

# Statistical Learning Project

Filippo Santin, Gurjeet Singh, Francesca Zen

18/5/2021

## 1. Introduction

In the following report we present the analysis computed on a particular disease, stroke, and its correlation to other factors such as smoking, glucose level, bmi and so on.

“Stroke” is the medical term for damage to brain tissue or the death of a portion of it, due to insufficient blood supply to an area of the brain.

Our aim is to see if and how the variables we are dealing with are related to each other, in order to predict which individual is more probable to get a stroke.

The symptoms of stroke vary from patient to patient, depending on the severity of the condition, the affected brain area, causes, type of stroke, etc.

Stroke is characterized by sudden onset and for this reason it involves the need for immediate therapeutic intervention and adapted to the needs of the patient. In this sense, looking for relation between features may help to prevent it.

In order to have a guide for the interpretation of the data we underline the following information:

- The normal values of glucose level are between 60 and 110 mg/dl and with a value greater than 126 mg/dl a person is considered diabetic;
- a body mass index (BMI) between 18.5-24.9 indicates a normal/healthy weight, below 18.5 indicates underweight, 25.0-29.9 indicates overweight and above 30.0 indicates obese person

## 2. Exploiting the Dataset

The dataset we used was provided by kaggle <sup>1</sup> and it is composed of 5,110 entries with a total of 12 columns: `id`, `gender`, `age`, `hypertension`, `heart_disease`, `ever_married`, `work_type`, `Residence_type`, `avg_glucose_level`, `bmi`, `smoking_status`, `stroke`.

```
stroke_data <- read.csv('healthcare-dataset-stroke-data.csv')
attach(stroke_data)
```

The preliminary part of the analysis focuses on the study of the dataset: we looked at the `id` column and verified that all the data collected was referring to different people, so no recidivist status were involved.

```
stroke_data<-stroke_data[,-1]
```

In order to use the variables through the analysis we then transformed the categorical variables into factors:

```
stroke_data$gender<- as.factor(gender)
stroke_data$ever_married<-as.factor(ever_married)
stroke_data$work_type<-as.factor(work_type)
```

---

<sup>1</sup><https://www.kaggle.com/fedesoriano/stroke-prediction-dataset>

```
stroke_data$Residence_type<-as.factor(Residence_type)
stroke_data$smoking_status<-as.factor(smoking_status)
```

What's more, the variable `bmi` was not numeric because of the presence of a string "N/A" which identifies the lack of the value, and so we transformed it using the command `stroke_data$bmi <- as.numeric(bmi)`; and then removed those NA values: `stroke_data<- na.omit(stroke_data)`.

## Warning: NA introdotti per coercizione

We ended up having 4,909 entries and 11 total columns. Here we give a quick overview of the main information about the dataset:

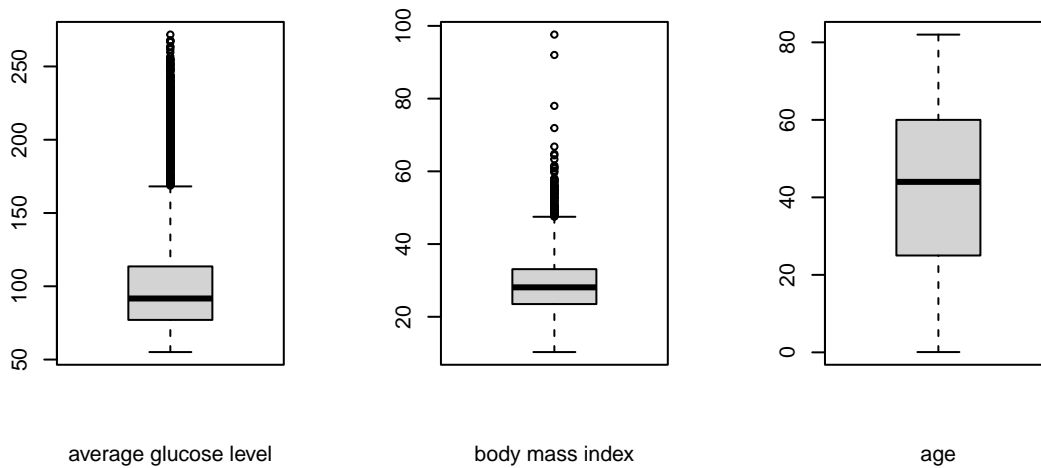
```
summary(stroke_data)
```

```
##      gender      age      hypertension      heart_disease      ever_married
## Female:2897  Min.   : 0.08  Min.   :0.00000  Min.   :0.0000  No :1705
## Male   :2011  1st Qu.:25.00  1st Qu.:0.00000  1st Qu.:0.0000  Yes:3204
## Other  :    1  Median :44.00  Median :0.00000  Median :0.0000
##              Mean   :42.87  Mean   :0.09187  Mean   :0.0495
##              3rd Qu.:60.00  3rd Qu.:0.00000  3rd Qu.:0.0000
##              Max.   :82.00  Max.   :1.00000  Max.   :1.0000
##      work_type  Residence_type avg_glucose_level      bmi
## children      : 671  Rural:2419  Min.   : 55.12  Min.   :10.30
## Govt_job       : 630  Urban:2490  1st Qu.: 77.07  1st Qu.:23.50
## Never_worked   : 22              Median : 91.68  Median :28.10
## Private        :2811              Mean   :105.31  Mean   :28.89
## Self-employed: 775              3rd Qu.:113.57  3rd Qu.:33.10
##              Max.   :271.74  Max.   :97.60
##      smoking_status      stroke
## formerly smoked: 837  Min.   :0.00000
## never smoked    :1852  1st Qu.:0.00000
## smokes          : 737  Median :0.00000
## Unknown         :1483  Mean   :0.04257
##              3rd Qu.:0.00000
##              Max.   :1.00000
```

```
attach(stroke_data)
```

