

# Exploratory Data Analysis (for) Stroke Data Set

Pascal Visser, 410729

2022-09-14

## Contents

<b>1. Intro</b>	<b>2</b>
<b>2. Dataset</b>	<b>2</b>
<b>3 Dataset modifying</b>	<b>5</b>
<b>4. Exploratory Data analysis</b>	<b>7</b>
4.1 Summary . . . . .	7
4.2 Distribution of classifier . . . . .	8
4.3 Age and gender distribution . . . . .	9
4.4 health related variables . . . . .	9
4.4.1 Smoke . . . . .	9
4.4.2 BMI and Glucose levels . . . . .	10
4.4.3 Cardiovascular diseases . . . . .	12
4.5 Lifestyle variables . . . . .	13
4.5.1 Marriage status . . . . .	13
4.5.2 Type of work . . . . .	14
4.5.2 Residence type . . . . .	15
<b>5. Comparing variables</b>	<b>16</b>
<b>5.1 health</b>	<b>16</b>
5.2 lifestyle . . . . .	18
<b>6. EDA conclusion</b>	<b>20</b>

<b>7. Machine learning</b>	<b>20</b>
7.1 Dataset changing . . . . .	21
7.1.1 Methods . . . . .	21
7.1.2 SMOTE . . . . .	22
7.2 Testing algorithms . . . . .	23
7.3 Choosing Alogrithm . . . . .	25
7.3.1 Baseline . . . . .	25
7.3.2 Advanced algorithms . . . . .	26
7.3.3 Ensemble learners . . . . .	26
7.4 Parameter optimization . . . . .	27
7.5 ML Verdict . . . . .	31

```
# load libraries
library(tidyverse)
library(naniar)
library(readr)
library(ggplot2)
library(lemon)
library(knitr)
library(pander)
library(cowplot)
library(dplyr)
library(DMwR)
library(gridExtra)
library(farff)
```

## 1. Intro

A stroke is a life-threatening medical condition, where there is no or poor blood flow to (parts of) the brain. This causes parts of the brain to stop functioning and die off in minutes time. A stroke or brain attack is labelled as the second leading cause of death worldwide and is responsible for an annual mortality of about 5.5 million.

This dataset consists of people that had a stroke and people that have not. The variables are values that indicate their lifestyle, health and environment factors. The goal of this research is to calculate the chance of a stroke based of their lifestyle. This can be accomplished by machine learning. This gives the following research question:

*Is it possible to produce an accurate algorithm with machine learning, to calculate the risk of a stroke based on lifestyle variables of people with stroke history and people that never had a stroke?*

## 2. Dataset

The data is as said, lifestyle information of people with stroke history or not. In this section the data will be talked trough and the variables explained. First the data is loaded in:

```
# load in the data
strokedata <- read.csv("../Data/Stroke_dataset.csv", header = T)

# visualize
kable(strokedata[1:6, 1:6], format = 'simple', caption = "Head of data")
```

Table 1: Head of data

id	gender	age	hypertension	heart_disease	ever_married
9046	Male	67	0	1	Yes
51676	Female	61	0	0	Yes
31112	Male	80	0	1	Yes
60182	Female	49	0	0	Yes
1665	Female	79	1	0	Yes
56669	Male	81	0	0	Yes

```
cat("The dataset is", dim(strokedata)[1], "rows long, and has", dim(strokedata)[2], "columns")
```

```
## The dataset is 5110 rows long, and has 12 columns
```

Above, the first six rows of the data are shown. Also is given how many records there are. The codebook will clarify the variables and their values. The codebook will be loaded in:

```
# load codebook
codebook <- read.table("../Data/codebook.txt", header = T, sep = ';')
kable(codebook, format = 'pipe', caption = "Codebook of stroke data")
```

Table 2: Codebook of stroke data

Column	Description	value	datatype
id	unique identifier of patient	number	int
gender	male of female	sex	chr
age	age of patient	number	dbl
hypertension	suffering from hypertension?	binary	int
heart_disease	suffering from heart disease?	binary	int
ever_married	married or married in the past?	binary	chr
work_type	sort of work	string	chr
Residence_type	type of residence	string	chr
avg_glucose_level	average level of glucose	number	dbl
bmi	body mass index	number	chr
smoking_status	smoking history	string	chr
stroke	stroke history	binary	int

The codebook gives a detailed description the columns content about the value/datatypes. This is necessary to understand the data.

Now that the data is loaded in, there can be looked at the negatives of the data. The biggest flaws are the datatypes of the columns that are not equal. For example: the columns age is the type double instead of

integer. The column BMI is character, while a double datatype is much more logical. Also, the columns hypertension and hearth disease are a yes or no question, in these columns the yes or no is represented by 0 or 1. The columns of marriage history is a yes or no question too, but here yes or no is represented by text. This also needs modifying.

The character variables are can be converted to factor classes.

Additionally, it's good to know how many NA's are present in the dataset. These missing values can screw the data. Also, values like 'unknown' can be declared as a missing value

```
# check unique values of character values

cat("Gender:", unique(strokedata$gender), "\n",
    "Married:", unique(strokedata$ever_married), "\n",
    "Work type:", unique(strokedata$work_type), "\n",
    "Residence type:", unique(strokedata$Residence_type), "\n",
    "Smoking:", unique(strokedata$smoking_status))

## Gender: Male Female Other
## Married: Yes No
## Work type: Private Self-employed Govt_job children Never_worked
## Residence type: Urban Rural
## Smoking: formerly smoked never smoked smokes Unknown
```

Above reveals the Smoking status variable has a value "Unknown". This need to be addressed in the modifying section. Also, the column Gender has one other record. For the N/A the scan function of the 'naniar' package is used.

```
# scan for missing values
nas <- miss_scan_count(data = strokedata, search = list("N/A", "Unknown"))
kable(nas, format = 'pipe', caption = "Number of missing values per variable")
```

Table 3: Number of missing values per variable

Variable	n
id	0
gender	0
age	0
hypertension	0
heart_disease	0
ever_married	0
work_type	0
Residence_type	0
avg_glucose_level	0
bmi	201
smoking_status	1544
stroke	0

The column BMI is responsible for 201 missing values, these missing values are likely responsible for the column to be parsed as character datatype. the column smoking\_status has 1544 missing values. 1544 values is a big number and need to be dealt with in the modifying.

### 3 Dataset modifying

As talked in the previous section, some column types makes no sense, these types need to be changed to more sensible types. The columns that will be converted are:

- age, dbl -> int
- bmi, chr -> dbl

And additionally, the contents of the column ever\_married will be converted to binary, and the smoking\_status unknown will be cleared.

First the column age will be cast to integer, because a double type with decimals makes no sense in age, secondly, the column bmi will be cast to double, to keep the decimal precision:

```
# cast columns to another type
strokedata <- transform(strokedata, age = as.integer(age))
strokedata <- transform(strokedata, bmi = as.double(bmi))
```

Now the ever\_married column need to be changed to binary and a integer datatype:

```
# replace yes and no by 1 or 0
strokedata$ever_married[strokedata$ever_married == "Yes"] <- "1"
strokedata$ever_married[strokedata$ever_married == "No"] <- "0"
strokedata <- transform(strokedata, ever_married = as.integer(ever_married))
```

The character classes will be cast to factor:

```
# Cast the character columns to factor
strokedata$Residence_type <- as.factor(strokedata$Residence_type)
strokedata$gender <- as.factor(strokedata$gender)
strokedata$work_type <- as.factor(strokedata$work_type)
```

The record with gender ‘Other’ will be replaced with an gender value

```
# seek the most present value inside the column
table(strokedata$gender)
```

```
##
## Female    Male    Other
##   2994    2115         1
```

```
# Transform record where gender is equal to other
strokedata$gender[strokedata$gender == 'Other'] <- "Male"
```

It has been chosen to replace the other value with male, because of the balance inside the variable.

Now what is left are the missing values, BMI has N/A and smoking status has “unknown”. For BMI the Na’s will be replaced with the mean of the column:

```
# replace the missing values with the column mean
strokedata$bmi[is.na(strokedata$bmi)] <- mean(strokedata$bmi, na.rm = T)

# keep the digit precision as before
strokedata$bmi <- round(strokedata$bmi, digits = 1)
```

The smoking status 'unknown' is harder to replace because the variables are character, first we need to know how many records there are of the column:

```
# Count the unique variables in the gender column
table(strokedata$smoking_status)
```

```
##
## formerly smoked    never smoked      smokes      Unknown
##           885           1892           789           1544
```

There are 885, 1892 and 789 records of not unknowns. Based on these numbers, we can calculate the probability of occurring in the smoking status variable.

```
# Calculate the probability of formerly smoker, current smokers and non-smokers given that there's only

prob.Formerly <- 885 / (885 + 1892 + 789)
prob.Never <- 1892 / (885 + 1892 + 789)
prob.Smoke <- 789 / (885 + 1892 + 789)
```

With the probabilities, we can replace the unknowns based on their weight in the column.

```
# Replacing 'Unknown' in smoking_status by the other 3 variables according to their weightage
strokedata$rand <- runif(nrow(strokedata))

strokedata <- strokedata%>%mutate(Probability = ifelse(rand <= prob.Formerly, "formerly smoked", ifelse
strokedata <- strokedata%>%mutate(smoking.status = ifelse(smoking_status == "Unknown", Probability, smo

# View the new Smoking Status column's unique values and their counts
table(strokedata$smoking.status)
```

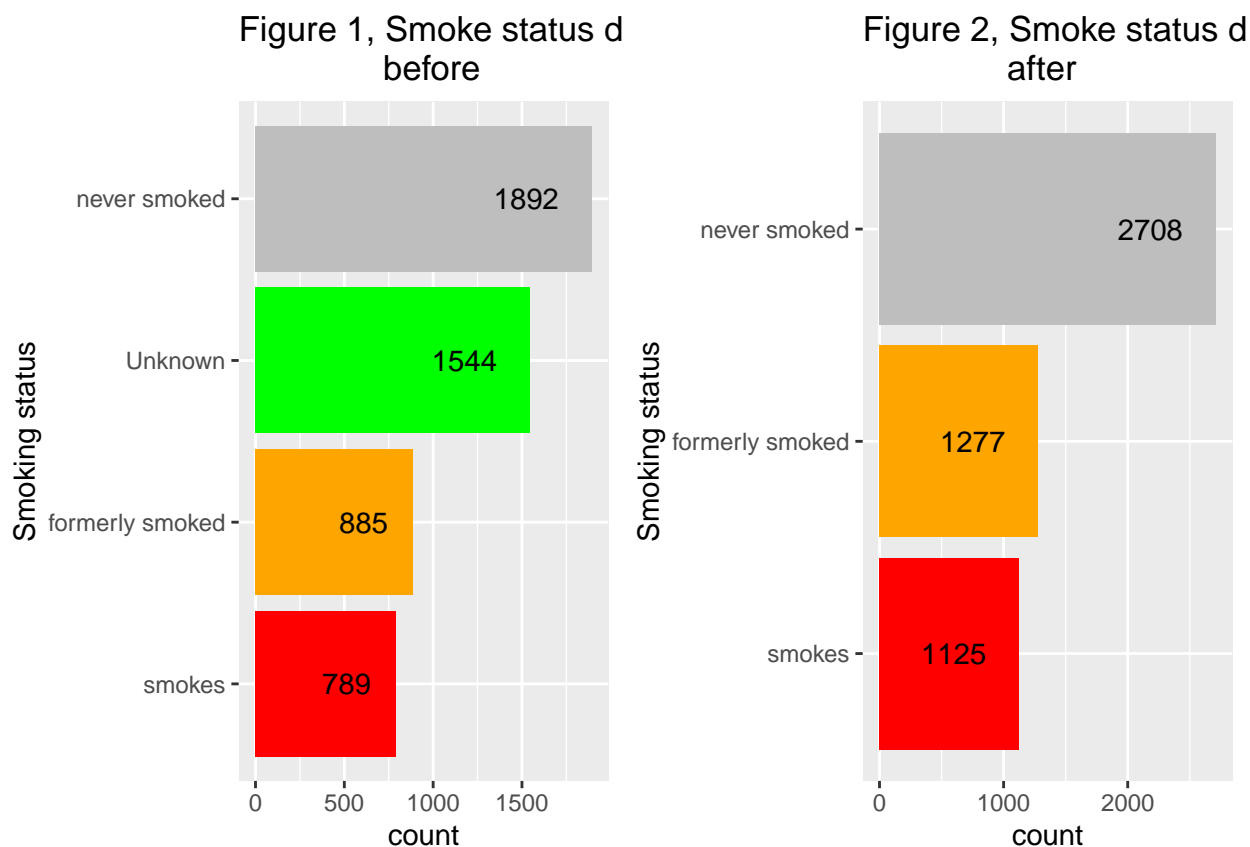
```
##
## formerly smoked    never smoked      smokes
##           1277           2708           1125
```

```
# Remove columns that are not needed
health <- subset(strokedata, select = -c(rand,Probability,smoking_status))

# revise the column name of smoking status
colnames(health)[12] <- "smoking_status"

# 'health' is the final modified dataset which will be used for the EDA section below.
```

Now we have a column without the unknowns, but with increased other variables based on their weight. Which is shown in the plots below:



In the plots, the effect is shown of replacing the unknown value by other values based on weight. With this method, the distribution of values is not affected.

## 4. Exploratory Data analysis

In this part, the modified data is visualised into useful graphs. These graphs will tell about the distribution of variables and the relations to each other.

