

# EDA Jorick Baron

Jorick Baron

```
library(tidyr)
library(dplyr)
library(ggplot2)
library(gridExtra)
library(stringr)
library(knitr)
library(kableExtra)
library(e1071)
```

## Research question

How accurate can a model be trained to detect the difficult to diagnose pancreatic cancer utilising a patient's urine sample?

## Codebook

In this EDA we will explore the data downloaded from [here](#). For future reference we will describe the data in the codebook below.

```
codebook <- read.delim("Data/codebook.csv", sep = ",")
kable(codebook, caption = "Table 1: Codebook", align = "lcccr", booktabs = T) %>%
  kable_styling(latex_options = c("scale_down"))
```

Table 1: Table 1: Codebook

Name	Fullname	Description	Type	Unit
sample_id	Sample ID	Unique string identifying each subject	string	NA
patient_cohort	Patient's Cohort	Cohort 1, previously used samples; Cohort 2, newly added samples	string	NA
sample_origin	Sample Origin	BPTB: Barts Pancreas Tissue Bank; ESP: Spanish National Cancer Research Centre; LIV: Liverpool University; UCL: University College	string	NA
age	Age	Age in years	int	years
sex	Sex	M = male, F = female	char	NA
diagnosis	Diagnosis (1=Control, 2=Benign, 3=PDAC)	1 = control, 2 = benign hepatobiliary disease; 3 = Pancreatic ductal adenocarcinoma, i.e. pancreatic cancer	int	NA
stage	Stage	For those with pancreatic cancer, what stage was it? One of IA, IB, IIA, IIIB, III, IV	string	NA
benign_sample_diagnosis	Benign Samples Diagnosis	For those with a benign, non-cancerous diagnosis, what was the diagnosis?	string	NA
plasma_CA19_9	Plasma CA19-9 U/ml	Blood plasma levels of CA 19-9 monoclonal antibody that is often elevated in patients with pancreatic cancer.	float	plasma units/milliliter
creatinine	Creatinine mg/ml	Urinary biomarker of kidney function	float	mg/ml
LYVE1	LYVE1 ng/ml	Urinary levels of Lymphatic vessel endothelial hyaluronan receptor 1, a protein that may play a role in tumor metastasis	float	ng/ml
REG1B	REG1B ng/ml	Urinary levels of a protein that may be associated with pancreas regeneration.	float	ng/ml
TFF1	TFF1 ng/ml	Urinary levels of Trefoil Factor 1, which may be related to regeneration and repair of the urinary tract	float	ng/ml
REG1A	REG1A ng/ml	Urinary levels of a protein that may be associated with pancreas regeneration.	float	ng/ml

## Loading data

First we will load in the data and to check if it has loaded in properly we look at the structure of the loaded data.

```
data <- read.table(file="Data/source/Debernardi_et_al_2020_data.csv", sep = ",",
                  header = T, na.strings = "")
```

```
str(data)
```

```
## 'data.frame': 590 obs. of 14 variables:
## $ sample_id      : chr "S1" "S10" "S100" "S101" ...
## $ patient_cohort : chr "Cohort1" "Cohort1" "Cohort2" "Cohort2" ...
## $ sample_origin  : chr "BPTB" "BPTB" "BPTB" "BPTB" ...
## $ age            : int 33 81 51 61 62 53 70 58 59 56 ...
## $ sex            : chr "F" "F" "M" "M" ...
## $ diagnosis       : int 1 1 1 1 1 1 1 1 1 1 ...
## $ stage          : chr NA NA NA NA ...
## $ benign_sample_diagnosis: chr NA NA NA NA ...
## $ plasma_CA19_9   : num 11.7 NA 7 8 9 NA NA 11 NA 24 ...
## $ creatinine      : num 1.832 0.973 0.78 0.701 0.215 ...
## $ LYVE1           : num 0.89322 2.03758 0.14559 0.0028 0.00086 ...
## $ REG1B           : num 52.9 94.5 102.4 60.6 65.5 ...
## $ TFF1            : num 654.3 209.5 461.1 142.9 41.1 ...
## $ REG1A           : num 1262 228 NA NA NA ...
```

Thus far it seems to have loaded correctly.

We will also check the first few records to maybe catch some possible errors.

```
head(data)
```

```
## sample_id patient_cohort sample_origin age sex diagnosis stage
## 1 S1 Cohort1 BPTB 33 F 1 <NA>
## 2 S10 Cohort1 BPTB 81 F 1 <NA>
## 3 S100 Cohort2 BPTB 51 M 1 <NA>
## 4 S101 Cohort2 BPTB 61 M 1 <NA>
## 5 S102 Cohort2 BPTB 62 M 1 <NA>
## 6 S103 Cohort2 BPTB 53 M 1 <NA>
## benign_sample_diagnosis plasma_CA19_9 creatinine LYVE1 REG1B
## 1 <NA> 11.7 1.83222 0.89321920 52.94884
## 2 <NA> NA 0.97266 2.03758500 94.46703
## 3 <NA> 7.0 0.78039 0.14558890 102.36600
## 4 <NA> 8.0 0.70122 0.00280488 60.57900
## 5 <NA> 9.0 0.21489 0.00085956 65.54000
## 6 <NA> NA 0.84825 0.00339300 62.12600
## TFF1 REG1A
## 1 654.2822 1262.000
## 2 209.4882 228.407
## 3 461.1410 NA
## 4 142.9500 NA
## 5 41.0880 NA
## 6 59.7930 NA
```

The data seems to have quite a few NAs, reading further into the description most NAs would be expected i.e. no stage if there is no cancer thus an NA.