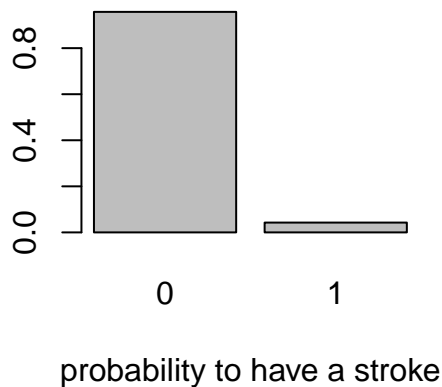
A visual transformation of these values is provided in the following boxplots:

```
par(mfrow=c(1,3))
boxplot(avg_glucose_level, xlab= 'average glucose level' )
boxplot(bmi, xlab = 'body mass index')
boxplot(age, xlab = 'age', pch=20)
```



Through the analysis on the Stroke dataset, we discovered that it was strongly bias, in the sense that 209 people on a total of 4909 get a stroke:

```
barplot(table(stroke)/dim(stroke_data)[1],
        xlab='probability to have a stroke')
```



This value is representative of the real situation in which there are not many stroke cases compared with the whole population. In Italy, for example, we have 200,000 cases over 59.226.539 people, i.e. 0.33%.

### 3. Searching for Relationships

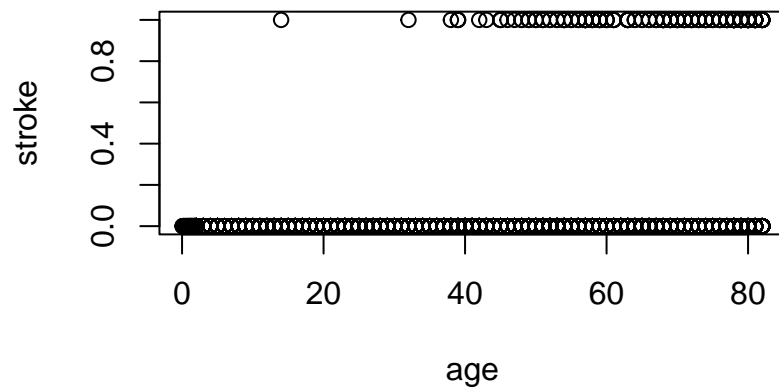
At this point we can ask some questions:

- Is it possible to prevent ictus?
- Which factors are the most related to it?
- How strong are the relations between the features?
- Are the given variables enough to predict a good accuracy of some possible person affected by ictus?

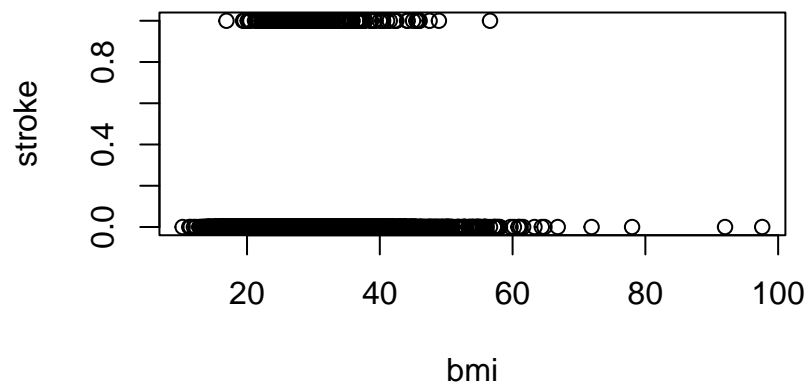
We will explore the data trying to answer them.

First of all we look at some intuitive relation between `stroke` and `age`, `bmi` and `avg_glucose_level`:

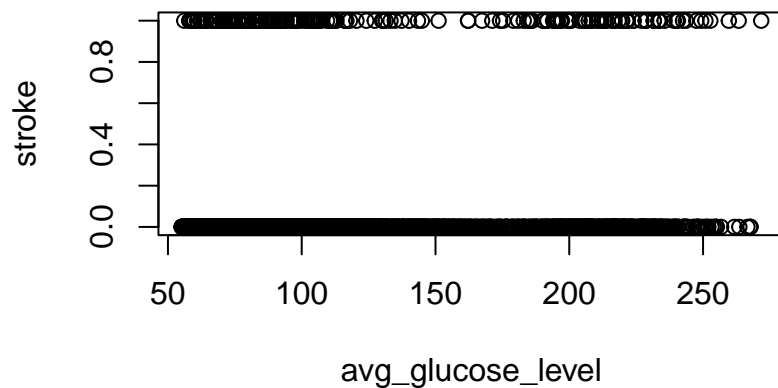
```
plot(stroke~age)
```



```
plot(stroke~bmi)
```



