

Thema09_EDA

Contents

Introduction	1
EDA	1
Codebook	1
Data exploration	4
Visualization	4

Introduction

Pancreatic ductal adenocarcinoma (PDAC) is one of the deadliest cancers. The chances of survival are increased when diagnosed in an early stage. However, PDAC shows symptoms when it already spread throughout the body. Most of the time, it's too late by then. There may be a way to detect PCAD in an early stage with a simple urine test, with the use of the following biomarkers: creatinine (Urinary biomarker of kidney function) LYVE1 (Urinary levels of Lymphatic vessel endothelial hyaluronan receptor 1, a protein that may play a role in tumor metastasis), REG1A and REG1B (Urinary levels of a protein that may be associated with pancreas regeneration.), and TFF1 (Urinary levels of Trefoil Factor 1, which may be related to regeneration and repair of the urinary tract)

the attributes in the data interesting for this research are the biomarker values mentioned before. There is also an attribute called diagnosis, in which the diagnosis of the sample is stated, where 1 means no PDAC, 2 means benign hepatobiliary disease (non cancerous, non harmful pancreatic condition), and 3 means that the sample has PDAC.

Research question: Is it possible to detect pancreatic cancer using values of the urinary biomarkers?

EDA

Codebook

```
myData <- read.csv("Data/Debernardi et al 2020 data.csv")

columns <- colnames(myData)
type <- c("character", "character", "character", "double", "character", "double", "logical", "logical",
unit <- c(NA, NA, NA, "years", "F/M", NA, NA, NA, "U/ml", "mg/ml", "ng/ml", "ng/ml", "ng/ml", "ng/ml")
descriptions = c("Unique string identifying each subject", "Cohort 1, previously used samples; Cohort 2
codebook <- data.frame(columns, type, unit, descriptions)
write.csv(codebook, "Codebook.csv", row.names = FALSE)
```

```
knitr::kable(codebook, caption="Table 1: The Codebook", format = 'latex') %>%  
  kable_styling(full_width = F) %>%  
  column_spec(1, bold = T) %>%  
  column_spec(4, width = "22em")
```

Table 1: Table 1: The Codebook

columns	type	unit	descriptions
sample_id	character	NA	Unique string identifying each subject
patient_cohort	character	NA	Cohort 1, previously used samples; Cohort 2, newly added samples
sample_origin	character	NA	BPTB: Barts Pancreas Tissue Bank, London, UK; ESP: Spanish National Cancer Research Centre, Madrid, Spain; LIV: Liverpool University, UK; UCL: University College London, UK
age	double	years	Age in years
sex	character	F/M	M = male, F = female
diagnosis	double	NA	1 = control (no pancreatic disease), 2 = benign hepatobiliary disease (119 of which are chronic pancreatitis); 3 = Pancreatic ductal adenocarcinoma, i.e. pancreatic cancer
stage	logical	NA	For those with pancreatic cancer, what stage was it? One of IA, IB, IIA, IIIB, III, IV
benign_sample_diagnosis	logical	NA	For those with a benign, non-cancerous diagnosis, what was the diagnosis?
plasma_CA19_9	double	U/ml	Blood plasma levels of CA 19-9 monoclonal antibody that is often elevated in patients with pancreatic cancer. Only assessed in 350 patients (one goal of the study was to compare various CA 19-9 cutpoints from a blood sample to the model developed using urinary samples).
creatinine	double	mg/ml	Urinary biomarker of kidney function
LYVE1	double	ng/ml	Urinary levels of Lymphatic vessel endothelial hyaluronan receptor 1, a protein that may play a role in tumor metastasis
REG1B	double	ng/ml	Urinary levels of a protein that may be associated with pancreas regeneration.
TFF1	double	ng/ml	Urinary levels of Trefoil Factor 1, which may be related to regeneration and repair of the urinary tract
REG1A	double	ng/ml	Urinary levels of a protein that may be associated with pancreas regeneration. Only assessed in 306 patients (one goal of the study was to assess REG1B vs REG1A)