## NAs

Let's check that nothing went wrong with those two anyway.

```
healthy <- subset(data, diagnosis == 1, select = c(diagnosis, stage, benign_sample_diagnosis))
cancerfree <- subset(data, diagnosis == 2, select = c(diagnosis, stage, benign_sample_diagnosis))
cancerous <- subset(data, diagnosis == 3, select = c(diagnosis, stage, benign_sample_diagnosis))

stage_na_count <- sum(is.na(data$stage))
bsd_na_count <- sum(is.na(data$benign_sample_diagnosis))
paste("all these numbers should be the same number",
      stage_na_count - nrow(cancerfree), bsd_na_count - nrow(cancerous), nrow(healthy))

## [1] "all these numbers should be the same number 183 183 183"
```

Those numbers lined up to expectations.

The NAs in columns "plasma\_CA19\_9" and "REG1A" are supposed to be there because not every patient had been fully tested:

"REG1A ... Only assessed in 306 patients", "plasma\_CA19\_9 ... Only assessed in 350 patients" see Debernardi et al 2020 documentation.csv in the source files.

However to make sure everything is correct these numbers will be tested.

```
n_plasma_CA19_9 <- nrow(data) - sum(is.na(data$plasma_CA19_9))
n_REG1A <- nrow(data) - sum(is.na(data$REG1A))
paste("REG1A:", n_REG1A, "plasma_CA19_9:", n_plasma_CA19_9)

## [1] "REG1A: 306 plasma_CA19_9: 350"
```

These numbers are correct.

Are there more NAs?

```
sum(is.na(data[, c(1:6, 10:13)]))

## [1] 0
```

0 NAs remaining.

## Data exploration

### Distribution

Class label checking the different diagnoses should be in similar number to each diagnosis.

```
paste("Amount of patients with diagnosis 1:", nrow(subset(data, diagnosis == 1)))

## [1] "Amount of patients with diagnosis 1: 183"
paste("Amount of patients with diagnosis 2:", nrow(subset(data, diagnosis == 2)))

## [1] "Amount of patients with diagnosis 2: 208"
paste("Amount of patients with diagnosis 3:", nrow(subset(data, diagnosis == 3)))

## [1] "Amount of patients with diagnosis 3: 199"
```

These are quite balanced and should not influence statistics.

Let's look at a summary of the data for a quick overview of the distributions.

```
summary(data[,c(4, 9:14)])
```

```
##      age      plasma_CA19_9      creatinine      LYVE1
## Min.   :26.00 Min.   :  0.0 Min.   :0.05655 Min.   : 0.000129
## 1st Qu.:50.00 1st Qu.:  8.0 1st Qu.:0.37323 1st Qu.: 0.167179
## Median :60.00 Median : 26.5 Median :0.72384 Median : 1.649862
## Mean   :59.08 Mean   : 654.0 Mean   :0.85538 Mean   : 3.063530
## 3rd Qu.:69.00 3rd Qu.: 294.0 3rd Qu.:1.13948 3rd Qu.: 5.205037
## Max.   :89.00 Max.   :31000.0 Max.   :4.11684 Max.   :23.890323
##
##      NA's      :240
##      REG1B      TFF1      REG1A
## Min.   :  0.0011 Min.   :  0.005 Min.   :  0.00
## 1st Qu.: 10.7572 1st Qu.:  43.961 1st Qu.:  80.69
## Median : 34.3034 Median : 259.874 Median : 208.54
## Mean   : 111.7741 Mean   : 597.869 Mean   : 735.28
## 3rd Qu.: 122.7410 3rd Qu.: 742.736 3rd Qu.: 649.00
## Max.   :1403.8976 Max.   :13344.300 Max.   :13200.00
##
##      NA's      :284
```

Much of the data seems to be imbalanced with outliers.