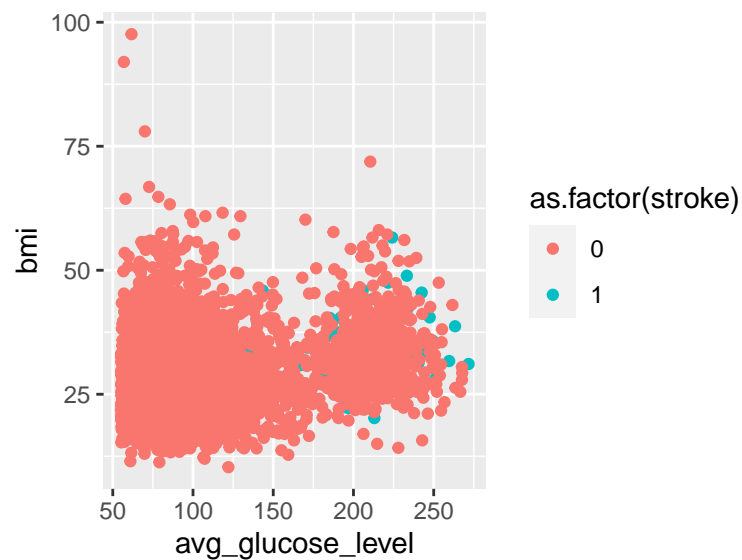
```
plot(stroke~avg_glucose_level)
```



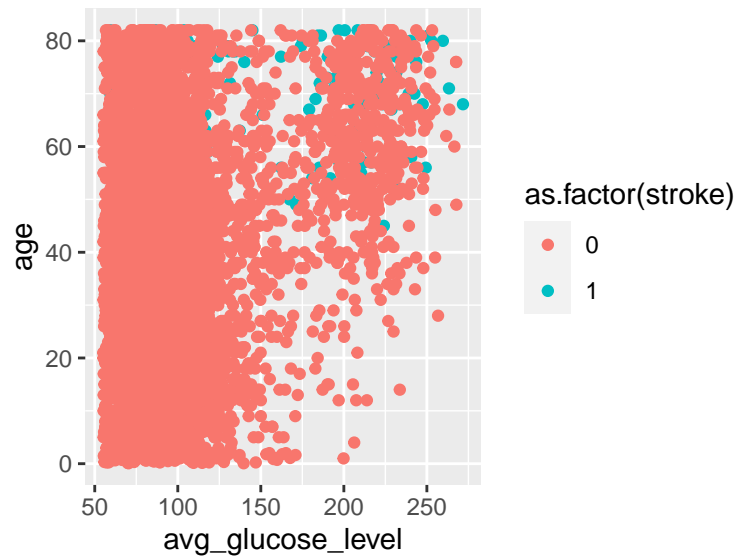
Looking to these plots we can see that the incidence of the virus increases progressively with age and that if we sum also the information about `avg_glucose_level` we may wonder if diabetic people are more probable to get a stroke or not. In addition there is no apparent relation of `stroke` with `bmi`.

We now highlight other visual relationship thanks to the scatter plots:

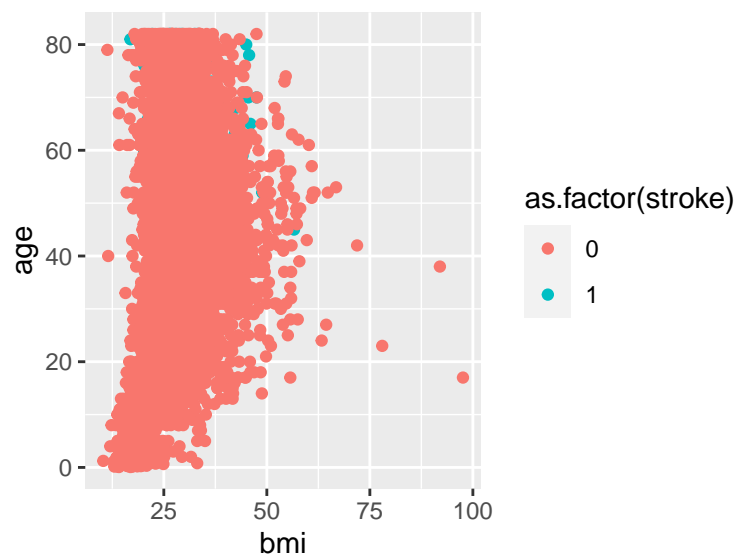
```
library(ggplot2)
par(mfrow=c(1,3))
ggplot(stroke_data, aes(x = avg_glucose_level, y = bmi,
                        col = as.factor(stroke))) + geom_point()
```



```
ggplot(stroke_data, aes(x = avg_glucose_level, y = age,
                        col = as.factor(stroke))) + geom_point()
```



```
ggplot(stroke_data, aes(x = bmi, y = age,
                        col = as.factor(stroke))) + geom_point()
```



```
par(mfrow=c(1,1))
```

discussione

## 4. Analysis on models

While going on with the classification using logistic regression we could meet the following problems: non-linearity of the data, correlation of error terms, heteroscedasticity, outliers, leverage point and collinearity.

