

```
[1]: df <- read.csv("master.csv")
head(df)
```

A data frame: 6 x 12											
d<country	year	sex	age	suicides.no	population	suicides.100k.pop	country.year	hdi.for.year	gdp.for.year	gdp.per.capita...	generation
<chr>	<int>	<chr>	<chr>	<int>	<int>	<dbl>	<chr>	<dbl>	<chr>	<int>	<chr>
1	Albania	1987	male	15-24 years	21	312900					
2	Albania	1987	male	35-54 years	16	302900	6.71	Albania1987	NA	2,156,624,900	796
3	Albania	1987	female	15-24 years	14	289000	4.83	Albania1987	NA	2,156,624,900	796
4	Albania	1987	male	75+ years	1	21800	5.19	Albania1987	NA	2,156,624,900	796
5	Albania	1987	male	25-34 years	9	274300	3.28	Albania1987	NA	2,156,624,900	796
6	Albania	1987	female	75+ years	1	35600	2.81	Albania1987	NA	2,156,624,900	796
G.I Generation											

```
In [2]: #Napraviću novu data frame koji će predstavljati sreden prethodni data frame
reworked.df <- df

#Potrešno je preimenovali kolone
names(reworked.df)[1] <- "country"

#Prvo obrisati suviševe tačke na krajevima naziva kolona
dots.delete <- function(name){
  return(sub("\\.(2,)*", "", name))
}

names(reworked.df) <- sapply(names(reworked.df), dots.delete)

#Svi su potrebne preimenovali "a" u "a", "b" u "b" i imali iste notacije u kolonama
colnnotation <- function(coln){
  return(sub("-", "", coln))
}

names(reworked.df) <- sapply(names(reworked.df), colnnotation)

#Menjamo "hdi" u malo "hdi"
names(reworked.df) <- sapply(names(reworked.df), tolower)

#gdp.for.year iz string u numeric
reworked.df$gdp.for.year <- sapply(reworked.df$gdp.for.year, function(x) {gsub(" ", "", x)})
reworked.df$gdp.for.year <- sapply(reworked.df$gdp.for.year, as.numeric)

#Potrešno napraviti faktore
reworked.dffsex <- as.factor(reworked.dffsex)
reworked.dffage <- factor(reworked.dffage, levels=c(
  "5-14 years", "15-24 years", "25-34 years", "35-54 years", "55-74 years", "75+ years"), ordered = TRUE)

#unique(reworked.dffgeneration)

reworked.dffgeneration <- factor(reworked.dffgeneration, levels=c(
  "G.I Generation", "Silent", "Boomers", "Generation X", "Millenials", "Generation Z"))

head(reworked.df)
```

A data frame: 6 x 12											
country	year	sex	age	suicides.no	population	suicides.100k.pop	country.year	hdi.for.year	gdp.for.year	gdp.per.capita	generation
<chr>	<int>	<chr>	<chr>	<int>	<int>	<dbl>	<chr>	<dbl>	<dbl>	<int>	<chr>
1	Albania	1987	male	15-24 years	21	312900					Generation X
2	Albania	1987	male	35-54 years	16	302900	6.71	Albania1987	NA	2156624900	796
3	Albania	1987	female	15-24 years	14	289700	4.83	Albania1987	NA	2156624900	796
4	Albania	1987	male	75+ years	1	21800	4.58	Albania1987	NA	2156624900	796
5	Albania	1987	male	25-34 years	9	274300	3.28	Albania1987	NA	2156624900	796
6	Albania	1987	female	75+ years	1	35600	2.81	Albania1987	NA	2156624900	796
G.I Generation											

```
In [3]: apply(reworked.df, 2, function(x) sum(is.na(x))))

country: 0 year: 0 sex: 0 age: 0 suicides.no: 0 population: 0 suicides.100k.pop: 0 country.year: 0 hdi.for.year: 19450 gdp.for.year: 0 gdp.per.capita: 0 generation: 0
```

```
In [26]: #Prvo izmisliti sta su suikidi i stavu obezbediti podatci za iste godine
country.year <- aggregate(reworked.df$year, ) , list(country = reworked.dffcountry), unique)

head(country.year)
```

