

UBS - Code Assessment

Exercise 2 - Outlier detection and statistics

```
In [4]: library(data.table)
library(stringr)
library(ggplot2)
```

2.A) Load the data of the csv file into a data format of your choice

I) Data Gathering

```
In [10]: raw_data <- read.csv("~/Users/alvarosanchezfernandez/Documents/Programming/Raw Data/data.csv")
raw_data
```

Description	X199501	X199502	X199503	X199504	X199505	X199506	X199507	X199508
GDP Level	NA	NA	179700.0	NA	NA	177070.000	NA	1.7e+08
LIBOR 1M %	3.8125	3.625	3.5	3.4375	3.25	3.125	2.625	2.875
Equities Index 1 (index points)	3017.3000	3041.200	3143.1	3220.4000	3340.60	3323.700	3449.900	3509.400
Equities Index 2 (index points)	NA	NA	NA	NA	NA	NA	NA	100.000
Equities Index 3 (index points)	NA	NA	NA	NA	NA	NA	NA	NA

II) Data Cleaning

```
In [11]: # Create Dataframe Organized Vertically
df <- transpose(raw_data)

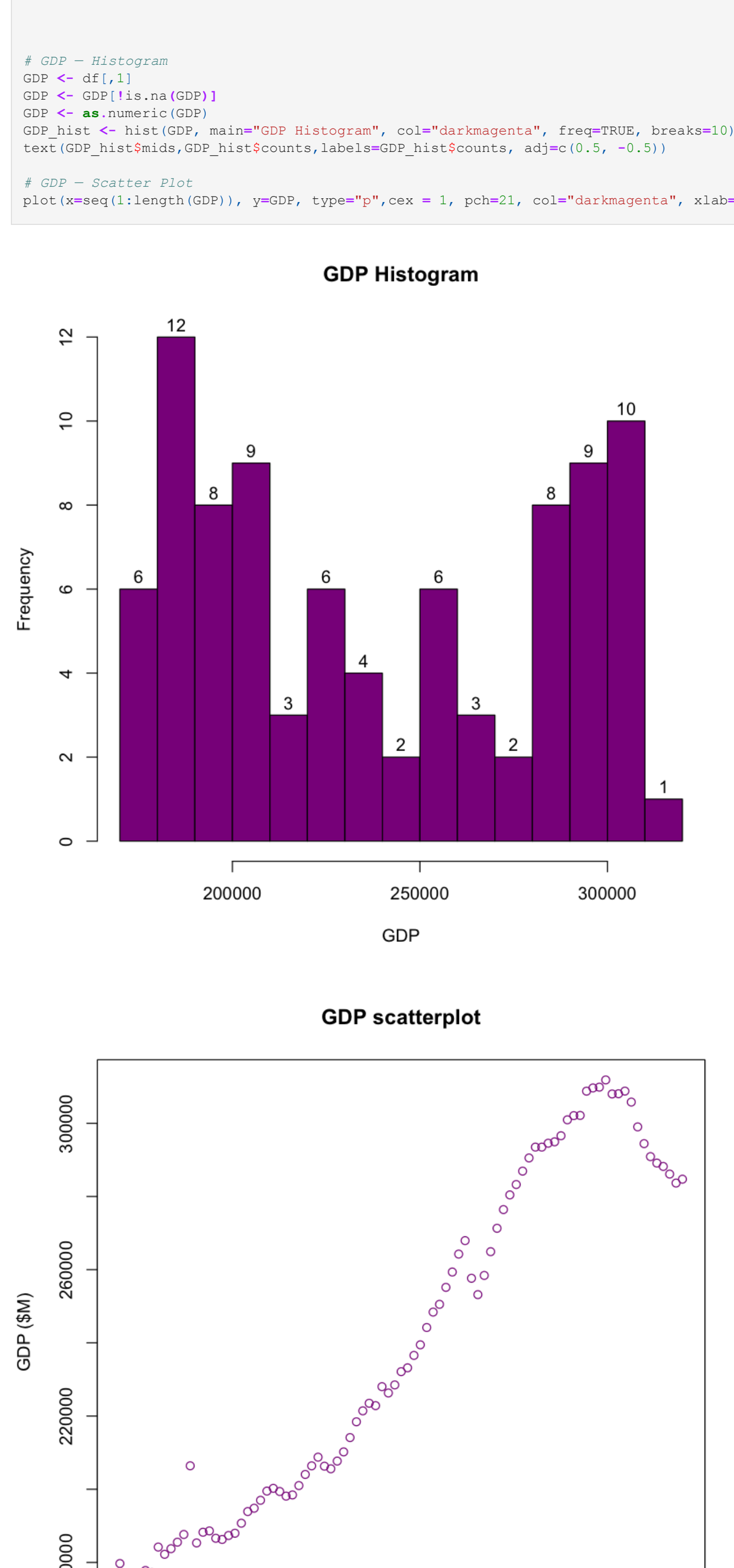
# Set Row Names
rownames(df) <- colnames(raw_data)

# Set Column Names
colnames(df) <- df[1,]
df <- df[-1, , ]
df

# Eliminate NA ?

# GDP - Histogram
GDP <- df[,2]
GDP <- na.omit(GDP)
GDP <- as.numeric(GDP)
GDP_hist <- hist(GDP, main="GDP Histogram", col="darkmagenta", freq=TRUE, breaks=10)
text(GDP_hist$mid$ids,GDP_hist$counts,labels=GDP_hist$counts,adj=c(0.5, -0.5))

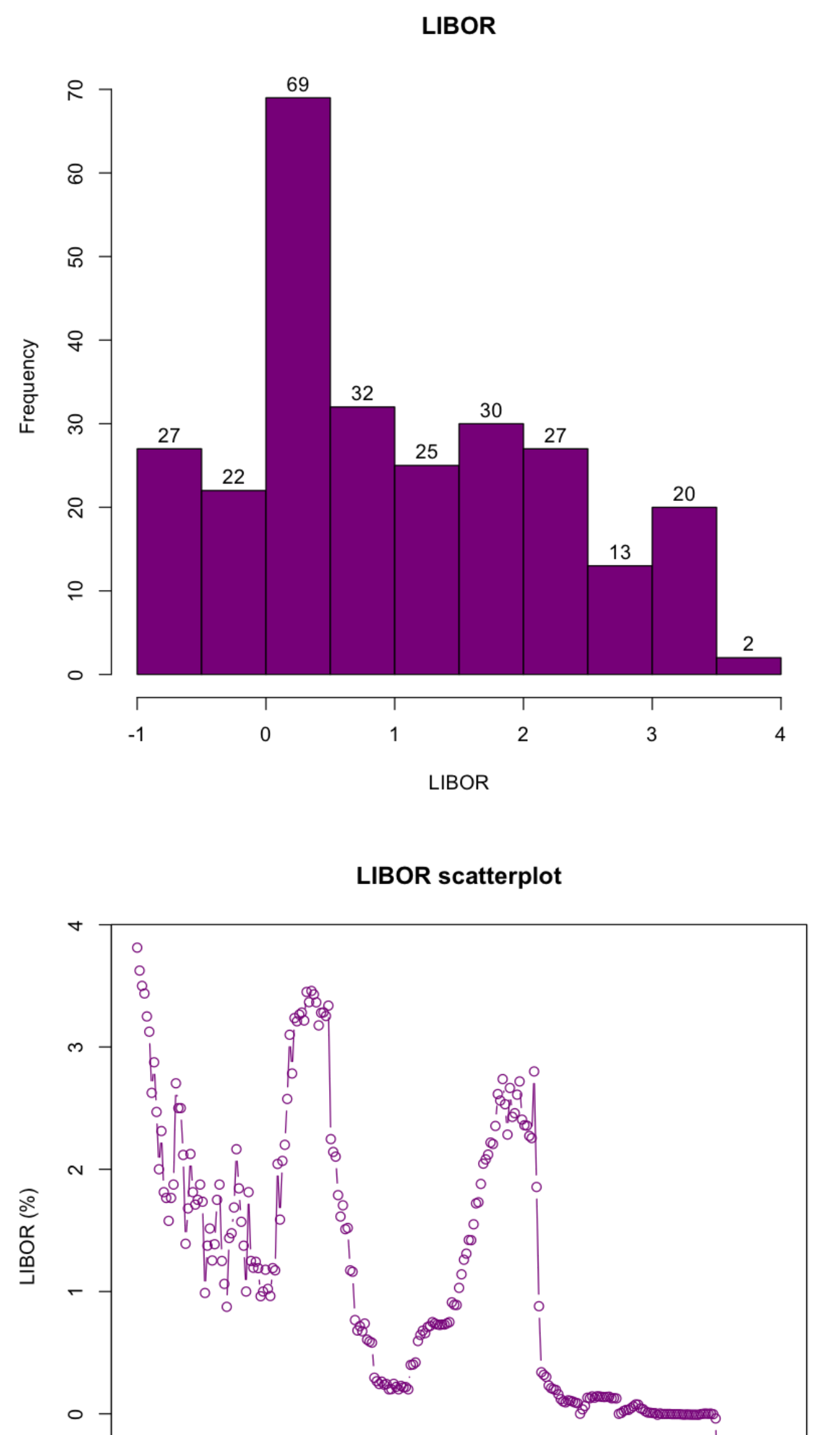
# GDP - Scatter Plot
plot(x=seq(1:length(GDP)), y=GDP, type="p", cex = 1, pch=21, col="darkmagenta", xlab="Observations", ylab="GDP ($M)")
```



