci$CrdY, napadaci$CrdY, alt = "greater", var.equal = TRUE)
```

```
##  
## Two Sample t-test  
##  
## data: veznjaci$CrdY and napadaci$CrdY  
## t = 1.7863, df = 110, p-value = 0.03841  
## alternative hypothesis: true difference in means is greater than 0  
## 95 percent confidence interval:  
## 0.06828648 Inf  
## sample estimates:  
## mean of x mean of y  
## 4.723077 3.765957
```

*#izveo sam i za two.sided, outcome je isti*

Budući da je p-value značajno malen, možemo odbaciti H0 u korist H1.

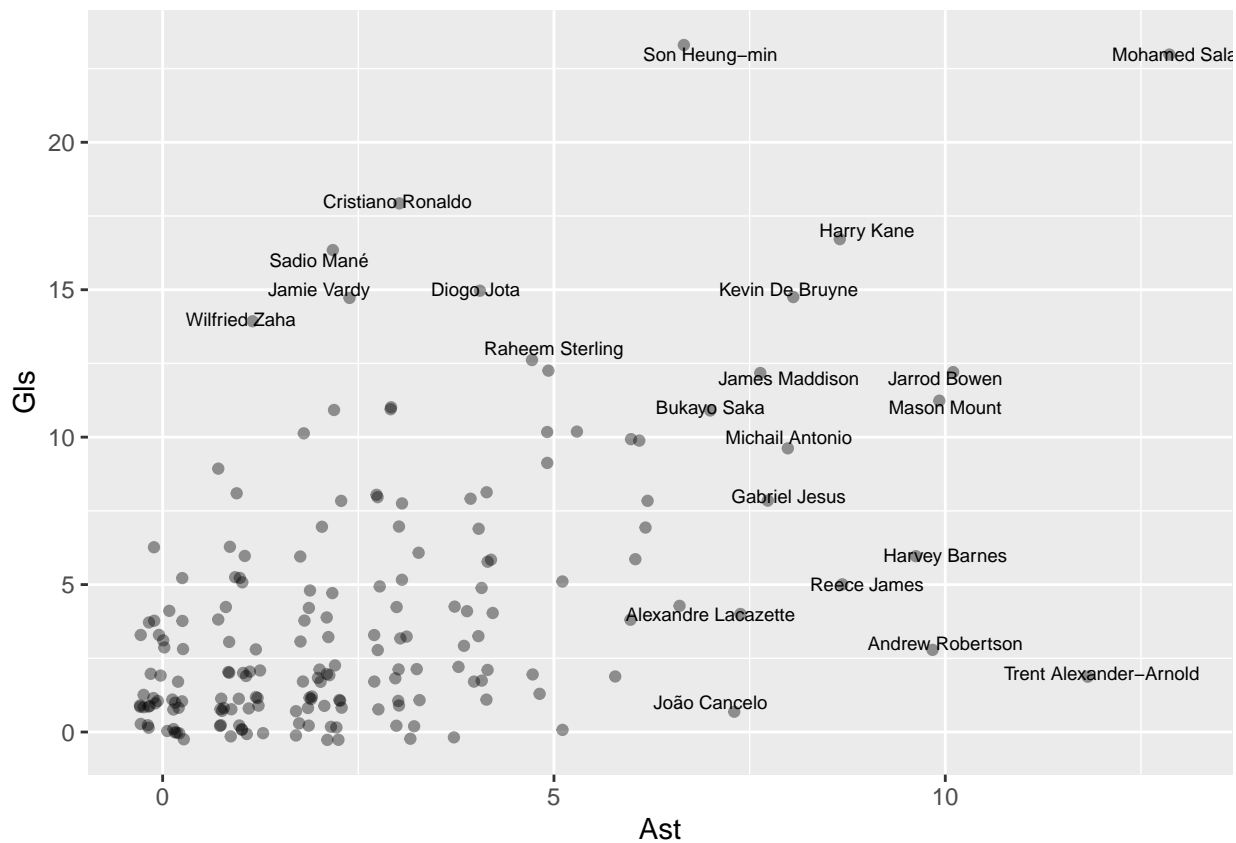
**Zaključci:** Na razini pouzdanosti od 95% odbacujemo H0 u korist H1 odnosno zaključujemo da je broj žutih kartona kod veznjaka veći od onog kod napadača.

**6. Možete li na temelju zadanih parametara odrediti uspješnost pojedinog igrača?** Što je zapravo uspješnost igrača? To je pitanje kojim smo se prvotno morali baviti i secirati što čini dobrog igrača ovisno o pozicijama.

Kao mjere uspješnosti igrača na raspolaganju imamo broj golova i broj asistencija. Naravno nije objektivno usporedjivat obrambene vezne i napadače prema broju golova tako da za neke pozicije sljedeća analiza nije najpogodnija.

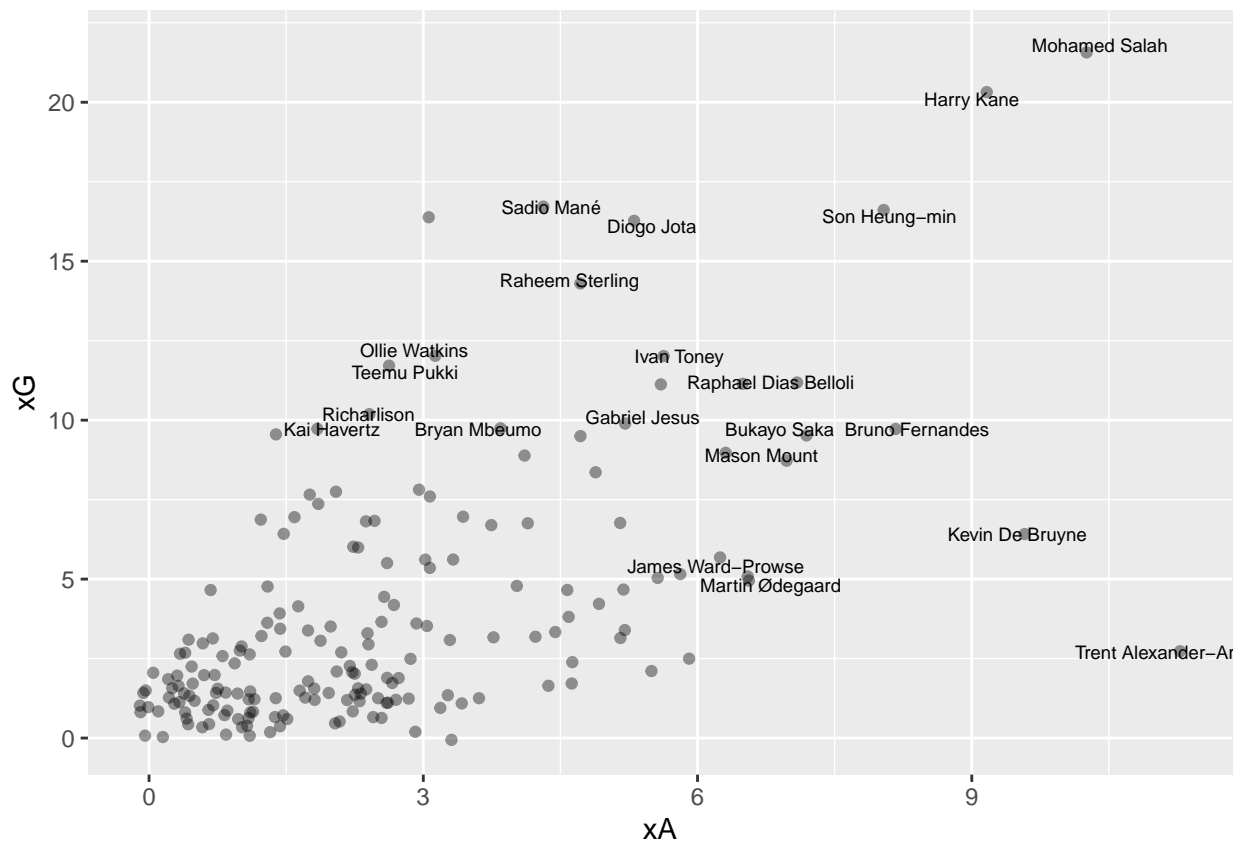
G/A bez golmana