### 4.1 Summary

First the summary of the values numeric values, we have three columns with important numeric values, age, glucose level and BMI. The summary of these variables tells about the distribution of values and possible outliers

```
#summary of subset data
pander(summary(health[c(3,9:10)],), caption = "Summary of numeric values")
```

Table 4: Summary of numeric values

age	avg_glucose_level	bmi
Min. : 0.00	Min. : 55.12	Min. :10.30
1st Qu.:25.00	1st Qu.: 77.25	1st Qu.:23.80

age	avg_glucose_level	bmi
Median :45.00	Median : 91.89	Median :28.40
Mean :43.22	Mean :106.15	Mean :28.89
3rd Qu.:61.00	3rd Qu.:114.09	3rd Qu.:32.80
Max. :82.00	Max. :271.74	Max. :97.60

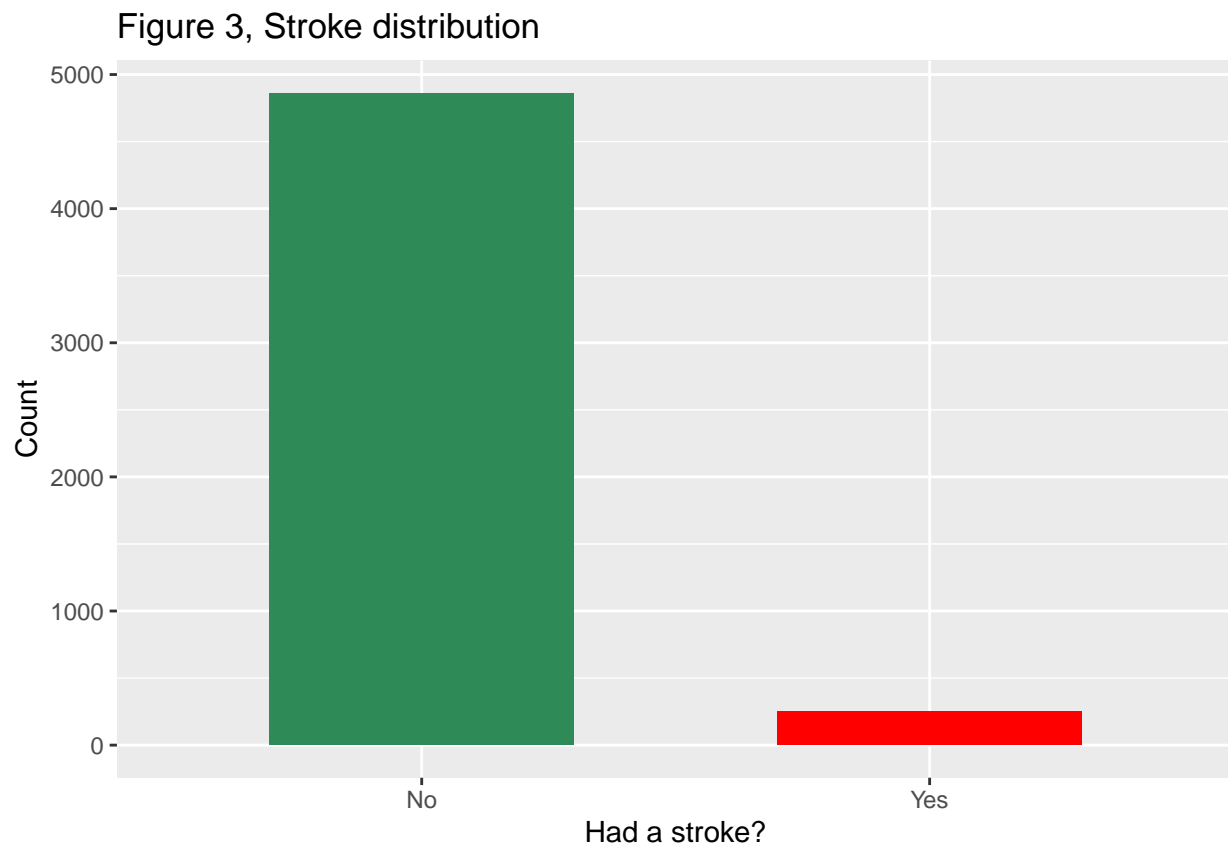
Looking at age, the minimum value is 0, which is a bit odd. Maybe this is an unknown value. Glucose looks normal, but BMI has a very low minimum, which is also odd.

## 4.2 Distribution of classifier

The research is about classifying patients chances of stroke based on there lifestyle and health. The figure below is the distribution of the stroke incidents in the data.

```
# Plot stroke distribution
fig1 <- ggplot(health, aes(x=factor(stroke))) +
  geom_bar(width = 0.6, fill = c("seagreen", "red")) +
  labs(title = "Figure 3, Stroke distribution", x = "Had a stroke?", y = "Count")

# plot the graph
fig1 + scale_x_discrete(breaks=c("0", "1"), labels=c("No", "Yes"))
```



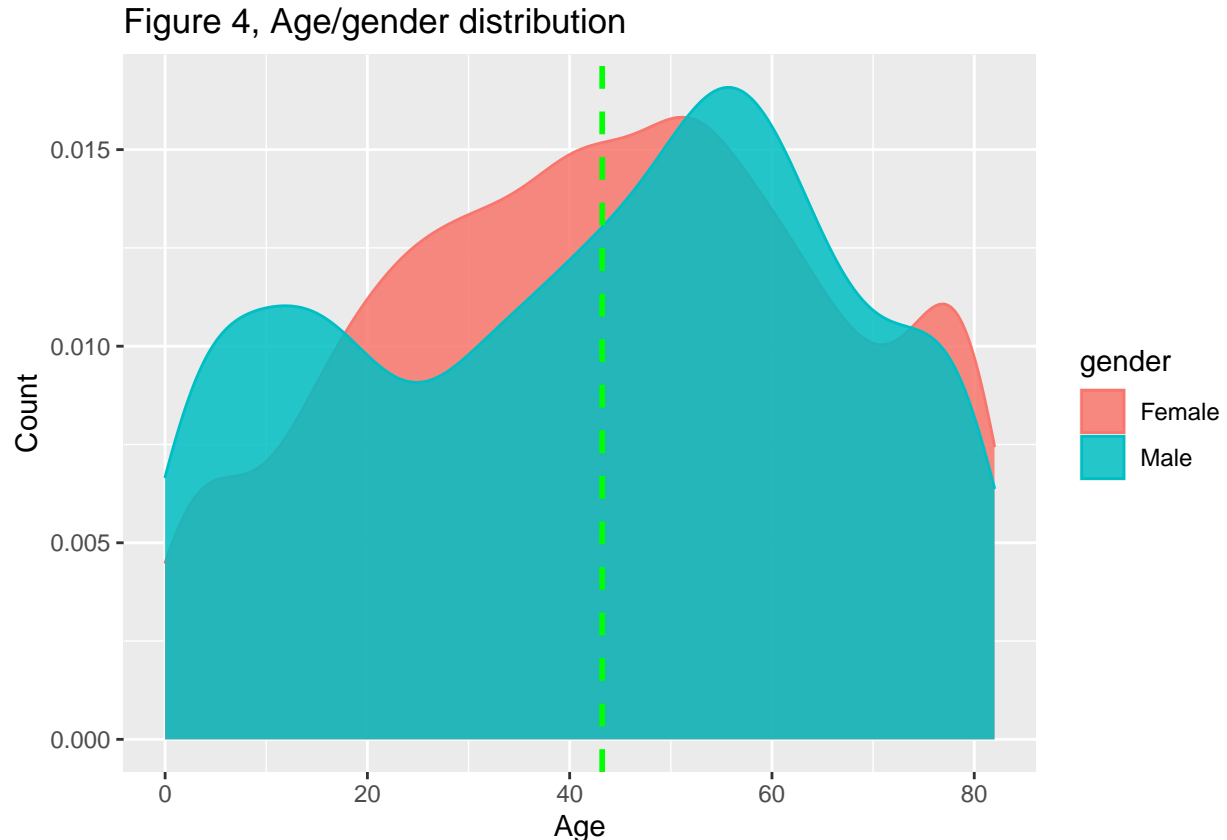
The number of stroke in the dataset is very low. The vast majority of the data didn't have a stroke. This needs to be balanced later.



### 4.3 Age and gender distribution

There are 5110 people in the dataset from which there is data, the gender and age distribution is as follows:

```
# make age and gender plot
ggplot(health, aes(x=age, fill=gender, color=gender)) +
  geom_density(alpha = .85) +
  labs(title = "Figure 4, Age/gender distribution", x = "Age", y = "Count") +
  geom_vline(aes(xintercept=mean(age)), color = 'green', lwd = 1, linetype = 'dashed')
```



Most people are around 40 - 60 years old, with some more females round the 20- 40 mark.

### 4.4 health related variables

This section tells about the health related variables such as, smoking habits, BMI, hypertension, heart disease and glucose levels

#### 4.4.1 Smoke

Smoking can affect your health in a bad way, therefore it is good to know how many people (formerly) are smoking.

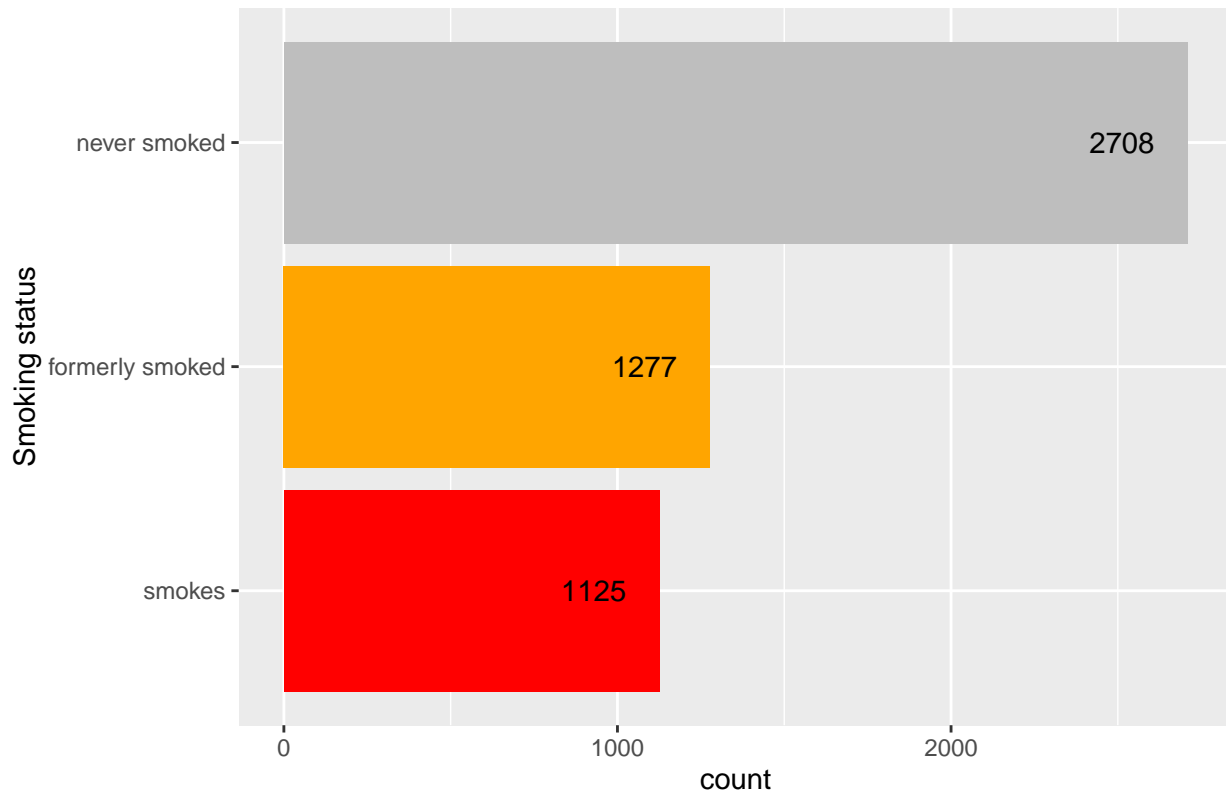
```
# Plot smoke distribution
health %>%
  group_by(smoking_status) %>%
```

```

summarise(count = length(smoking_status)) %>%
  mutate(smoking_status = factor(smoking_status)) %>%
  ggplot(aes(x = fct_reorder(smoking_status, count), y = count)) +
  geom_col(fill = c("Orange", "gray", "red")) +
  geom_text(aes(label = count, x = smoking_status, y = count), size = 4, hjust = 1.5) +
  coord_flip() +
  labs(x = "Smoking status", title = "Figure 5, Smoke status distribution")

```

Figure 5, Smoke status distribution



Most of the people never smoked, but a fairly amount of the people smokes or formerly smoked. This can indicate a bad overall health, which possible can contribute towards strokes. The term smokes is a bit too general, because it says nothing about how much someone smokes on a daily basis. The same applies to formerly smoking people.

