

Uvod

Ovaj rad će koristeći podatke o stopama suicida u svetu po godinama na osnovu čega će se pokušati odgovoriti na određena istraživačka pitanja koja su definisana dalje u tekstu. Na početku, dataset će se nedelno razdeliti na test i treninge kako bi se videlo, potom će se preispitati njegovom strukturu na bi li se bolje upoznao sa podacima. Nakon toga sledi istraživačka pitanja i na samom kraju rada, kad bi se odgovori na pitanja, pristupaće se modelovanju podataka da bi se videlo da li možemo da predviđamo stopu suicida na osnovu dostupnih podataka u ovom datasetu.

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

read.csv("waster.csv", header = FALSE)

A data.frame: 6 x 12												
d.country	year	sex	age	suicides_no	population	suicides.100k.pop	country	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	6.71	Albania1987	NA	2.156.624.900	796	Generation X
2	Albania	1987	male	35-54 years	16	308000	5.19	Albania1987	NA	2.156.624.900	796	Silent
3	Albania	1987	female	15-24 years	14	289700	4.83	Albania1987	NA	2.156.624.900	796	Generation X
4	Albania	1987	male	75+ years	1	21800	4.59	Albania1987	NA	2.156.624.900	796	G.I. Generation
5	Albania	1987	male	25-34 years	9	274300	3.28	Albania1987	NA	2.156.624.900	796	Boomers
6	Albania	1987	female	75+ years	1	35600	2.81	Albania1987	NA	2.156.624.900	796	G.I. Generation

#Napravicu nov dataframe koji ce predstavljati sreden prethodni dataframe
reworked.df <- df

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

#Prvo obrisati suvisne tacke na krajevima naziva kolona
dots.delete <- function(name){
 return(sub("\\.", "", name))
}

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

#Sad je potrebno promeniti "." u "-" da bi imali iste notacije u kolonama
columnnotation <- function(column){
 return(gsub(".", "-", column))
}

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

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

#gdp.for.year 12 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)

#Potrebno napravitii faktore
reworked.df\$sex <- as.factor(reworked.df\$sex)
reworked.df\$age <- factor(reworked.df\$age, levels=c("15-24 years", "25-34 years", "35-54 years", "75+ years"), ordered = TRUE)

unique(reworked.df\$generation)

reworked.df\$generation <- factor(reworked.df\$generation, levels=c("G.I. Generation", "Silent", "Boomers", "Generation X", "Millennials", "Generation Z"))

head(reworked.df)

A data.frame: 6 x 12												
country	year	sex	age	suicides_no	population	suicides.100k.pop	country	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	6.71	Albania1987	NA	2156624900	796	Generation X
2	Albania	1987	male	35-54 years	16	308000	5.19	Albania1987	NA	2156624900	796	Silent
3	Albania	1987	female	15-24 years	14	289700	4.83	Albania1987	NA	2156624900	796	Generation X
4	Albania	1987	male	75+ years	1	21800	4.59	Albania1987	NA	2156624900	796	G.I. Generation
5	Albania	1987	male	25-34 years	9	274300	3.28	Albania1987	NA	2156624900	796	Boomers
6	Albania	1987	female	75+ years	1	35600	2.81	Albania1987	NA	2156624900	796	G.I. Generation

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: 19456 gdp.for.year: 0 gdp.per.capita: 0 generation: 0

#Proverimo da li su za svaku drzavu ohshebdeni podaci za iste godine
country.year <- aggregate(reworked.df["year"], list(country = reworked.df\$country), unique)

head(country.year)

country

A data.frame: 6 x 2

year

```
In [2]: #Izavršiti nov datframe koji će predstavljati sveden prethodni datframe
reworked.df <- df
```

#Potrebno je preimenovati kolone

```
names(reworked.df)[1] <- "country"
```

#Prvo obrisati starije tačke na krajevima naziva kolona

```
dots.delete <- function(name){
  return(sub("\\.", "", name))
}
```

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

#Sad je potrebno promeniti " " u " " da bi imali iste notacije u kolonama

```
columnnotation <- function(column){
  return(sub("-", "", column))
}
```

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

#Najveći "id" u nali "id" u nali

```
names(reworked.df) <- sapply(names(reworked.df), tolower)
```

#Prvo for year i2 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)
```

#Potrebno napraviti faktore

```
reworked.df$sex <- as.factor(reworked.df$sex)
reworked.df$age <- factor(reworked.df$age, levels=c(
  "15-14 years", "15-24 years", "25-34 years", "35-54 years", "55-74 years", "75+ years", ordered = TRUE))
unique(reworked.df$generation)
```

```
reworked.df$generation <- factor(reworked.df$generation, levels=c(
  "G.I. Generation", "Silent", "Boomers", "Generation X", "Millennials", "Generation Z"))
```