Now let's take a closer look at the data itself using box-plots.

```
p1<- ggplot(data=data)+
  geom_boxplot(mapping = aes(x = "",
                             y = age))+
  ylab("age in years") +
  xlab(NULL)

p2<- ggplot(data=data)+
  geom_boxplot(mapping = aes(x = "",
                             y = plasma_CA19_9))+
  ylab("plasma_CA19_9 in U/ml") +
  xlab(NULL)+
  ylim(0,500)

p3<- ggplot(data=data)+
  geom_boxplot(mapping = aes(x = "",
                             y = creatinine))+
  ylab("creatinine in mg/ml") +
  xlab(NULL)

p4<- ggplot(data=data)+
  geom_boxplot(mapping = aes(x = "",
                             y = LYVE1))+
  ylab("LYVE1 in ng/ml") +
  xlab(NULL)+
  ylim(0,17)

p5<- ggplot(data=data)+
  geom_boxplot(mapping = aes(x = "",
                             y = REG1B))+
  ylab("REG1B in ng/ml") +
  xlab(NULL)+
  ylim(0,600)
```

```

p6<- ggplot(data=data)+
  geom_boxplot(mapping = aes(x = "",
                             y = TFF1))+
  ylab("TFF1 in ng/ml") +
  xlab(NULL)+
  ylim(0,5000)

p7<- ggplot(data=data)+
  geom_boxplot(mapping = aes(x = "",
                             y = REG1A))+
  ylab("REG1A in ng/ml") +
  xlab(NULL) +
  ylim(0,5000)

grid.arrange(p1, p2, p3, p4, p5, p6, p7, nrow = 2)

```

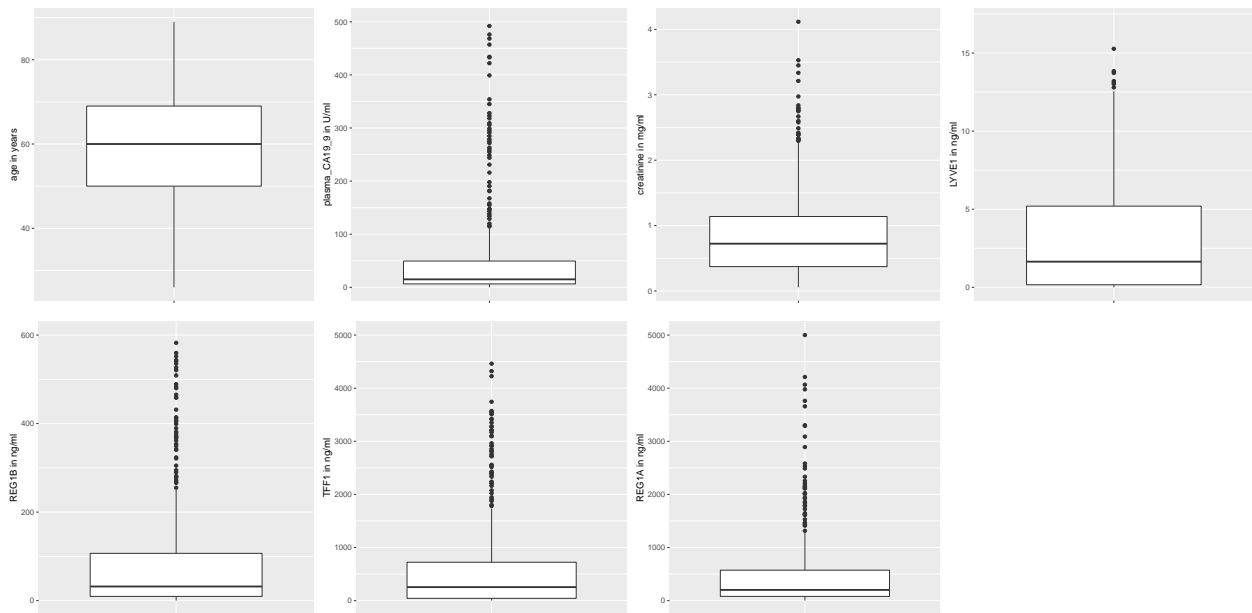


Figure 1: boxplots of different values

There are many outliers to take a good look at the whiskers y-limits are in place. Still it's a lot, maybe adding another dimension can correct this.

To add this extra dimension let's look at the difference in diagnoses. To properly do this we will also assign levels to the diagnosis column in the dataframe.

```

data$diagnosis <- factor(data$diagnosis)

gp1 <- ggplot(data = data)+
  geom_boxplot(mapping = aes(x = diagnosis,
                             y = age,
                             group=diagnosis,
                             fill=diagnosis))+
  scale_x_discrete(labels=c("control",
                           "benign hepatobiliary disease",