II) LIBOR - Detect and Treat Outliers

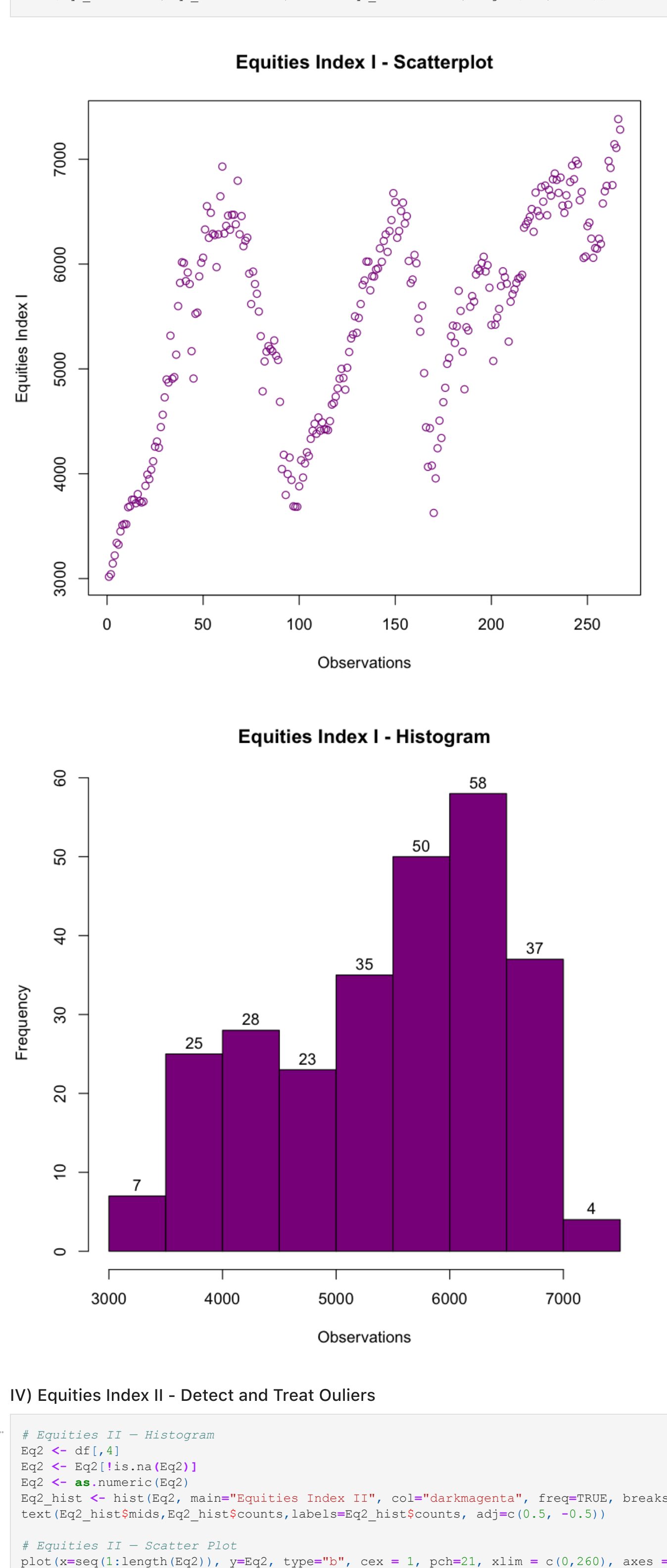
```
In [12]: # LIBOR - Histogram
LIBOR <- df[,3]
LIBOR <- na.omit(LIBOR)
LIBOR <- as.numeric(LIBOR)
LIBOR_hist <- hist(LIBOR, main="LIBOR", col="darkmagenta", freq=TRUE, breaks=10)
text(LIBOR_hist$mid$ids,LIBOR_hist$counts,labels=LIBOR_hist$counts,adj=c(0.5, -0.5))

# LIBOR - Scatter Plot
plot(x=seq(1:length(LIBOR)), y=LIBOR, type="p", cex = 1, pch=21, col="darkmagenta", xlab="Observations", ylab="LIBOR (%)")
```



