# Data Preparation

1. Loaded file using *pandas.read_csv*. Separator was ";". I used quoting=3 because there were some Unicode errors with quotes in the file which was causing certain elements to be in the wrong column. This way when separating elements *pandas* does not care about quotes. Later used *converters* to remove the quotes around the strings.
2. Data had appropriate types, no steps necessary
3. Found typos using *dataframe.value_counts().* Fixed them using *dataframe.mask().*
4. Removed whitespaces using *converter* during the loading file step
5. Casted data to lowercase using *converter* during the loading file step
6. For sanity checks given the description of the data, I assumed that age < 120, if duration=0 then y=no, and if pdays=999 then poutcome!=success. 2 rows were removed for the age check.
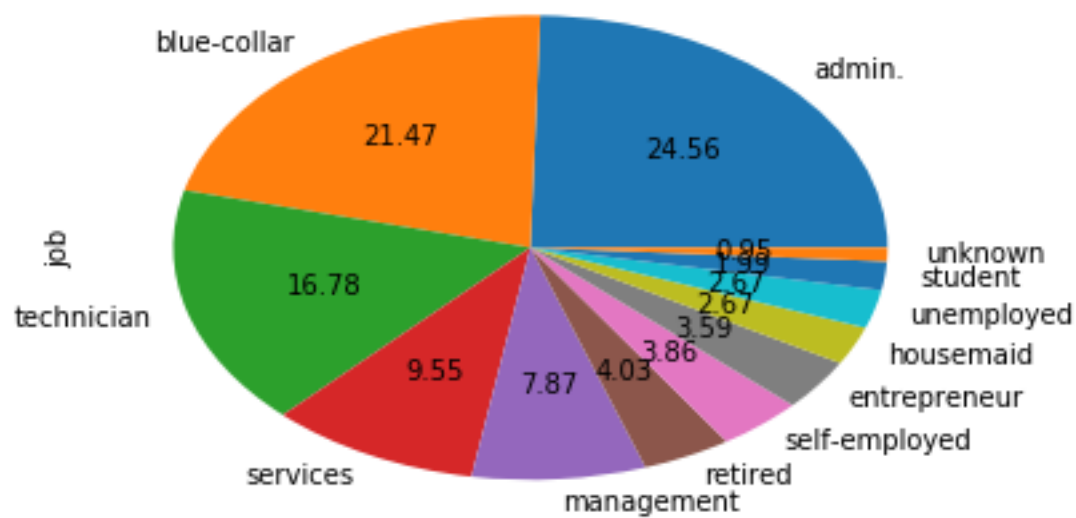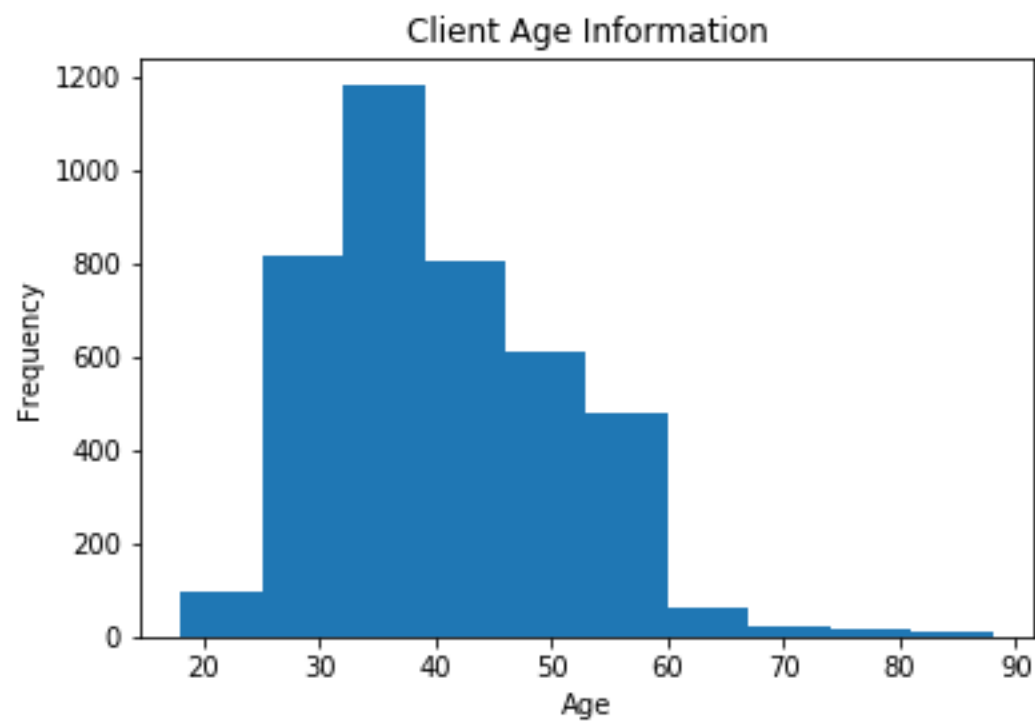7. Filled missing values using *dataframe.fillna()*
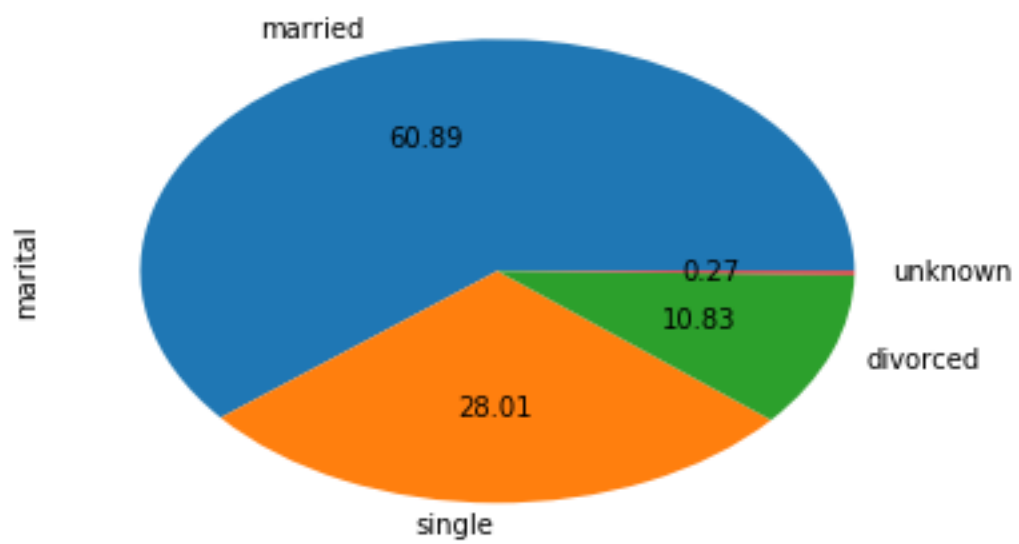
# Data Exploration