#### 4.4.2 BMI and Glucose levels

Body mass index (BMI) and glucose levels are also a good indication of general health. High BMI and glucose levels are considered bad for health.

```

# Plot BMI distribution
fig2 <- ggplot(health, aes(x = bmi)) +
  geom_histogram(aes(y = ..density..), bins = 60, color = 1, fill = "gray") +
  geom_density(lwd = 1, linetype = 2, color = 'red') +
  labs(title = "Figure 6, Bmi distribution", x = "Body Mass Index", y = "Density") +
  geom_vline(aes(xintercept=mean(bmi)), color = 'green', lwd = 1) +
  geom_vline(xintercept = 35, color = 'blue', lwd = 1)

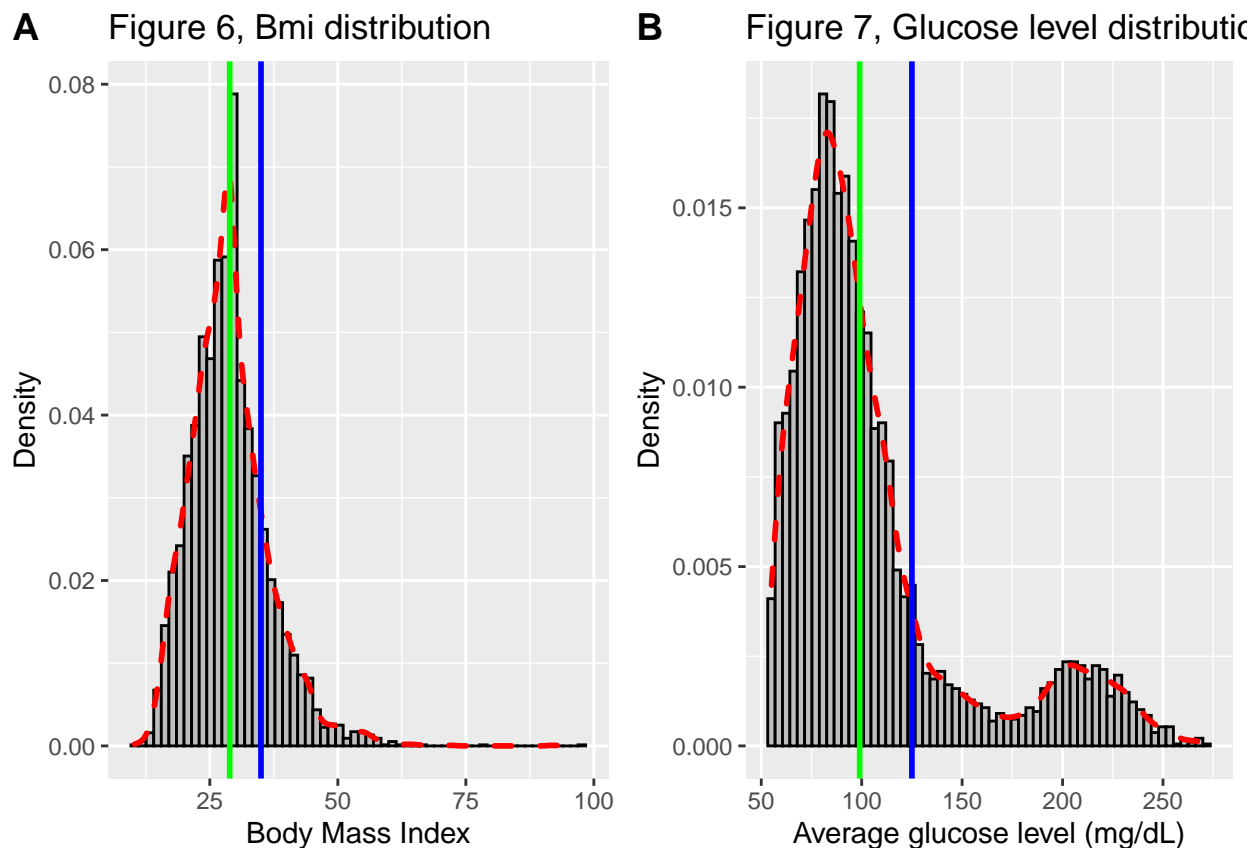
```

```

# Plot Glucose Distribution
fig3 <- ggplot(health, aes(x = avg_glucose_level)) +
  geom_histogram(aes(y = ..density..), bins = 60, color = 1, fill = "gray") +
  geom_density(lwd = 1, linetype = 2, color = 'red') +
  labs(title = "Figure 7, Glucose level distribution", x = "Average glucose level (mg/dL)", y = "Density") +
  geom_vline(xintercept = 99, color = 'green', lwd = 1) +
  geom_vline(xintercept = 125, color = 'blue', lwd = 1)

# Plot distributions side-by-side
plot_grid(fig2, fig3, labels = "AUTO", rel_widths = c(1,1))

```



BMI has a sort of normal distribution, with a mean of ~ 28. A BMI between 25 - 30 means overweight and 30+ means obese. There is a fairly amount of people with a BMI higher than 35 (the blue line). This means that a reasonable amount of people in the dataset is serious overweight.

For glucose level, the distribution is a little bit right skewed with a bulge toward the end of the x axis. The mean of 106 is high. 99 mg/dL or lower is normal, 100 to 125 mg/dL indicates you have pre-diabetes, and 126 mg/dL or higher indicates you have diabetes. The green line represents 99 mg/dL and the blue line 125 mg/dL. These borders show that, according to the above information. All the people left of the green line are healthy and all the people right of the blue line have diabetes. So, there are some serious unhealthy people according to these numbers.

```

# show number of people with high glucose and BMI
length(health$id[health$bmi > 35 & health$avg_glucose_level > 125])

```

```
## [1] 265
```

It tells that with the above information 265 people are serious unhealthy.

#### 4.4.3 Cardiovascular diseases

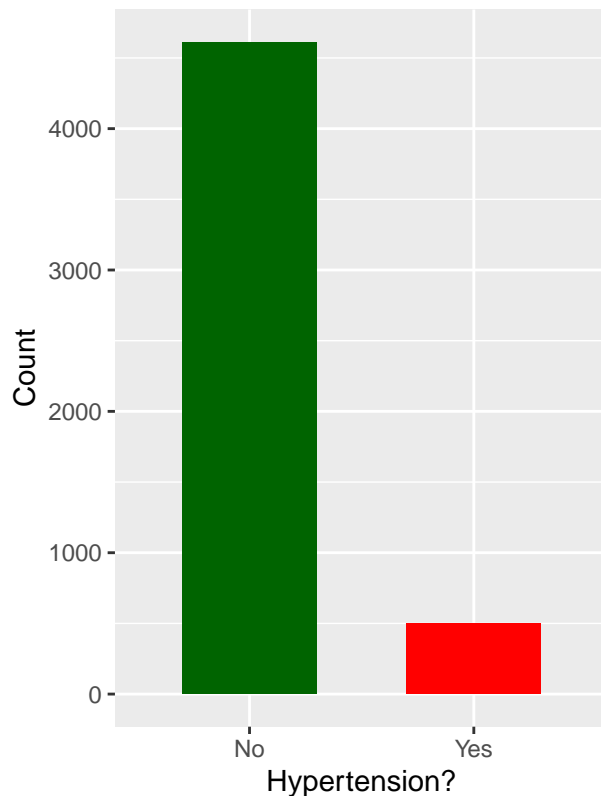
There are two cardiovascular diseases in the data, which are heart disease and hypertension.

```
# Show hypertension
fig4 <- ggplot(health, aes(x=factor(hypertension))) +
  geom_bar(width = 0.6, fill = c("darkgreen", "red")) +
  labs(title = "Figure 8, Hypertension", x = "Hypertension?", y = "Count") +
  scale_x_discrete(breaks=c("0", "1"), labels=c("No", "Yes"))

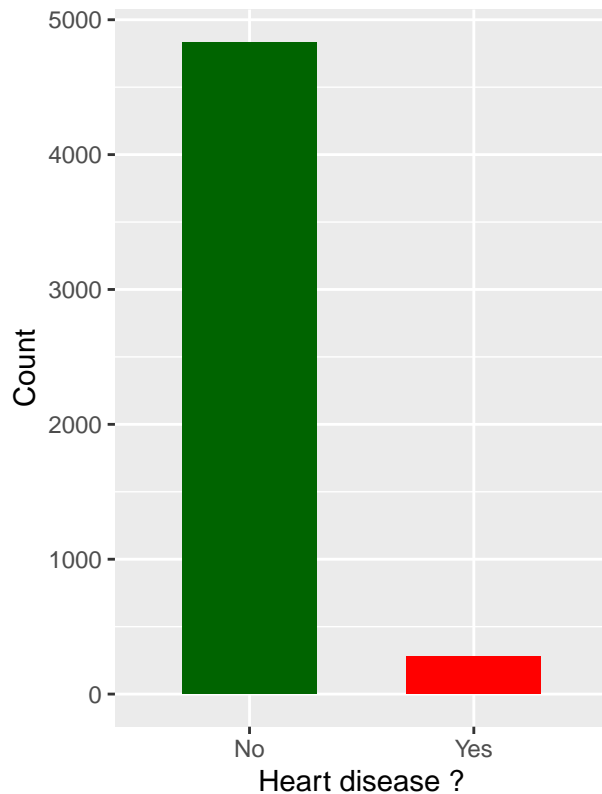
# Show heart_disease
fig5 <- ggplot(health, aes(x=factor(heart_disease))) +
  geom_bar(width = 0.6, fill = c("darkgreen", "red")) +
  labs(title = "Figure 9, Heart_disease", x = "Heart disease ?", y = "Count") +
  scale_x_discrete(breaks=c("0", "1"), labels=c("No", "Yes"))

# Plot distributions side-by-side
plot_grid(fig4, fig5, labels = "AUTO")
```

**A** Figure 8, Hypertension



**B** Figure 9, Heart\_disease



The cardiovascular diseases are both overwhelming no, this is also an imbalance problem.

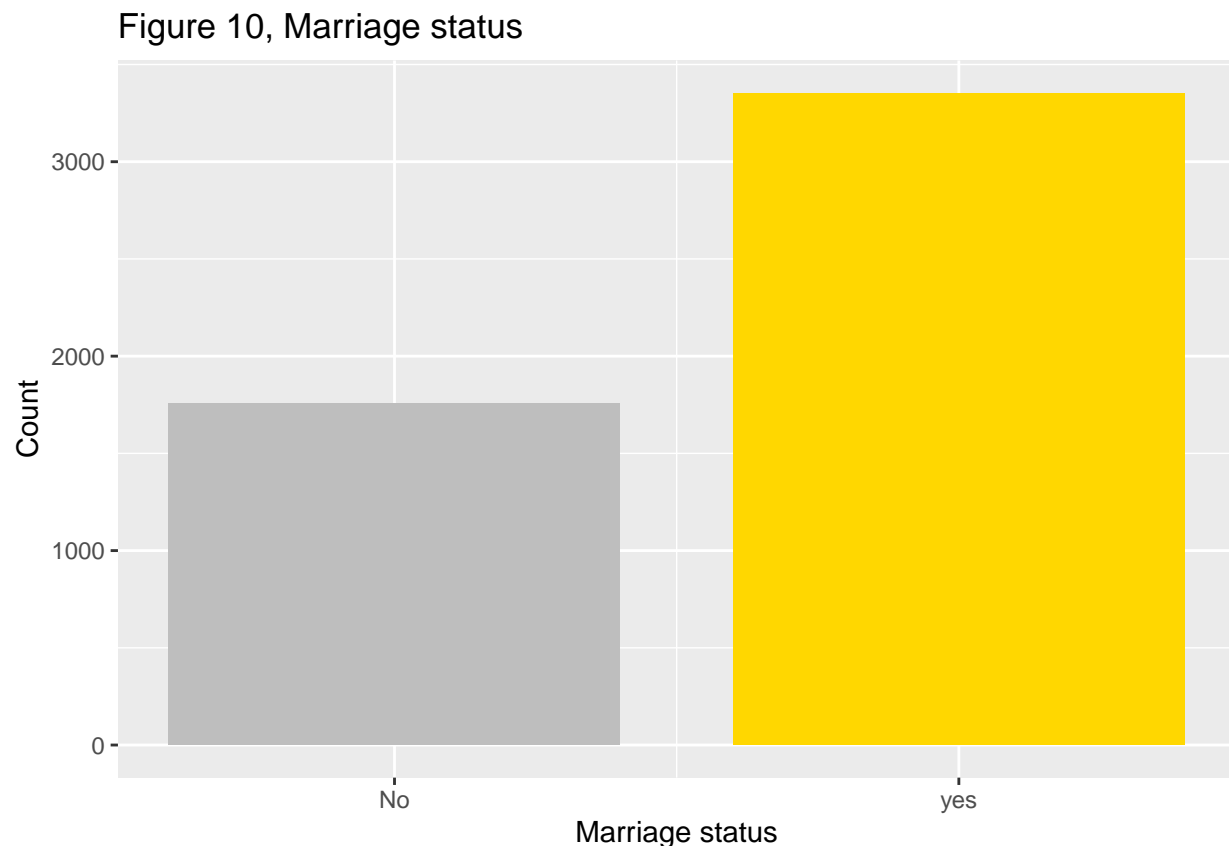
## 4.5 Lifestyle variables

There are three lifestyle variables: married?, work type and residence type. These variables tell about the person's possible stress level. All these factors can contribute toward stress and possible health affection.

