Practical Data Science

Assignment 1

Chi Ting, Low

S3611774

# Table of Contents

## Data Preparation

The aim of the current report is to conduct data cleaning and data exploration from the given dataset using Python 2.7 (Python Software Foundation, 2018). The dataset used is related to direct marketing campaigns of a Portuguese banking institution. This campaign was conducted based on phone calls in order to access a client will subscribed to their product or not. Attributes information:

- age – client age (numeric)

- job – type of job (categorical)

- marital – marital status (categorical)

- education – education status (categorical)

- default – credit default history (categorical)

- housing – has housing loan (categorical)

- contact – communication type (categorical)

- month – month of last contact (categorical)

- day_of_week -day of the week of last contact (categorical)

- duration – contact duration (numeric)

- campaign – number of contact performed (numeric)

- pdays – number of day passed after last contact (numeric)

- previous – number of contact before the campaign (numeric)

- poutcome – outcome of previous campaign (categorical)

- emp.var.rate – employment variation rate, quarterly indicator (numeric)

- cons.price.idx – consumer price index, monthly indicator (numeric)

- cons.conf.idx – consumer confidence index, monthly indicator (numeric)

- euribor3m – euribor 3 month rate, daily indicator (numeric)

- nr.employed – number of employees, quarterly indicator (numeric)

- y – has the client subscribed a term deposit (categorical)

Task 1.1

To read the data, pandas packages (McKinney, 2010) is called and pd.read_csv function is used along with arguments (sep, decimal, header, names and na_values) to load the data into python environment. The argument *sep* is used to separate stacked values in the data set. The argument *decimal* is used to recognize the values into decimal values. The argument *header* is used to set the column names. In current report, the *header* is set to 0 which allows the column names to be replaced (references). The *name* argument is used to rename the columns names. Lastly, the *na_values* argument is used to set the missing values to string NaN.

Task 1.2

After loading dataset into python environment, *.dtypes* method is used to access the types of variables assigned to each variable (Boschetti & Massaron, 2015). As observed, some variables are assigned as object. This indicated that the variables do not assigned properly into the data frame except age, duration, campaign, cons_price_idx and cons_conf_idx which are assigned into float64 (numeric with double decimal point). To convert the variables to the data types presented in the attribute information above two methods is called (*.to_numeric* and *.astype*). The *.to_numeric* is used to convert the object variable into numerical data type. The *.astype* is used to convert the variable into categorical data types.

Task 1.3

To examine the typos in the data due to imputation error method *.value_counts*() is used. The types are found in the marital, education, default, housing, loan, contact, month,

day_of_week, pdays, previous, poutcome, emp_var_rate, cons_price_idx, cons_conf_idx, euribor3m, and nr_employed variables. These variables contain strings that is not assigned or contain no meaningful information for the variables. By using masks method, these invalid strings are set into NaN.

Task 1.4

To examine whether there are instances of extra whitespaces in the data, method *str.strip* is called on the categorical variables. For example, *data_copy["job"].str.strip()* which is used to omit the whitespace in the job column in the data frame.

Task 1.5

To set the text data into lower case method *.str.lower()* is used to set the text data into lower-case. For example in *data_copy["marital"].str.lower(),* this allows the method to be called and search the marital variables columns and set all the row into lower case.

Task 1.6

To check whether there are impossible values for each variable in the data set, *mean* method is used to perform a series of sanity checks to filter out impossible values. In the variable age, the mean age is 40.112. However, in the row 385 (age = 143) and 394 (age = 250) the age is above mean age by large margin and the chance of an individual to live longer than 143 and 250 is slim.

In addition, in the campaign variables which indicated the number of contact to client during campaign have the mean of 3.042. This indicated that on average each client will receive 3 calls. But, in row 569, 650 and 678 shows that a person has received 999 and 93.994 call which indicated this is impossible. Similar approach is conducted to the variables; previous (mean = 0.191), emp.var.rate (mean = 0.085), cons.price.idx (mean = 93.558), cons.conf.idx (mean = -

40.499), euribor3m (mean = 4.859) and nr.employed (mean = 5166.517) which shows that there are values beyond the mean or lower than mean. All these values is replaced as missing values.