### 4.1 Full and Reduced Models

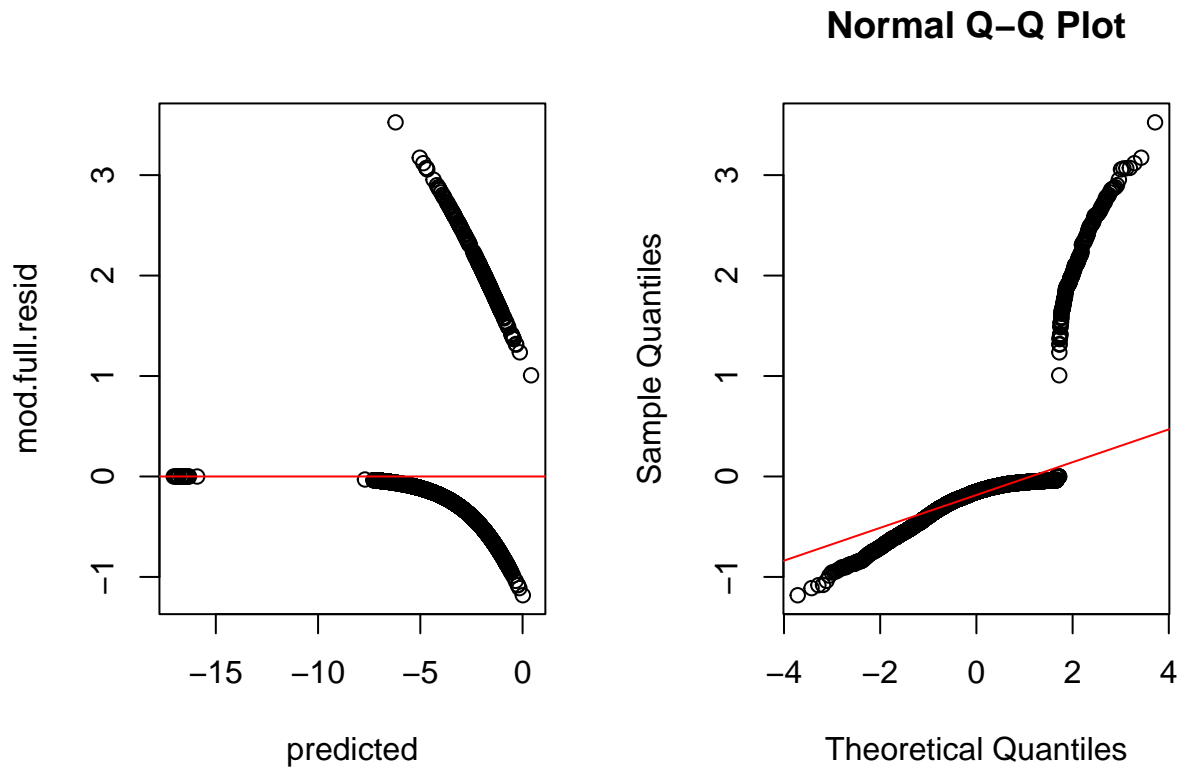
```
mod.full <- glm(stroke~., data=stroke_data, family = binomial)
summary(mod.full)
```

```
##
## Call:
## glm(formula = stroke ~ ., family = binomial, data = stroke_data)
##
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
## -1.1823  -0.2947  -0.1524  -0.0744   3.5251
##
## Coefficients:
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept)    -7.360e+00  1.067e+00  -6.895 5.37e-12 ***
## genderMale      -1.463e-02  1.544e-01  -0.095 0.924525
## genderOther    -1.135e+01  2.400e+03  -0.005 0.996225
## age             7.348e-02  6.347e-03  11.578 < 2e-16 ***
## hypertension    5.249e-01  1.750e-01   2.999 0.002711 **
## heart_disease    3.488e-01  2.072e-01   1.683 0.092381 .
## ever_marriedYes -1.152e-01  2.473e-01  -0.466 0.641394
## work_typeGovt_job -6.817e-01  1.114e+00  -0.612 0.540660
## work_typeNever_worked -1.082e+01  5.090e+02  -0.021 0.983036
## work_typePrivate  -5.208e-01  1.100e+00  -0.473 0.635943
## work_typeSelf-employed -9.459e-01  1.119e+00  -0.845 0.397906
## Residence_typeUrban  4.514e-03  1.500e-01   0.030 0.975990
## avg_glucose_level  4.652e-03  1.294e-03   3.595 0.000324 ***
## bmi             4.062e-03  1.188e-02   0.342 0.732387
## smoking_statusnever smoked -6.722e-02  1.886e-01  -0.356 0.721556
## smoking_statussmokes  3.139e-01  2.295e-01   1.368 0.171310
## smoking_statusUnknown -2.753e-01  2.471e-01  -1.114 0.265193
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
##      Null deviance: 1728.4  on 4908  degrees of freedom
## Residual deviance: 1363.2  on 4892  degrees of freedom
## AIC: 1397.2
##
## Number of Fisher Scoring iterations: 15
```

Here we see that `age` and `hypertension` are the variables most related to `stroke`.

Let's use the residual plots to get more information about this model:

```
mod.full.resid <- residuals(mod.full, type="deviance") # because we have a binary response
predicted <- predict(mod.full, type = "link")
par(mfrow=c(1,2))
plot(mod.full.resid-predicted)
abline(h=0, col='red')
qqnorm(mod.full.resid)
qqline(mod.full.resid, col='red')
```



The residual plots are not satisfactory. From the right plot we can see that the data are not normal.