```

```

                                "pancreatic cancer"))+
xlab("diagnosis")+
ylab("age in years")

gp2 <- ggplot(data = data)+
  geom_boxplot(mapping = aes(x = diagnosis,
                             y = plasma_CA19_9,
                             group=diagnosis,
                             fill=diagnosis))+
  scale_x_discrete(labels=c("control",
                             "benign hepatobiliary disease",
                             "pancreatic cancer"))+
  xlab("diagnosis")+
  ylab("plasma_CA19_9 in U/ml")+
  ylim(0,500)

gp3 <- ggplot(data = data)+
  geom_boxplot(mapping = aes(x = diagnosis,
                             y = creatinine,
                             group=diagnosis,
                             fill=diagnosis))+
  scale_x_discrete(labels=c("control",
                             "benign hepatobiliary disease",
                             "pancreatic cancer"))+
  xlab("diagnosis")+
  ylab("creatinine in mg/ml")

gp4 <- ggplot(data = data)+
  geom_boxplot(mapping = aes(x = diagnosis,
                             y = LYVE1,
                             group=diagnosis,
                             fill=diagnosis))+
  scale_x_discrete(labels=c("control",
                             "benign hepatobiliary disease",
                             "pancreatic cancer"))+
  xlab("diagnosis")+
  ylab("LYVE1 in ng/ml")+
  ylim(0,17)

gp5 <- ggplot(data = data)+
  geom_boxplot(mapping = aes(x = diagnosis,
                             y = REG1B,
                             group=diagnosis,
                             fill=diagnosis))+
  scale_x_discrete(labels=c("control",
                             "benign hepatobiliary disease",
                             "pancreatic cancer"))+
  xlab("diagnosis")+
  ylab("REG1B in ng/ml")+
  ylim(0,600)

gp6 <- ggplot(data = data)+
  geom_boxplot(mapping = aes(x = diagnosis,

```

```

y = TFF1,
group=diagnosis,
fill=diagnosis))+
scale_x_discrete(labels=c("control",
                           "benign hepatobiliary disease",
                           "pancreatic cancer"))+

xlab("diagnosis")+
ylab("TFF1 in ng/ml")+
ylim(0,5000)

gp7 <- ggplot(data = data)+
  geom_boxplot(mapping = aes(x = diagnosis,
                             y = REG1A,
                             group=diagnosis,
                             fill=diagnosis))+
  scale_x_discrete(labels=c("control",
                           "benign hepatobiliary disease",
                           "pancreatic cancer"))+

xlab("diagnosis")+
ylab("REG1A in ng/ml")+
ylim(0,5000)

grid.arrange(gp1, gp2, gp3, gp4, gp5, gp6, gp7, nrow = 3)

```

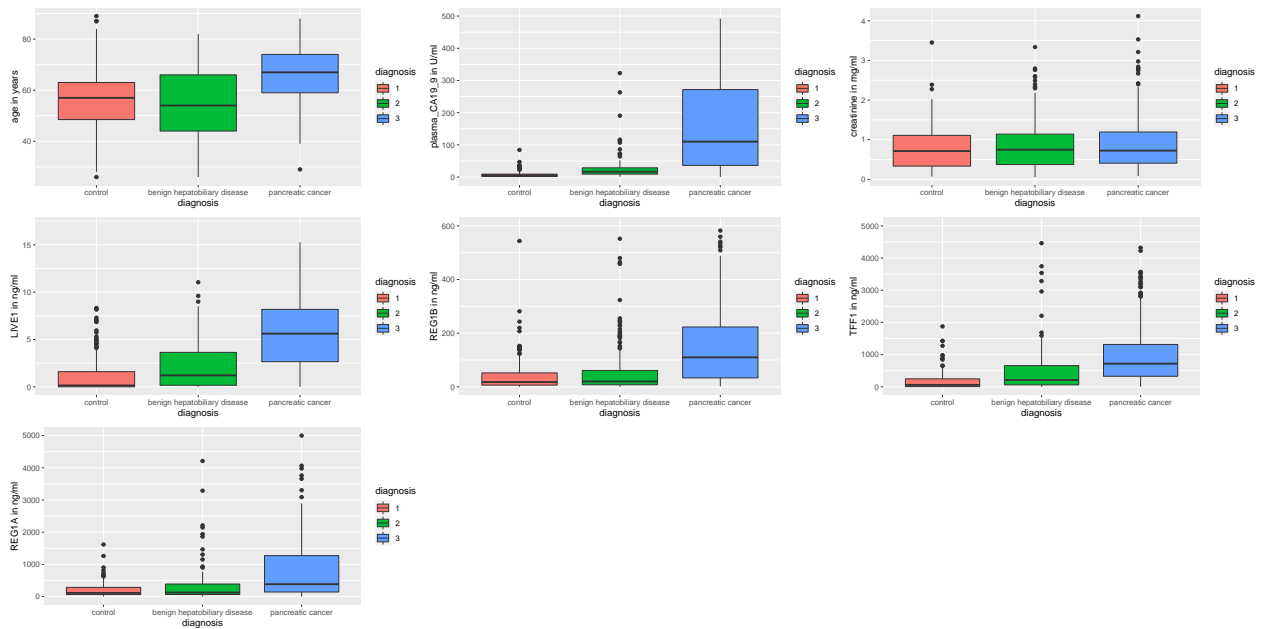


Figure 2: boxplots with added dimension (diagnosis)

The data still has many outliers but by many columns a pattern does emerge.