```
nog <- nogometasi %>% filter(Pos != "GK") %>% filter(X90s >= 19)  
dobri <- nog %>% filter(Ast > 6 | Gls > 12)  
losi <- nog %>% filter(Ast <= 6 & Gls <= 12)  
ggplot(dobri, aes(x = Ast, y = Gls)) + geom_jitter(width = 0.4, height = 0.4, alpha = 0.4) + geom_text(
```



xG/xA bez golmana

```
dobrixG <- nog %>% filter(xG > 9.5 | xA > 6)
losixG <- nog %>% filter(xG <= 9.5 & xA <= 6)
ggplot(dobrixG, aes(x = xA, y = xG)) + geom_jitter(width = 0.3, height = 0.3, alpha = 0.4) + geom_text(
```



\*\*napomena: u gornja dva grafa dodan je *jitter* efekt kako bi se stekao bolji dojam količine točaka jer se koriste diskretni podaci\*\*

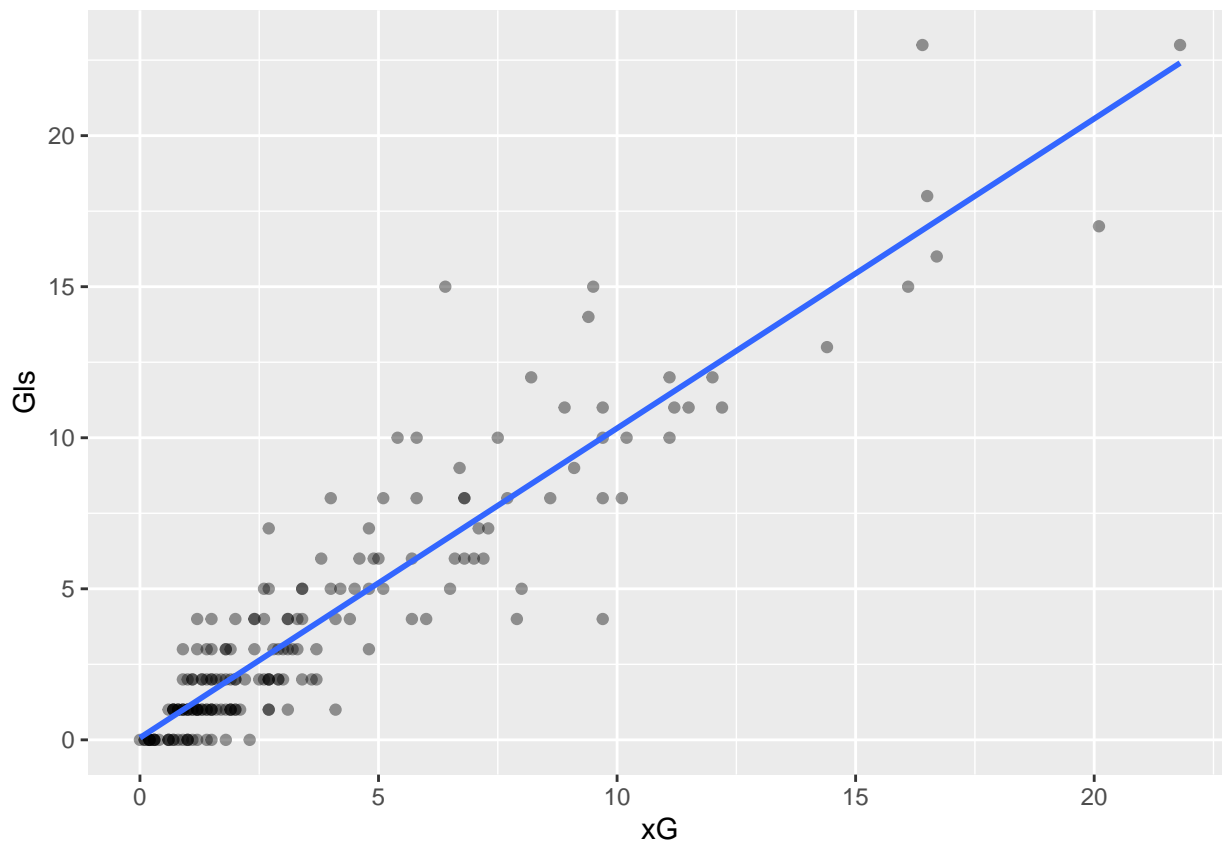
Gper90/Aper90 bez golmana

### Određivanje uspješnosti po broju golova preko mjere očekivanih golova.

Osobi koja ne prati nogomet pojam očekivanih golova (xG) je možda nepoznat pa ćemo napomenuti da se radi o mjeri koja pokazuje procjenu vjerojatnosti u kojima neka prilika završi zgoditkom.

Gls/xG

```
ggplot(nog, aes(x = xG, y = GlS)) + geom_point(alpha = 0.4) + stat_smooth(method = lm, formula = y~x, s
```



Prema grafu se da naslutiti da postoji jasna linearna veza između golova i očekivanih golova što daje motivaciju za daljnje istraživanje.