## 4.2 Comparison between Reduced Models

```
# Reduced model 1 : We remove at least all the features that have collinearity between
# each other (work_type, ever_married) and the residence type
mod.red1 <- glm(stroke ~ age + bmi + avg_glucose_level + hypertension + smoking_status +
               gender, family=binomial)
summary(mod.red1)
```

```
##
## Call:
## glm(formula = stroke ~ age + bmi + avg_glucose_level + hypertension +
##      smoking_status + gender, family = binomial)
##
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
## -1.0138  -0.2991  -0.1569  -0.0715   3.6908
##
## Coefficients:
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept)   -7.915289    0.587035 -13.484  < 2e-16 ***
## age             0.070996    0.005764  12.318  < 2e-16 ***
## bmi             0.002695    0.011699   0.230  0.817837
## avg_glucose_level 0.004915    0.001282   3.834  0.000126 ***
## hypertension    0.527714    0.173926   3.034  0.002412 **
## smoking_statusnever smoked -0.070206    0.187586  -0.374  0.708209
```



```

## smoking_statussmokes      0.343917  0.227595  1.511 0.130766
## smoking_statusUnknown     -0.256575  0.244806 -1.048 0.294603
## genderMale                 0.017558  0.152828  0.115 0.908533
## genderOther                -7.261569 324.743861 -0.022 0.982160
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
## Null deviance: 1728.4 on 4908 degrees of freedom
## Residual deviance: 1372.5 on 4899 degrees of freedom
## AIC: 1392.5
##
## Number of Fisher Scoring iterations: 11
# Reduced model 2: remove gender
mod.red2 <- glm(stroke ~ age + bmi + avg_glucose_level + hypertension + smoking_status, family=binomial)
summary(mod.red2)

##
## Call:
## glm(formula = stroke ~ age + bmi + avg_glucose_level + hypertension +
##      smoking_status, family = binomial)
##
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
## -1.0108  -0.2983  -0.1571  -0.0716   3.6890
##
## Coefficients:
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept)    -7.906885   0.580996 -13.609 < 2e-16 ***
## age              0.070991   0.005760  12.325 < 2e-16 ***
## bmi              0.002663   0.011686   0.228 0.819754
## avg_glucose_level 0.004928   0.001277   3.859 0.000114 ***
## hypertension    0.527864   0.173922   3.035 0.002405 **
## smoking_statusnever smoked -0.073031  0.185814  -0.393 0.694294
## smoking_statussmokes  0.343591   0.227529   1.510 0.131019
## smoking_statusUnknown -0.258084   0.244363  -1.056 0.290900
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
## Null deviance: 1728.4 on 4908 degrees of freedom
## Residual deviance: 1372.5 on 4901 degrees of freedom
## AIC: 1388.5
##
## Number of Fisher Scoring iterations: 7
# Final Reduced Model
mod.red <- glm(stroke~age + bmi + avg_glucose_level+ hypertension, data=stroke_data, family = binomial)
summary(mod.red)

##
## Call:
## glm(formula = stroke ~ age + bmi + avg_glucose_level + hypertension,

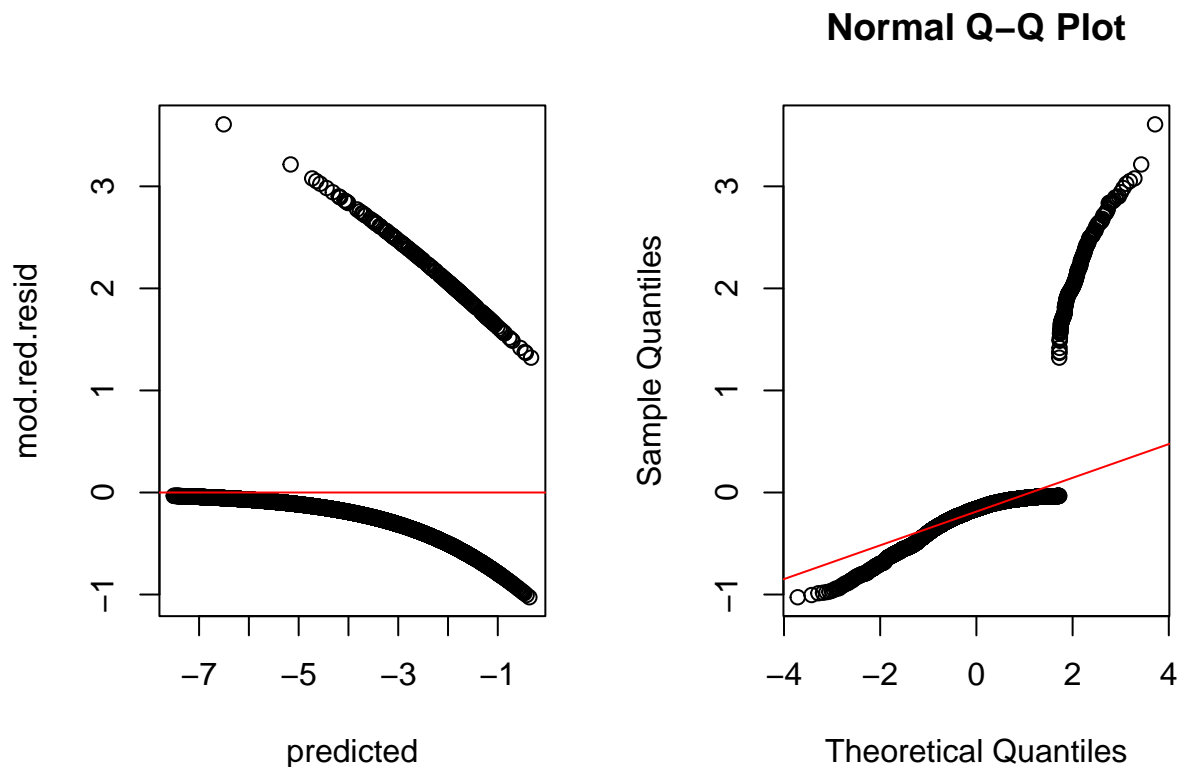
```

```
##      family = binomial, data = stroke_data)
##
## Deviance Residuals:
##      Min        1Q      Median        3Q        Max
## -1.0265   -0.2986   -0.1600   -0.0755    3.6075
##
## Coefficients:
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept)   -7.852291    0.541579  -14.499 < 2e-16 ***
## age            0.069793    0.005593   12.479 < 2e-16 ***
## bmi           0.002621    0.011598    0.226  0.82121
## avg_glucose_level 0.004984    0.001276    3.905 9.41e-05 ***
## hypertension  0.543399    0.173304    3.136 0.00172 **
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
##      Null deviance: 1728.4  on 4908  degrees of freedom
## Residual deviance: 1378.4  on 4904  degrees of freedom
## AIC: 1388.4
##
## Number of Fisher Scoring iterations: 7
```

The mod.red has age, hypertension, avg\_glucose\_level, bmi as the variables with the highest level of significance.

```
mod.red.resid <- residuals(mod.red, type="deviance")
predicted <- predict(mod.red, type = "link")

par(mfrow=c(1,2))
plot(mod.red.resid~predicted)
abline(h=0, col='red')
qqnorm(mod.red.resid)
qqline(mod.red.resid, col='red')
```



Residual plot:

```
#par(mfrow=c(2,2))
#plot(mod.red)
#par(mfrow=c(1,1))
```

We can see from the residual vs predicted values the presence of high non-linearity in the dataset. In the qqplot instead we see that residuals do not follow a normal distribution.