Table 2: Table 2: Values of biomarkers and the diagnosis

age	diagnosis	plasma_CA19_9	creatinine	LYVE1	REG1B	TFF1	REG1A
33	1	11.7	1.83222	0.8932192	52.94884	654.2822	1262.000
81	1	NA	0.97266	2.0375850	94.46703	209.4882	228.407
51	1	7.0	0.78039	0.1455889	102.36600	461.1410	NA
61	1	8.0	0.70122	0.0028049	60.57900	142.9500	NA
62	1	9.0	0.21489	0.0008596	65.54000	41.0880	NA
53	1	NA	0.84825	0.0033930	62.12600	59.7930	NA
70	1	NA	0.62205	0.1743808	152.27700	117.5160	NA
58	1	11.0	0.89349	0.0035740	3.73000	40.2940	NA
59	1	NA	0.48633	0.0019453	7.02100	26.7820	NA
56	1	24.0	0.61074	0.2787785	83.92800	19.1850	NA
77	1	NA	0.29406	0.0011762	6.21800	28.2970	NA
71	1	23.0	1.05183	0.8603368	243.08200	608.2840	NA
49	1	NA	0.85956	1.4163140	151.83077	74.1899	505.571
53	1	7.0	1.91139	1.5167730	150.89000	590.6860	NA
56	1	12.0	0.91611	0.5996449	93.81100	93.5760	NA

Data exploration

Visualization

To look at the data and any possible missing data, here's a table containing only the relevant columns for this research.

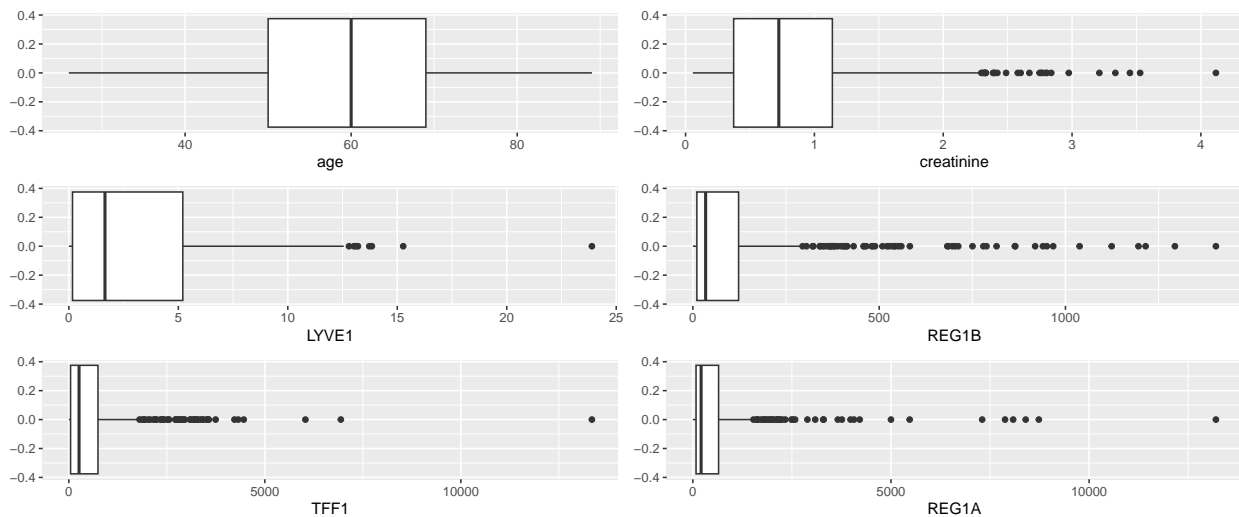
```
relevant <- myData[c(4, 6, 9:14)]
knitr::kable(head(relevant, n=15), caption="Table 2: Values of biomarkers and the diagnosis") %>%
  kable_styling(position = 'center')
```

As is it obvious to see, the REG1A column contains a lot of missing values. The question is exactly how is this possible and what does it mean for the rest of the data? The first and most simple solution is that we don't need this column at all, since the REG1B is similar in function and containing all of the data. Another solution that's not very logical, is to use only the rows where there is a value in the REG1A column. The most logical solution in my opinion is to look at every column separately, and then remove the missing data. When the cleaned data then is plotted, it's possible to look at all the columns together and see if there is any correlation or trend to figure out.