```
head(reworked.df)
```

3	Argentina	1985, 1986, 1987, 1988, 1989, 1990, 1991, 1992, 1993, 1994, 1995, 1996, 1997, 1998, 1999, 2000, 2001, 2002, 2003, 2004, 2005, 2006, 2007, 2008, 2009, 2010, 2011, 2012, 2013, 2014, 2015
4	Armenia	1990, 1991, 1992, 1993, 1994, 1995, 1996, 1997, 1998, 1999, 2000, 2001, 2002, 2003, 2006, 2007, 2008, 2009, 2010, 2011, 2012, 2013, 2014, 2015, 2016

```
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: 13456 gdp.for.year: 0 gdp.per.capita: 0 generation: 0

```
In [26]: #Proverimo da li su za svaku državu obezbeđeni podaci za iste godine
country.year <- aggregate(reworked.df["year"], list(country = reworked.df$country), unique)
```

```
head(country.year)
```

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

Istraživačka pitanja:

Bilo je ista da ovaj dataset ima dosta nedostajućih vrednosti za varijablu hdi.for.year. Nema svaka država merenja za svaku godinu već npr. Argentina merenja za 1985, dok Albanija ima merenja tek od 1987. Takođe postoji nedostajućih podataka kao što su hdi.for.year, gdp.for.year i gdp.per.capita koji se ponavljaju čak 12 puta (6 starijih grupa za muškarce i 6 starijih grupa za ženske pa za dati godinu).

*Varijable suicides.100k.pop predstavlja količinu varijabli suicides.no i population

Istraživačka pitanja:

- Koje države imaju najveću stopu samoubistva?
 - Koje starije grupe su nasleđene suicida?
 - Da li postoji značajno veća stopa kod određenog pola?
 - Da li je u nekoj generaciji suicid zastupljeniji?
 - Od čega zavisi ta stopa?
1. Koje države imaju najveću stopu suicida

Za odgovor na ovo pitanje koristiću godinu sa najveće opservacija. Razlog je taj što nemaju sve dve podatke za sve godine. Iako bi to uo ignorisali, dobili bi takve podatke da su neke više pozicionirane samo zbog toga što imaju više opservacija (godina). Zbog toga smatram da se ovo pitanje, kao i nekoliko sledećih, treba vezati za određenu godinu.

```
In [5]: which.max(table(reworked.df$year))
```

2009: 25

```
In [6]: year2009 <- as.data.frame(reworked.df[reworked.df$year == 2009, ])
year2009.suicides.no <- aggregate(year2009[, c(
  "suicides.no", "population")], list(country = year2009$country), sum)
```

year2009.suicides.aggsuicides.100k.pop <- year2009.suicides.aggsuicides.no / year2009.suicides.aggppopulation

year2009.suicides.aggsuicides.100k.pop <- year2009.suicides.aggsuicides.no / year2009.suicides.aggppopulation

order(year2009.suicides.aggsuicides.100k.pop, decreasing = TRUE,)

```
head(year2009.suicides.aggs)
```

A dataset: 6 x 4											
country		suicides.no	population	suicides.100k.pop							
<chr>	<chr>	<int>	<int>	<dbl>							
48	Lithuania	1138	3016497	0.0003772588							
65	Republic of Korea	15402	47380517	0.0003259238							
11	Belarus	2743	9169969	0.0002991296							
67	Russian Federation	51620	138058423	0.0002769863							
37	Guatemala	131	470448	0.0002707560							
44	Kazakhstan	3838	14370635	0.0002670724							

1. Koje starije grupe su sklonije suicidu?

Računice se za sve države. Za radiku od prošlog primera neću računati nove vrednosti suicides.100k.pop, već ću sabirati postojeće. Suma tumirana vrednost nije ista kao i ona koja se dobija kad suicides.no i population ali smatram da odlikava isti pojavu.

```
In [7]: year2009.age.aggs <- aggregate(year2009[, c("suicides.100k.pop")], list(age=year2009$age), sum)
```

```
#install.packages("RColorBrewer")
library(RColorBrewer)
```

ggplot(year2009.age.aggs, aes(x=age, y=suicides.100k.pop, fill=sex)) +

geom_bar(stat="identity", color="black") +

labs(x="Starosne grupe", y="Sumirana stopa suicida") +

scale_fill_brewer(palette="Blues")



1. Kojoj polji je skloniji suicidu?

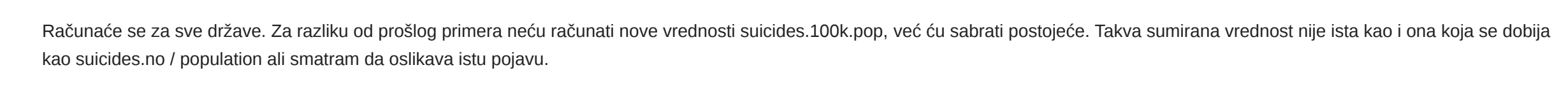
```
In [8]: year2009.sex.aggs <- aggregate(year2009[, "suicides.100k.pop"], list(sex = year2009$sex), sum)
```

ggplot(year2009.sex.aggs, aes(x=sex, y=suicides.100k.pop, fill=sex)) +

geom_bar(stat="identity", color="black") +

labs(x="Pol", y="Sumirana stopa suicida") +

scale_fill_brewer(palette="Blues")