```
fit.gls = lm(Gls~xG,data=nog)
```

Potrebno je provjeriti jesu li narušene osnovne pretpostavke o rezidualima prije nego nastavimo dalje. Pretpostavke reziduala su normalnost i homogenost varijance.

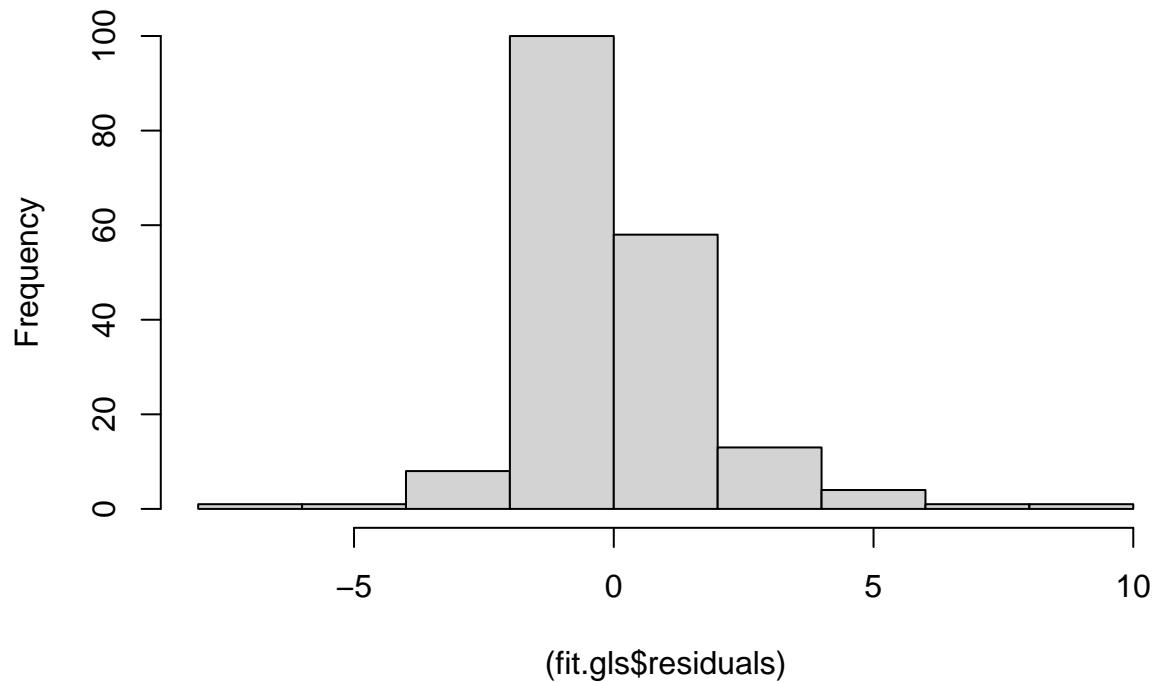
### Normalnost

Normalnost možemo provjeriti grafički pomoću histograma.

```
hist((fit.gls$residuals))
```



## Histogram of (fit.gls\$residuals)



Statistički ju možemo provjeriti pomoću Kolmogorov-Smirnovljevog testa.

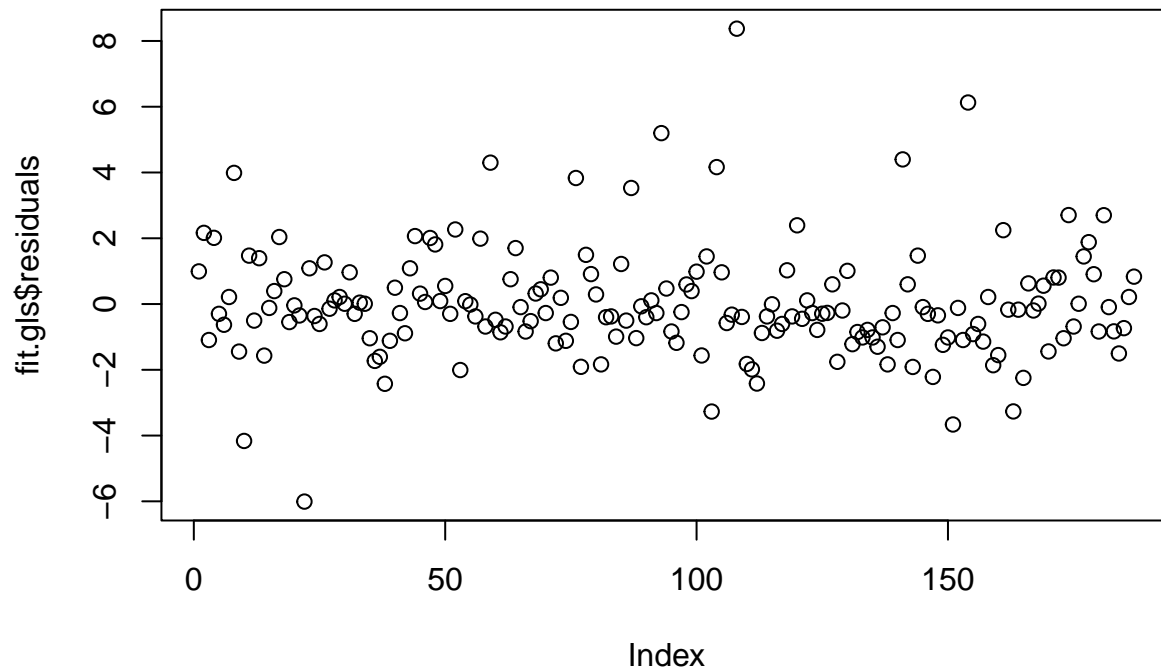
```
require(nortest)
lillie.test(fit.gls$residuals)
```