knowing this, let's have a look at the important raw data with the use of boxplots.

```
p1 <- ggplot(myData, aes(x=age)) +
  geom_boxplot()
p2 <- ggplot(myData, aes(x=creatinine)) +
  geom_boxplot()
p3 <- ggplot(myData, aes(x=LYVE1)) +
  geom_boxplot()
p4 <- ggplot(myData, aes(x=REG1B)) +
  geom_boxplot()
p5 <- ggplot(myData, aes(x=TFF1)) +
  geom_boxplot()
p6 <- ggplot(myData, aes(x=REG1A)) +
  geom_boxplot()

grid.arrange(p1, p2, p3, p4, p5, p6, nrow = 3)
```



as you can see, there are also a lot of outliers which we need to check. If we remove the outliers and clean up the data, without removing the missing values of REG1A, we'll get something like this.

```
Q <- quantile(myData$creatinine, probs=c(.25, .75), na.rm = FALSE)
iqr <- IQR(myData$creatinine)
up <- Q[2]+1.5*iqr
low <- Q[1]-1.5*iqr

myDataNO <- subset(myData, myData$creatinine > low & myData$creatinine < up)

Q <- quantile(myDataNO$LYVE1, probs=c(.25, .75), na.rm = FALSE)
iqr <- IQR(myDataNO$LYVE1)
```

```

up <- Q[2]+1.5*iqr
low <- Q[1]-1.5*iqr

myDataNO <- subset(myDataNO, myDataNO$LYVE1 > low & myDataNO$LYVE1 < up)

Q <- quantile(myDataNO$REG1B, probs=c(.25, .75), na.rm = FALSE)
iqr <- IQR(myDataNO$REG1B)
up <- Q[2]+1.5*iqr
low <- Q[1]-1.5*iqr

myDataNO <- subset(myDataNO, myDataNO$REG1B > low & myDataNO$REG1B < up)

Q <- quantile(myDataNO$TFF1, probs=c(.25, .75), na.rm = FALSE)
iqr <- IQR(myDataNO$TFF1)
up <- Q[2]+1.5*iqr
low <- Q[1]-1.5*iqr

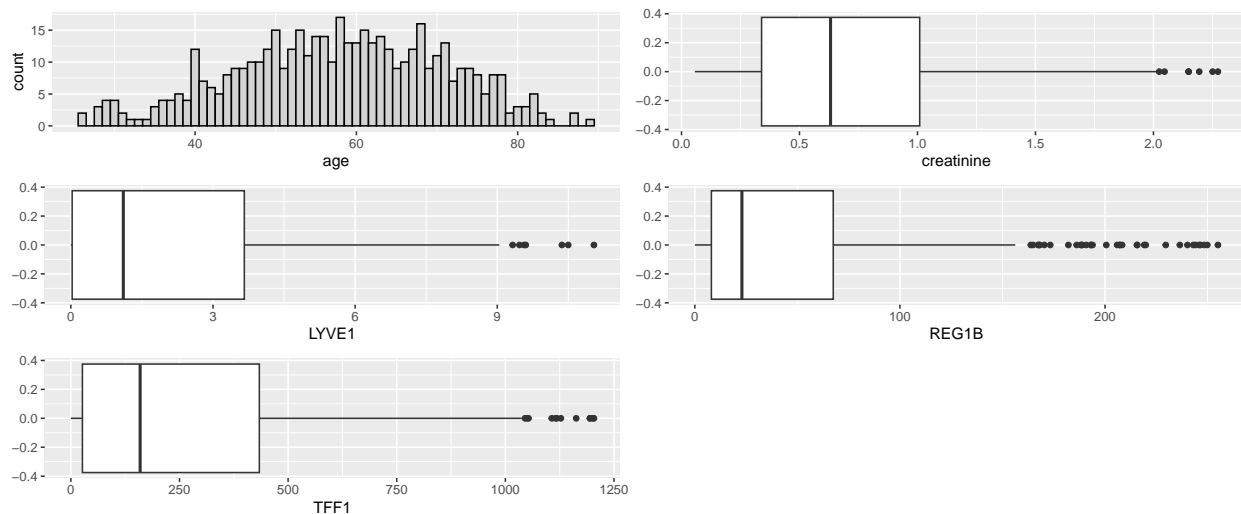
myDataNO <- subset(myDataNO, myDataNO$TFF1 > low & myDataNO$TFF1 < up)

p1 <- ggplot(myDataNO, aes(x=age)) +
  geom_histogram(fill = "lightgrey", col = "black", binwidth = 1)

p2 <- ggplot(myDataNO, aes(x=creatinine)) +
  geom_boxplot()
p3 <- ggplot(myDataNO, aes(x=LYVE1)) +
  geom_boxplot()
p4 <- ggplot(myDataNO, aes(x=REG1B)) +
  geom_boxplot()
p5 <- ggplot(myDataNO, aes(x=TFF1)) +
  geom_boxplot()

grid.arrange(p1, p2, p3, p4, p5, nrow = 3)