A data frame: 6 x 2											
country	year										
<chr>	<int>										
1	Albania										
2	Antigua and Barbuda										
3	Argentina										
4	Armenia										
5	Austria										
6	Australia										

Bilo je ista da ovaj dataset ima dosta nedostajajki vrednosti za varijablu hdi.for.year. Nema svaka država merenja za ovaku godinu već npr. Argentina ima merenja za 1986. dok Albanija ima merenja tek od 1987. Takođe postoji ponavljanje podataka kao što su hdi.for.year, gdp.for.year i gdp.per.capita koji se ponavljaju čak 12 puta (6 starosnih grupa za muškar i 6 starosnih grupa za ženski pol za dati godinu).

*Varijabla suicides.100k.pop predstavlja koeficijent varijabli suicides.no i population

**Iskazačka piramida:

1. Koje države imaju najveću stopu samoubistava?
2. Koje starosne grupe su nasklonije suicidu?
3. Da li postoji značajno veća stopa kod obojenog pola?
4. Da li u nekoj generaciji suicid zastupljeniji?
5. Od čega zavisi ta stopa?

```
In [5]: #1. Koje države imaju najveću stopu suicida
#2a odgovor na ovo pitanje koristeći godinu sa najviše opservacija
which.max(table(reworked.dffyear))
```

2009: 25

```
In [6]: year2009 <- as.data.frame(reworked.df[reworked.dffyear == 2009, ])

year2009.suicides.agg <- aggregate(year2009[, c(
  "suicides.no", "population")], list(country = year2009.country), sum)

year2009.suicides.agg$suicides.100k.pop <- year2009.suicides.agg$suicides.no / year2009.suicides.agg$population

year2009.suicides.agg <- year2009.suicides.agg[
  order(year2009.suicides.agg$suicides.100k.pop, decreasing = TRUE), ]

head(year2009.suicides.agg)
```

A data frame: 6 x 4											
country	suicides.no	population	suicides.100k.pop								
<chr>	<int>	<int>	<dbl>								
48	Lithuania	1128	3016497	0.000377358							
55	Republic of Armenia	15402	47390517	0.003250238							
11	Belarus	2743	9169969	0.0002991286							
67	Russian Federation	37408	134085433	0.0002789863							
77	Suriname	131	470448	0.00002784580							
44	Kazakhstan	3838	14370635	0.0002670724							