III) Equities Index I - Detect and Treat Outliers

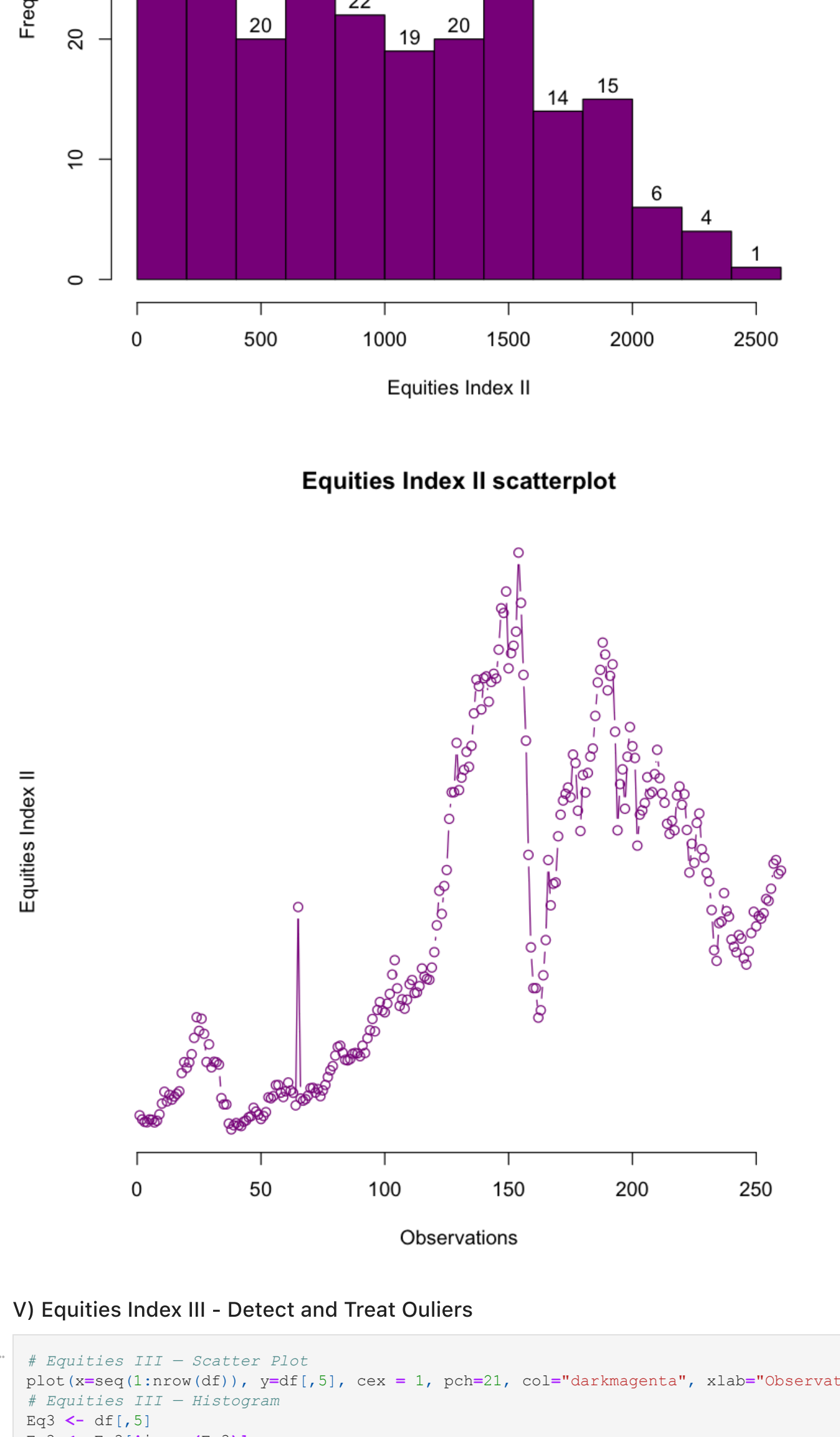
```
In [23]: # Equities I - Scatter Plot
plot(x=seq(1:nrow(df)), y=df[,3], cex = 1, pch=21, col="darkmagenta", xlab="Observations", ylab="Equities Index I")
# Equities I - Histogram
Eq1 <- df[,3]
Eq1 <- na.omit(Eq1)
Eq1 <- as.numeric(Eq1)
Eq1_hist <- hist(Eq1, main="Equities Index I - Histogram", col="darkmagenta", freq=TRUE, breaks=10)
text(Eq1_hist$mid$ids,Eq1_hist$counts,labels=Eq1_hist$counts,adj=c(0.5, -0.5))
```