## Log transformation

Let's use statistical tests to test the skewness to see how imbalanced the data is.

```
skewness(data$age)
```

```
## [1] -0.2157312
```

```
skewness(data$plasma_CA19_9, na.rm = T)
```

```
## [1] 7.950382
```

```
skewness(data$creatinine)
```

```
## [1] 1.458965
```

```
skewness(data$LYVE1)
```

```
## [1] 1.386933
```

```
skewness(data$REG1B)
```

```
## [1] 3.316992
```

```
skewness(data$LYVE1)
```

```
## [1] 1.386933
```

```
skewness(data$REG1A, na.rm = T)
```

```
## [1] 4.425404
```

Here we see that everything is skewed greatly except age.

A way of dealing with this skewness is to apply a log transformation on the data due to the high positively skewed data.

```
trans <- as.data.frame(log(data$plasma_CA19_9))
names(trans) <- "plasma_CA19_9"
trans$creatinine <- log(data$creatinine)
trans$LYVE1 <- log(data$LYVE1)
trans$REG1B <- log(data$REG1B)
trans$TFF1 <- log(data$TFF1)
trans$REG1A <- log(data$REG1A + 1)
trans$diagnosis <- data$diagnosis
head(trans)
```

```
##   plasma_CA19_9 creatinine    LYVE1    REG1B    TFF1    REG1A diagnosis
## 1      2.459589  0.60552835 -0.1129233  3.969326  6.483539  7.141245         1
## 2           NA -0.02772069  0.7117653  4.548251  5.344668  5.435498         1
## 3      1.945910 -0.24796148 -1.9269684  4.628555  6.133704         NA         1
## 4      2.079442 -0.35493360 -5.8763945  4.103948  4.962495         NA         1
## 5      2.197225 -1.53762901 -7.0590899  4.182661  3.715716         NA         1
## 6           NA -0.16457988 -5.6860408  4.129165  4.090889         NA         1
```

Now having transformed the data lets see how this influences the distribution.

```
tgpl <- ggplot(data = trans)+
  geom_boxplot(mapping = aes(x = diagnosis,
                             y = plasma_CA19_9,
                             group=diagnosis,
                             fill=diagnosis))+
  scale_x_discrete(labels=c("control",
                             "benign hepatobiliary disease",
                             "pancreatic cancer"))+
```



```

xlab("diagnosis")+
ylab("plasma_CA19_9 in U/ml (log transformed)")

tgp2 <- ggplot(data = trans)+
  geom_boxplot(mapping = aes(x = diagnosis,
                             y = creatinine,
                             group=diagnosis,
                             fill=diagnosis))+
  scale_x_discrete(labels=c("control",
                             "benign hepatobiliary disease",
                             "pancreatic cancer"))+
  xlab("diagnosis")+
  ylab("creatinine in mg/ml (log transformed)")

tgp3 <- ggplot(data = trans)+
  geom_boxplot(mapping = aes(x = diagnosis,
                             y = LYVE1,
                             group=diagnosis,
                             fill=diagnosis))+
  scale_x_discrete(labels=c("control",
                             "benign hepatobiliary disease",
                             "pancreatic cancer"))+
  xlab("diagnosis")+
  ylab("LIVE1 in ng/ml (log transformed)")

tgp4 <- ggplot(data = trans)+
  geom_boxplot(mapping = aes(x = diagnosis,
                             y = REG1B,
                             group=diagnosis,
                             fill=diagnosis))+
  scale_x_discrete(labels=c("control",
                             "benign hepatobiliary disease",
                             "pancreatic cancer"))+
  xlab("diagnosis")+
  ylab("REG1B in ng/ml (log transformed)")

tgp5 <- ggplot(data = trans)+
  geom_boxplot(mapping = aes(x = diagnosis,
                             y = TFF1,
                             group=diagnosis,
                             fill=diagnosis))+
  scale_x_discrete(labels=c("control",
                             "benign hepatobiliary disease",
                             "pancreatic cancer"))+
  xlab("diagnosis")+
  ylab("TFF1 in ng/ml (log transformed)")

tgp6 <- ggplot(data = trans)+
  geom_boxplot(mapping = aes(x = diagnosis,
                             y = REG1A,
                             group=diagnosis,
                             fill=diagnosis))+
  scale_x_discrete(labels=c("control",

```

```

    "benign hepatobiliary disease",
    "pancreatic cancer"))+
  xlab("diagnosis")+
  ylab("REG1A in ng/ml (log transformed)")

grid.arrange(tgp1, tgp2, tgp3, tgp4, tgp5, tgp6, nrow = 3)