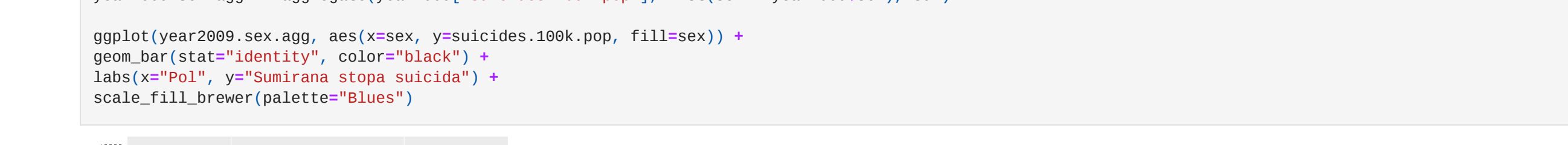
```
In [7]: #2. Koje starosne grupe su sklonije suicidu?
#Načinac se za sve države, za razliku od prelog primera neću računati nove vrednosti suicides.100k.pop,
#već ću sabrati postojeće. Takvo sumirana vrednost nije ista kao i ona
#koja se dobija suicides.no/population ali samratm da oslikava istu pojavu

year2009.age <- aggregate(year2009[,"suicides.100k.pop"], list(age=year2009.age), sum)

#install.packages("ggcolorbrewer")
```

```
library(ggplot2)
library(ggcolorbrewer)

ggplot(year2009.age.agg, aes(x=age, y=suicides.100k.pop, fill=age)) +
  geom_bar(stat="identity", color="black") +
  labs(x="Starosne grupe", y="Sumirana stopa suicida") +
  scale_fill_brewer(palette="blues")
```



```
In [8]: #3. Koja pol je skloniji suicidu?

year2009.sex.agg <- aggregate(year2009[,"suicides.100k.pop"], list(sex = year2009$sex), sum)

ggplot(year2009.sex.agg, aes(x=sex, y=suicides.100k.pop, fill=sex)) +
  geom_bar(stat="identity", color="black", position = "dodge") +
  labs(x="Pol", y="Sumirana stopa suicida") +
  scale_fill_brewer(palette="blues")
```



```
In [9]: sex.age.agg <- aggregate(year2009[,"suicides.100k.pop"], list(sex = year2009$sex, age = year2009$age), sum)

ggplot(sex.age.agg, aes(x = age, y = suicides.100k.pop, fill = sex)) +
  geom_bar(stat="identity", color="black", position = "dodge") +
  labs(x="Starosne grupe", y="Sumirana stopa suicida", title="Stopa suicida po starosnim grupama i polu")
  scale_fill_brewer(palette="blues")
```



```
In [10]: #4. U kojoj generaciji je suicid zastupljeniji?
#2a odgovor na ovo pitanje koristeći godinu sa najviše opservacija
which.max(table(reworked.dffyear))

generation
```

A data frame: 5 x 2											
generation	suicides.100k.pop										
<chr>	<dbl>										
Silent	6177.07										
Boomers	2506.18										
Generation X	1857.17										
Millenials	1429.18										
Generation Z	106.44										

```
In [11]: #5. Od čega zavisi stopa suicida?
#2a odgovor na ovo pitanje moram koristiti i varijablu sa ROI vrednost.
#Početna hipoteza je da stopa zavisi od ROI - ja i GDP - ja

hdi.df <- as.data.frame(reworked.df[is.na(reworked.dffhdi.for.year) == FALSE, ])
which.max(table(hdi.dffyear))
```

2010: 6

```
In [12]: #izdeta je godina 2010. zbog najvećeg broja ROI vrednosti

year2010 <- as.data.frame(hdi.df[hdi.dffyear == 2010, ])

head(year2010)
```

A data frame: 6 x 12											
country	year	sex	age	suicides.no	population	suicides.100k.pop	country.year	hdi.for.year	gdp.for.year	gdp.per.capita	generation
<chr>	<int>	<chr>	<chr>	<int>	<int>	<dbl>	<chr>	<dbl>	<chr>	<int>	<chr>
253	Albania	2010	male	55-74 years	20	241852	0.87	Albania2010	0.722	11926953259	4359
254	Albania	2010	male	35-54 years	20	371611	5.38	Albania2010	0.722	11926953259	4359
255	Albania	2010	male	25-34 years	9	179720	5.01	Albania2010	0.722	11926953259	4359
256	Albania	2010	male	75+ years	2	50767	3.84	Albania2010	0.722	11926953259	4359
257	Albania	2010	male	15-24 years	10	279588	3.58	Albania2010	0.722	11926953259	4359
258	Albania	2010	female	25-34 years	6	183579	3.27	Albania2010	0.722	11926953259	4359
Generation X											

```
In [13]: year2010.agg <- aggregate(
  year2010[, c("hdi.for.year", "gdp.for.year", "gdp.per.capita")], list(country = year2010$country), unique)

year2010.agg2 <- aggregate(year2010[, c("suicides.no", "population")], list(country = year2010$country), sum)
year2010.agg2$suicides.100k.pop <- year2010.agg2$suicides.no / year2010.agg2$population
year2010.agg[c("suicides.no", "population", "suicides.100k.pop")] <- year2010.agg2[
  c("suicides.no", "population", "suicides.100k.pop")]

head(year2010.agg)
```