IV) Equities Index II - Detect and Treat Outliers

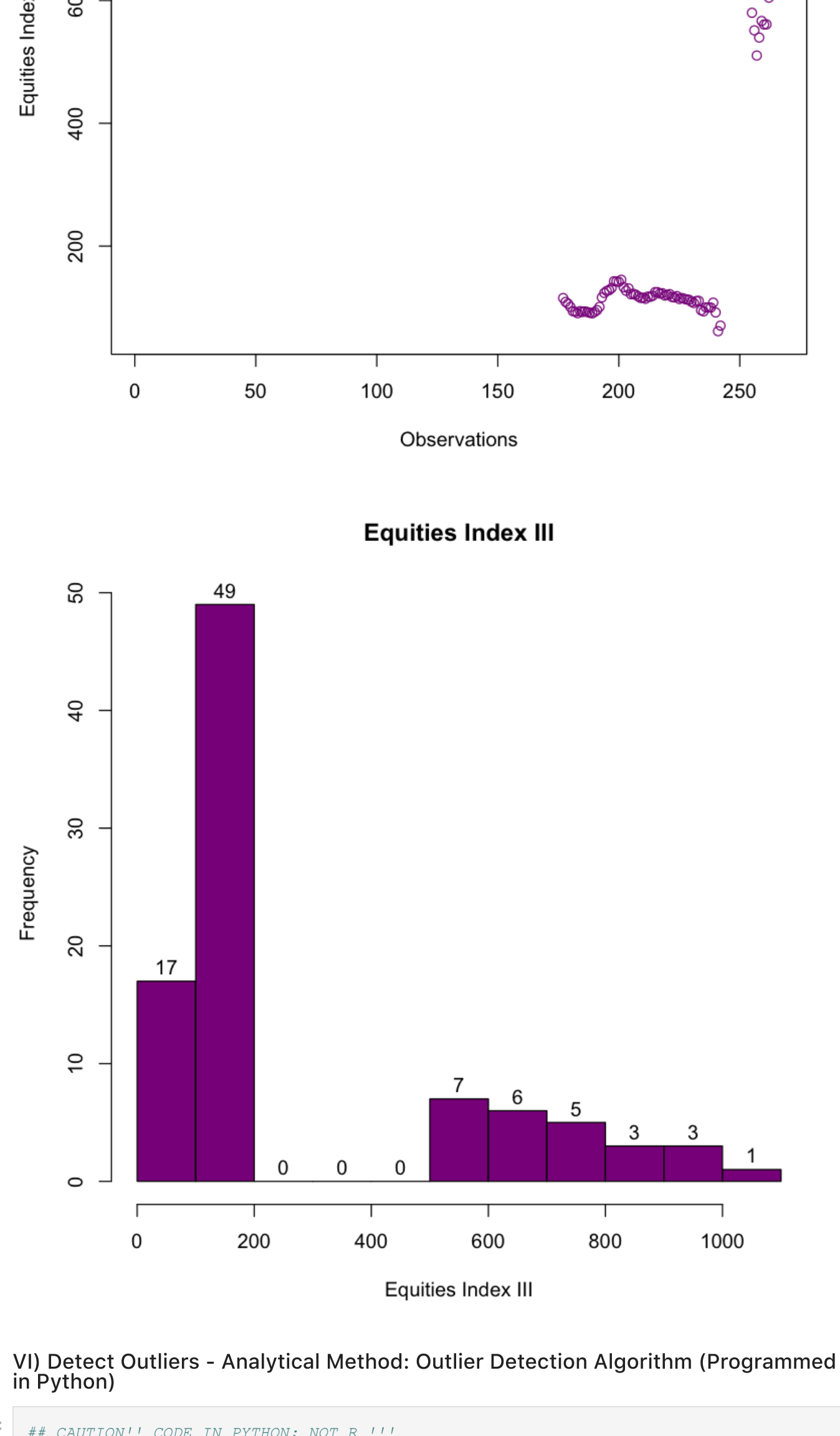
```
In [148]: # Equities II - Histogram
Eq2 <- df[,4]
Eq2 <- na.omit(Eq2)
Eq2 <- as.numeric(Eq2)
Eq2_hist <- hist(Eq2, main="Equities Index II", col="darkmagenta", freq=TRUE, breaks=10)
text(Eq2_hist$mid$ids,Eq2_hist$counts,labels=Eq2_hist$counts,adj=c(0.5, -0.5))

# Equities II - Scatter Plot
plot(x=seq(1:length(Eq2)), y=Eq2, type="p", cex = 1, pch=21, xlim = c(0,260), axes = FALSE, las = 1)
```