```

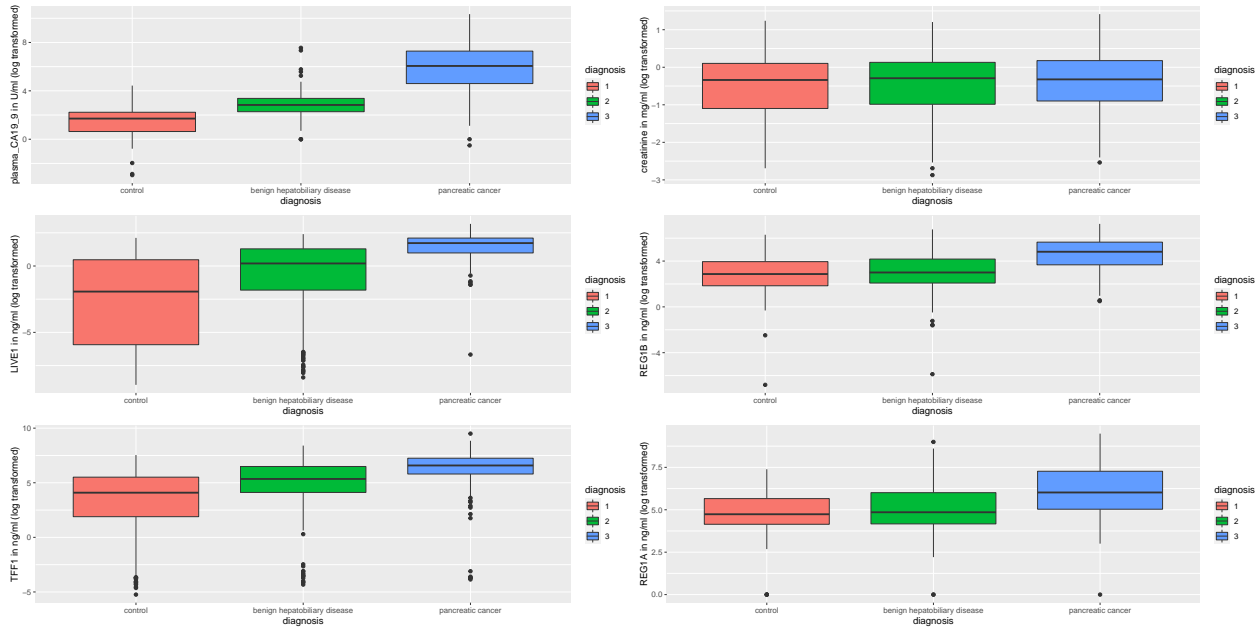


Figure 3: boxplots of transformed values

This data looks more normalized than before.

However it's good practise to test normality after transformations.

```

shapiro.test(trans$plasma_CA19_9)

##
## Shapiro-Wilk normality test
##
## data: trans$plasma_CA19_9
## W = NaN, p-value = NA

qq1 <- ggplot(trans, aes(sample = plasma_CA19_9, colour = diagnosis)) +
  stat_qq() +
  stat_qq_line() +
  ggtitle("plasma_CA19_9")

shapiro.test(trans$creatinine)

##
## Shapiro-Wilk normality test
##
## data: trans$creatinine
## W = 0.98211, p-value = 1.254e-06

```

```
qq2 <- ggplot(trans, aes(sample = creatinine, colour = diagnosis)) +
  stat_qq() +
  stat_qq_line() +
  ggtitle("creatinine")
```

```
shapiro.test(trans$LYVE1)
```

```
##
## Shapiro-Wilk normality test
##
## data: trans$LYVE1
## W = 0.81496, p-value < 2.2e-16
```

```
qq3 <- ggplot(trans, aes(sample = LYVE1, colour = diagnosis)) +
  stat_qq() +
  stat_qq_line() +
  ggtitle("LYVE1")
```

```
shapiro.test(trans$REG1B)
```

```
##
## Shapiro-Wilk normality test
##
## data: trans$REG1B
## W = 0.9695, p-value = 9.888e-10
```

```
qq4 <- ggplot(trans, aes(sample = REG1B, colour = diagnosis)) +
  stat_qq() +
  stat_qq_line() +
  ggtitle("REG1B")
```

```
shapiro.test(trans$TFF1)
```

```
##
## Shapiro-Wilk normality test
##
## data: trans$TFF1
## W = 0.82765, p-value < 2.2e-16
```

```
qq5 <- ggplot(trans, aes(sample = TFF1, colour = diagnosis)) +
  stat_qq() +
  stat_qq_line() +
  ggtitle("TFF1")
```

```
shapiro.test(trans$REG1A)
```

```
##
## Shapiro-Wilk normality test
##
## data: trans$REG1A
## W = 0.97007, p-value = 5.564e-06
```

```
qq6 <- ggplot(trans, aes(sample = REG1A, colour = diagnosis)) +
  stat_qq() +
  stat_qq_line() +
  ggtitle("REG1A")
```

```
grid.arrange(qq1, qq2, qq3, qq4, qq5, qq6, nrow = 2)
```