Task 1.7

Lastly, to examine the missing values in the data set *.isnull().sum()* method is called to examine the total missing values in each variable. It shows that marital, education, default, and housing have one missing values; pdays has two missing values; loan, contact, month, day_of_week and campaign have 3 missing values; previous, poutcome, emp.var.rate, cons.conf.idx and y have 4 missing values; euribor3m and age have 5 missing values; cons.price.idx and nr_employed have 6 missing values after the process of data cleaning and replace all the typos and impossible values to missing values. To replace all the missing values method, *.fillna* and *.mean* to fill the missing value by using mean of each correspondent variables.

**Extension**

To replace all the missing values, three other approaches are used; median, using fix-value and dropping the missing values. Same approach is used to fill the missing values by using median of the correspondent variables. Lastly, using the method *.fillna* to fill all the missing values using fix value 9999 which is demonstrated in Task1.ipynb in line 96, 98, and 100.

The last approach to deal with missing values is used *.dropna* method. This method is used to remove the entire rows that contain missing values from data frame. As this approach applied, the number of observation in the data set is drop from 4119 to 4103 as demonstrated in Task1.ipynb in line 119.
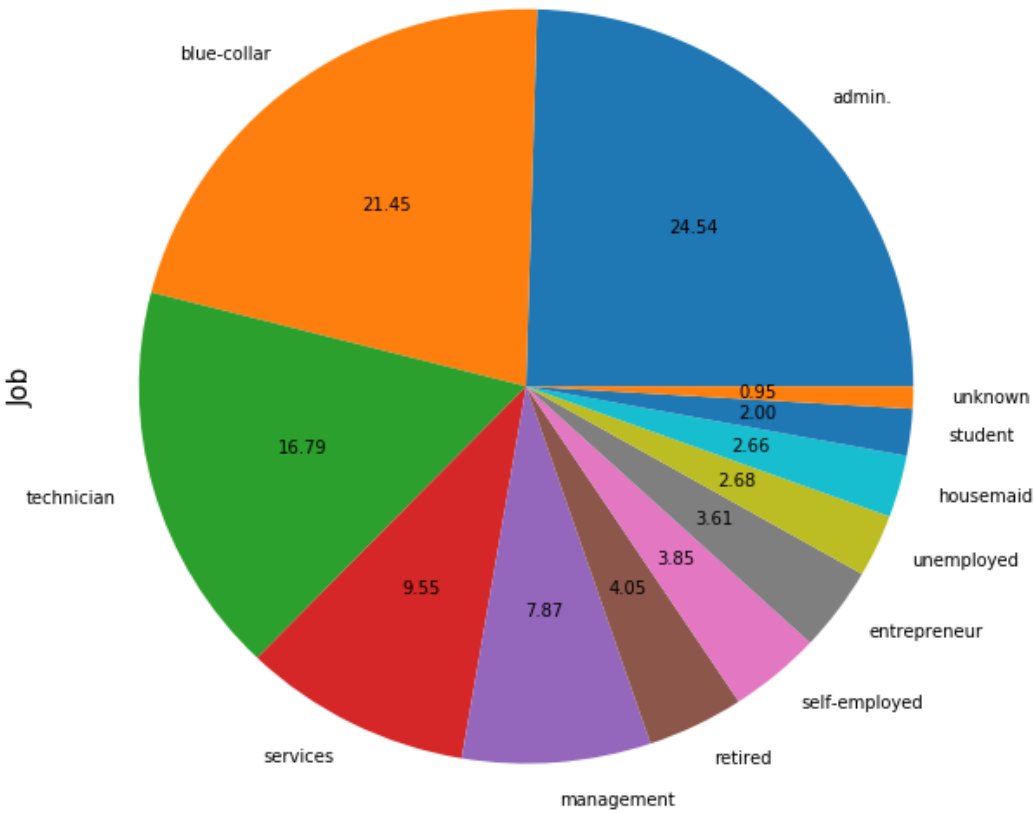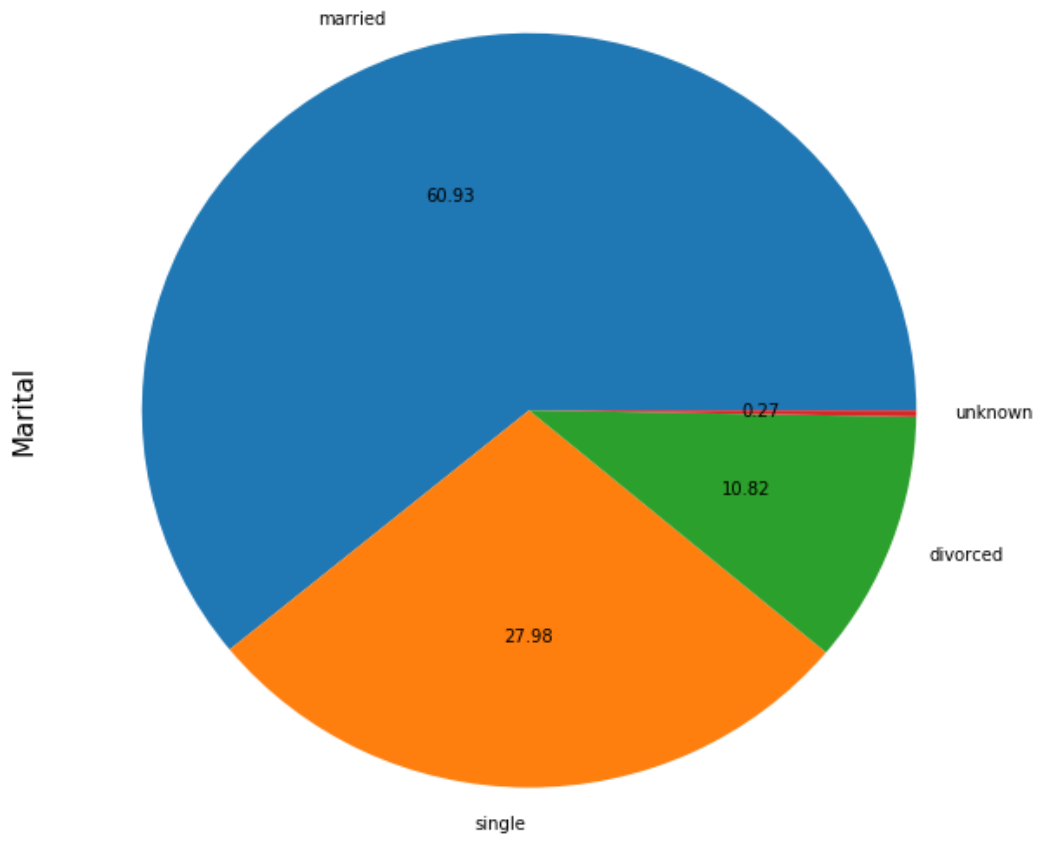
Data Exploration