```
##
##  Lilliefors (Kolmogorov-Smirnov) normality test
##
## data:  fit.gls$residuals
## D = 0.12357, p-value = 2.428e-07
```

Budući da reziduali ne pokazuju preveliko odstupanje od normalnosti (u smislu zakrivljenosti ili drugih razlika u distribuciji) te je poznato da je t-test robustan na (ne)normalnost - u analizi podataka se u ovakvim slučajevima i dalje mogu donositi statistički zaključci iz regresijskih modela.

Homogenost varijance provjerit ćemo grafički prikazom reziduala. Bitno nam je da se reziduali ne šire povećanjem y.

```
plot(fit.gls$residuals)
```



Pogledajmo rezultat analize...

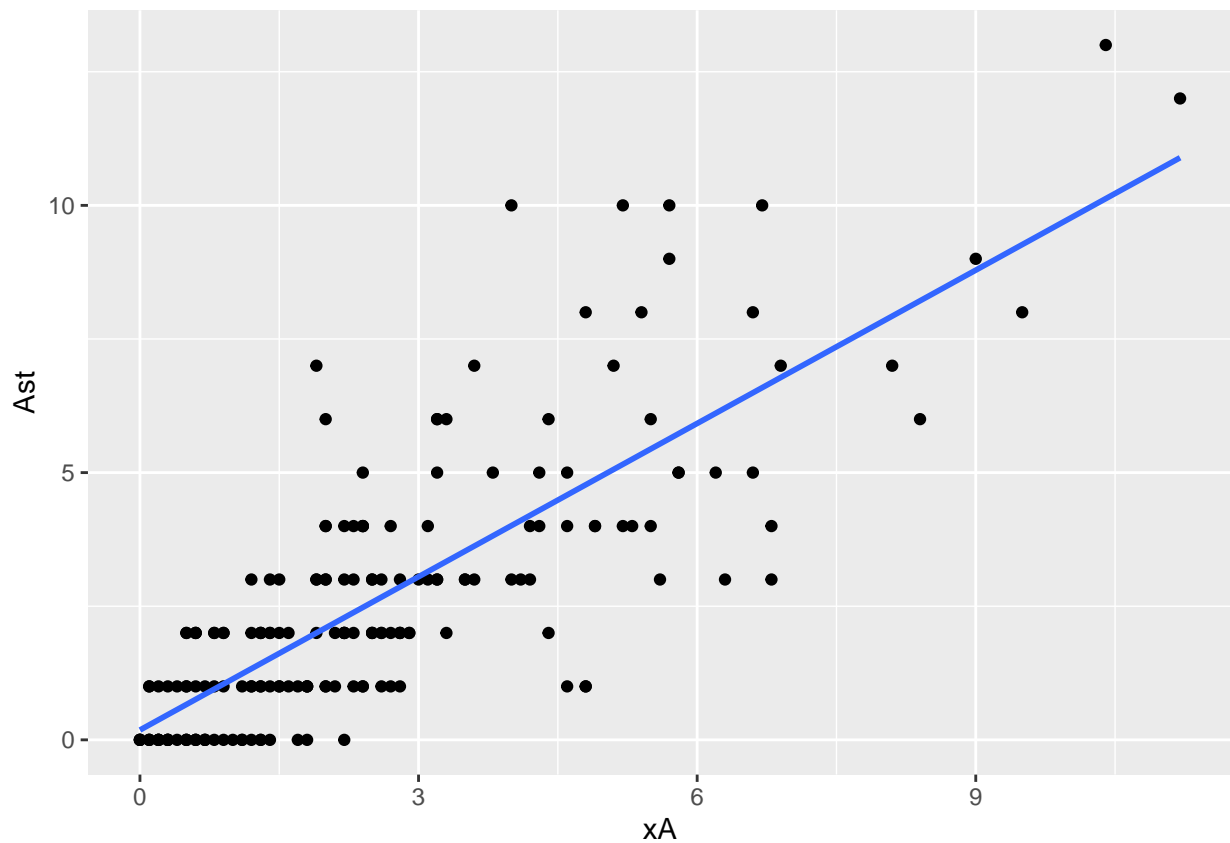
```
summary(fit.gls)
```