A data frame: 6 x 7											
country	hdi.for.year	gdp.for.year	gdp.per.capita	suicides.no	population	suicides.100k.pop					
<chr>	<dbl>	<dbl>	<dbl>	<int>	<int>	<dbl>					
1	Albania	0.722	1192695e+10	4359	96	2738025	3.508740e-05				
2	Argentina	0.811	4.22677e+11	11273	2943	3757854	7.823015e-05				
3	Austria	0.721	9.12009e+10	3460	70	267025	2.77272e-05				
4	Australia	0.827	1.14203e+12	54887	2420	2084747	1.160808e-04				
5	Austria	0.879	3.918027e+11	49181	2764	7696421	1.580202e-04				
6	Bahamas	0.774	1.009578e+10	30239	10	333869	2.986174e-05				

```
In [14]: Shapiro.test(year2010.agg$suicides.100k.pop)
Shapiro.test(year2010.agg$hdi.for.year)

ggplot(year2010.agg, aes(hdi.for.year, suicides.100k.pop)) + geom_point()
cor.test(year2010.agg$hdi.for.year, year2010.agg$suicides.100k.pop, method="spearman")
```

Shapiro-Wilk normality test

data: year2010.agg\$suicides.100k.pop

W = 0.93355, p-value = 0.0802438

Shapiro-Wilk normality test

data: year2010.agg\$hdi.for.year

W = 0.96029, p-value = 0.0002197

Warning message in cor.test.default(year2010.agg\$hdi.for.year, year2010.agg\$suicides.100k.pop) :
"Cannot compute exact p-value with ties"

Spearman's rank correlation rho

data: year2010.agg\$hdi.for.year and year2010.agg\$suicides.100k.pop

S = 59529, p-value = 0.0081826

alternative hypothesis: true rho is not equal to 0

sample estimates:

rho

0.3972966



Varijabla suicides.100k.pop nije normalno raspoređena. Varijabla hdi.for.year nije normalno raspoređena. Korelacija između ove dve varijable statistički značajna ali je slabija.

```
In [15]: Shapiro.test(year2010.agg$gdp.per.capita)
cor.test(year2010.agg$gdp.per.capita, year2010.agg$suicides.100k.pop, method="spearman")

ggplot(year2010.agg, aes(gdp.per.capita, suicides.100k.pop)) + geom_point()
```

Shapiro-Wilk normality test

data: year2010.agg\$gdp.per.capita

W = 0.83498, p-value = 2.966e-08

Spearman's rank correlation rho

data: year2010.agg\$gdp.per.capita and year2010.agg\$suicides.100k.pop

S = 74468, p-value = 0.02394

alternative hypothesis: true rho is not equal to 0

sample estimates:

rho

0.2460538



Korelacija između ove dve varijable nije statistički značajna.

```
In [16]: #2a decision tree cu uzorke isti ovaj df jer sadrži najviševe hdi vrednosti
#install.packages(c("rpart", "rpart.plot"))

library(rpart)
library(rpart.plot)
```

train.data <- as.data.frame(year2010.agg)

third.q <- quantile(train.data\$suicides.100k.pop, 0.75)

train.data\$high.rate <- ifelse(train.data\$suicides.100k.pop > third.q, "yes", "no")

train.data\$high.rate <- as.factor(train.data\$high.rate)

train.data[c("country", "suicides.100k.pop")] <- NULL

head(train.data)

The scatter plot displays the relationship between 'suicides.1000.gd' (y-axis) and 'population' (x-axis). The y-axis ranges from 0 to 30, and the x-axis ranges from 0 to 100. The plot shows a positive correlation, with a dense cluster of points at lower population values (0-50) and a more dispersed set of points at higher population values (50-100).