```

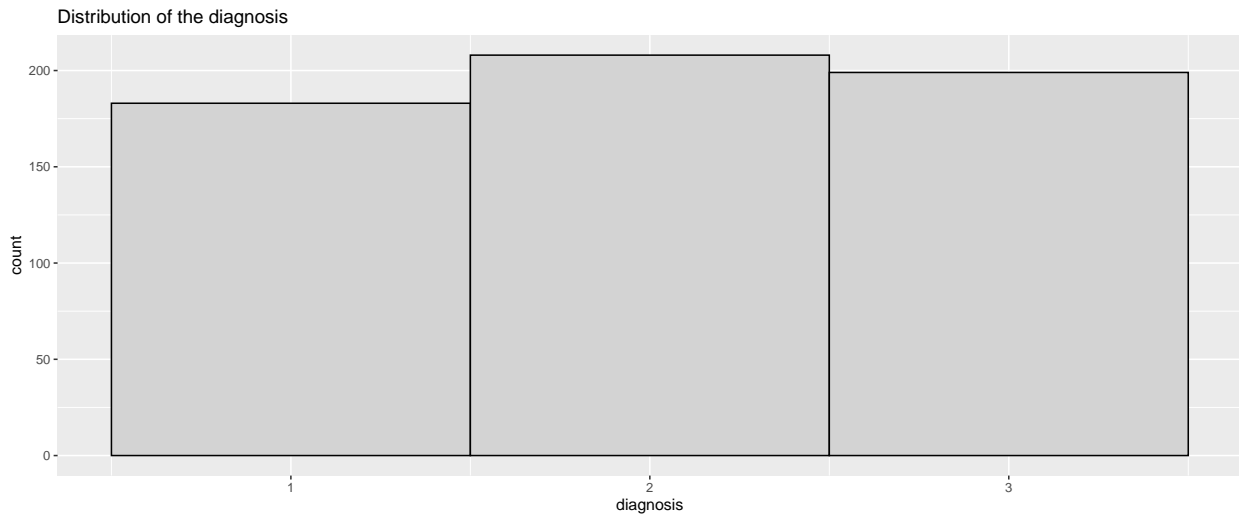


This looks a lot cleaner, however it is not wise to remove these outliers as they may contain important data in a later stage.

Before proceeding, it is quite helpful to look at the distribution of the diagnosis column, to see if they're