```
##
## Call:
## lm(formula = Gls ~ xG, data = nog)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -6.0076 -0.8704 -0.2742  0.6911  8.3735
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  0.06929    0.17266   0.401   0.689
## xG           1.02456    0.03110  32.941 <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 1.696 on 185 degrees of freedom
## Multiple R-squared:  0.8543, Adjusted R-squared:  0.8536
## F-statistic: 1085 on 1 and 185 DF, p-value: < 2.2e-16
```

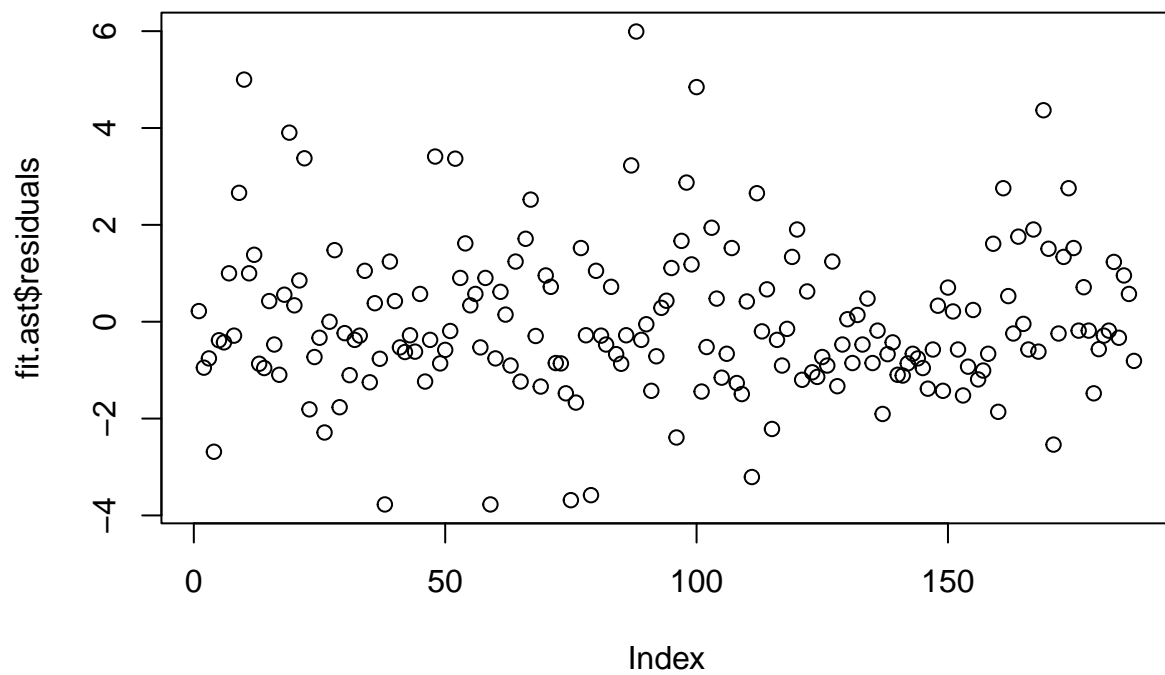
Kao mjeru valjanosti linearne veze razmatramo varijablu  $R^2$ . Ona iznosi 0.854 što je dovoljno dobro za reći da mjerom xG relativno dobro možemo odrediti uspjehnost igrača.

A/xA

```
fit.ast = lm(Ast~xA,data=nog)
ggplot(nog, aes(x = xA, y = Ast)) + geom_point() + stat_smooth(method = lm, formula = y~x, se = F)
```



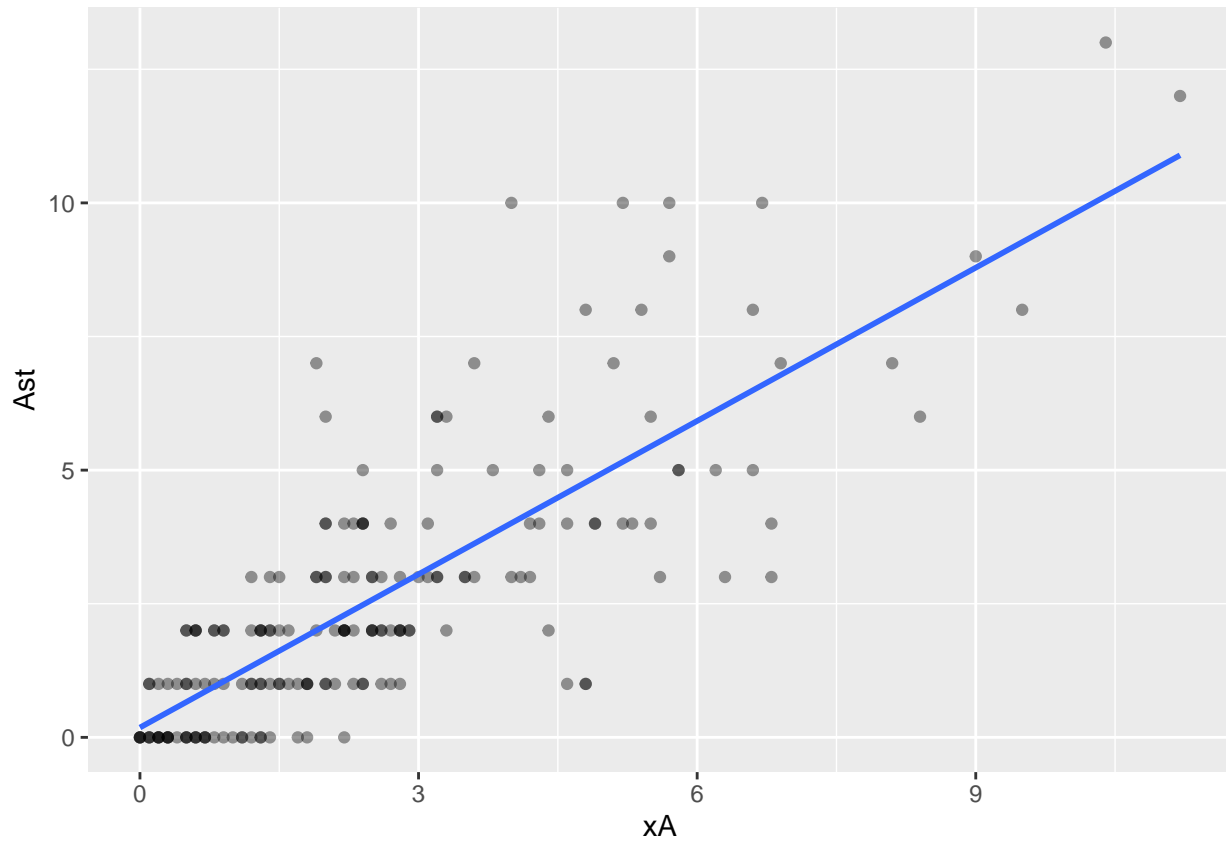
```
plot(fit.ast$residuals)
```



Određivanje uspešnosti po broju asistencija preko mjere očekivanih asistencija.

Prikažimo raspodjelu na grafu...

```
ggplot(nog, aes(x = xA, y = Ast)) + geom_point(alpha = 0.4) + stat_smooth(method = lm, formula = y~x, s
```



Prema grafu se da naslutiti da postoji linearna veza između te dvije mjere, no uvjerljivo kao golovi. Svejedno, provest ćemo istraživanje i vidjeti rezultate.