### 4.5.1 Marriage status

The below plot shows the distribution of marriage. A married life could be stressful, even as a divorce. The 'yes' in ever married tells not the whole story, because a person could be married in the past, but is now divorced. These previously married people also are categorised into married. Which makes the 'yes' records a bit iffy.

```
# Plot marriage status
ggplot(health, aes(x = ever_married)) +
  geom_bar(width = 0.8, fill = c("gray", "gold")) +
  scale_x_continuous(breaks = c(0,1), labels = c("No", "yes")) +
  labs(x = "Marriage status", y = "Count", title = "Figure 10, Marriage status")
```



As is shown, most people have been married at some point in their lives. The not-married group is a sizeable group with people that have not been married. This maybe could be young people, which would be logic. The code below tests this hypothesis:

```
# Show the mean of married or not married
cat("Mean age of not married: ", round(mean(health$age[health$ever_married == "0"]), 2), '\n')
```

```
## Mean age of not married: 21.98
```

```
cat("Mean age of ever married: ", round(mean(health$age[health$ever_married == "1"]), 2))
```

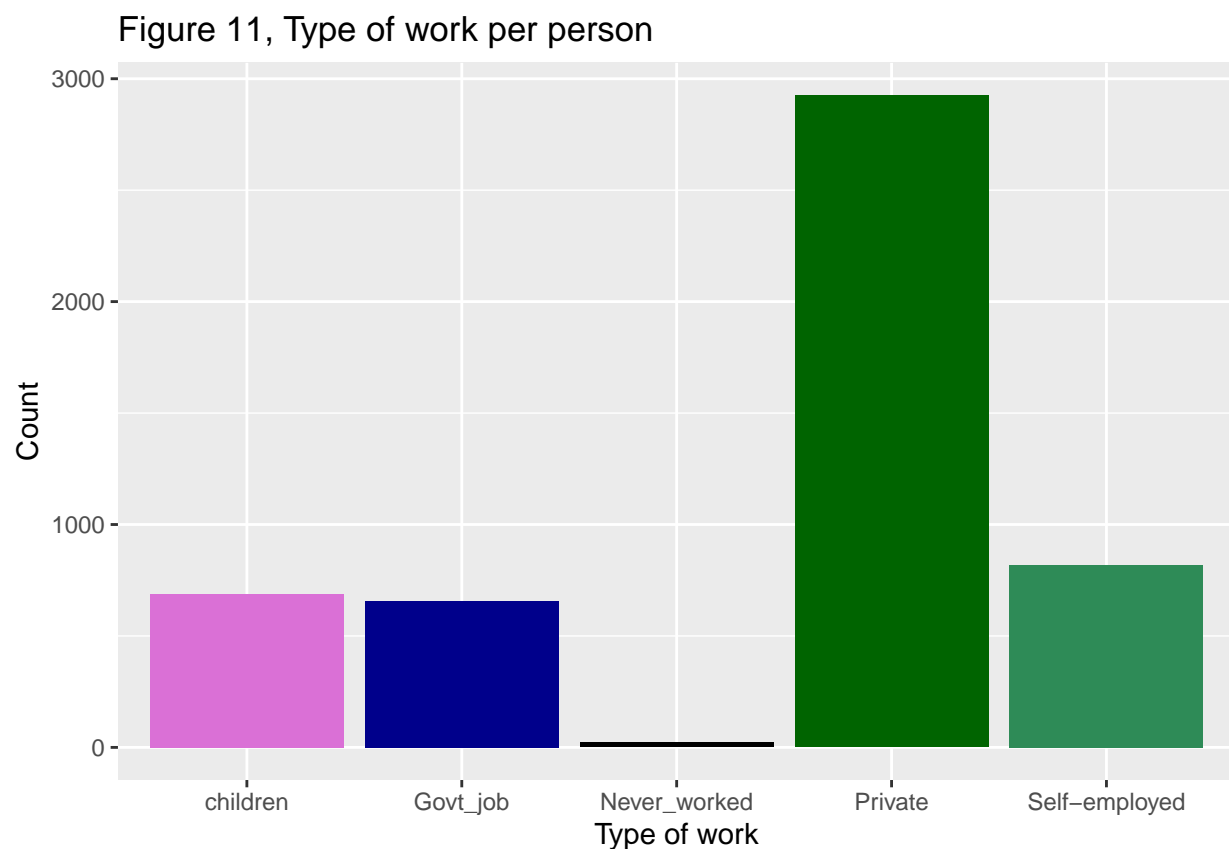
```
## Mean age of ever married: 54.34
```

As thought, the mean age is ~ 22 years old. Which means that this group mostly are young people.

#### 4.5.2 Type of work

The type of work someone performs can contribute toward higher stress levels. The plot beneath displays the distribution of work types:

```
# show distribution of work type
ggplot(health, aes(x = work_type)) +
  geom_bar(position = 'identity', fill = c("orchid", "darkblue", "black", "darkgreen", "seagreen")) +
  labs(x = "Type of work", y = "Count", title = "Figure 11, Type of work per person")
```



looking at the graph, most people work for a private company. Children, government job and self employed are somewhat balanced. Children would be mostly women is the guess.

```
# Show genders with children
table(health$gender[health$work_type == 'children'])
```

```
##
## Female   Male   Other
##    326    361     0
```

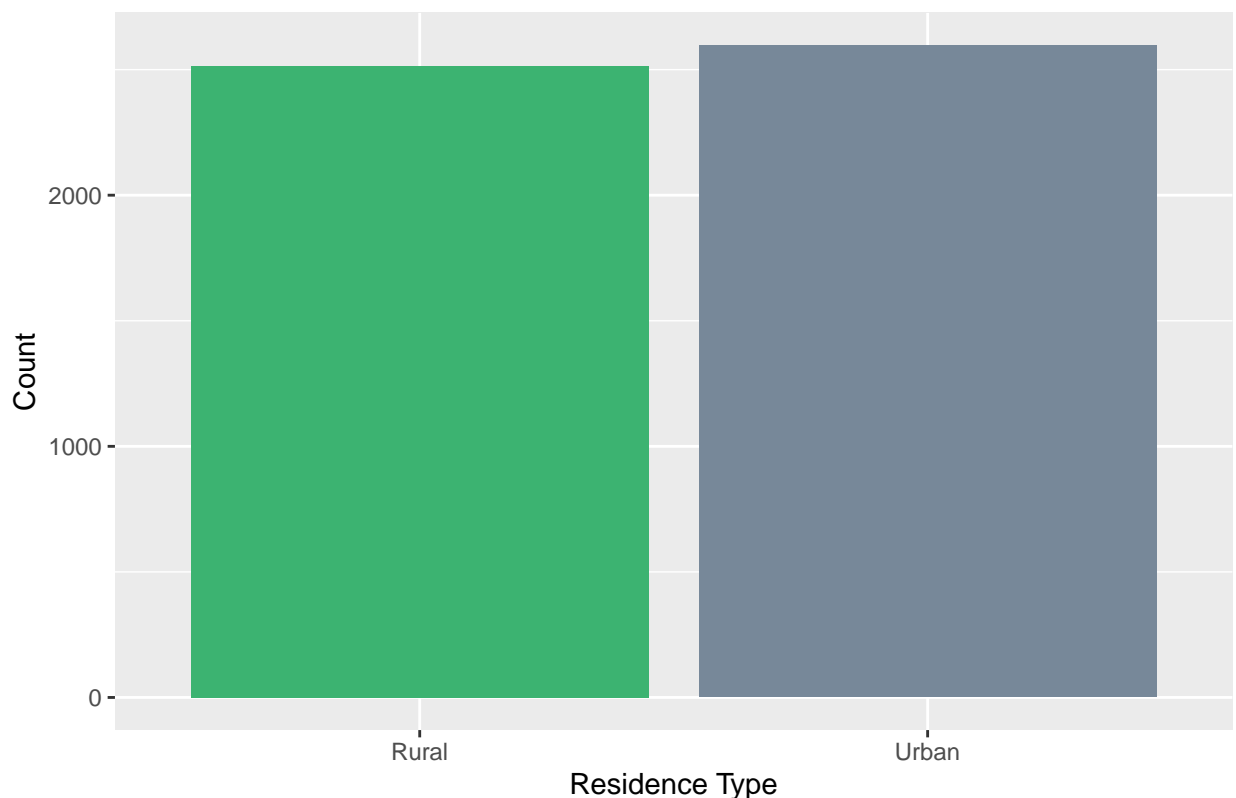
Surprisingly, children is balanced between male and female. so the thought of extra stress on women with children (which possibly can lead to higher stroke chances) can be rejected.

#### 4.5.2 Residence type

The type of living area can have influence on someone's life, living in a busy urban area maybe contributes to stress. The plot of residence type below shows the spread of living areas.

```
# Distribution of living area
ggplot(health, aes(x = Residence_type)) +
  geom_bar(position = 'identity', fill = c("mediumseagreen", "lightslategray")) +
  labs(x = "Residence Type", y = "Count", title = "Figure 12, Living area per person")
```

Figure 12, Living area per person



The bar plot of residence type shows a balance between the two types. There can be assumed that mostly young people live in urban areas and elderly people leave the busy areas.

```
# Show mean age of living area
cat("Mean age of urban: ", round(mean(health$age[health$Residence_type == "Urban"]), 2), '\n')
```

```
## Mean age of urban: 43.53
```

```
cat("Mean age of rural: ", round(mean(health$age[health$Residence_type == "Rural"]), 2))
```

```
## Mean age of rural: 42.89
```

The mean ages did not show any significant age difference between living rural or urban

## 5. Comparing variables

Now that all variables are shown, as standalone. It is time to compare variables with each other. Looking at the variables, there are two types, lifestyle and health. Obviously, health variables are the first thing to look at, because stroke is a cardiovascular disease. For this reason, health variables will be compared against stroke events.

### 5.1 health

In the dataset there are multiple people with heart disease and hypertension, could this combination be related towards stroke events? We take a look at hypertension at stroke and heart disease at stroke, also we look at heart disease and hypertension together.

```
# total cases of heart diseases and hypertension
cat(" Number of cases with stroke: ", sum(health$stroke),"\n","Number of cases with heart disease: ", sum(health$heart_disease),"\n","Number of cases with hypertension: ", sum(health$hypertension))

## Number of cases with stroke: 249
## Number of cases with heart disease: 276
## Number of cases with hypertension: 498
```

Inspecting the numbers above, there are more heart disease and hypertension cases than stroke events. So, not all heart disease and hypertension cases could be linked to stroke events.

So, how many stroke cases did also have underlying conditions such as heart disease and hypertension?

```
# stroke/heart disease
cat("The number of case with stroke and heart disease are: ",sum(health$stroke[health$heart_disease == 1]))

## The number of case with stroke and heart disease are: 47
```

```
# stroke/hypertension
cat("The number of case with stroke and hypertension are: ",sum(health$stroke[health$hypertension == 1]))

## The number of case with stroke and hypertension are: 66
```

```
# stroke/ heart disease and hypertension
cat("The number of case with stroke and heart disease & hypertension are: ",sum(health$stroke[health$heart_disease == 1 & health$hypertension == 1]))

## The number of case with stroke and heart disease & hypertension are: 13
```

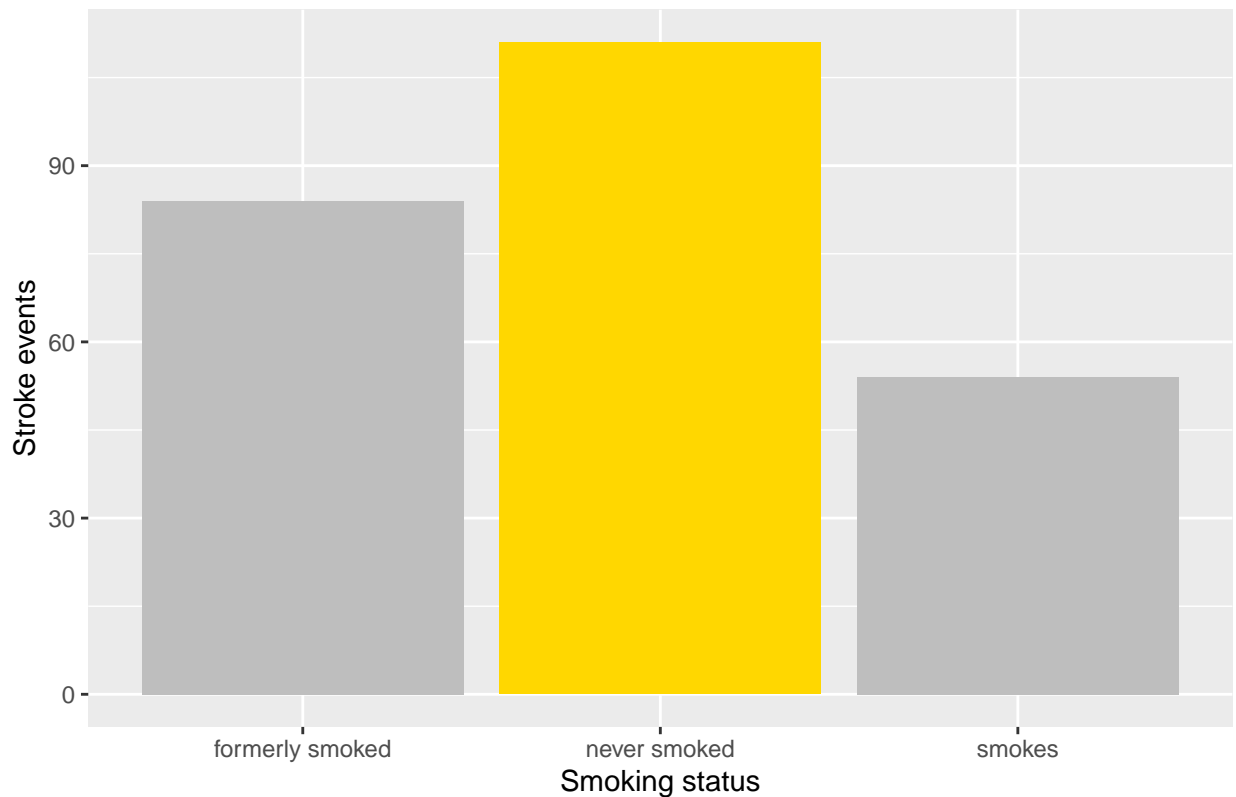
The relationship between stroke and heart disease & hypertension is not big. Only 13 cases of the 249 stroke events have two other cardiovascular diseases. This means that, according to this dataset, the relationship between hypertension/heart disease and stroke is not great.