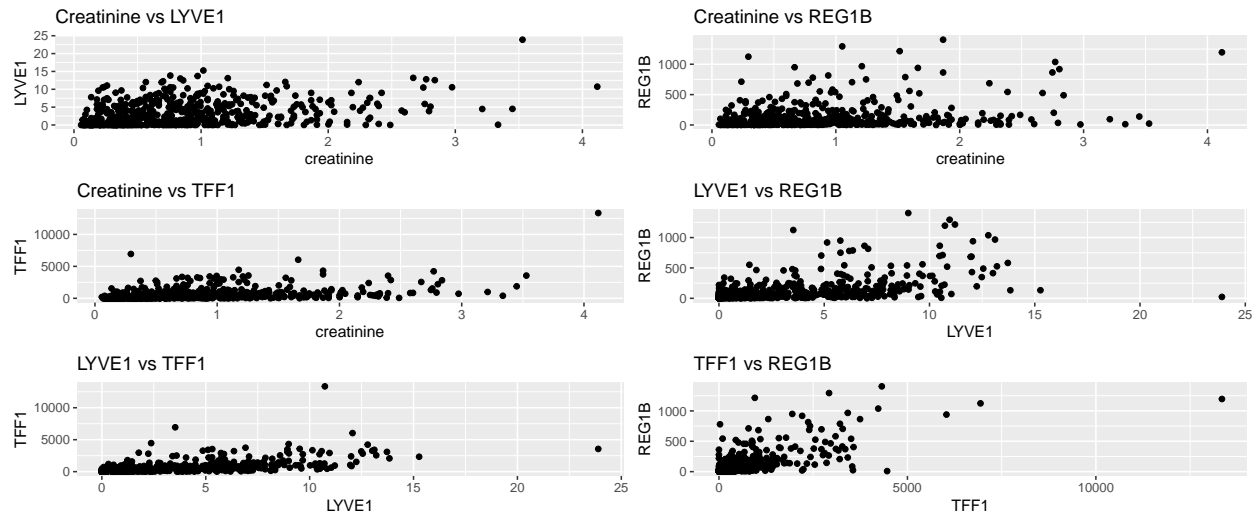
evenly distributed or not.

```
ggplot(myData, aes(x=diagnosis)) +  
  geom_histogram(fill = "lightgrey", col = "black", binwidth = 1) +  
  ggtitle("Distribution of the diagnosis")
```



Now we can check if the different biomarkers correlate with each other in any way.

```
p1 <- ggplot(myData, aes(x=creatinine, y=LYVE1)) +  
  geom_point() +  
  ggtitle("Creatinine vs LYVE1")  
p2 <- ggplot(myData, aes(x=creatinine, y=REG1B)) +  
  geom_point() +  
  ggtitle("Creatinine vs REG1B")  
p3 <- ggplot(myData, aes(x=creatinine, y=TFF1)) +  
  geom_point() +  
  ggtitle("Creatinine vs TFF1")  
p4 <- ggplot(myData, aes(x=LYVE1, y=REG1B)) +  
  geom_point() +  
  ggtitle("LYVE1 vs REG1B")  
p5 <- ggplot(myData, aes(x=LYVE1, y=TFF1)) +  
  geom_point() +  
  ggtitle("LYVE1 vs TFF1")  
p6 <- ggplot(myData, aes(x=TFF1, y=REG1B)) +  
  geom_point() +  
  ggtitle("TFF1 vs REG1B")  
  
grid.arrange(p1, p2, p3, p4, p5, p6, nrow=3)
```



I think it's safe to assume that none of these correlate with each other in any way. One thing to notice however, is that LYVE1, REG1B and TFF1 have a lot of values close to zero. We can figure out if this is important or not in a future stage of the research.

To ensure that all the values are completely independant we can perform a principal component analysis.

What we should do now, is look at the values of those urinary levels with each diagnosis and see if there are any obvious differences.