```
fit.ast = lm(Ast~xA,data=nog)
```

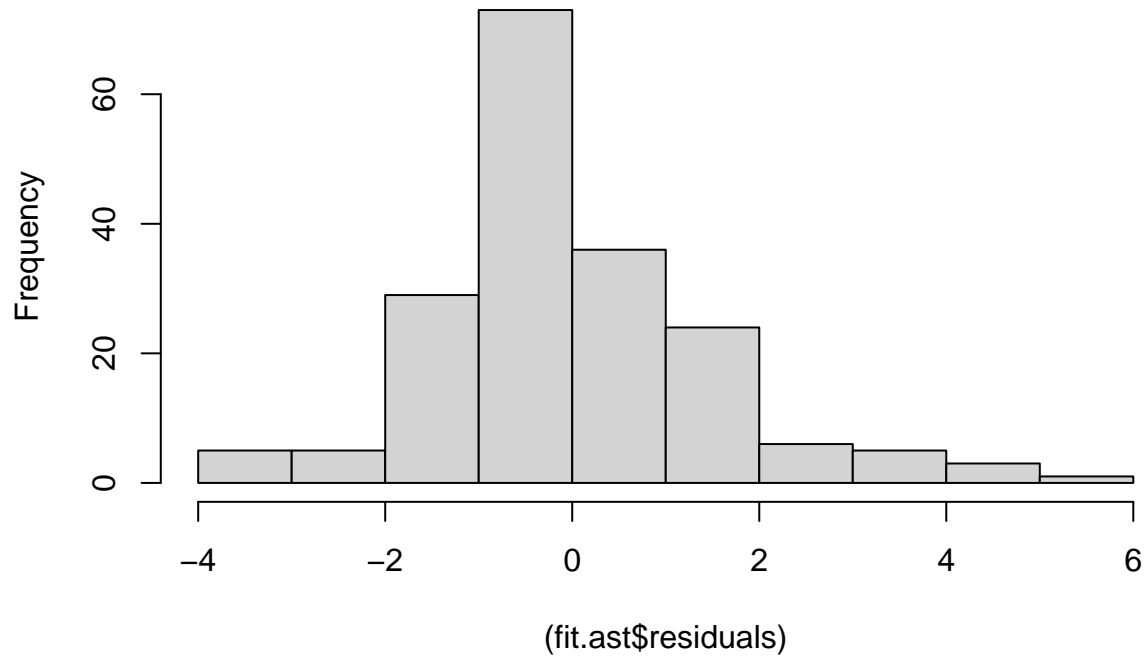
Potrebno je provjeriti jesu li narušene osnovne pretpostavke o rezidualima prije nego nastavimo dalje. Pretpostavke reziduala su normalnost i homogenost varijance.

### Normalnost

Normalnost možemo provjeriti grafički pomoću histograma.

```
hist((fit.ast$residuals))
```

## Histogram of (fit.ast\$residuals)



Statistički ju možemo provjeriti pomoću Kolmogorov-Smirnovljevog testa.

```
require(nortest)
lillie.test(fit.ast$residuals)
```

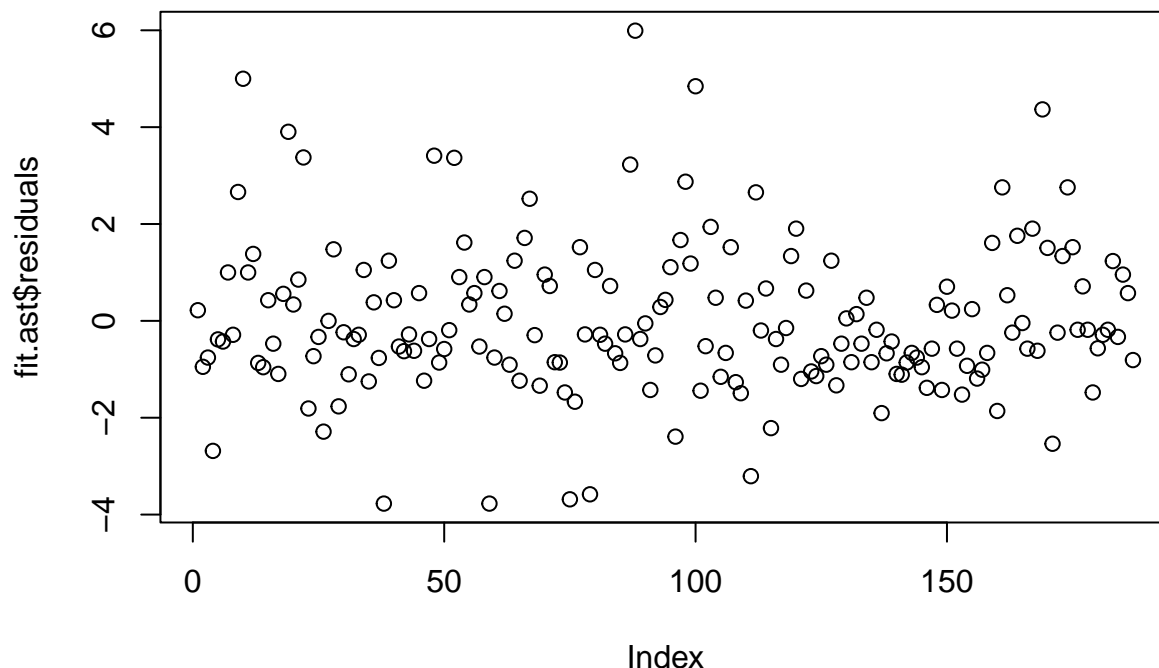
```
##
##  Lilliefors (Kolmogorov-Smirnov) normality test
##
## data:  fit.ast$residuals
## D = 0.12465, p-value = 1.774e-07
```

Budući da reziduali ne pokazuju preveliko odstupanje od normalnosti (u smislu zakrivljenosti ili drugih razlika u distribuciji) te je poznato da je t-test robustan na (ne)normalnost - u analizi podataka se u ovakvim slučajevima i dalje mogu donositi statistički zaključci iz regresijskih modela.

### Homogenost

Homogenost varijance provjerit ćemo grafički prikazom reziduala. Bitno nam je da se reziduali ne šire povećanjem y.

```
plot(fit.ast$residuals)
```



Pogledajmo rezultat analize...