Instead in the standard deviation vs predicted we can see that homoscedasticity does not hold since the line of the standard residual is not flat, hence even by standardizing the residual we end up having high variance among residuals.

In the end by looking at the leverage plot, we see the presence of some sample with high leverage values (bottom right), which could influence the prediction of the model.

*But I don't know in which range of leverage value is considered to change a lot the prediction of the model. Furthermore R does not show the index of the sample with high leverage, I guess because a lot of values could change the prediction. Some outliers with high variance are: idx: 119, 183, 246.*

Anova computation:

```
anova(mod.full, mod.red, test="Chisq")
```

```
## Analysis of Deviance Table
```

```
##
```

```
## Model 1: stroke ~ gender + age + hypertension + heart_disease + ever_married +
```

```
##      work_type + Residence_type + avg_glucose_level + bmi + smoking_status
```

```
## Model 2: stroke ~ age + bmi + avg_glucose_level + hypertension
```

```
##   Resid. Df Resid. Dev  Df Deviance Pr(>Chi)
```

```
## 1      4892      1363.2
## 2      4904      1378.4 -12  -15.127   0.2346
```

As expected from the anova test rejects that the complex model is more significant than the reduced one, since the p-value is not less than 5%. Hence the full model does not help with our prediction.

Outliers:

```
View(stroke_data[c("119", "183", "246"), ])
```

### 4.3 Mixed Approach for Models

In order to see which variables were relevant on our research we tested various models using the mixed approach: we started by the null model and then added variables and removed those with p-value over a threshold of 0.1 (which indicates low relationship with the `stroke` variable). In addition, we used the F-statistic to determine those variables that were in good relation, so they may contribute positively to the models:

```
##
## F test to compare two variances
##
## data:  age and avg_glucose_level
## F = 0.25778, num df = 4908, denom df = 4908, p-value < 2.2e-16
## alternative hypothesis: true ratio of variances is not equal to 1
## 95 percent confidence interval:
##  0.2437502 0.2726159
## sample estimates:
## ratio of variances
##      0.2577793
##
## F test to compare two variances
##
## data:  age and hypertension
## F = 6096.4, num df = 4908, denom df = 4908, p-value < 2.2e-16
## alternative hypothesis: true ratio of variances is not equal to 1
## 95 percent confidence interval:
##  5764.585 6447.244
## sample estimates:
## ratio of variances
##      6096.367
##
## F test to compare two variances
##
## data:  hypertension and avg_glucose_level
## F = 4.2284e-05, num df = 4908, denom df = 4908, p-value < 2.2e-16
## alternative hypothesis: true ratio of variances is not equal to 1
## 95 percent confidence interval:
##  3.998287e-05 4.471776e-05
## sample estimates:
## ratio of variances
##      4.228409e-05
##
## F test to compare two variances
##
```

```
## data: age and heart_disease
## F = 10810, num df = 4908, denom df = 4908, p-value < 2.2e-16
## alternative hypothesis: true ratio of variances is not equal to 1
## 95 percent confidence interval:
## 10221.95 11432.46
## sample estimates:
## ratio of variances
## 10810.27

##
## F test to compare two variances
##
## data: avg_glucose_level and heart_disease
## F = 41936, num df = 4908, denom df = 4908, p-value < 2.2e-16
## alternative hypothesis: true ratio of variances is not equal to 1
## 95 percent confidence interval:
## 39653.86 44349.79
## sample estimates:
## ratio of variances
## 41936.15

##
## F test to compare two variances
##
## data: age and bmi
## F = 8.2471, num df = 4908, denom df = 4908, p-value < 2.2e-16
## alternative hypothesis: true ratio of variances is not equal to 1
## 95 percent confidence interval:
## 7.798263 8.721756
## sample estimates:
## ratio of variances
## 8.247093
```

All the p-values of the tests were above  $2.2e-16$  and so each of them were good relationships. Only the value of the F-statistic varied.

We start with a simple model with the only iteration between **stroke** and **age**.

```
mod1 <- glm(stroke~age, family=binomial)
summary(mod1)
```

```
##
## Call:
## glm(formula = stroke ~ age, family = binomial)
##
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
## -0.7162  -0.3073  -0.1639  -0.0776   3.5579
##
## Coefficients:
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept) -7.377231   0.362383  -20.36  <2e-16 ***
## age          0.074969   0.005318   14.10  <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
```

```
##
## Null deviance: 1728.4 on 4908 degrees of freedom
## Residual deviance: 1407.7 on 4907 degrees of freedom
## AIC: 1411.7
##
## Number of Fisher Scoring iterations: 7
```

We then add the avg\_glucose\_level variable and also its interaction with age:

```
mod2 <- glm(stroke~age+avg_glucose_level+age*avg_glucose_level, family = binomial)
summary(mod2)
```

avg\_glucose\_level\*age and avg\_glucose\_level have p-value over the threshold, so we remove them, then add hypertension:

```
mod3 <-glm(stroke ~ age + hypertension, family = binomial)
summary(mod3)
```

All the p-values makes the variables relevant so we keep all of them and we add the interaction between hypertension and glucose\_level and age:

```
##
## Call:
## glm(formula = stroke ~ age + hypertension + hypertension * avg_glucose_level +
## hypertension * age, family = binomial)
##
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
## -0.9157  -0.3025  -0.1536  -0.0703   3.6637
##
## Coefficients:
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept)    -8.057914    0.434123  -18.561 < 2e-16 ***
## age              0.072805    0.006162   11.815 < 2e-16 ***
## hypertension    2.298823    1.028218    2.236 0.025369 *
## avg_glucose_level  0.005670    0.001476    3.842 0.000122 ***
## hypertension:avg_glucose_level -0.002234    0.002691   -0.830 0.406447
## age:hypertension  -0.020893    0.014248   -1.466 0.142539
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
## Null deviance: 1728.4 on 4908 degrees of freedom
## Residual deviance: 1375.5 on 4903 degrees of freedom
## AIC: 1387.5
##
## Number of Fisher Scoring iterations: 8
```