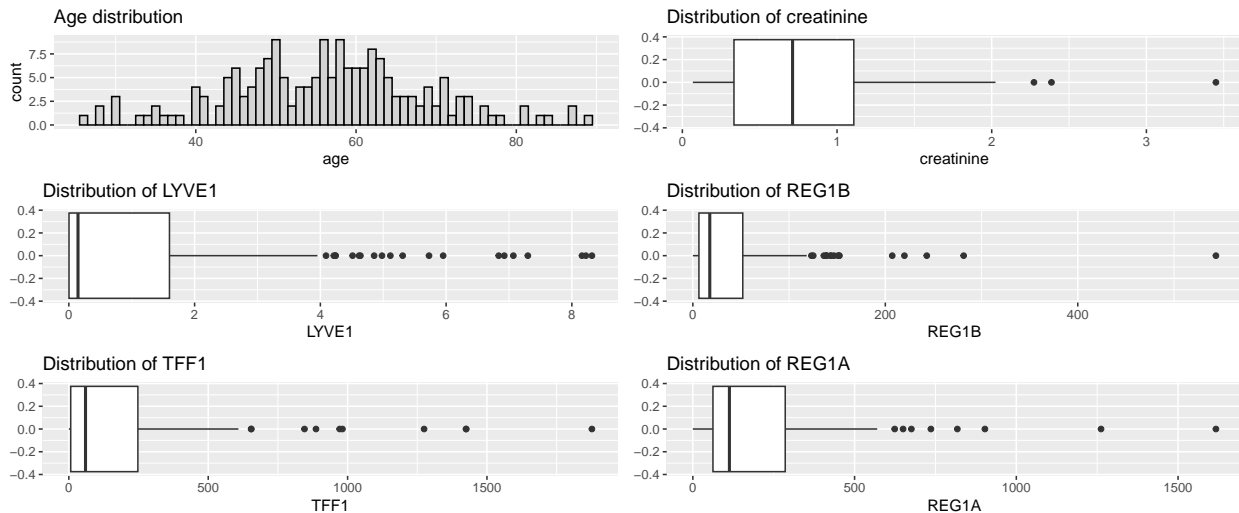
Let's start with the samples that do not have PDAC

```
myData1 <- subset(myData, myData$diagnosis == 1)

p1 <- ggplot(myData1, aes(x=age)) +
  geom_histogram(fill = "lightgrey", col = "black", binwidth = 1) +
  ggtitle("Age distribution")

p2 <- ggplot(myData1, aes(x=creatinine)) +
  geom_boxplot() +
  ggtitle("Distribution of creatinine")
p3 <- ggplot(myData1, aes(x=LYVE1)) +
  geom_boxplot() +
  ggtitle("Distribution of LYVE1")
p4 <- ggplot(myData1, aes(x=REG1B)) +
  geom_boxplot() +
  ggtitle("Distribution of REG1B")
p5 <- ggplot(myData1, aes(x=TFF1)) +
  geom_boxplot() +
  ggtitle("Distribution of TFF1")
p6 <- ggplot(myData1, aes(x=REG1A)) +
  geom_boxplot() +
  ggtitle("Distribution of REG1A")

grid.arrange(p1, p2, p3, p4, p5, p6, nrow = 3)
```

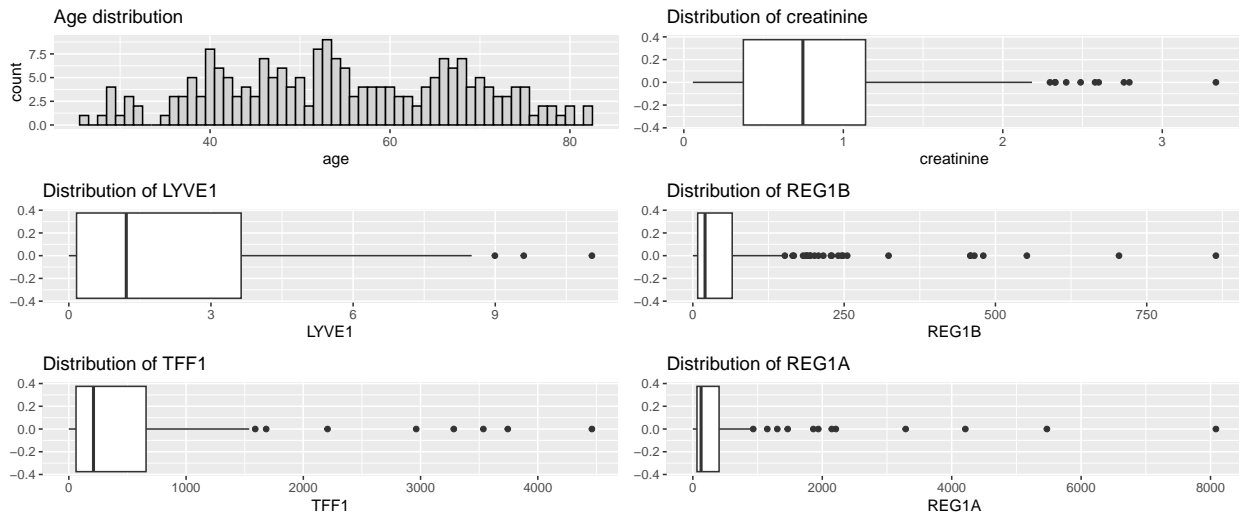
Here are the values of the samples with non-cancerous pancreatic conditions.

```
myData2 <- subset(myData, myData$diagnosis == 2)

p1 <- ggplot(myData2, aes(x=age)) +
  geom_histogram(fill = "lightgrey", col = "black", binwidth = 1) +
  ggtitle("Age distribution")

p2 <- ggplot(myData2, aes(x=creatinine)) +
  geom_boxplot() +
  ggtitle("Distribution of creatinine")
p3 <- ggplot(myData2, aes(x=LYVE1)) +
  geom_boxplot() +
  ggtitle("Distribution of LYVE1")
p4 <- ggplot(myData2, aes(x=REG1B)) +
  geom_boxplot() +
  ggtitle("Distribution of REG1B")
p5 <- ggplot(myData2, aes(x=TFF1)) +
  geom_boxplot() +
  ggtitle("Distribution of TFF1")
p6 <- ggplot(myData2, aes(x=REG1A)) +
  geom_boxplot() +
  ggtitle("Distribution of REG1A")

grid.arrange(p1, p2, p3, p4, p5, p6, nrow = 3)
```