Next, looking at smoke. There is an expectation that smoking is related to cardiovascular diseases. So, expected is that the most stroke events occurred with smokers or formerly smokers.

```
# plot smoking status with stroke events
ggplot(health[health$stroke == T,], aes(x=smoking_status)) + geom_bar(position = 'dodge', fill = c("gray", "red", "blue")) +
  labs(x = "Smoking status", y = "Stroke events", title = "Figure 13, Relation of Smoking and stroke events")
```



Figure 13, Relation of Smoking and stroke events

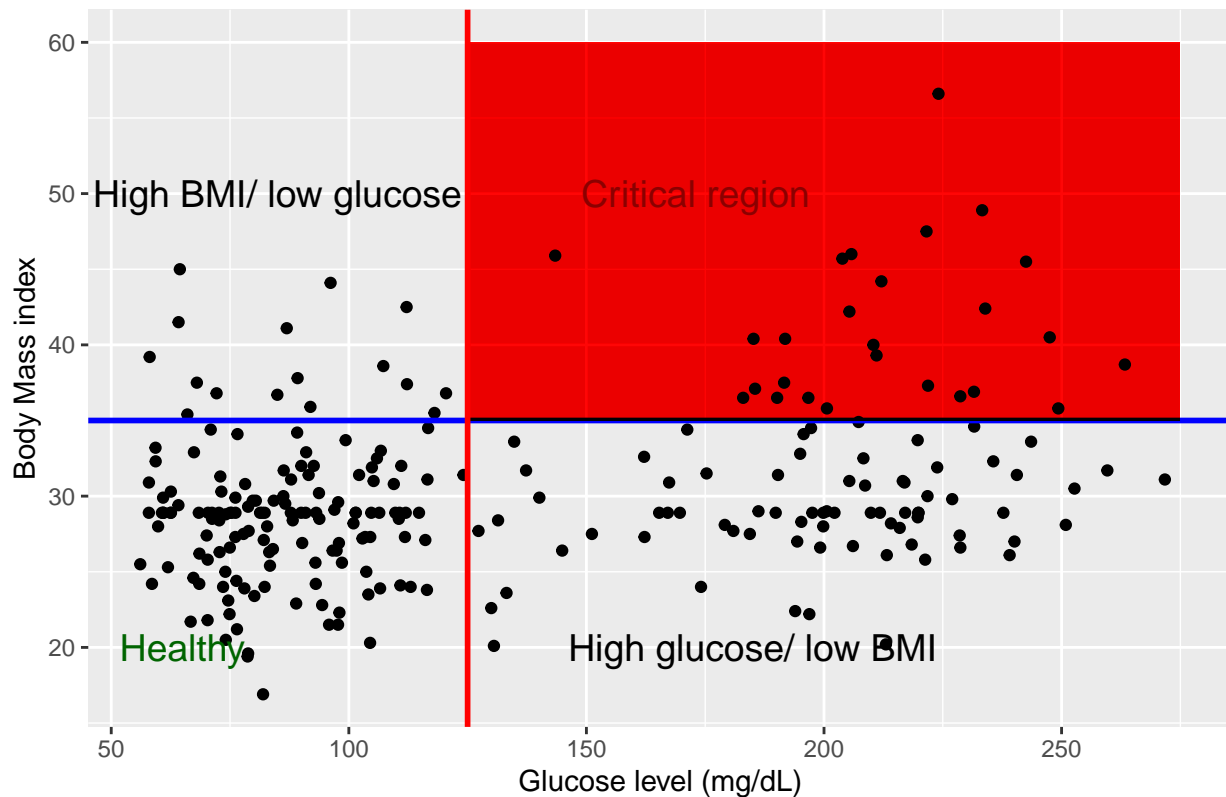


The graph above shows that the Group with the most stroke events are the non-smokers. Something that wasn't expected, but take into mind that the non smokers group is also the biggest group in the dataset. Outnumbering the other smoking groups by at least 1000 records.

high BMI and glucose levels can also contribute towards cardiovascular diseases. The graph beneath checks if it is true that high BMI and glucose levels increase the chance of a stroke.

```
# plot BMI/glucose of people with stroke history
ggplot(health[health$stroke == T,], aes(x = avg_glucose_level, y = bmi), ) + geom_point(show.legend = T) +
  labs(x = "Glucose level (mg/dL)", y = "Body Mass index", title = "Figure 14, Stroke events, BMI/glucose") +
  geom_hline(yintercept = 35, color = "blue", lwd = 1) +
  geom_vline(xintercept = 125, color = "red", lwd = 1) +
  annotate("text", x= 173, y = 50, label="Critical region", color= "darkred", lwd = 5) +
  annotate("text", x= 185, y = 20, label="High glucose/ low BMI", color= "black", lwd = 5) +
  annotate("text", x= 85, y = 50, label="High BMI/ low glucose", color= "black", lwd = 5) +
  annotate("text", x= 65, y = 20, label="Healthy", color= "darkgreen", lwd = 5) +
  geom_rect(aes(xmin= 125, xmax = 275, ymin = 35, ymax = 60), fill = "red", alpha = 0.002)
```

Figure 14, Stroke events, BMI/glucose level values



All the records in the graph above are people that had a stroke. Remembering the section about BMI en glucose levels, a BMI of 35 is considered to high and a glucose level of 125 mg/dl is also to high. The region in red is called the ‘critical region’. This region is where there are records of people with a BMI high than 35 and a glucose level high than 125 mg/dL. So, the critical region is where the very unhealthy people are. The expectation is that the most stroke events happens by these people.

But, the most of the records are outside of the critical region, most people reside in the ‘healthy’ region. So the assumption of the contribution of high BMI and high glucose levels towards strokes can be rejected based on this dataset.

## 5.2 lifestyle

The lifestyle variables tells about the living, working and marriage status of a person. A successfully stress free life would reduce the chances of stroke, won’t it? To test this hypothesis the lifestyle variables are compared against the classifier stroke.

The first look is at Work type, the type of work someone does

```
# probability table of working type, in percentages
round(prop.table(table(health$work_type[health$stroke == T])) * 100, 1)
```

```
##
##      children      Govt_job  Never_worked      Private Self-employed
##           0.8          13.3           0.0          59.8          26.1
```

The table overhead tells that the most stroke events happens by people that are working for a private company 59,8%. To recall the plot of work types earlier, the most people work at private companies, so it is not unexpected that the most stroke events happens by people that work for a private company.

Let's take a look a living area. The living area's are rural or urban. Earlier is said that the mean age of living is not that different. So, the assumption that young people, that mostly live in urban regions, decreases the percentage of strokes events in that living region can be rejected.

```
# probability table living area, in percentages
round(prop.table(table(health$Residence_type[health$stroke == T]))) * 100, 1)
```

```
##
## Rural Urban
## 45.8 54.2
```

Viewing the percentages above, the most stroke events happens at urban region. With not a big difference between the percentages. Which leave marriage status as the last variable.

The variable ever married indicate if a person is married of have been married at some point in their life. The percentages beneath shows the stroke events by married status.

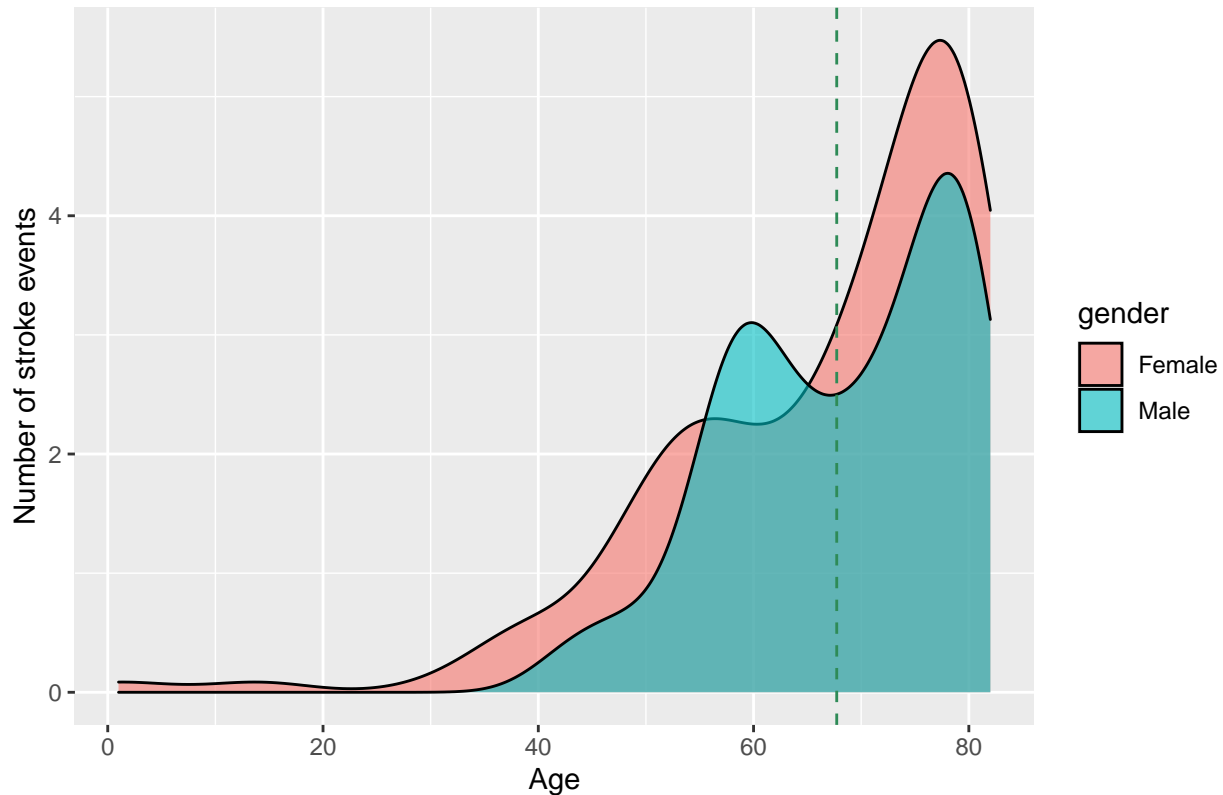
```
# Probability table of marriage status, in percentages
mar_prop <- round(prop.table(table(health$ever_married[health$stroke == T])) * 100, 1)
names(mar_prop) <- c("Not married", "Married")
mar_prop
```

```
## Not married    Married
##          11.6      88.4
```

Interesting is that The majority of the stroke events happens by people that are or have been married, 88,4%. But, take into mind that the most people that are married are older then the not married people. Speaking of age, what is the age where most stroke events take place?

```
# plot age and gender of people with stroke history
ggplot(health[health$stroke == T,], aes(x = age, fill = gender)) +
  geom_density(aes(y = stat(count)), alpha=0.6) +
  geom_vline(xintercept = mean(health$age[health$stroke == T]), linetype = 'dashed', color = 'seagreen')
labs(x = "Age", y = "Number of stroke events", title = "Figure 15, Age and gender in relation to stroke events")
```

Figure 15, Age and gender in relation to stroke events



In the distribution above is clearly seen that strokes mostly affects older people, which is expected. The mean age of stroke is 67 years old. It is also more frequent with female than with male.

## 6. EDA conclusion

The data is in a good condition with several variables to work with for machine learning. The removal of the Na's made the quality of the data better, mainly through the conversion of the term 'unknown' by smoking status. Also the casting of types helps a lot. The variables do not show in the EDA a very big difference in relation with stroke. But the variables age and marriage status shows some hope. But the data is still having one big flaw.

The plots show that the distribution of certain variables is skewed, namely: stroke, hypertension and heart\_disease. These variables have a very big imbalance inside. This imbalance can cause problems with machine learning because saying 'No' at these variables can give you a 95%+ accuracy. This imbalance needs to be tackled before moving on into machine learning. By doing this, the dataset can be workable for machine learning.

## 7. Machine learning

The goal in the machine learning project is to predict a chance of a stroke event based on the health and lifestyle variables. This part covers the cleaning of the data for machine learning, determining quality metrics and investigating the performance of machine learning algorithms.

For the machine learning (ML) a program named Weka will be used. Weka is a collection of machine learning algorithms for data mining tasks. It contains tools for data preparation, classification, regression, clustering, association rules mining, and visualization. In this analysis Weka version 3.9.6 is used.

In the section below the word accuracy is many times being said. In this experiment **accuracy means the precision of the ML algorithm in which it can rightly predict the chance of a stroke label**. So, an accuracy of 70% means that ,based on the inputted values of that person, there is a 70% chance that a persons stroke history is right. In other words if the output is true, you're likely to get a stroke in the near future.

## 7.1 Dataset changing

As talked over in the EDA, the dataset is okay but infected with imbalanced data. Namely the variables hypertension, heart disease and the classifier stroke. By leaving the imbalance in the data, the ML algorithm will bow towards the 'No', because this automatic will give a 95% accuracy. Any sort of algorithm will tent to high accuracy, if that high percentage can be achieved by taking the majority, it will. So, it is necessary to tackle this imbalance.

### 7.1.1 Methods

With imbalanced data sets, an algorithm doesn't get the necessary information about the minority class to make an accurate prediction. Which will lead to misleading accuracies.

There are a few ways to fix imbalanced dataset:

1. Undersampling
2. Oversampling
3. Synthetic Data Generation
4. Cost Sensitive Learning

*Undersampling* is reducing the number of observations from the majority class to achieve a balance. Undersampling can be split into two types: random and informative. Random undersampling means, taking random records from the majority and eliminate these to achieve a balance. Informative undersampling means eliminate records based on a pre-defined condition. undersampling seems good, but the danger lies in deleting (unwanted) important information for the algorithms. Not the best choice.

*Oversampling* is replicating records from the minority class. This can also be random or informative. the advantage with this method is that there is no information loss. But, with extra records from the minority, it can lead to overfitting.

*Synthetic Data Generation* is a method to, in stead of replication, add observations to the minority. It is a type of oversampling. These new observations are artificial generated. The prefect method for this is SMOTE (synthetic minority oversampling technique). The SMOTE algorithm creates artificial data based on feature space. similarities from minority samples. It creates a random set of minority class observations to shift the classifier learning bias toward the minority class. This method looks really good from this dataset.

*Cost Sensitive Learning* is the last technique. It puts weight on misclassified observations. There are penalties applied at false positives and false negatives. By doing this, the algorithm knows which classifications are 'heavier' to classify. It look good to use on the dataset.

Inspecting the four methods above for fixing imbalanced dataset. Synthetic Data Generation with SMOTE looks the best. This method generates new data for the minority class without losing important information or oversample the data with duplicate observations.