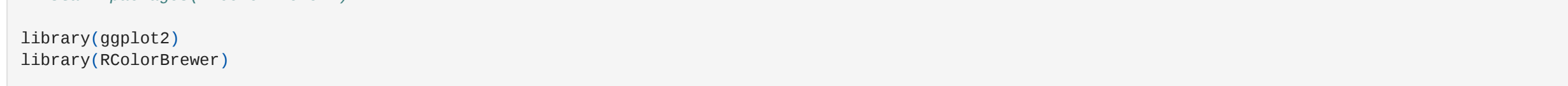
```
In [9]: sex.age.aggs <- aggregate(year2009[, "suicides.100k.pop"], list(sex = year2009$sex, age = year2009$age), sum)
```

ggplot(sex.age.aggs, aes(x = age, y = suicides.100k.pop, fill = sex)) +

geom_bar(stat="identity", color="black") +

scale_fill_brewer(palette="paired") +

labs(x = "Starosne grupe", y = "Sumirana stopa suicida", title="Stopa suicida po starosnim grupama i polu")



1. U kojoj generaciji je suicid zastupljeniji?

```
In [10]: generation <- aggregate(year2009[, "suicides.100k.pop"], list(generation=year2009$generation), sum)
```

```
#install.packages("ggplot2")
library(ggplot2)
```

ggplot(generation, aes(x=generation, y=suicides.100k.pop, fill=generation)) +

geom_bar(stat="identity", color="black") +

labs(x="Generacija", y="Sumirana stopa suicida") +

scale_fill_brewer(palette="Blues")

A dataset: 6 x 2											
country		suicides.100k.pop									
<chr>	<chr>	<dbl>									
	Silent	6177.07									
	Boomers	2506.18									
	Generation X	1957.17									
	Millennials	1429.38									
	Generation Z	106.44									

1. Od čega zavisi stopa suicida? Za odgovor na ovo pitanje moram koristiti i varijablu sa HDI vrednostima. Početna hipoteza je da stopa zavisi od HDI-a i GDP-a. Koristiće se godina sa najvećim brojem HDI opservacija.

```
In [11]: hdi.df <- as.data.frame(reworked.df[is.na(reworked.df$hdi.for.year) == FALSE, ])
which.max(table(hdi.df$year))
```

2010: 6

```
In [12]: year2010 <- as.data.frame(hdi.df[hdi.df$year == 2010, ])
head(year2010)
```

A dataset: 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>	<chr>	<int>	<int>	<dbl>	<chr>	<dbl>	<chr>	<int>	<chr>
253	Albania	2010	male	55-74 years	20	241852	8.27	Albania2010	0.722	11509953259	4359	Silent
254	Albania	2010	male	35-54 years	20	371611	5.38	Albania2010	0.722	11509953259	4359	Generation X
255	Albania	2010	male	25-34 years	9	179720	5.01	Albania2010	0.722	11509953259	4359	Generation X
256	Albania	2010	male	75+ years	2	50767	3.94	Albania2010	0.722	11509953259	4359	Silent
257	Albania	2010	male	15-24 years	10	276225	3.94	Albania2010	0.722	11509953259	4359	Generation X
258	Albania	2010	female	25-34 years	6	183579	3.27	Albania2010	0.722	11509953259	4359	Generation X

```
In [13]: year2010.aggs <- aggregate(year2010[, c("hdi.for.year", "gdp.for.year", "gdp.per.capita")], list(sex = year2010$sex, age = year2010$age), sum)
```

year2010.aggs2 <- aggregate(year2010[, c("suicides.no", "population")], list(country = year2010\$country), sum)

year2010.aggs2suicides.100k.pop <- year2010.aggs2suicides.no / year2010.aggs2population

year2010.aggs2suicides.100k.pop <- year2010.aggs2suicides.no / year2010.aggs2population

order(year2010.aggs2suicides.100k.pop, decreasing = TRUE,)

```
head(year2010.aggs2)
```

A dataset: 6 x 7											
country		hdi.for.year	gdp.for.year	gdp.per.capita	suicides.no	population	suicides.100k.pop				
<chr>	<chr>	<dbl>	<dbl>	<dbl>	<int>	<int>	<dbl>				
1	Albania	0.722	1.150995e+10	4359	96	2736025	3.508740e-05				
2	Argentina	0.811	4.236274e+11	11273	2943	37578454	7.831615e-05				
3	Armenia	0.721	9.260208e+09	3460	73	2676225	2.727723e-05				
4	Australia	0.927	1.144261e+12	54887	2420	20847547	1.160006e-04				
5	Austria	0.879	3.918927e+11	49181	1264	768421	1.586262e-04				
6	Bahamas	0.774	1.009576e+10	30239	10	333869	2.995317e-05				