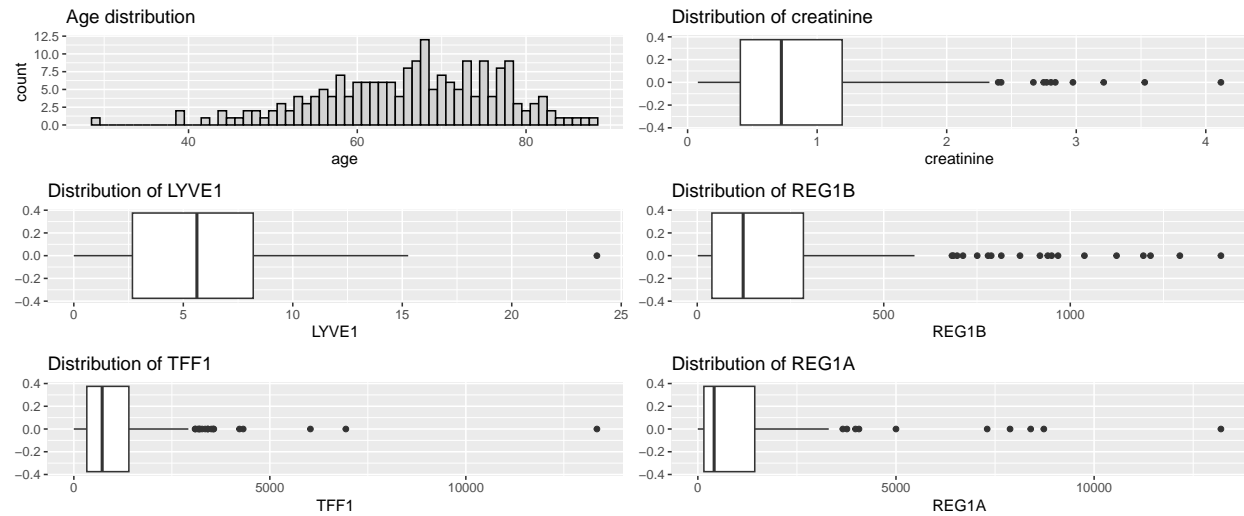
And finally the samples that do have PDAC.

```
myData3 <- subset(myData, myData$diagnosis == 3)

p1 <- ggplot(myData3, aes(x=age)) +
  geom_histogram(fill = "lightgrey", col = "black", binwidth = 1) +
  ggtitle("Age distribution")

p2 <- ggplot(myData3, aes(x=creatinine)) +
  geom_boxplot() +
  ggtitle("Distribution of creatinine")
p3 <- ggplot(myData3, aes(x=LYVE1)) +
  geom_boxplot() +
  ggtitle("Distribution of LYVE1")
p4 <- ggplot(myData3, aes(x=REG1B)) +
  geom_boxplot() +
  ggtitle("Distribution of REG1B")
p5 <- ggplot(myData3, aes(x=TFF1)) +
  geom_boxplot() +
  ggtitle("Distribution of TFF1")
p6 <- ggplot(myData3, aes(x=REG1A)) +
  geom_boxplot() +
  ggtitle("Distribution of REG1A")

grid.arrange(p1, p2, p3, p4, p5, p6, nrow = 3)
```



At first glance, there are notable differences but we can further investigate this in future steps of this research.

```
write.csv(myData, "Data/DataCleaned.csv")
```