We remove hypertension\*avg\_glucose\_level and add heart\_disease:

```
mod5 <- glm(stroke~age + hypertension +heart_disease, family = binomial)
summary(mod5)
```

We add heart\_disease \*avg\_glucose\_level + heart\_disease\*hypertension + heart\_disease\*age

```
##
## Call:
```

```
## glm(formula = stroke ~ age + hypertension + heart_disease + heart_disease *
##     avg_glucose_level + heart_disease * hypertension + heart_disease *
##     age, family = binomial)
##
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
## -1.0105  -0.2975  -0.1555  -0.0739   3.6108
##
## Coefficients:
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept)    -7.704998   0.410975 -18.748 < 2e-16 ***
## age              0.070089   0.005905  11.868 < 2e-16 ***
## hypertension    0.672166   0.190121   3.535 0.000407 ***
## heart_disease   1.974750   1.515413   1.303 0.192537
## avg_glucose_level 0.003563   0.001430   2.492 0.012714 *
## heart_disease:avg_glucose_level 0.005404   0.003173   1.703 0.088520 .
## hypertension:heart_disease -0.677822   0.451136  -1.502 0.132973
## age:heart_disease -0.030265   0.020012  -1.512 0.130445
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
##      Null deviance: 1728.4  on 4908  degrees of freedom
## Residual deviance: 1367.3  on 4901  degrees of freedom
## AIC: 1383.3
##
## Number of Fisher Scoring iterations: 7
```

We can remove heart\_disease because it has a p-value of 0.1925 but also heart\_disease\*hypertension and heart\_disease\*age. We go on adding bmi:

```
##
## Call:
## glm(formula = stroke ~ age + hypertension + heart_disease * avg_glucose_level +
##     bmi, family = binomial)
##
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
## -1.2178  -0.2959  -0.1607  -0.0769   3.5813
##
## Coefficients:
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept)    -7.728286   0.545569 -14.166 < 2e-16 ***
## age              0.068598   0.005696  12.043 < 2e-16 ***
## hypertension    0.540689   0.174307   3.102 0.00192 **
## heart_disease   -0.359185   0.528041  -0.680 0.49636
## avg_glucose_level 0.003656   0.001456   2.512 0.01200 *
## bmi              0.004684   0.011691   0.401 0.68868
## heart_disease:avg_glucose_level 0.005156   0.003190   1.616 0.10605
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
```

```
## Null deviance: 1728.4 on 4908 degrees of freedom
## Residual deviance: 1371.9 on 4902 degrees of freedom
## AIC: 1385.9
##
## Number of Fisher Scoring iterations: 7
```

We remove bmi because it has a p-value of 0.68868 and also heart\_disease\*avg\_glucose\_level and add its relation with age and hypertension because they are the most relevant variables:

```
##
## Call:
## glm(formula = stroke ~ age + hypertension + heart_disease * hypertension +
## heart_disease * age + bmi * age + bmi * hypertension, family = binomial)
##
## Deviance Residuals:
## Min 1Q Median 3Q Max
## -0.9454 -0.3007 -0.1597 -0.0679 3.5726
##
## Coefficients:
## Estimate Std. Error z value Pr(>|z|)
## (Intercept) -9.4992570 1.3689325 -6.939 3.94e-12 ***
## age 0.0977302 0.0210542 4.642 3.45e-06 ***
## hypertension 1.3441611 0.7537776 1.783 0.0745 .
## heart_disease 3.0023732 1.3624036 2.204 0.0275 *
## bmi 0.0686198 0.0417573 1.643 0.1003
## hypertension:heart_disease -0.6077302 0.4395576 -1.383 0.1668
## age:heart_disease -0.0325363 0.0191242 -1.701 0.0889 .
## age:bmi -0.0008542 0.0006618 -1.291 0.1968
## hypertension:bmi -0.0198876 0.0234357 -0.849 0.3961
## ---
## Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
## Null deviance: 1728.4 on 4908 degrees of freedom
## Residual deviance: 1380.5 on 4900 degrees of freedom
## AIC: 1398.5
##
## Number of Fisher Scoring iterations: 8
```

We remove them plus and left the model with:

```
mod9 <- glm(stroke~ age + avg_glucose_level + hypertension+ heart_disease*age,
             family = binomial)
summary(mod9)
```

```
##
## Call:
## glm(formula = stroke ~ age + avg_glucose_level + hypertension +
## heart_disease * age, family = binomial)
##
## Deviance Residuals:
## Min 1Q Median 3Q Max
## -0.9897 -0.2980 -0.1557 -0.0737 3.6232
##
## Coefficients:
```



```
##               Estimate Std. Error z value Pr(>|z|)
## (Intercept)    -7.816578   0.407175 -19.197 < 2e-16 ***
## age            0.070133   0.005889  11.908 < 2e-16 ***
## avg_glucose_level 0.004702   0.001253   3.752 0.000176 ***
## hypertension    0.536550   0.172602   3.109 0.001880 **
## heart_disease    2.765299   1.396557   1.980 0.047694 *
## age:heart_disease -0.032872  0.019486  -1.687 0.091604 .
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
##    Null deviance: 1728.4  on 4908  degrees of freedom
## Residual deviance: 1372.0  on 4903  degrees of freedom
## AIC: 1384
##
## Number of Fisher Scoring iterations: 7
```

We can also try with the categorical variables but actually the results do not change, so we keep this mod9.