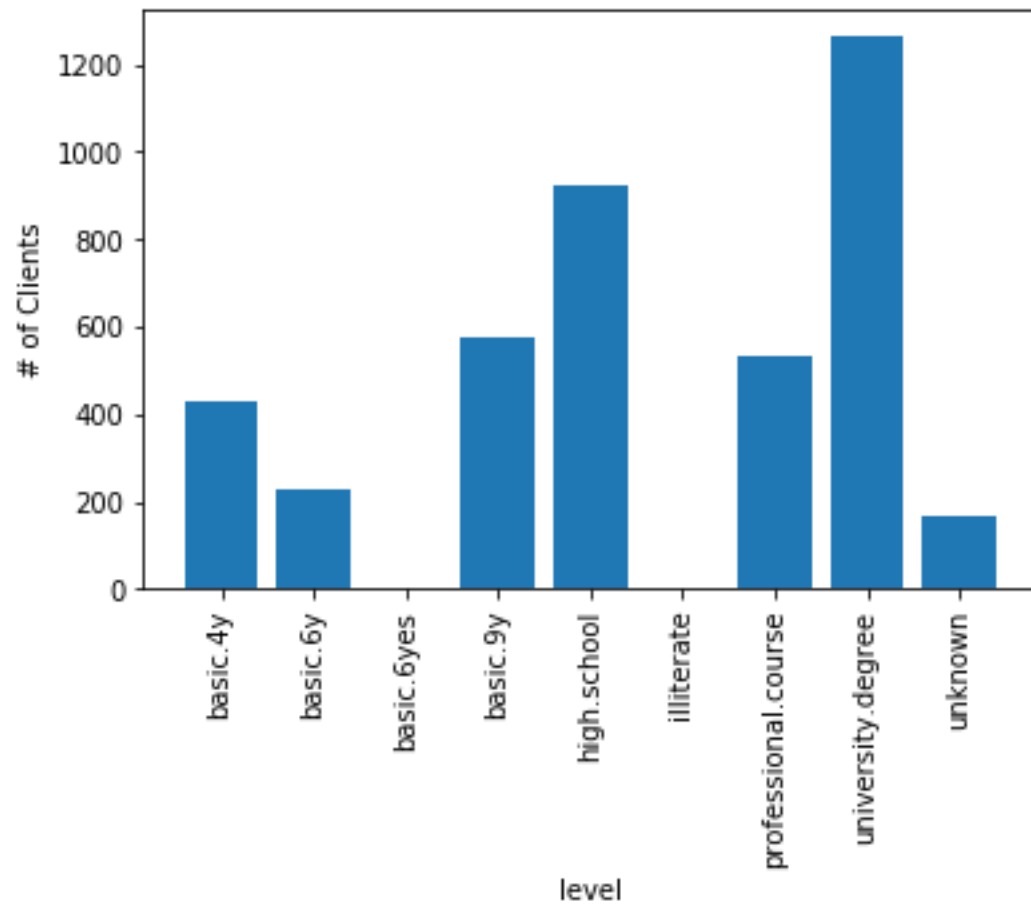## Subsection 1

**Why I chose a certain type of graph:**

- Pie Graph: used to show proportional relationship of values (Eg: proportion of people who are single vs married). Used for categorical