### 7.1.2 SMOTE

The SMOTE technique creates a new dataset by oversampling observations from the minority class. Which is excellent for this imbalanced dataset. The easiest way is to use the SMOTE() function from the 'DMwR' package. But, there is a problem with using SMOTE, there are categorical variables and numeric variables. SMOTE works best with numeric data. Luckily, the categorical variables are not big. The biggest is work\_type with four labels. The code below convert the variables to numeric, there are four variables that need changing:

```
# gender
health$gender <- as.character(health$gender)

health$gender[health$gender == "Male"] <- 1
health$gender[health$gender == "Female"] <- 0
names(health)[names(health) == "gender"] <- "is_Male?"

health$`is_Male?` <- as.integer(health$`is_Male?`)

#work_type
health$work_type <- as.character(health$work_type)

health$work_type[health$work_type == "Never_worked"] <- 0
health$work_type[health$work_type == "children"] <- 1
health$work_type[health$work_type == "Self-employed"] <- 2
health$work_type[health$work_type == "Govt_job"] <- 3
health$work_type[health$work_type == "Private"] <- 4

health$work_type <- as.integer(health$work_type)

#Residence_type
health$Residence_type <- as.character(health$Residence_type)

health$Residence_type[health$Residence_type == "Rural"] <- 0
health$Residence_type[health$Residence_type == "Urban"] <- 1
names(health)[names(health) == "Residence_type"] <- "living_urban?"

health$`living_urban?` <- as.integer(health$`living_urban?`)

#smoking_status
health$smoking_status <- as.character(health$smoking_status)

health$smoking_status[health$smoking_status == "never smoked"] <- 0
health$smoking_status[health$smoking_status == "formerly smoked"] <- 1
health$smoking_status[health$smoking_status == "smokes"] <- 2

health$smoking_status <- as.integer(health$smoking_status)

# relocate classifier
health <- health %>% relocate(stroke, .after = last_col())

head(health)

##      id is_Male? age hypertension heart_disease ever_married work_type
## 1  9046      1  67              0              1              1         4
```

## 2	51676	0	61	0	0	1	2
## 3	31112	1	80	0	1	1	4
## 4	60182	0	49	0	0	1	4
## 5	1665	0	79	1	0	1	2
## 6	56669	1	81	0	0	1	4

	living_urban?	avg_glucose_level	bmi	smoking_status	stroke
## 1	1	228.69	36.6	1	1
## 2	0	202.21	28.9	0	1
## 3	0	105.92	32.5	0	1
## 4	1	171.23	34.4	2	1
## 5	0	174.12	24.0	0	1
## 6	1	186.21	29.0	1	1

Now all the variables are numeric and the classifier is relocated toward the last column. The code below preforms the SMOTE balancing on the dataset. note that SMOTE tent to calculate decimal values of variables with 1 - 4. therefore all values are converted to factors to tackle this problem. Also the classifier is slightly modified to grow more towards SMOTE's preferences.

```
# changes certain variables to factor to prevent problems
health$smoking_status <- as.factor(health$smoking_status)
health$work_type <- as.factor(health$work_type)
health$`is_Male?` <- as.factor(health$`is_Male?`)
health$`living_urban?` <- as.factor(health$`living_urban?`)
health$ever_married <- as.factor(health$ever_married)
health$age <- as.factor(health$age)
health$avg_glucose_level <- as.factor(health$avg_glucose_level)
health$bmi <- as.factor(health$bmi)
health$hypertension <- as.factor(health$hypertension)
health$heart_disease <- as.factor(health$heart_disease)

# change classifier to factor with boolean labels
health$stroke <- as.numeric(health$stroke)
health$stroke[health$stroke == 0] <- "FALSE"
health$stroke[health$stroke == 1] <- "TRUE"
health$stroke <- as.factor(health$stroke)

# write pre-SMOTE data to csv
write.csv(health, "../Data/health.csv")

# preform SMOTE on data
x <- SMOTE(form = stroke ~ ., data = health, perc.over = 3000, perc.under = 150)

# Save SMOTE data to csv
write.csv(x = x, file = "../Data/SMOTEdata.csv")
```

The result are two dataset, one without SMOTE and one with SMOTE applied. The SMOTE data has a perc.over of 3000 and a perc.under of 150. Which means that the minority class is synthetic oversampled by 3000% and the majority class synthetic oversampled is by 150%. This creates a ratio of 60/40 instead of 95/5

## 7.2 Testing algorithms

We ended up with two datasets. One with the imbalance and one with the imbalanced SMOTE'd out. These two dataset will be tested on accuracy to determine which dataset is most suitable to proceed with. Also,

With a too high accuracy overfitting might occur. Overfitting is getting a very high accuracy(> 95%) from an algorithm which is too precise. For example, we speak of overfitting when a tree algorithm creates a tree of 457 nodes and gets a accuracy of 98,9%. It works excellent on the handed data, and only on that data. So, a extremely high accuracy on the data in not necessarily good.

Weka self offers a handful of filters and tools. One of this tools is ‘ClassBalancer’. It reweights the instances in the data so that each class has the same total weight. This form of balancing will be applied on the imbalanced data, Which makes the datasets to test a total of three. Normal, Normal with classbalancer and SMOTE modified data.

there will be eight algorithms tested:

- ZeroR
- OneR
- J48
- Random Forest
- NaiveBayes
- Simple Logistic
- Logistic
- Sequential minimal optimization(SMO)

The results are tabled below accuracies, all algorithms are executed on default options, and with 10-fold cross validation.

**All algorithm tests that are executed are with 10-fold cross validation**

```
# table different accuracies of different dataset
algo_accuracies <- read.table("../Data/accuracies.txt", sep = ",", header = T)
colnames(algo_accuracies) <- c("Type of algorithm", "Normal", "Normal (classBalancer)", "SMOTE")
kable(algo_accuracies, format = "simple", caption = "Accuracies of different algorithms")
```

Table 5: Accuracies of different algorithms

Type of algorithm	Normal	Normal (classBalancer)	SMOTE
ZeroR	95.1	49.8	59.2
OneR	95.1	-	97.7
J48	95.1	60.7	90.2
RandomForest	94.4	54.6	95.6
NaiveBayes	89.1	76.2	80.1
SimpleLogistic	95.1	75.1	81.7
Logistic	95.1	76.4	81.8
SMO	95.1	76.3	82.0

In the table is seen that the normal imbalanced dataset almost always has an accuracy of 95%. This is because the imbalance in the dataset is equal to 95/5. So, the algorithm deals with this by saying that almost everything is false. In 95% of the cases, that is correct. As earlier concluded, the dataset with the imbalance is not good.

Looking at the dataset with the ClassBalance filter. The accuracies are reasonable balanced, with the bayes and function algorithms getting +70% accuracies. But the way how the ClassBalance filter acts is questionable. To put it simply, records from the minority are duplicated to achieve a 50/50 ratio and random duplication is not desirable.

The SMOTE dataset has a couple of overfitted accuracies, but the logistics algorithms and bayes preforms well. The way that created the extra records is also more desirable that the ClassBalance filter. So, The SMOTE dataset is the one to go with in the next section of the ML phase.



## 7.3 Choosing Alogrithm

Choosing the right algorithm is a crucial aspect to proceed. Beforehand it is important to know, what makes algorithm ‘good’? What are the characteristics of a desirable algorithm. For this question some quality metrics are determined.

Example of different quality metrics:

- accuracy
- speed
- size of tree\*
- Confusion matrix

Obviously a high accuracy is desired by all means. Then the question is, how accurate is good enough? As base, a accuracy of 90% is considered good. with diseases like strokes, 80% is not good enough. Which means that one in five is classified wrong, which can have serious consequences. Also, How is this accuracy achieved? Look out for overfitting.

Speed of the algorithm is important by big usage. Implementing the model in a local environment, 1 to 2 minutes is okay. As it would not calculate more than 10 times a day. But in big scale usage, an fast algorithm is a demand.

In case of a tree, the size of a Decision tree is a important factor. With a large tree with many leave, overfitting could be possible. As the algorithm is to specific on the trainings set. A small, but accurate tree is desirable.

A confusion matrix, is a specific table that allows visualization of the performance of an algorithm. Each row of the matrix represents the instances in an actual class while each column represents the instances in a predicted class. This indicates also the rate of true positives, true negatives, false positives and false negatives. The errors in the matrix will weight heavy on the model performance.

In conclusion: accuracy, confusion matrix and size of tree will weight the heaviest in selecting a good algorithm.

### 7.3.1 Baseline

With a proper dataset and determined quality metrics,the ‘SMOTEdata.csv’ file is loaded into Weka. The first two rows are removed, the blank attribute and ID. In the edit viewer, the Attributes are quickly visual scanned, to see if all the labels says numeric,except the classifier.

In The test options, cross-validation is set to 10 Folds (default). First a baseline performance is set (ZeroR).

TP = True positive rate FP = False positive rate

```
# Show the performance of ZeroR in table
baseline <- read.csv("../Data/zeroR_weka.csv", sep = ";", header = TRUE)
colnames(baseline)[1] <- "Algorithm ->"
pander(baseline, caption = "ZeroR mesurement")
```

Table 6: ZeroR mesurement

Algorithm ->	ZeroR
Accuracy (%)	59,2
Speed (sec)	0,01
T TP (%)	0

Algorithm ->	ZeroR
T FP (%)	0
F TP (%)	100
F FP (%)	100

ZeroR is a very simple algorithm that looks only at the classifier, in this case there are 11205 instances False, and 7719 True. ZeroR takes the majority and uses that for its accuracy, in this case 59,2%. SO, it is known that the baseline is 59,2%. This is the Zero measurement.

### 7.3.2 Advanced algorithms

Now that the baseline is determined, it is time for testing more complex algorithms. Earlier for comparison between normal and SMOTE data a couple of algorithms already were tested on default settings. This time speed, TP and FP also added.

```
# Show the performance of advanced algorithms in table form
advanced_algo <- read.csv("../Data/advanced_weka.csv", sep = ";", header = T)
colnames(advanced_algo)[1] <- "Algorithms ->"
kable(advanced_algo, caption = "Advanced algorithms")
```

Table 7: Advanced algorithms

Algorithms ->	OneR	J48	RandomForest	Naive.Bayes	Simple.logistic	Logistic	SMO
Accuracy (%)	97,7	90,2	95,6	80,1	81,7	81,8	82
Speed (sec)	0.02	0.43	3.6	0.04	1.14	0.2	18.5
T TP (%)	100	85,2	93,1	77,4	77,1	77,5	78,3
T FP (%)	3,9	6,2	2,5	18	15	15,2	15,3
F TP (%)	96,1	93,8	97,5	82	85	84,8	84,7
F FP (%)	0	14,8	6,9	22,6	22,9	22,5	21,7
Size of tree	.	925	.	.	.	.	.

Looking at the table above, there are three algorithms with accuracy >90% and four algorithms with accuracies around the 80%. Most algorithms are fast, with the exception of SMO. Especially SMO is slow. The high scoring algorithms have a low FP rate by True and False. The lower scoring algorithms have a more divided spread of TP and FP rates.

On paper, there are three excellent algorithms, but they are likely infected with (a form of) overfitting. The size of the J48 tree says enough. Next is it logic to look into parameters to reduce the amount of overfitting on the algorithms.

### 7.3.3 Ensemble learners

While the advanced algorithms look good, there are also ensemble learners. The main purpose of these is using an ensemble model to use a group of weak learners and form a strong learner. Random forest is an example of an ensemble learner. It uses bagging (which is explained later) and feature randomness when building each individual tree to try to create an uncorrelated forest of trees whose prediction by committee is more accurate than that of any individual tree.

There are 3 ensemble learners that will be tested: Bagging, AdaBoostM1 and stacking.

- Bagging is a technique for reducing prediction variance by producing additional data for training from a dataset by combining repetitions with combinations to create multi-sets of the original data.
- AdaBoostM1 or adaptive boosting is a ensemble learner that combines weak learners and weights their output into a weighted sum. This sum is used to tweak the algorithms to create a accurate output.
- Stacking involves training a learning algorithm to combine the predictions of several other learning algorithms.

The algorithms that are used by stacking are: ZeroR, OneR, J48, RandomForest, NaiveBayes, SimpleLogistic, Logistic and SMO. Especially SMO makes the algorithm very slow

```
# Show the performance of ensemble algorithms in table form
ensemble <- read.csv("../Data/ensemble_weka.csv", sep = ";", header = TRUE)
colnames(ensemble)[1] <- "Algorithm ->"
pander(ensemble, caption = "Ensemble learners")
```

Table 8: Ensemble learners

Algorithm ->	Bagging	AdaBoostM1	Stacking
Accuracy (%)	91,5	81,5	59,2
Speed (sec)	1.09	0.32	215.84
T TP %	89,2	75,7	0
T FP %	6,9	14,4	0
F TP %	93,1	85,6	100
F FP%	10,8	24,3	100

Bagging produces a fast and high accuracy output, which seem promising. AdaBosstM1 creates are even faster algorithm, but it is not as high as bagging. Stacking is just awful. It took nearly 45 minutes to complete the cross validation and it produces 59,2%...

So, in this case bagging seems great.

## 7.4 Parameter optimization

With the above tests, there are three algorithms that at default are already good: J48, Random forest and bagging. In this section these algorithms are optimized by tweaking the parameters.

*J48*

There are a number of parameters to play with. Two parameters (to tackle the overfitting) are the most important: confidenceFactor and minNumObj. This determines how aggressive the pruning process will be. The higher this value, the more ‘confident’ you are that the data you are learning from is a good representation of all possible events, and therefore the less pruning that will occur. Smaller values induce more pruning. The minNumObj determines the minimum observations that are allowed at each leaf of the tree. How higher this number, how more ‘strict’ the leafs are.