Task 2.1

To conduct simple exploration of the data set, a basic visualization is performed for each variable using histogram, pie charts, bar charts and density plot for the data set.
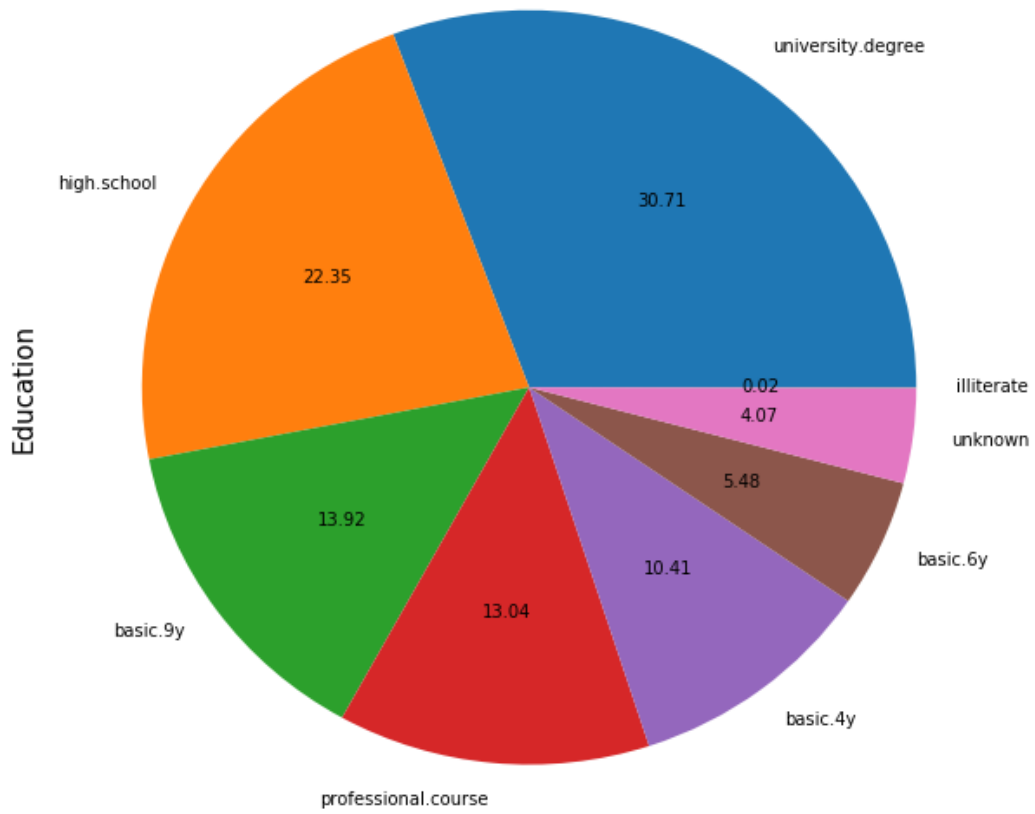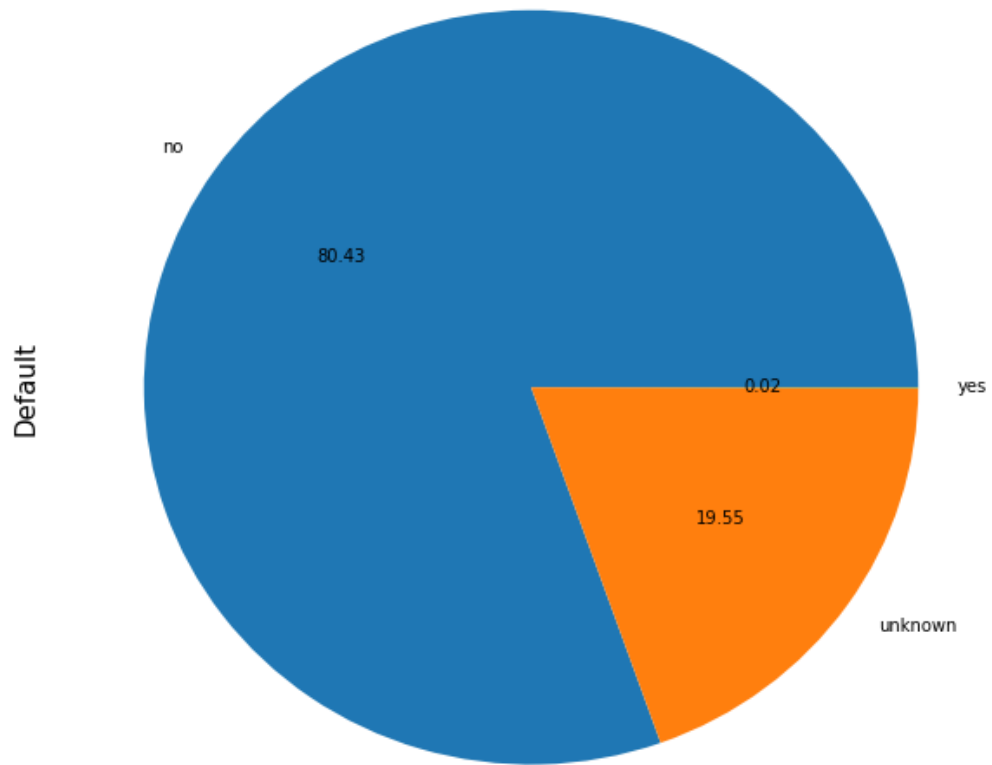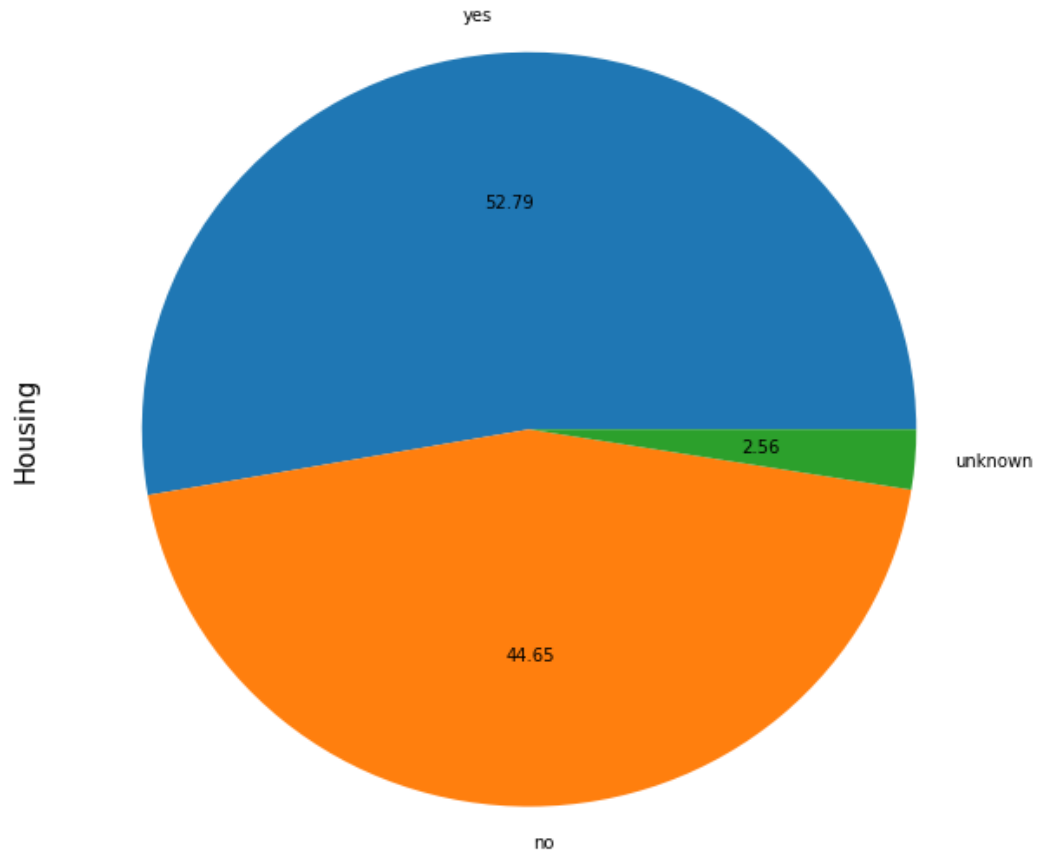