and binary data. Except for *education* because there were couple values where the proportion was too small to see.
- Box plot: used when data had many outliers, in this case, *duration and campaign* columns. Also helps to show the data was left skewed
- Histograms: to show distribution and when data was easily divisible into bins
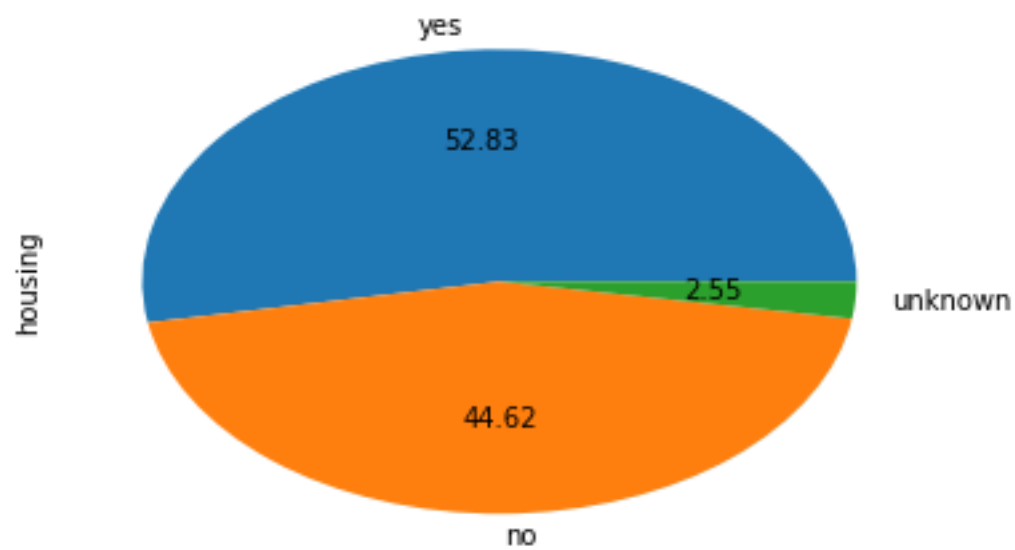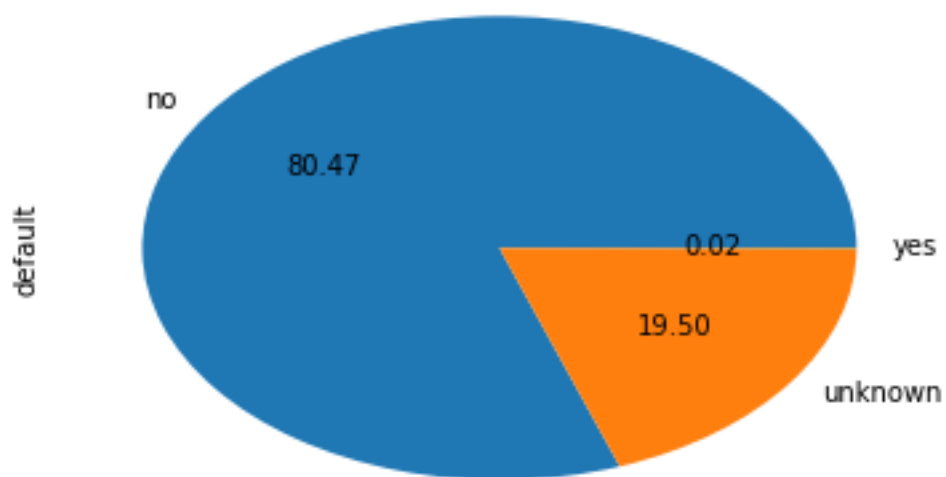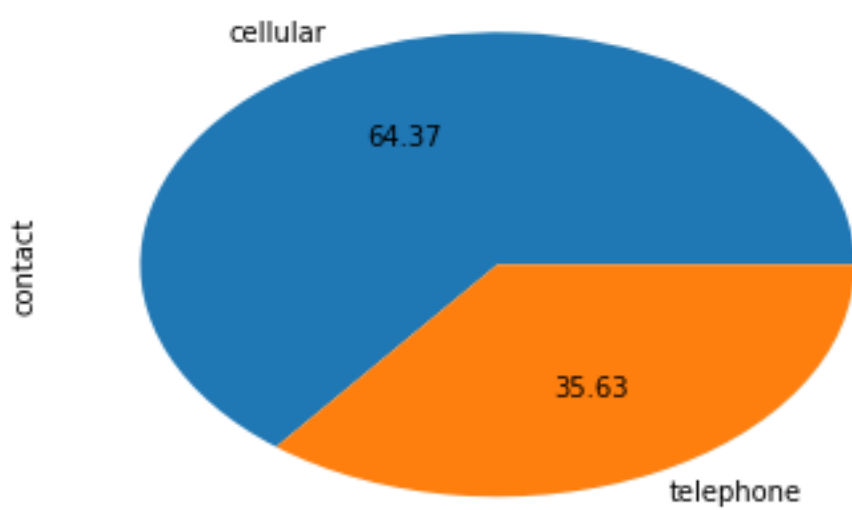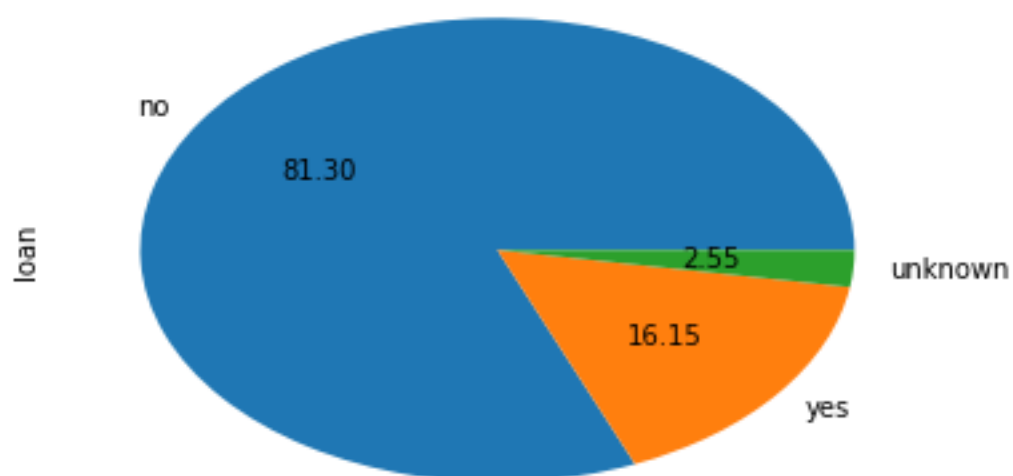