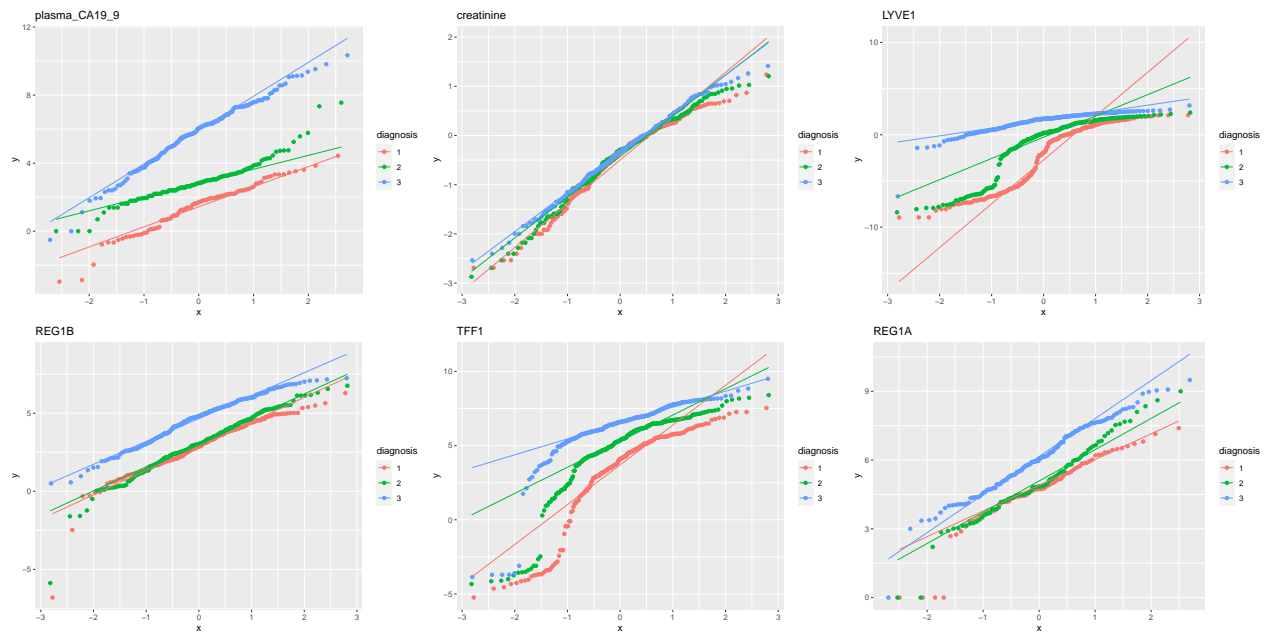


Figure 4: qqplots displaying normalcy

The data is despite the transformation still not fully normalised however we can still continue but this should be kept this in mind in case of future problems.

## Correlations

Now using the transformed data let's create a new dataframe.

```
new_data <- cbind(data[3:6], trans[2:6])
```

Using the new dataframe let's explore if the data is correlated.

```
matrix_data <- drop_na(new_data[,c(2, 5:9)])
cor_matrix <- cor(matrix_data)
heatmap(cor_matrix, scale = "column", col = heat.colors(5, rev = T))
legend(x="right", legend=c("full","strong", "medium", "minimal", "none"),fill=heat.colors(5))
title(main = "Heatmap depicting correlations")
```

REG1 A and B seem moderately correlated (0.7641084), otherwise no real strong correlation is observed.

Now we also should check if any variable is seemingly influential for the diagnosis so we can see later if the machine learning picks up on this.

```
new_data$has_cancer <- ifelse(new_data$diagnosis == 3, 1, 0)
new_data$has_cancer <- factor(new_data$has_cancer)
```

```
t.test(new_data$age ~ new_data$has_cancer)
```

```
##
## Welch Two Sample t-test
```