So, With a low confidence factor and a low minimum number of objects the tree can go its own way. It can split on decimal values many times and possible overfit enormously, like at default. The graph below shows the accuracy in relation to a bigger number of objects.

```

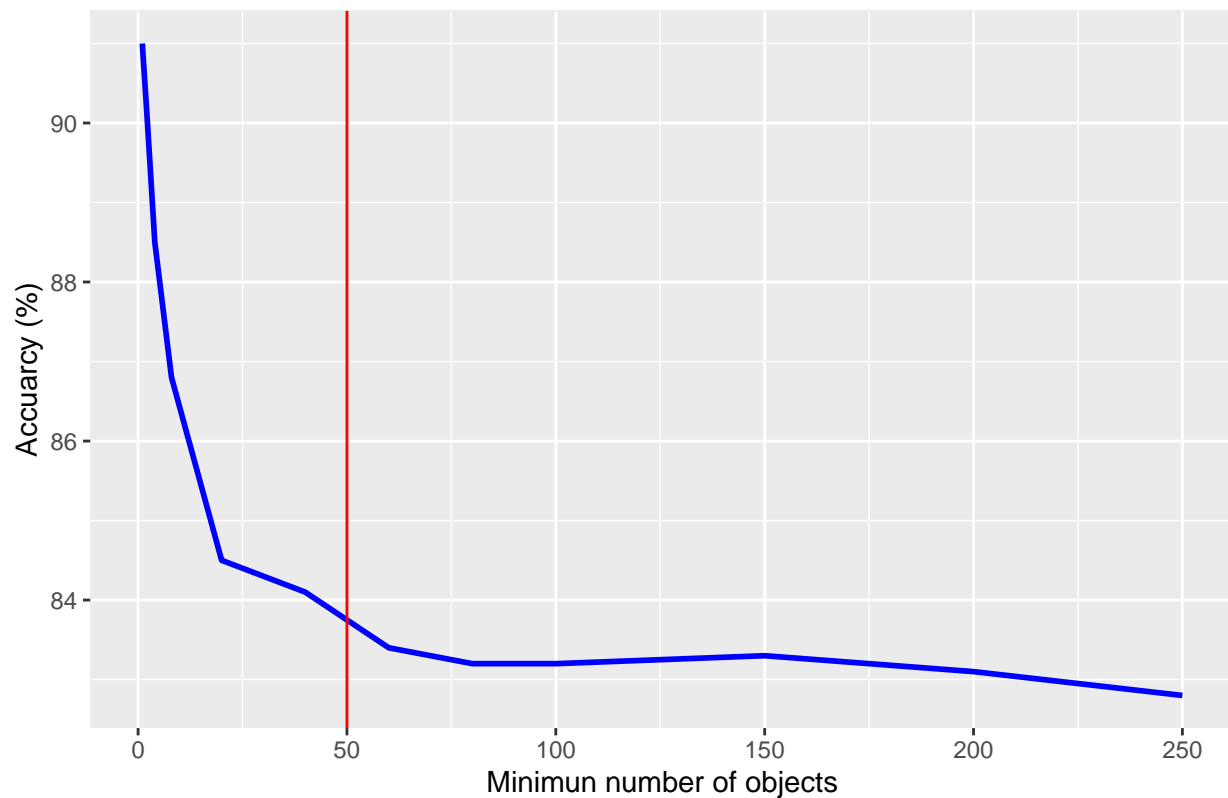
# Store data
accuracy <- c(91.0,90.2,88.5,86.8,84.5,84.1,83.4,83.2,83.2,83.3,83.1,82.8)
minNumObj <- c(1,2,4,8,20,40,60,80,100,150,200,250)

# Put data in dataframe
x <- as.data.frame(x = accuracy)
x$minNumObj <- minNumObj

# plot data
ggplot(x, aes(x = minNumObj, y = accuracy)) + geom_line(color = "blue", lwd = 1) + geom_vline(xintercept = 50)

```

Figure 16, Accuracy J48 in relation to minNumObj parameter



In the graph above, it is clearly seen that the minimal number of objects affects the accuracy in a way. The higher the minimum of objects, the lower the accuracy goes. But, after a minimum of 50 objects, the accuracy is more stable in relation to its minimum objects. So, 50 seems to be the sweet spot. Next, the confidence factor is tested in combination with minNumObj of 50.

```

# Store data in variables
confi_values <- c(0.01, 0.05, 0.10, 0.20, 0.30, 0.40, 0.50)
acc_values <- c(84.1, 84.3, 84.3, 84.1, 84.3, 84.3, 84.3)
size_of_tree <- c(35, 45, 45, 58, 58, 73, 77)

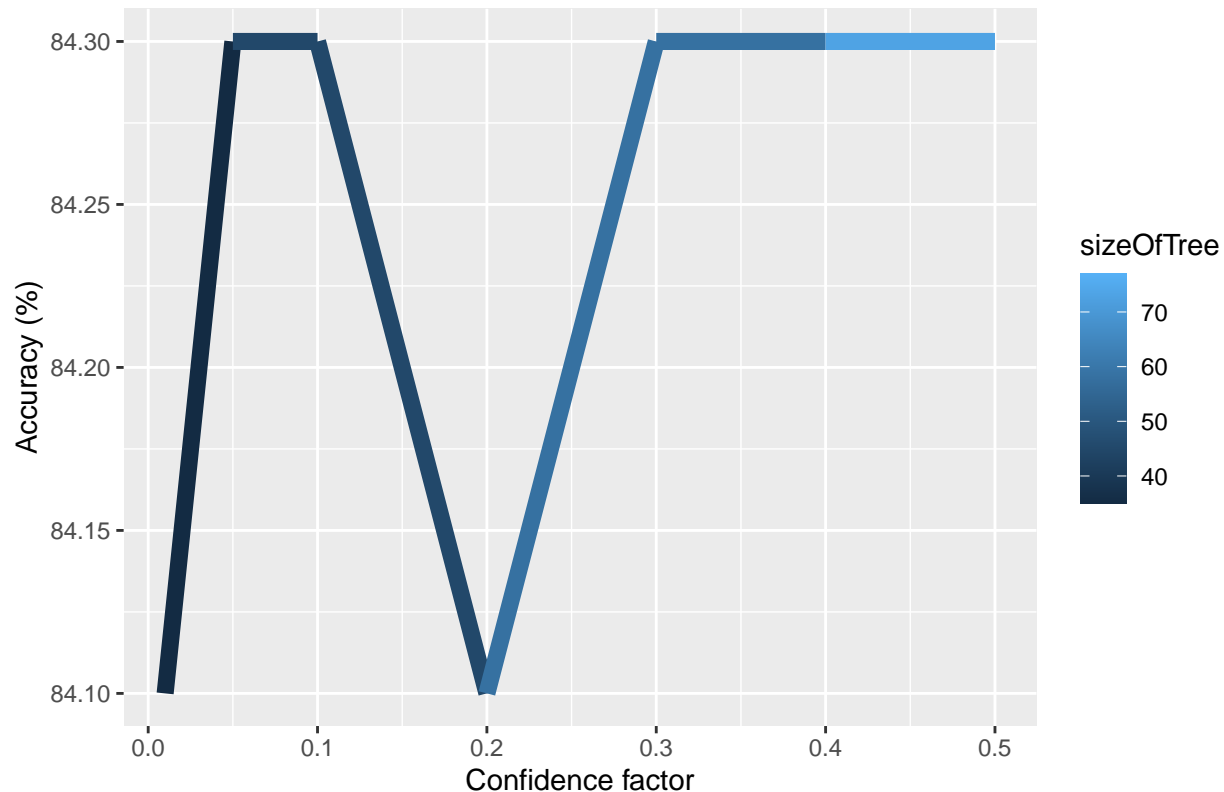
# put data in data frame
x <- as.data.frame(confi_values)
x$acc_values <- acc_values
x$sizeOfTree <- size_of_tree

```

```
# plot data
```

```
ggplot(x, aes(x = confi_values, y = acc_values)) + geom_line(aes(color=sizeOfTree), lwd = 3) + labs(y =
```

Figure 17, J48 confidence factor difference with minNumObj = 50



The line plot above tells that there is not many difference with a high or low confidence factor. Noticeable is that the size of the tree grows with a higher confidence factor. It is not desirable to have a big complex tree with no usable accuracy gain. So, the verdict is that an optimal tree has a minNumObj of 50, confidence factor of 0.01. by these parameter the size of the tree is 35. which is not to simple nor complex. Which gives an accuracy of **84.13%**

#### Random forest

The random forest gives a base accuracy of 95.6%, which is really high. But, is it high because it is good or because is it overfitted? To test this there are a couple of parameters to play with, maxDepth and numIterations. Especially maxDepth is important. This determines how 'deep' the tree may go. The default is 0, which means unlimited. This unlimited depth causes complex overfitted trees, which get high accuracies. So, by limiting its depth, the accuracy gets more reliable. the graph underneath, looks at the effect on accuracy with different maxDepth.

```
# Gather data
```

```
depth <- c(1,2,3,4,5,6,7,8,9,10,11,12,13,14,15,16,17,18,19,20)
```

```
accuracy <- c(78.7, 81.9, 82.5, 82.9, 83.3, 83.7, 84.6, 85.5, 86.8, 88.2, 89.7, 91.4, 92.6, 93.6, 94.5,
```

```
# enter data in dataframe
```

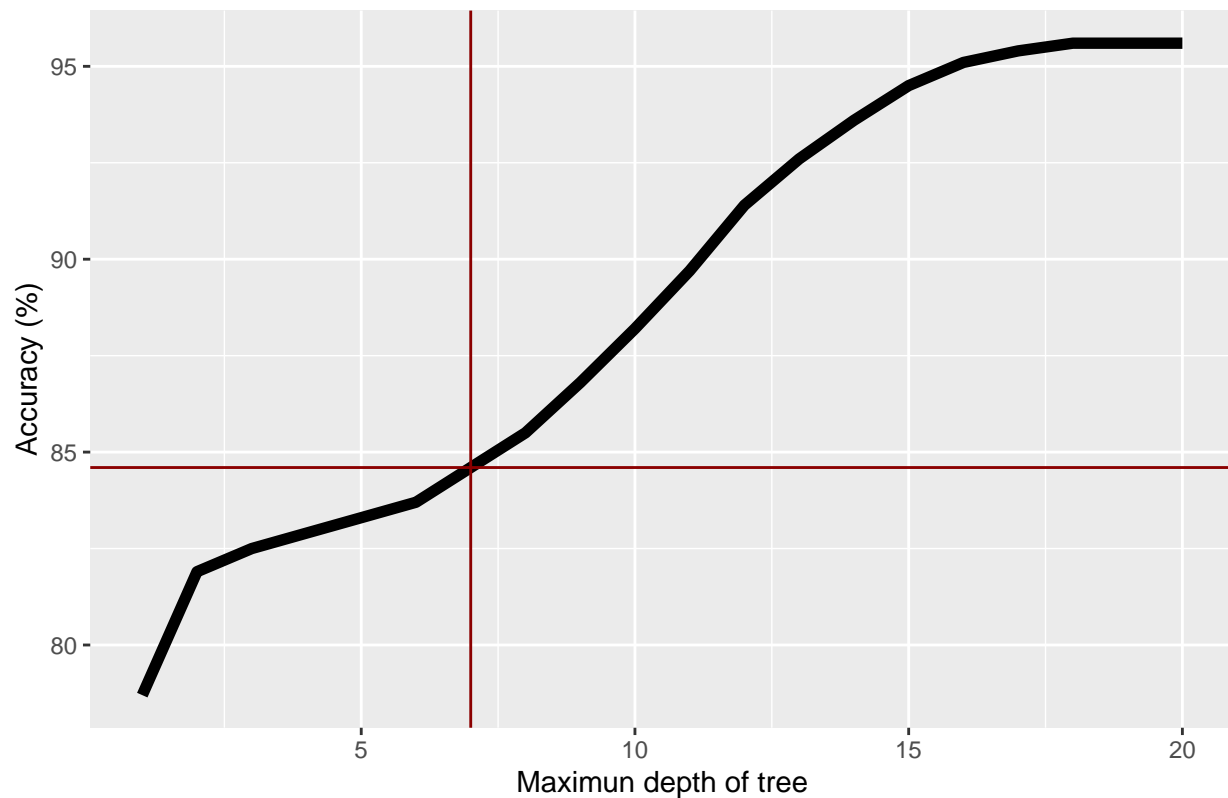
```
x <- as.data.frame(depth)
```

```
x$accuracy <- accuracy
```

```
# Plot data
```

```
ggplot(x, aes(x = depth, y = accuracy)) + geom_line(lwd = 2) + labs(x = "Maximun depth of tree", y = "A
```

Figure 18, RandomForest Effect of maxDepth parameter on Accuracy



The accuracy increases slowly with each higher depth level. After a depth of 15, it flattens out around the 95% mark. To tackle overfitting, a low depth is desirable, but a depth of 2 gives 'low' accuracies. but a too high depth creates overfitting problems. So, we take a depth of 7. Not too low and not too high. A depth of 7 gives a accuracy of **84.6%**

### Bagging

Bagging or **B**ootstrap **a**ggregating is a ensemble model to improve stability and accuracy of combined machine learning algorithms. It also reduces variance and helps avoid overfitting. The preventing of overfitting is great in this experiment. The Bagging model has a couple of parameter to tweak, but the most impacting is numIteration. The number of iterations the model preforms. The default is 10. Which earlier gave an accuracy of 91,5%. But more iteration means that it will be slower. So, the time to build a model will also be noted.

```
# Store data in variables
numOfIters <- c(1,10,20,30,40,50,60,70,80,90,100)
accuracy <- c(87.2, 91.4, 91.5, 92.0, 91.9, 92.0, 92.0, 92.0, 92.0, 92.0, 92.0)
timeModel <-c(0.08, 0.81, 1.61, 2.29, 3.04, 3.65, 4.46, 5.24, 6.09, 6.79, 7.51)

# Make dataframe for the data
x <- as.data.frame(numOfIters)
x$accuracy <- accuracy
x$timeModel <- timeModel

# make plot object for number of iterations
```

```

plot1 <- ggplot(x, aes(x = numOfIters, y = accuracy)) + geom_line(lwd = 1.5) + labs(x = "Number of Iterations", y = "Accuracy (%)")
#make plot object for time
plot2 <- ggplot(x, aes(x = numOfIters, y = timeModel)) + geom_line(lwd = 1.5) + labs(x = "Number of Iterations", y = "Speed (seconds)")
# plot side-by-side
grid.arrange(plot1, plot2, ncol=2)