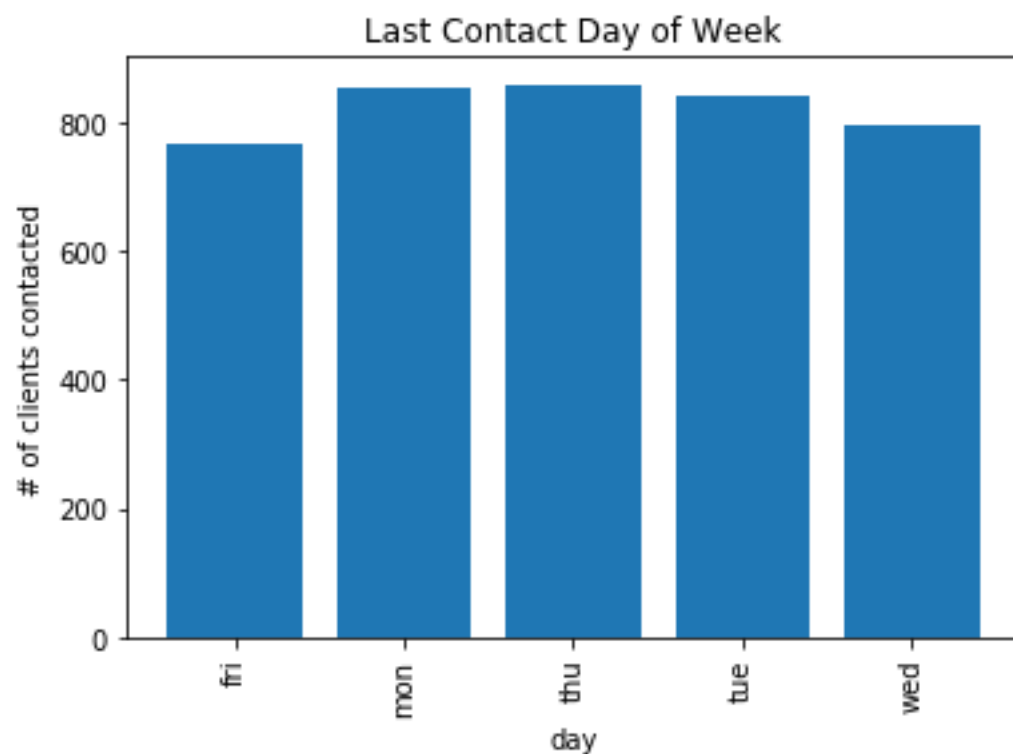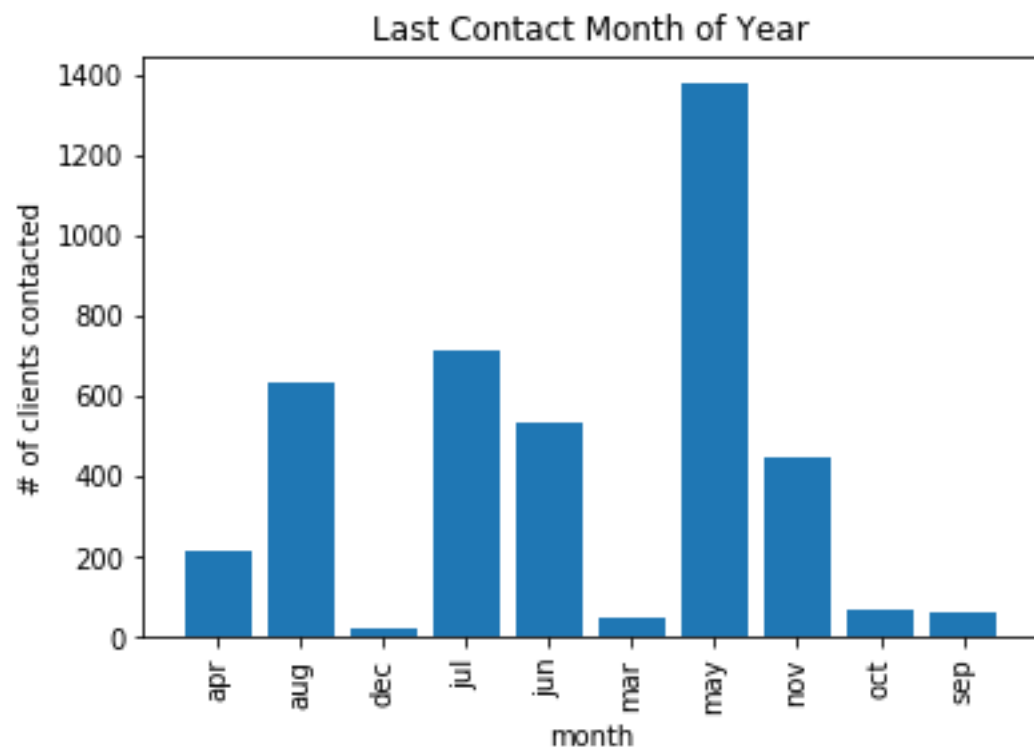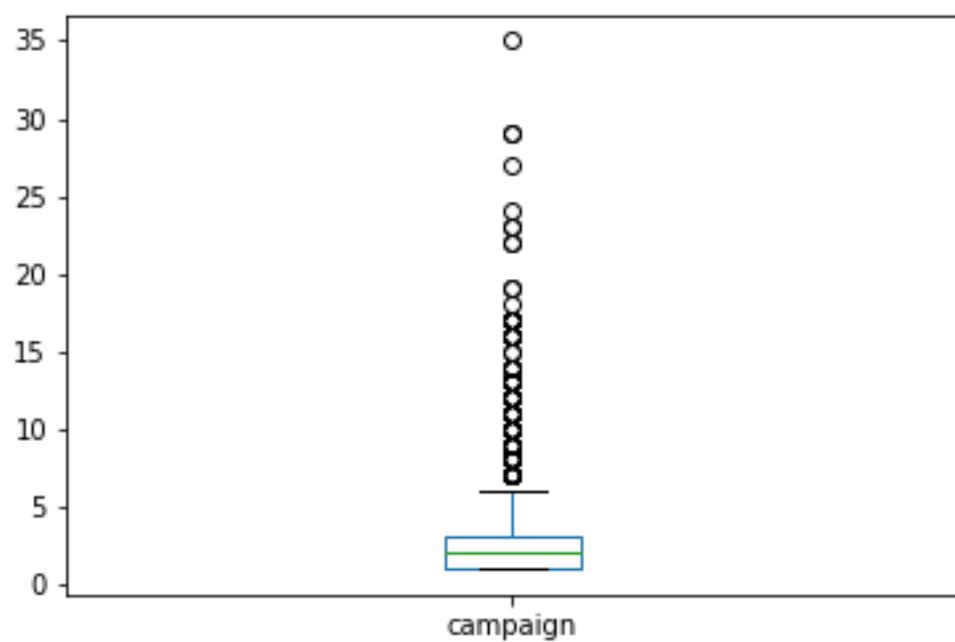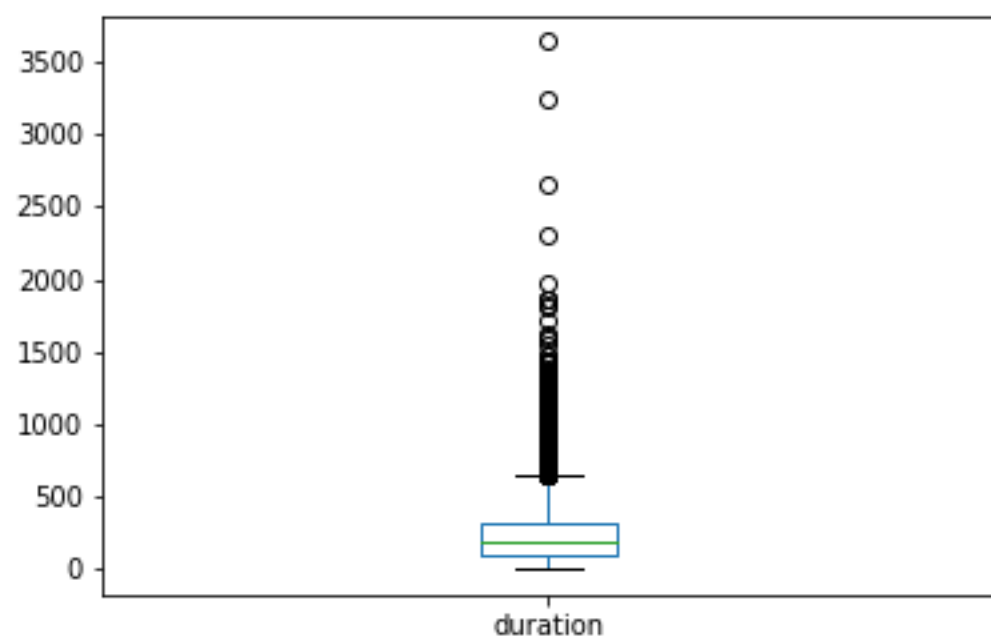- Bar graph: used when pie graphs would have been unreadable