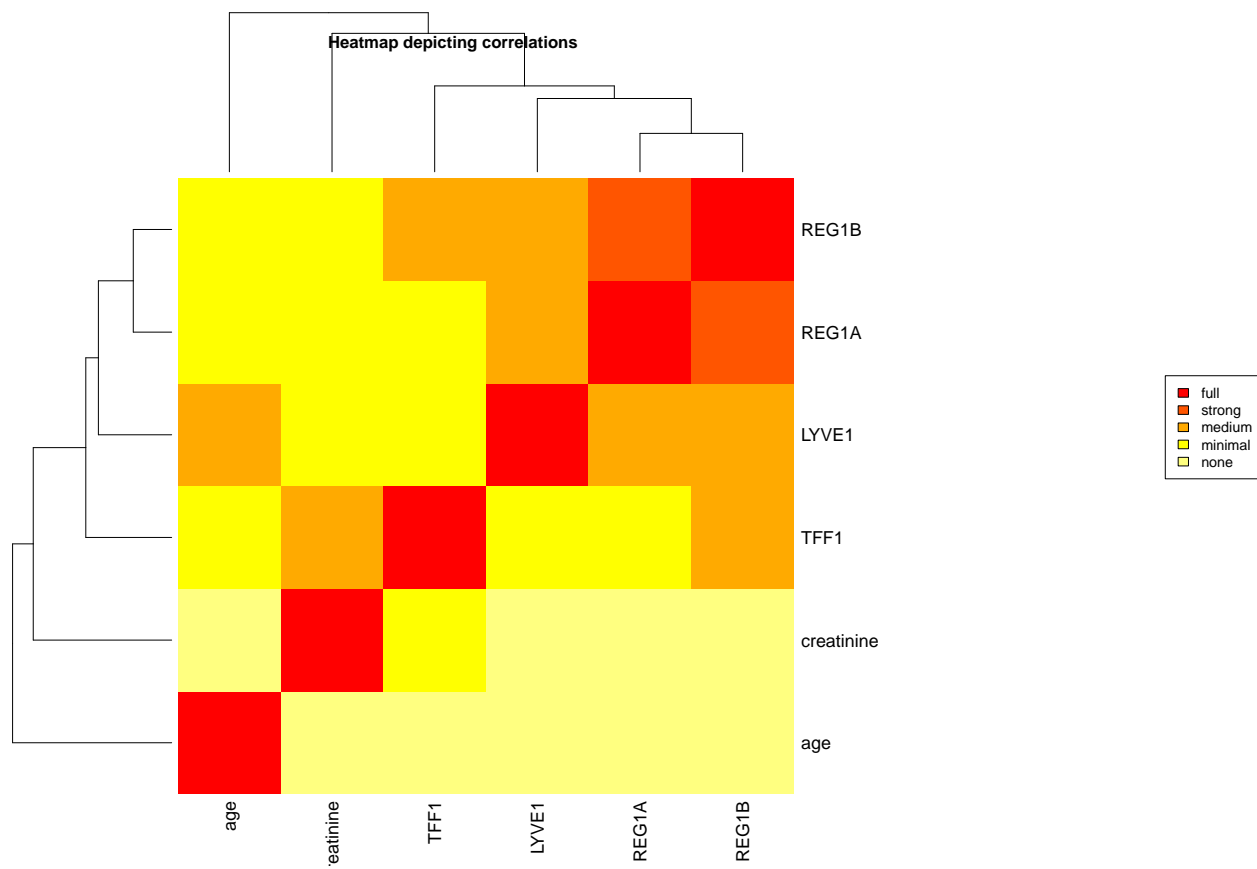


Figure 5: heatmap displaying correlation of values

```
##
## data: new_data$age by new_data$has_cancer
## t = -10.846, df = 473.94, p-value < 2.2e-16
## alternative hypothesis: true difference in means is not equal to 0
## 95 percent confidence interval:
## -12.656785 -8.774077
## sample estimates:
## mean in group 0 mean in group 1
## 55.46547 66.18090
t.test(new_data$creatinine ~ new_data$has_cancer)
```

```
##
## Welch Two Sample t-test
##
## data: new_data$creatinine by new_data$has_cancer
## t = -1.427, df = 411.11, p-value = 0.1543
## alternative hypothesis: true difference in means is not equal to 0
## 95 percent confidence interval:
## -0.2372607 0.0376788
## sample estimates:
## mean in group 0 mean in group 1
## -0.4793594 -0.3795684
t.test(new_data$LYVE1 ~ new_data$has_cancer)
```

```
##
## Welch Two Sample t-test
##
## data: new_data$LYVE1 by new_data$has_cancer
## t = -17.495, df = 520.79, p-value < 2.2e-16
## alternative hypothesis: true difference in means is not equal to 0
## 95 percent confidence interval:
## -3.536142 -2.822156
## sample estimates:
## mean in group 0 mean in group 1
## -1.746755 1.432394
t.test(new_data$REG1B ~ new_data$has_cancer)
```

```
##
## Welch Two Sample t-test
##
## data: new_data$REG1B by new_data$has_cancer
## t = -13.059, df = 456.65, p-value < 2.2e-16
## alternative hypothesis: true difference in means is not equal to 0
## 95 percent confidence interval:
## -1.982107 -1.463572
## sample estimates:
## mean in group 0 mean in group 1
## 2.910900 4.633739
t.test(new_data$TFF1 ~ new_data$has_cancer)
```

```
##
## Welch Two Sample t-test
##
```

```
## data: new_data$TFF1 by new_data$has_cancer
## t = -10.433, df = 532.25, p-value < 2.2e-16
## alternative hypothesis: true difference in means is not equal to 0
## 95 percent confidence interval:
## -2.679849 -1.830573
## sample estimates:
## mean in group 0 mean in group 1
## 3.910662 6.165872
t.test(new_data$REG1A ~ new_data$has_cancer)
```

```
##
## Welch Two Sample t-test
##
## data: new_data$REG1A by new_data$has_cancer
## t = -6.923, df = 300.15, p-value = 2.687e-11
## alternative hypothesis: true difference in means is not equal to 0
## 95 percent confidence interval:
## -1.5725494 -0.8764183
## sample estimates:
## mean in group 0 mean in group 1
## 4.850289 6.074772
```

No p-value except Creatinine seems to be small enough to not be statistically significant. We will expect to see this in the model.