V) Equities Index III - Detect and Treat Outliers

```
In [232]: # Equities III - Scatter Plot
plot(x=seq(1:nrow(df)), y=df[,5], cex = 1, pch=21, col="darkmagenta", xlab="Observations", ylab="Equities Index III")
# Equities III - Histogram
Eq3 <- df[,5]
Eq3 <- na.omit(Eq3)
Eq3 <- as.numeric(Eq3)
Eq3_hist <- hist(Eq3, main="Equities Index III", col="darkmagenta", freq=TRUE, breaks=10)
text(Eq3_hist$mid$ids,Eq3_hist$counts,labels=Eq3_hist$counts,adj=c(0.5, -0.5))
```



VI) Detect Outliers - Analytical Method: Outlier Detection Algorithm (Programmed in Python)

```
In [2]: ## CAUTION!! CODE IN PYTHON, NOT R !!!

import pandas as pd
LIBOR_outliers = df.values.tolist()

GDP=pd.read_csv("~/Users/alvarosanchezfernandez/Documents/Programming/Raw Data/data.csv")
GDP=GDP[GDP['LIBOR']!=0]
GDP_vect=[]
GDP_vect.append(GDP['GDP'])
GDP_vect.append(GDP['LIBOR'])
GDP_vect.append(GDP['Equities Index I'])
GDP_vect.append(GDP['Equities Index II'])
GDP_vect.append(GDP['Equities Index III'])

LIBOR=pd.read_csv("~/Users/alvarosanchezfernandez/Documents/Programming/Raw Data/data.csv")
LIBOR=LIBOR[LIBOR['LIBOR']!=0]
LIBOR_vect=[]
LIBOR_vect.append(LIBOR['LIBOR'])
LIBOR_vect.append(LIBOR['Equities Index I'])
LIBOR_vect.append(LIBOR['Equities Index II'])
LIBOR_vect.append(LIBOR['Equities Index III'])

Eq1=pd.read_csv("~/Users/alvarosanchezfernandez/Documents/Programming/Raw Data/data.csv")
Eq1=Eq1[Eq1['Equities Index I']!=0]
Eq1_vect=[]
Eq1_vect.append(Eq1['Equities Index I'])
Eq1_vect.append(Eq1['Equities Index II'])
Eq1_vect.append(Eq1['Equities Index III'])

Eq2=pd.read_csv("~/Users/alvarosanchezfernandez/Documents/Programming/Raw Data/data.csv")
Eq2=Eq2[Eq2['Equities Index II']!=0]
Eq2_vect=[]
Eq2_vect.append(Eq2['Equities Index II'])
Eq2_vect.append(Eq2['Equities Index III'])

Eq3=pd.read_csv("~/Users/alvarosanchezfernandez/Documents/Programming/Raw Data/data.csv")
Eq3=Eq3[Eq3['Equities Index III']!=0]
Eq3_vect=[]
Eq3_vect.append(Eq3['Equities Index III'])

Eq4=pd.read_csv("~/Users/alvarosanchezfernandez/Documents/Programming/Raw Data/data.csv")
Eq4=Eq4[Eq4['Equities Index III']!=0]
Eq4_vect=[]
Eq4_vect.append(Eq4['Equities Index III'])

# Outlier Detection Algorithm
def outlier_detection(dataset, name):
    # Configuration of x-axis scale: x=range(y)
    # between consecutive elements of the dataset
    dx = (max(dataset) - min(dataset)) / (len(dataset))

    # III) Calculation of difference (delta)
    # between consecutive elements of the dataset
    dataset_delta = []
    for i in range(len(dataset)-1):
        dataset_delta.append(abs(dataset[i]-dataset[i+1]))

    # III) Calculation of the mean and standard deviation
    # of the dataset differences
    dataset_delta_mean = np.mean(dataset_delta)
    dataset_delta_std = np.std(dataset_delta) ** (1/2)

    # IV) Calculation of Chebyshev k: n° of data that
    # contain at least 95% of observations
    k = (1/0.05) ** (1/2)

    # V) Calculation of the epsilon: min distance
    # between 2 points to be considered same cluster:
    eps = 2 * (dx ** 2 + 95 * CI of daily variability) ** (1/2)

    # VI) Determination of M: min 3 points within Eps. distance
    M = 3

    # VII) Clustering process
    # Creation of the Dataset vector: [y,i,dx]
    dataset_vect = []
    for i in range(len(dataset)):
        dataset_vect.append([dataset[i], i, dx])

    # Clustering = DBSCAN (eps=eps, min_samples=M, fit(dataset_vect))
    clustering = DBSCAN(eps=eps, min_samples=M, fit(dataset_vect))

    # VIII) Printing Results
    print("name: " + name)
    print(" - Delta mean: ", dataset_delta_mean)
    print(" - Delta std: ", dataset_delta_std)
    print(" - k: ", k)
    print(" - Epsilon: ", eps)
    print(" - M: ", M)
    print("Clustering: ", clustering)
    print()

    return clustering

cluster_GDP = outlier_detection(GDP, "GDP")
cluster_LIBOR = outlier_detection(LIBOR, "LIBOR")
cluster_Eq1 = outlier_detection(Eq1, "Equities I")
cluster_Eq2 = outlier_detection(Eq2, "Equities II")
cluster_Eq3 = outlier_detection(Eq3, "Equities III")

Outlier Detection Algorithm:
- Delta mean: 3080.26293181817
- Delta std: 56.4175059131117
- k: 4.47213595499958
- Epsilon: 7348.799001038649
- M: 3

Outlier Detection Algorithm:
- Delta mean: 0.1270601503759402
- Delta std: 0.4429045218447486
- k: 4.47213595499958
- Epsilon: 449.5432769403915
- M: 3

Outlier Detection Algorithm:
- Delta mean: 170.56627819548874
- Delta std: 11.98752391705578
- k: 4.47213595499958
- Epsilon: 7348.799001038649
- M: 3

Outlier Detection Algorithm:
- Delta mean: 77.20996138996139
- Delta std: 11.98752391705578
- k: 4.47213595499958
- Epsilon: 7348.799001038649
- M: 3
```



```
df_clean <- df[~72,] # Clean Dataset From Outliers of Equity II list

# Approach 3: Winsorize abnormal outliers (we set them equal to the x% percent)
# Approach 4: Do nothing (they hold statistical significance)
```