**Graphs:**

## Client Age Information

Client Education Information

## loan

no 81.30

2.55 unknown

16.15 yes

## contact

cellular 64.37

35.63 telephone

## Last Contact Month of Year



## Last Contact Day of Week

## Days Passed Since Contact From Prev Campaign

A bar chart titled "Days Passed Since Contact From Prev Campaign" with y-axis labeled "# of clients contacted" ranging from 0 to 4000, and x-axis labeled "# of days" with values 0, 4, 12, 16, 1, 5, 9, 13, 17, 21, 2, 6, 10, 14, 18, 3, 7, 11, 15, 19, 999. The bar at 999 reaches approximately 3950.

## Contacts Performed Before Campaign

A histogram titled "Contacts Performed Before Campaign" with y-axis labeled "# of clients contacted" ranging from 0 to 3500, and x-axis labeled "# of contacts performed" with values 0 through 6. The bar at 0 reaches approximately 3500, the bar at 1 approximately 500.

nonexistent

poutcome

85.52

3.45    success

11.03

failure

Employement Variation

2500

2000

1500

1000

500

0

# of clients

−3    −2    −1    0    1

Emp.var.rate

Consumer Price Index

Consumer Confidence Index

Euribor 3 Month Rate

Number of Employees

## Clients Subscribed



# Subsection 2
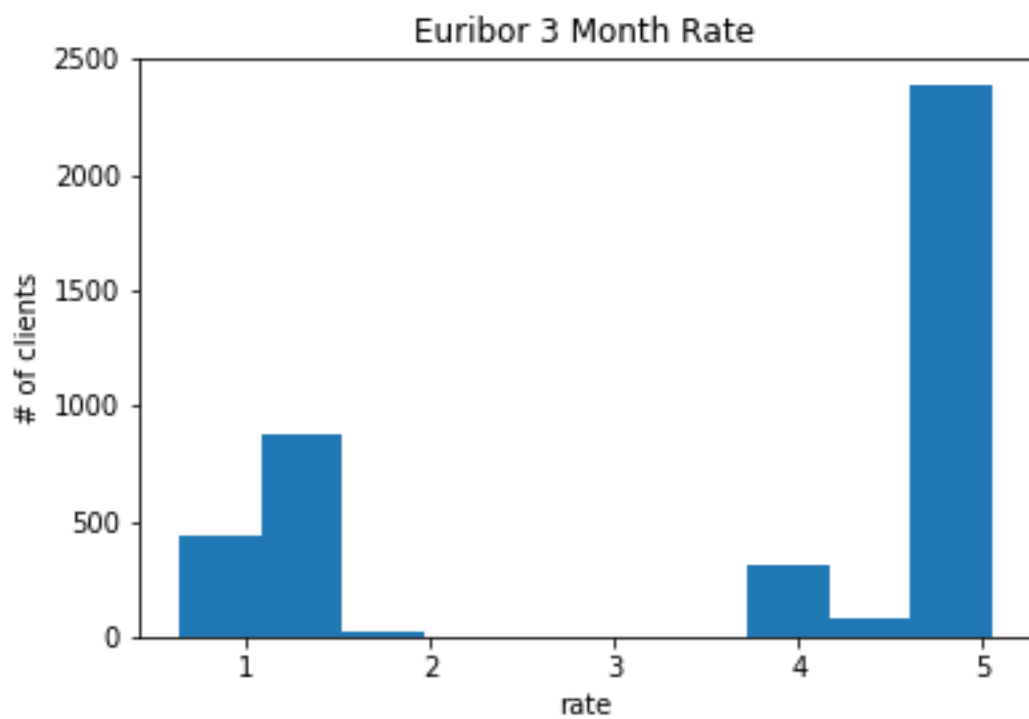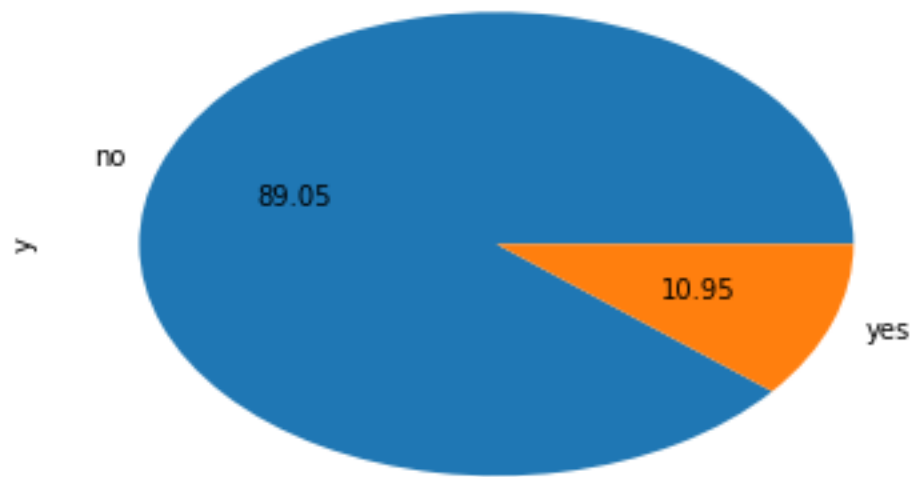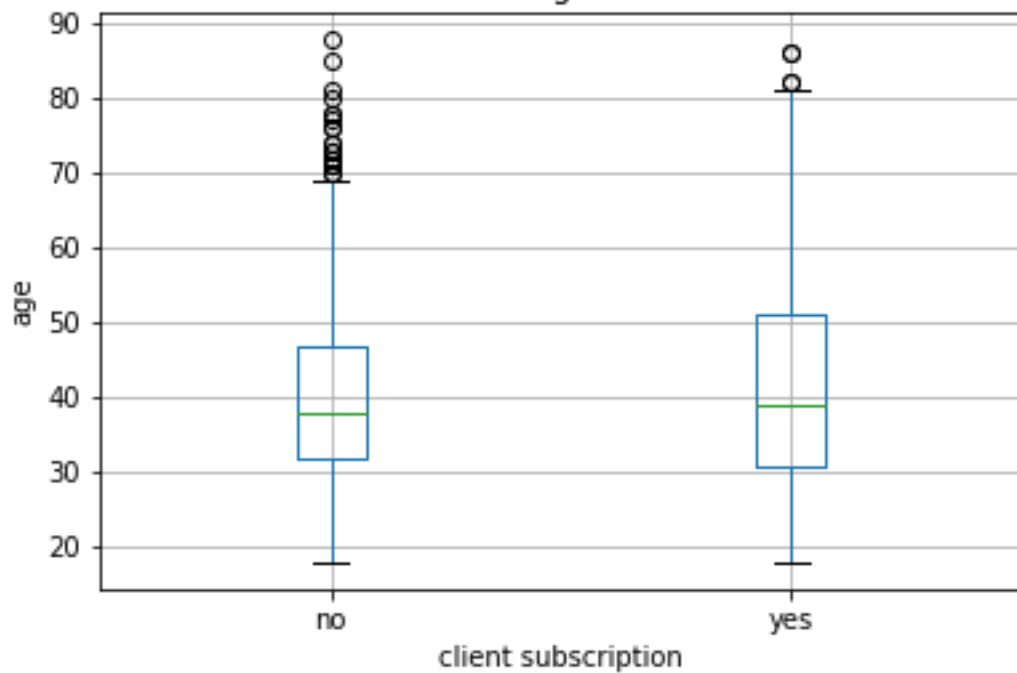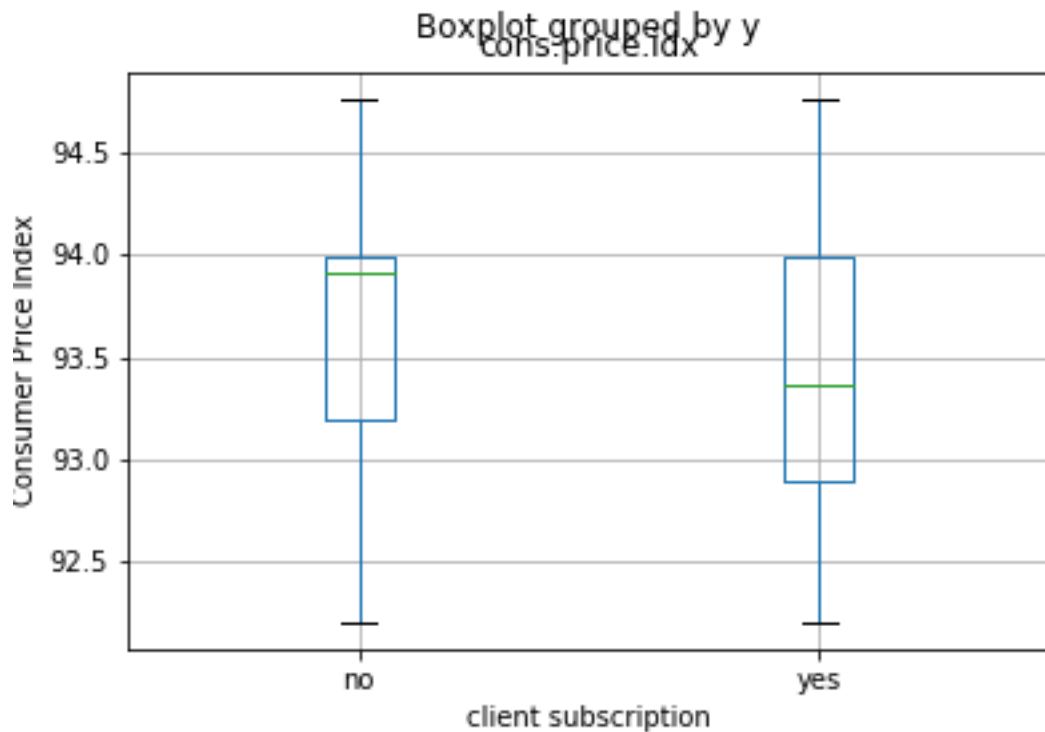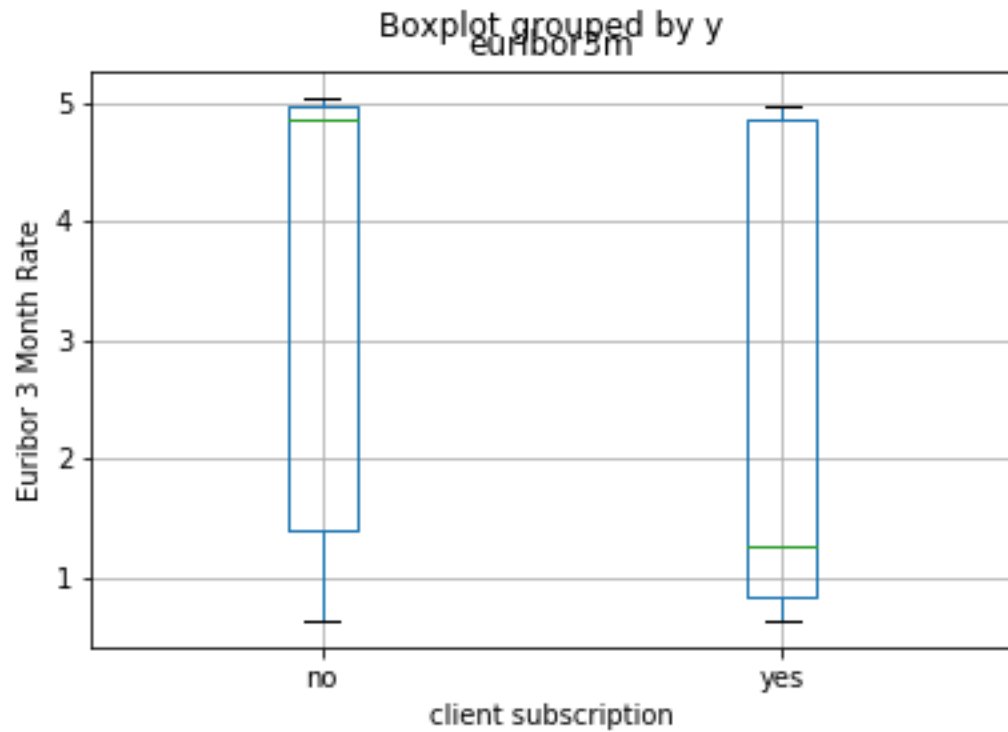
## Boxplot grouped by y

Hypothesis:  Older clients tend to subscribe rather than say no.
Why?: The boxplot for clients who did not subscribe has the upper whisker at age 70. This means that 75% of people who said no are under the age of 70. Also notice the upper whisker on boxplot for yes. It is at age 80. It seems that there is not relationship between middle age people and subscription shown in this graph because the median is around the same place. Also the IQR is roughly the same.

Boxplot grouped by y
cons:price.idx



Hypothesis: clients with high consumer price index will  not subscribe.
Why?: Notice the median for clients who said yes is much lower than other boxplot. Everything else is pretty much the same. We can see both graphs have similar distributions but the clients who said no are more densely populated around the values 93.2 to 94.0. Therefore it is possible that high consumer price index leads to clients saying no.

Boxplot grouped by y
euribor3m

Hypothesis: clients with high Euribor 3 month rate will not subscribe.
Why?: Notice the median for clients who said yes is much lower than other boxplot. Everything else is pretty much the same. We can see that clients who said no have a median of 4.8 which is near the max value. Thus it is possible that clients with a high rate tend to decline subscription.

# Subsection 3
## Scatter Matrix