```
summary(fit.ast)
```

```
##
## Call:
## lm(formula = Ast ~ xA, data = nog)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -3.7727 -0.8683 -0.2871  0.7205  5.9921
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  0.18391    0.17833   1.031   0.304
## xA           0.95600    0.05302  18.031 <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 1.558 on 185 degrees of freedom
## Multiple R-squared:  0.6373, Adjusted R-squared:  0.6354
## F-statistic: 325.1 on 1 and 185 DF, p-value: < 2.2e-16
```

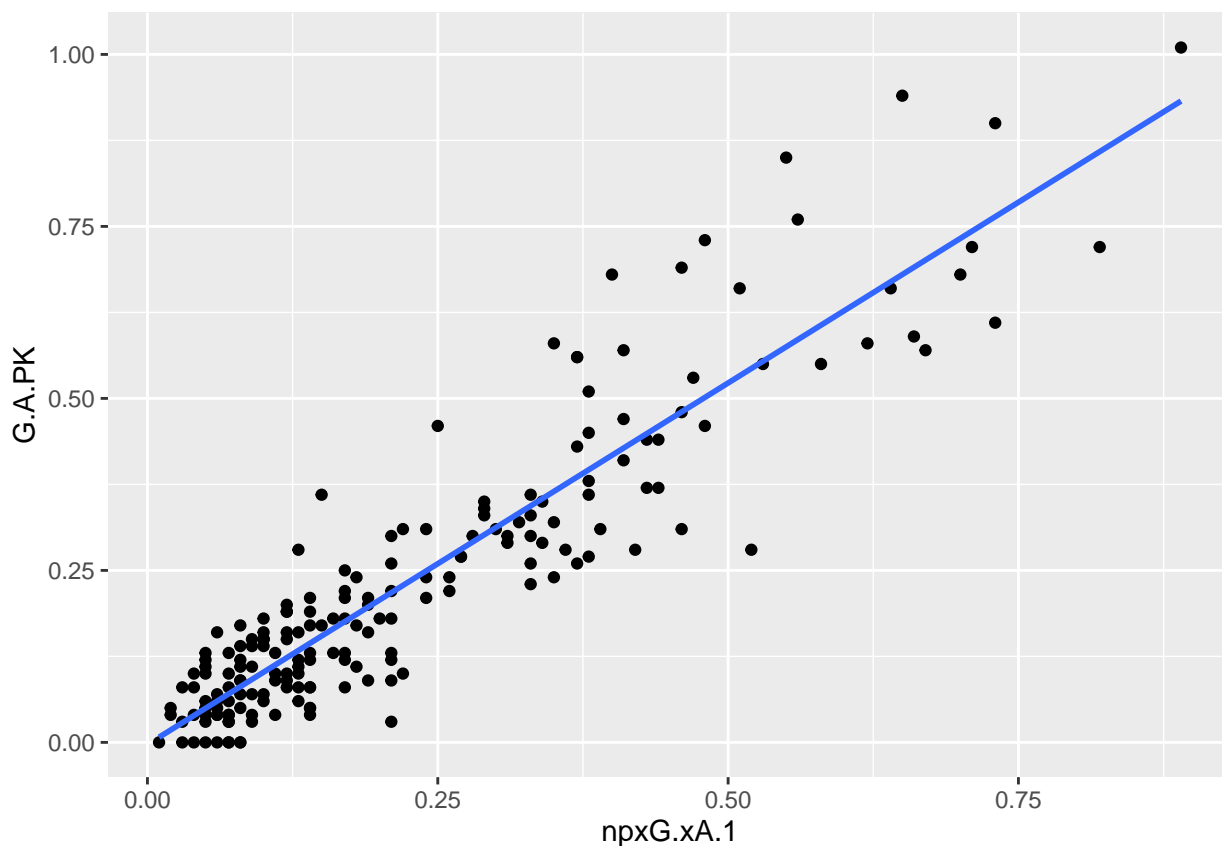
Kao mjeru valjanosti linearne veze razmatramo varijablu  $R^2$ . Ona iznosi 0.64 što nije dovoljno dobro kao prethodno istraživanje. Iako je  $R^2$  relativna mjera, smatramo da nije dovoljno dobra za reći da mjerom xA možemo precizno odrediti uspješnost igrača.

Ovaj rezultat prirodno ima smisla kada pogledamo nogometni kontekst očekivane asistencije koja ovisi o tome hoće li asistirani igrač zabiti.

**Određivanje uspješnosti po broju golova i asistencija bez kaznenih udaraca u po 90 min preko mjere očekivanih golova i asistencija bez kaznenih udaraca po 90 min.**

$G+A/npxG+xA$

```
ggplot(nog, aes(x = npxG.xA.1, y = G.A.PK)) + geom_point() + stat_smooth(method = lm, formula = y~x, se
```



Možemo opravdano naslutiti da postoji jaka linearna veza između ovih mjera.

```
fit.ga = lm(G.A.PK~npxG.xA.1,data=nog)
```

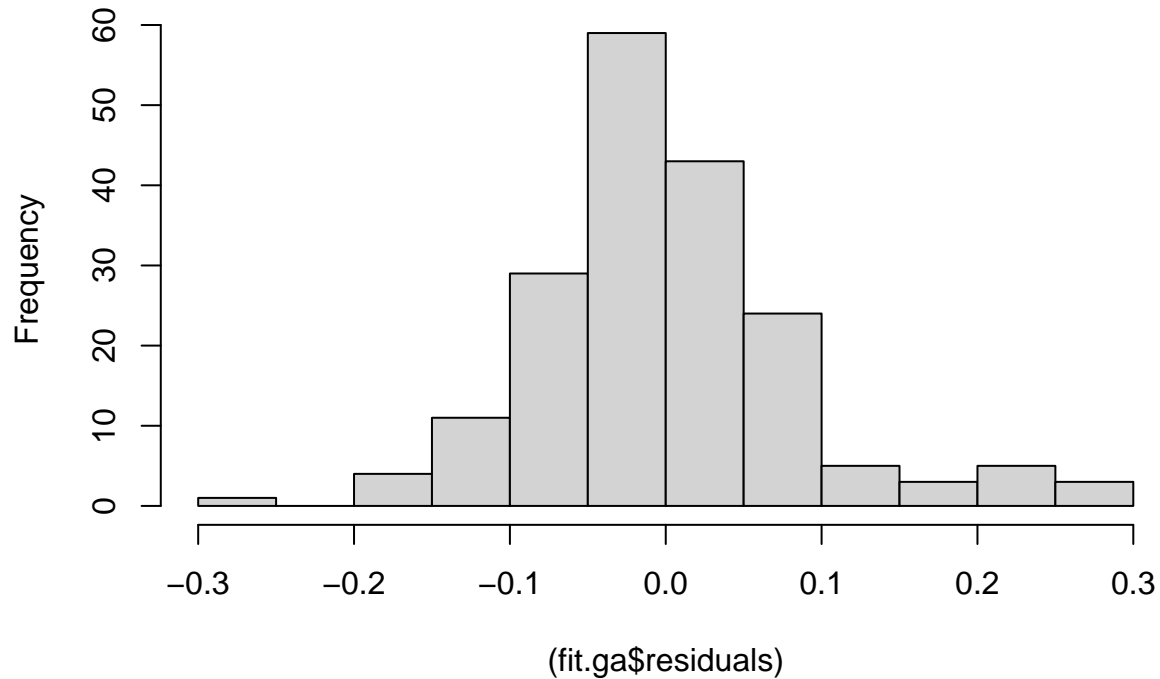
Potrebno je provjeriti jesu li narušene osnovne pretpostavke o rezidualima prije nego nastavimo dalje. Pretpostavke reziduala su normalnost i homogenost varijance.

### Normalnost

Normalnost možemo provjeriti grafički pomoću histograma.