```

Figure 19, Accuracy vs number of iterations

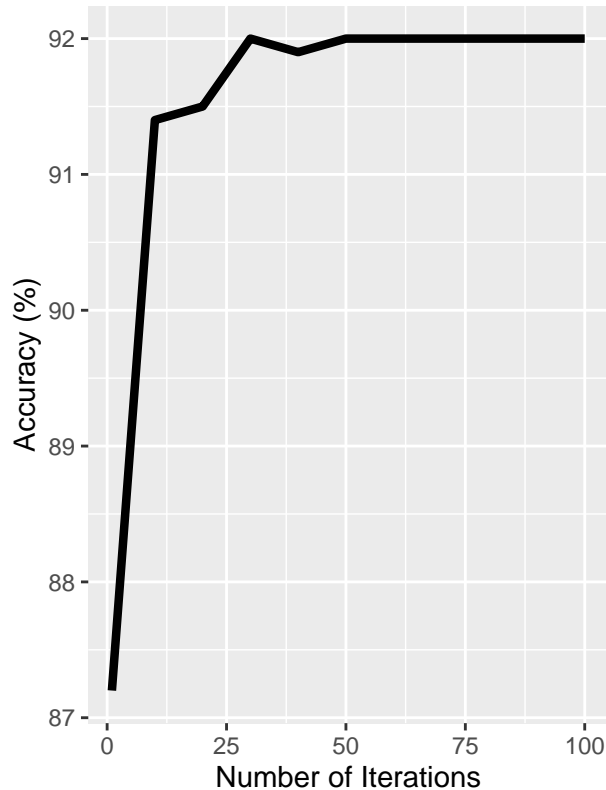
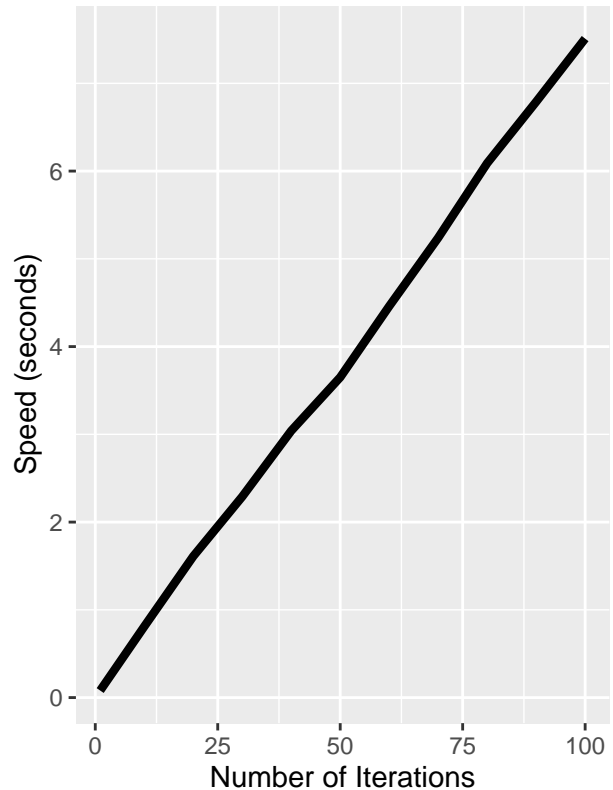


Figure 20, Time vs number of iterations



Looking at the graphs above, at +- 30 iterations it flattens out. There is no more gain after that. The down side is that, the more iteration the longer it takes to finish. Remember that after the model build it executes a 10-fold cross validation, which takes time.

The optimal number of iterations is around 30, where it peaks in accuracy and only takes 2.29 seconds to build a model. Fast and accurate. Which gives a accuracy of **92.0%**

## 7.5 ML Verdict

After parameter optimization, there are three accuracies of three models: J48 with 84.13%, RandomForest with 84.6% and Bagging with 92%.

```

# Store data
algoName <- c("ZeroR", "J48", "RandomForest", "Bagging")
acc <- c(59.2, 84.13, 84.6, 92.0)
timealgo <- c(0.01, 0.26, 1.12, 2.29)

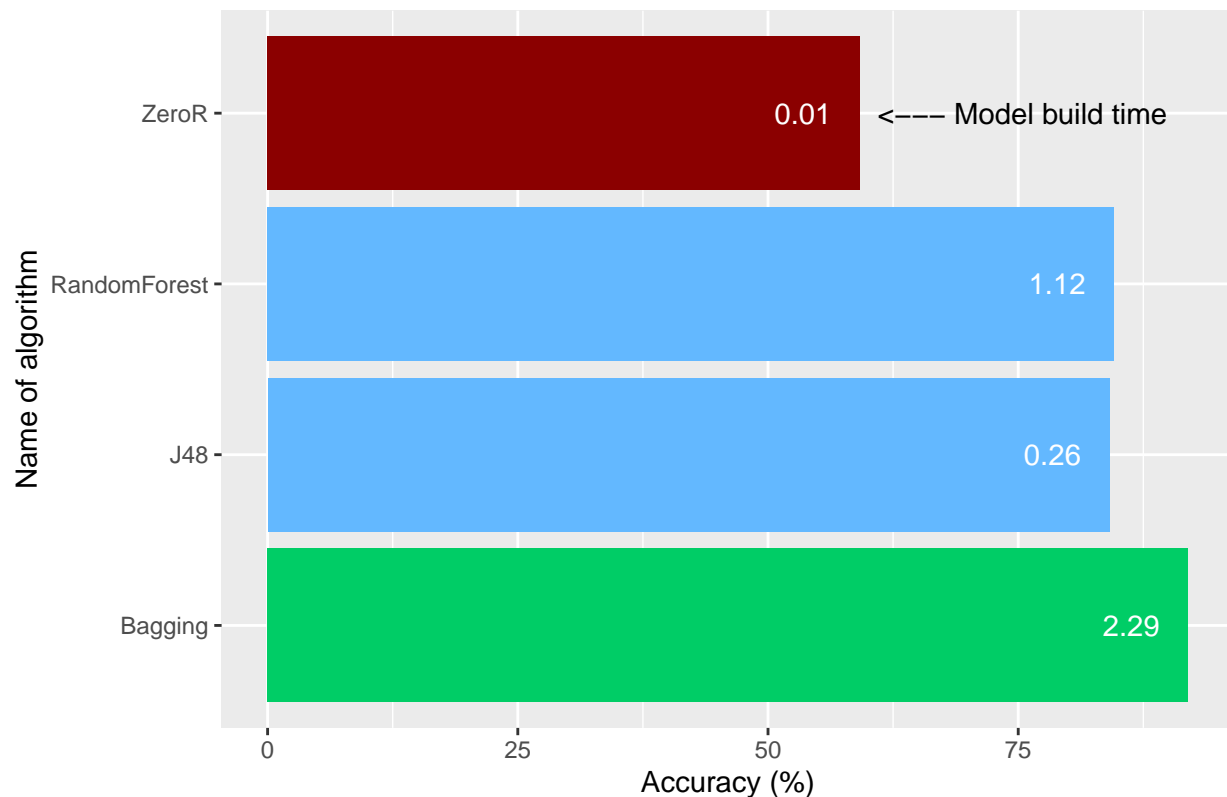
# dataframe data

```

```
x <- as.data.frame(algoName)
x$acc <- acc
x$timealgo <- timealgo

# plot data
ggplot(x, aes(x = acc, y = algoName)) + geom_bar(stat = 'identity', fill = c("darkred", "steelblue1", "steelblue1", "steelblue1"),
  geom_text(aes(label = timealgo), position = position_dodge(width = 0.9), hjust = 1.5, color = 'white'),
  labs(x = "Accuracy (%)", y = "Name of algorithm", title = "Figure 21, Performance selected algorithms"),
  annotate(x = 75, y = 4, geom = "text", label = " <--- Model build time")
```

Figure 21, Performance selected algorithms



In the graph above, Bagging preforms the overall best, it is to ‘slowest’ of the 4 but still fast enough for not heavy usage. If we recall the quality metrics: The demand was a high scoring algorithm that also was fast. A good confusion matrix was also desirable.

A confusion matrix is a table with predictions against actual values. It includes two dimensions. Below is the confusion matrix for the optimized bagging algorithm.

```
# make confusion matrix
confusion_matrix <- matrix(data = c(90.6, 6.9, 9.4, 93.1), nrow = 2, ncol = 2, dimnames = list(c(TRUE,F), c(TRUE,F)))
confusion_matrix

##          TRUE FALSE
## TRUE   90.6    9.4
## FALSE   6.9   93.1
```

This confusion matrix shows at the top the predicted value and the rows as actual values. The numbers are percentages. The matrix shows a solid 90+ percentage for The True positives of True and false. There is

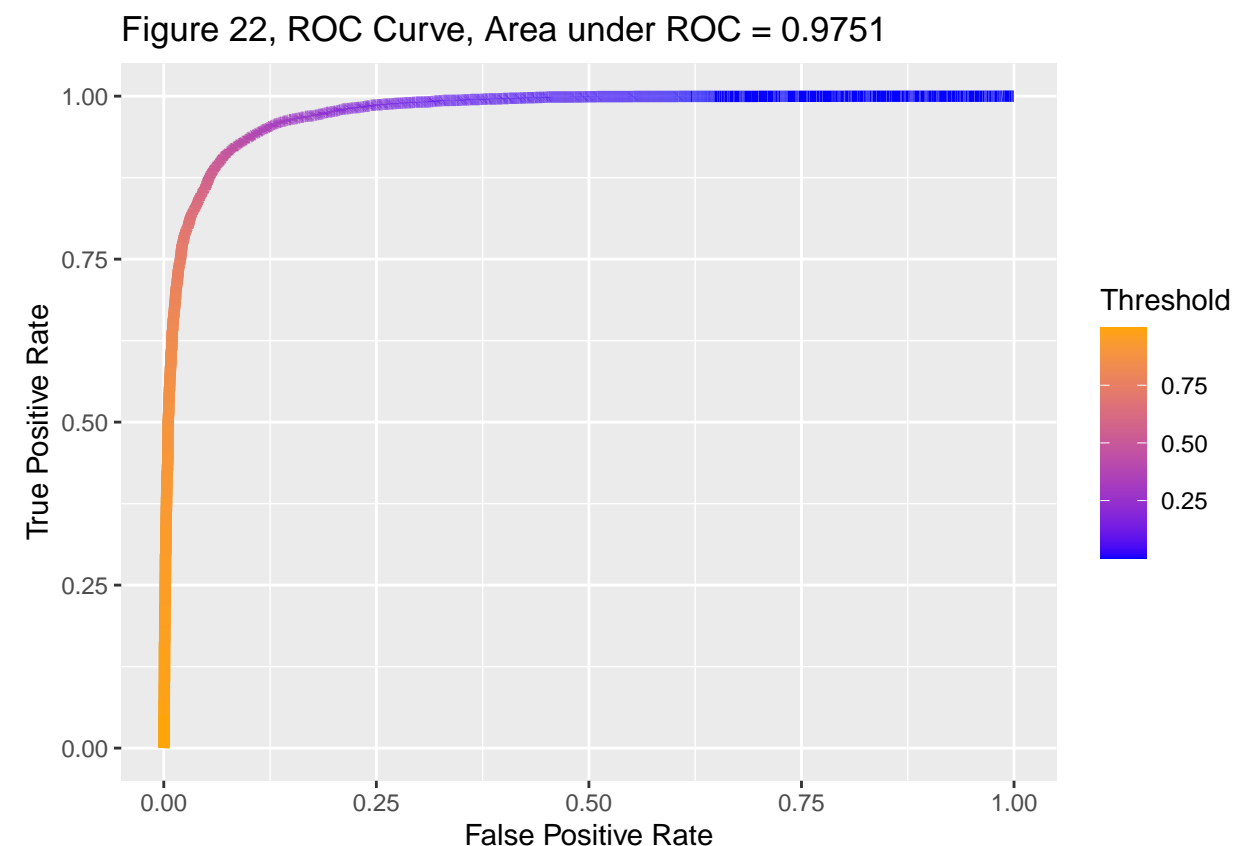


a 9.4% false positive rate, which is not great but not the worst. 6.9% is the false negative rate, which is a little bit too high. It means that roughly 7% of all the people are told that they don't have a higher change on stroke events, but they do.

An ROC curve is a plot that illustrates the diagnostic ability of a binary classifier system. It is a tool to control and select optimal models and discard worse models. It is related to a cost/benefit analysis. The Area under the curve (AUC), describes the performance. A perfect classification has an AUC of 1. Which means that the line lies perfectly along the axis. How more the bulge lies in the direction of the top left corner, the better the AUC is.

```
# load roc curve data for file
roc <- readARFF(path = "../Data/roc_curve.arff", show.info = F)

# plot curve
ggplot(roc, aes(y = `True Positive Rate`, x = `False Positive Rate`)) + geom_line(aes(color = Threshold))
```



Above is the ROC curve of the bagging algorithm, with an AUC of 0.9751, the model performs very well. The curve indicates a good model.

Overall the bagging algorithm is the best. High accuracy, fast and a good scoring confusion matrix. Also, the AUC is promising. Therefore is the bagging algorithm chosen to be used. To get the same results, simply follow the instructions below:

When the data is loaded in, remove the first empty row and the row with ID's. Next, apply the options below.

The Bagging algorithm with the following options:

- numIterations = 30

This algorithm will be featured in the wrapper, so that it can be executed as a command line application. The goal is that a user can input records or data and that the command line application with the bagging model, outputs a accuracy. The wrapper and it supporting files can be found here:

[https://github.com/PascalVisser/Thema\\_09\\_MLAnalysis/tree/main/JavaWrapper](https://github.com/PascalVisser/Thema_09_MLAnalysis/tree/main/JavaWrapper)