Let's now see some relevant information, such as outliers, leverage point and collinearity through some plots:

```
par(mfrow=c(2,2))
#plot(mod9)
par(mfrow=c(1,1))
```

## 4.4 Polynomial models

We try to use a polynomial model using the numerical variables with degree 2.

```
mod.red.poly1 <- glm(stroke~age + bmi + avg_glucose_level+ hypertension+
                     I(bmi^2), family = binomial)
summary(mod.red.poly1)
```

```
##
## Call:
## glm(formula = stroke ~ age + bmi + avg_glucose_level + hypertension +
##      I(bmi^2), family = binomial)
##
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
## -1.0399  -0.2986  -0.1603  -0.0753   3.5922
##
## Coefficients:
##               Estimate Std. Error z value Pr(>|z|)
## (Intercept)    -8.6441614   1.2126773  -7.128 1.02e-12 ***
## age            0.0694084   0.0056315  12.325 < 2e-16 ***
## bmi            0.0536596   0.0708178   0.758  0.44862
## avg_glucose_level 0.0049654   0.0012760   3.891 9.97e-05 ***
## hypertension    0.5453191   0.1732628   3.147  0.00165 **
## I(bmi^2)        -0.0007579   0.0010455  -0.725  0.46852
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
```

```
## Null deviance: 1728.4 on 4908 degrees of freedom
## Residual deviance: 1377.8 on 4903 degrees of freedom
## AIC: 1389.8
##
## Number of Fisher Scoring iterations: 7
mod.red.poly2 <- glm(stroke~age + bmi + avg_glucose_level+ hypertension
+ I(avg_glucose_level^2), family = binomial)
summary(mod.red.poly2)

##
## Call:
## glm(formula = stroke ~ age + bmi + avg_glucose_level + hypertension +
## I(avg_glucose_level^2), family = binomial)
##
## Deviance Residuals:
## Min 1Q Median 3Q Max
## -1.0306 -0.2984 -0.1604 -0.0758 3.5989
##
## Coefficients:
## Estimate Std. Error z value Pr(>|z|)
## (Intercept) -7.734e+00 7.828e-01 -9.880 <2e-16 ***
## age 6.971e-02 5.605e-03 12.438 <2e-16 ***
## bmi 2.600e-03 1.160e-02 0.224 0.8226
## avg_glucose_level 3.158e-03 8.844e-03 0.357 0.7210
## hypertension 5.438e-01 1.733e-01 3.138 0.0017 **
## I(avg_glucose_level^2) 6.078e-06 2.914e-05 0.209 0.8348
## ---
## Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
## Null deviance: 1728.4 on 4908 degrees of freedom
## Residual deviance: 1378.3 on 4903 degrees of freedom
## AIC: 1390.3
##
## Number of Fisher Scoring iterations: 7
mod.red.poly <- glm(stroke~age + bmi + avg_glucose_level+ hypertension+
I(bmi^2) + I(avg_glucose_level^2), family = binomial)
summary(mod.red.poly)

##
## Call:
## glm(formula = stroke ~ age + bmi + avg_glucose_level + hypertension +
## I(bmi^2) + I(avg_glucose_level^2), family = binomial)
##
## Deviance Residuals:
## Min 1Q Median 3Q Max
## -1.0443 -0.2991 -0.1602 -0.0755 3.5832
##
## Coefficients:
## Estimate Std. Error z value Pr(>|z|)
## (Intercept) -8.524e+00 1.334e+00 -6.388 1.68e-10 ***
## age 6.933e-02 5.644e-03 12.283 < 2e-16 ***
```

```
## bmi                5.377e-02  7.085e-02   0.759  0.44790
## avg_glucose_level  3.078e-03  8.840e-03   0.348  0.72768
## hypertension      5.457e-01  1.733e-01   3.149  0.00164 **
## I(bmi^2)          -7.598e-04  1.046e-03  -0.726  0.46761
## I(avg_glucose_level^2) 6.283e-06  2.912e-05   0.216  0.82918
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
##    Null deviance: 1728.4  on 4908  degrees of freedom
## Residual deviance: 1377.7  on 4902  degrees of freedom
## AIC: 1391.7
##
## Number of Fisher Scoring iterations: 7
```

Nothing interesting with polynomial models. No improvement in the results.

## 5. LDA

Assumption: samples are normally distributed and have same variance in every class => strong assumption.

```
library(MASS)
lda.fit <- lda(stroke~age+bmi+avg_glucose_level+hypertension+work_type+gender+smoking_status)
lda.pred <- predict(lda.fit)
table(lda.pred$class, stroke)
```

```
##      stroke
##         0    1
##    0 4674  200
##    1   26    9
```

## 6. QDA

Assumption: sample are normally distributed BUT NOT SAME variance among classes.

```
qda.fit <- qda(stroke~age+bmi+avg_glucose_level+hypertension+heart_disease+smoking_status, data = stroke_data)
# ERROR rank deficiency, i.e. some variables
# are collinear and one or more covariance matrices cannot be inverted to obtain the estimates in group
qda.pred <- predict(qda.fit, stroke_data)
#TODO: Calculate deviance and standardize deviance

table(qda.pred$class, stroke)
```

```
##      stroke
##         0    1
##    0 4260  123
##    1  440   86
```

## 7. ROC and RECALL-PRECISION CURVES

```
library(pROC)
```

```
## Type 'citation("pROC")' for a citation.
```

```
##  
## Caricamento pacchetto: 'pROC'  
## I seguenti oggetti sono mascherati da 'package:stats':  
##  
##      cov, smooth, var  
library(ROCR)
```

Si stima che la percentuale di persone che possono avere un ictus andrà via via crescendo dal momento che l'età media della popolazione è in costante crescita.