2.C) Provide a measure of correlation for the different variables at hand

```
# Change Dataset to numeric values
df_clean_num <- as.data.frame(apply(df_clean, 2, as.numeric))
sapply(df_clean_num, class)

# Calculate Correlation Matrix
cor(df_clean_num, function(x) !is.factor(x))) # Function useful to calculate
cor(df_clean_num, use="pairwise.complete.obs")
```

(index points)	Equities Index 2 (index points)				
'numeric'	Equities Index 1 (index points)				
'numeric'	Equities Index 2 (index points)				
	GDP Level	LIBOR 1M	Equities Index 1	Equities Index 2	Equities Index 3

(index points)				
Equities Index 3 (index points)	-0.06735564	-0.8880737	0.4542635	-0.6522365
				1.0000

2.D) Build a simple model that explains Equity 2 by one or more other variables provided. Provide an assessment of the out-of-sample performance of your model for the out-of-sample period 2015/01 to 2017/03

I) Model Generation

```
# Estimate dataset values from 2015 to 2017 (out-of-sample period)
df_regression <- df_clean_num[-c(1240:nrow(df_clean_num)) , ]
#df_regression

# Generate variable names
```

```
Eq2 <- df_regression[,4]
Eq3 <- df_regression[,5]

# Benchmark Regression ( $E_t Y = b_0 + b_1 \cdot GDP$ )
lm_case1 = lm(Eq2~GDP)
anova(lm_case1)
summary(lm_case1)

# Check whether adding Eq.1 improves regression results
```

```
# Check whether eliminating any parameter improves regression results
lm_case4 = lm(EQ2~GDP+EQ1+EQ3)
anova(lm_case4)
summary(lm_case4)

# Check whether adding LjBOR improves regression results
lm_case5 = lm(EQ2~GDP+LjBOR)
anova(lm_case5)
summary(lm_case5)

# Check whether adding EQ1 improves regression results
lm_case6 = lm(EQ2~GDP+LjBOR+EQ1)
anova(lm_case6)
summary(lm_case6)
```

```
lm(formula = EQ2 ~ GDP)

Residuals:
    Min       1Q   Median       3Q      Max
-960.9 -231.1 -129.0  111.5 1170.0

Coefficients:
            Estimate Std. Error t value Pr(>|t|)
(Intercept) -1.882e+03  2.545e+02  -7.396 1.67e-10 ***