```
hist((fit.ga$residuals))
```

## Histogram of (fit.ga\$residuals)



Statistički ju možemo provjeriti pomoću Kolmogorov-Smirnovljevog testa.

```
require(nortest)
lillie.test(fit.ga$residuals)
```

```
##
##  Lilliefors (Kolmogorov-Smirnov) normality test
##
## data:  fit.ga$residuals
## D = 0.079875, p-value = 0.005499
```

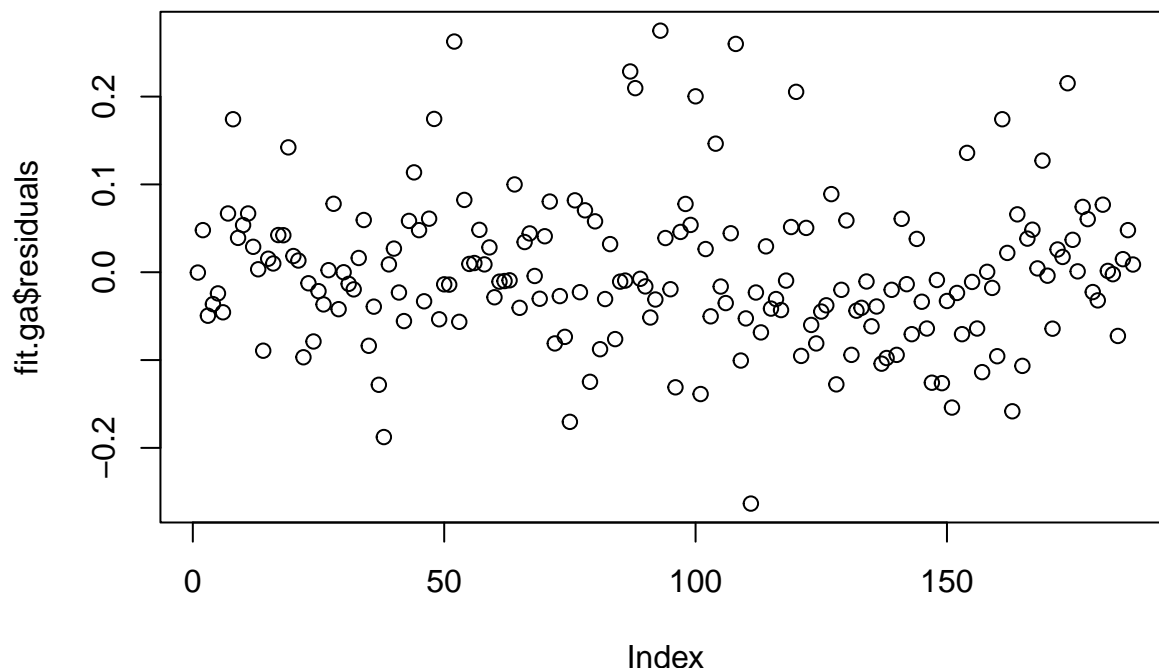
Budući da reziduali ne pokazuju preveliko odstupanje od normalnosti (u smislu zakrivljenosti ili drugih razlika u distribuciji) te je poznato da je t-test robustan na (ne)normalnost - u analizi podataka se u ovakvim slučajevima i dalje mogu donositi statistički zaključci iz regresijskih modela.

### Homogenost

Homogenost varijance provjerit ćemo grafički prikazom reziduala. Bitno nam je da se reziduali ne šire povećanjem y.

```
plot(fit.ga$residuals)
```





Pogledajmo rezultat analize...

```
summary(fit.ga)
```

```
##
## Call:
## lm(formula = G.A.PK ~ npxG.xA.1, data = nog)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -0.26343 -0.04521 -0.01004  0.04428  0.27504
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept) -0.003017   0.009901  -0.305   0.761
## npxG.xA.1    1.050857   0.033811  31.080 <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.08462 on 185 degrees of freedom
## Multiple R-squared:  0.8393, Adjusted R-squared:  0.8384
## F-statistic: 966 on 1 and 185 DF, p-value: < 2.2e-16
```

Kao mjeru valjanosti linearne veze razmatramo varijablu  $R^2$ . Ona iznosi 0.84 što opravdava naše izvorne pretpostavke.

## 7.

Svi koji prate nogomet malo detaljnije znaju čiji igrači se cijene. Brazilci su najbolji dribleri, Španjolci najbolji u tiki-taki, Hrvati najbolji u penalima, ali u Engleskoj su najbolji Englezi. Javnost to zove “*English tax*” i time se cilja na činjenicu kako engleski klubovi skuplje plaćaju i prodaju domaće igrače u odnosu na strane. Je li to opravdano, pokazat će nam naš xz-test. Koristit ćemo ga jer **razlog** i napokon ćemo saznati doprinose li oni sveukupnom uspjehu tima ili je to još jedna preuveličana engleska nogometna bajka. ...

**Zaključci:** ...

## **8. Rezultati**

...

...