```

Multiple R-squared: 0.2, Adjusted R-squared: 0.1965
 F-statistic: 57.5 on 1 and 230 DF, p-value: 8.275e-13

	Df	Sum Sq	Mean Sq	F value	Pr(>F)
GDP	1	21051082.4	21051082.4	125.555198	1.325578e-17
EQ1	1	429331.1	429331.1	2.560664	1.138135e-01
Residuals	74	14207133.4	167664.0	NA	NA

Call:
 lm(formula = EQ2 ~ GDP + EQ1)

Residuals:

Min	1Q	Median	3Q	Max
-1003.86	-276.76	-57.19	138.31	1090.27

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)
--	----------	------------	---------	----------

	Df	Sum Sq	Mean Sq	F value	Pr(>F)
GDP	1	125271.2	125271.2	2.040670	0.16937586
EQ3	1	249208.1	249208.1	4.059604	0.05829403
Residuals	19	1166358.7	61387.3	NA	NA

Call:
lm(formula = EQ3 ~ GDP + EQ3)

```

Eq3          6.875e+00  3.412e+00    2.015   0.0583 .
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 247.8 on 19 degrees of freedom
(217 observations deleted due to missingness)
Multiple R-squared:  0.243,    Adjusted R-squared:  0.1634 
F-statistic: 3.05 on 2 and 19 DF,  p-value: 0.07099


```

```
Residuals:
    Min       1Q   Median       3Q      Max
-478.81 -123.81   25.84  133.55  366.13

Coefficients:
              Estimate Std. Error t value Pr(>|t|)
(Intercept) 3967.64760 1245.07832   3.187  0.00511 **
GDP          -0.01996   0.00698    -2.860  0.01042 *
GDP2          0.00097    0.00027    3.689  0.00049 ***
---
Signif. levels:  0.0001 '***'  0.001 '**'  0.01 '*'  0.05 '.'  0.1 ' '  1.0 ''
```

		Eq2	GDP	Eq2	Libor	Residuals
GDP	1	21051082	21051082.4	147.4788	2.7265665e-19	
LIBOR	1	2273729	2273728.7	15.9292	1.532950e-04	
Residuals	74	10562736	142739.7	NA	NA	

Call:
lm(formula = Eq2 ~ GDP + LIBOR)

```

Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 377.8 on 74 degrees of freedom
(162 observations deleted due to missingness)
Multiple R-squared:  0.6893,    Adjusted R-squared:  0.6799
F-statistic: 81.7 on 2 and 74 Df,    p-value: < 2.2e-16


```

	Df	Sum Sq	Mean Sq	F value	Pr(>F)
GDP	1	2.105108e+07	2.105108e+07	1.454861e+02	4.748399e-19

```

Residuals:
    Min       1Q   Median       3Q      Max
-927.85  -213.94   -56.25   241.67   836.69

Coefficients:
            Estimate Std. Error t value Pr(>|t|)
(Intercept) -2.773e+03  3.324e+02  -8.343 1.39e-12 ***
GDP          1.458e-02  1.507e-03   9.671 1.03e-14 ***
LIBOR       2.097e+02  5.875e+01   3.570 0.000635 ***
EQ1         -6.132e+04  5.815e+02  -10.01 0.9916e-16 ***

```

```
# Preparing Data
df_clean_num_NA <- na.omit(df_clean_num)

# Defining x and y variables
y_real <- df_clean_num_NA[-c(1:22), 4]
y_prediction_1m = 2774 + 0.01457*df_clean_num_NA[-c(1:22), 1] + 209.5*df_clean_n
x_axis = seq(1:9)
```

169369.512907242

Equities Index II Prediction

000

Year	Equities Index I (Green)	Equities Index II (Blue)
1990	1050	1250
1991	1150	1300
1992	1100	1250
1993	1050	1200
1994	1000	1150
1995	950	1100
1996	900	1050
1997	850	1000
1998	800	950
1999	750	900
2000	700	850

Observations