Age Distribution

# Job Distribution



Job

- admin. 24.54
- blue-collar 21.45
- technician 16.79
- services 9.55
- management 7.87
- retired 4.05
- self-employed 3.85
- entrepreneur 3.61
- unemployed 2.68
- housemaid 2.66
- student 2.00
- unknown 0.95

# Marital Distribution

Marital

married

60.93

unknown

0.27

divorced

10.82

single

27.98

# Education Distribution

# Default in Credit Distribution



Default

no 80.43

0.02 yes

19.55

unknown

# Housing Distribution

yes

52.79

Housing

2.56

unknown

44.65

no

# Contact Distribution



Contact

cellular

64.37

35.63

telephone

# Month of Contact Distribution



| Month |
|-------|
| may 33.41 |
| jul 17.28 |
| aug 15.40 |
| jun 12.92 |
| nov 10.85 |
| apr 5.24 |
| oct 1.68 |
| sep 1.56 |
| mar 1.12 |
| dec 0.54 |

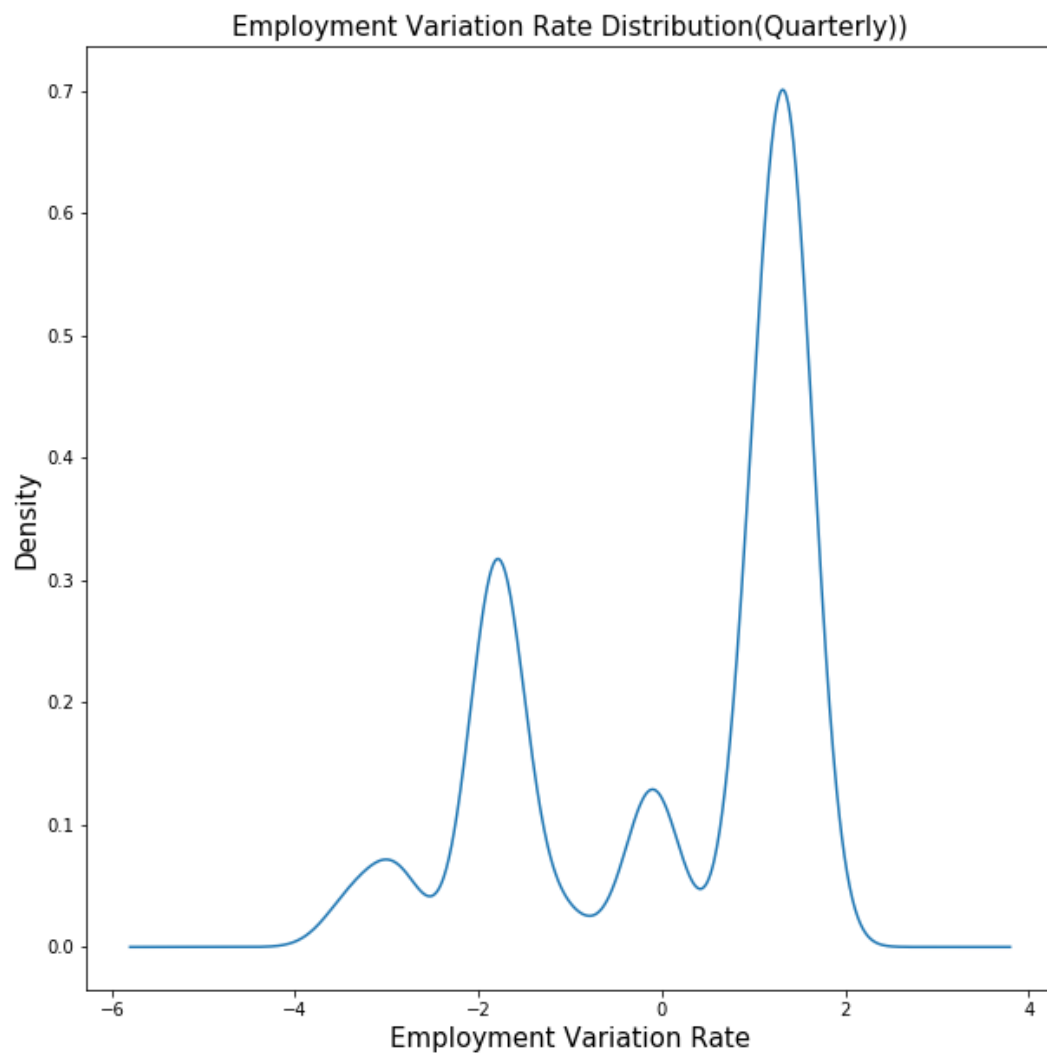# Last Contact From Week Distribution

Duration of Last Contact Distribution

Number of Contact During Campaign Distribution
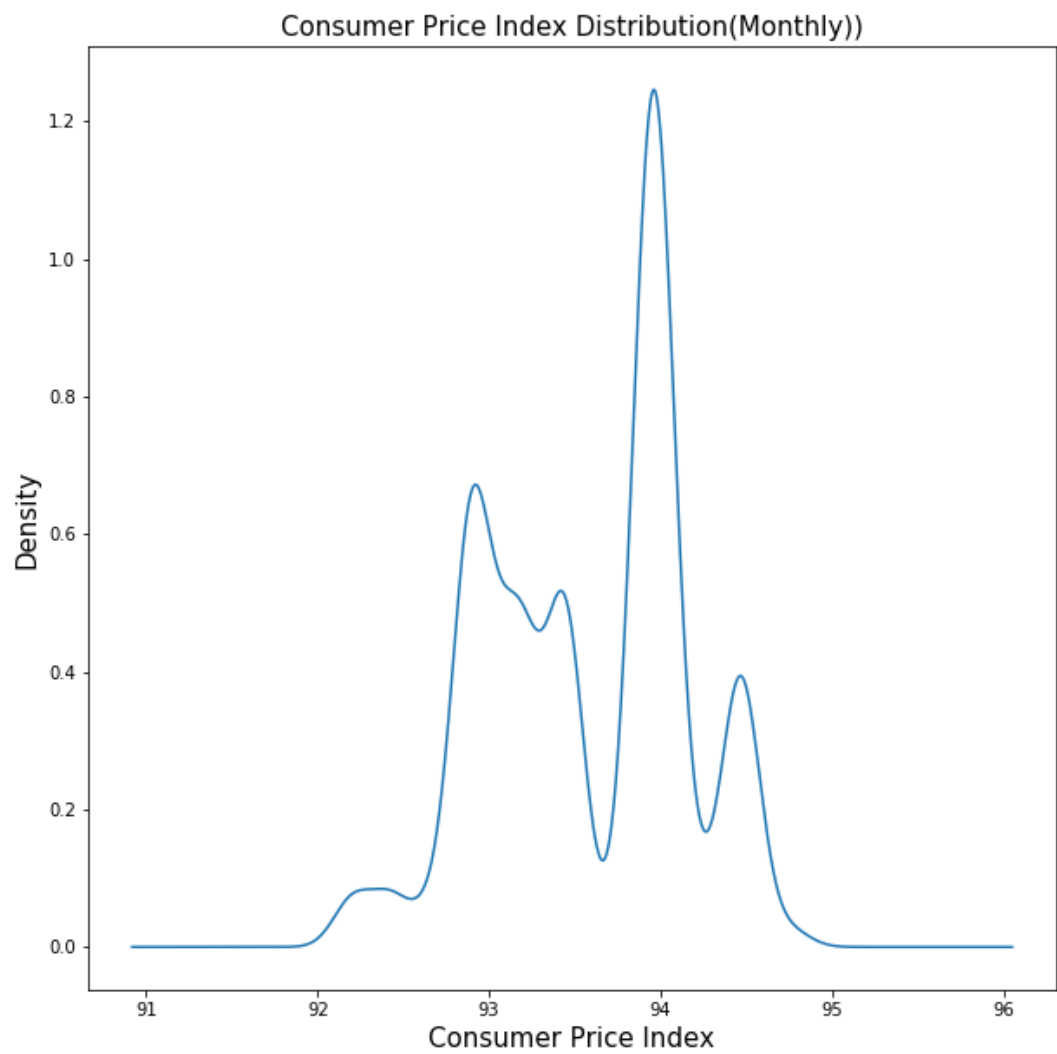
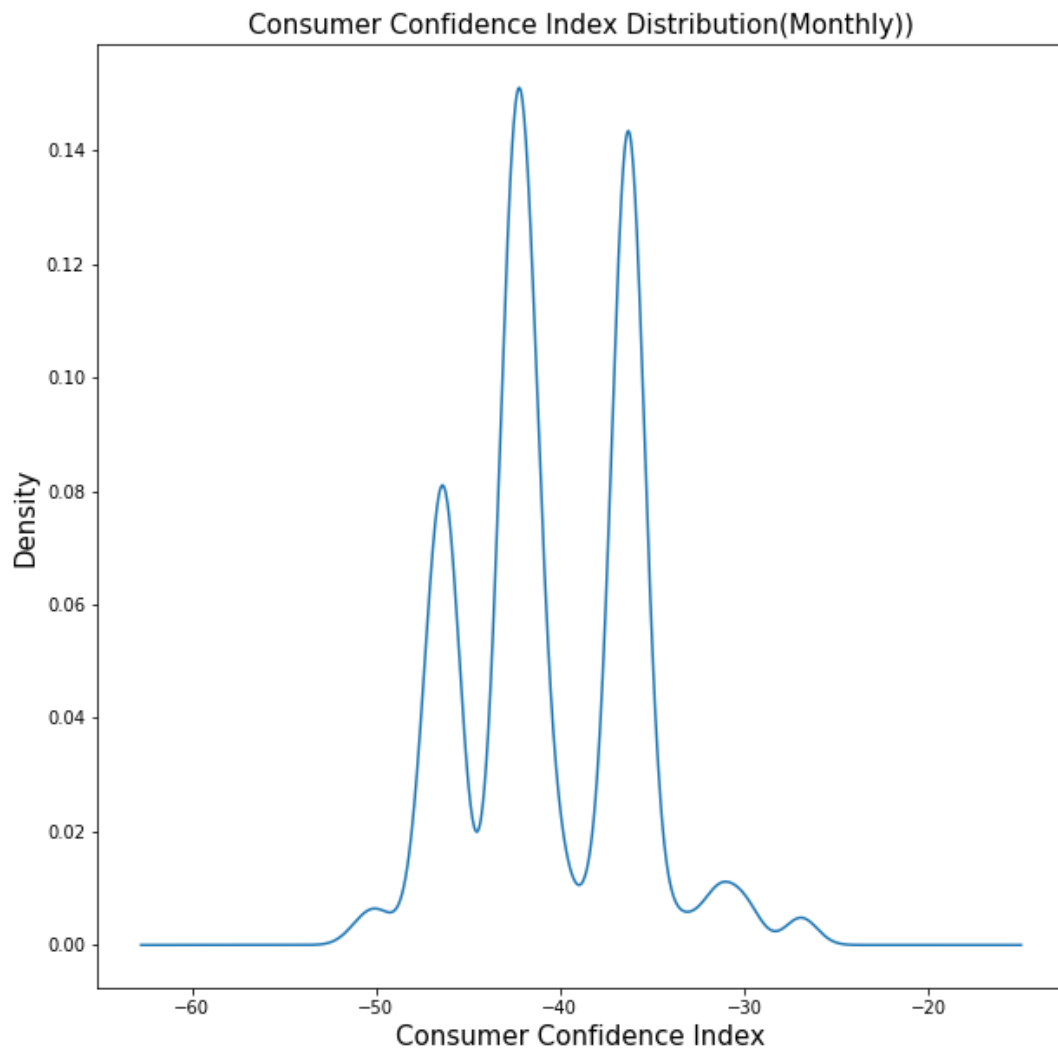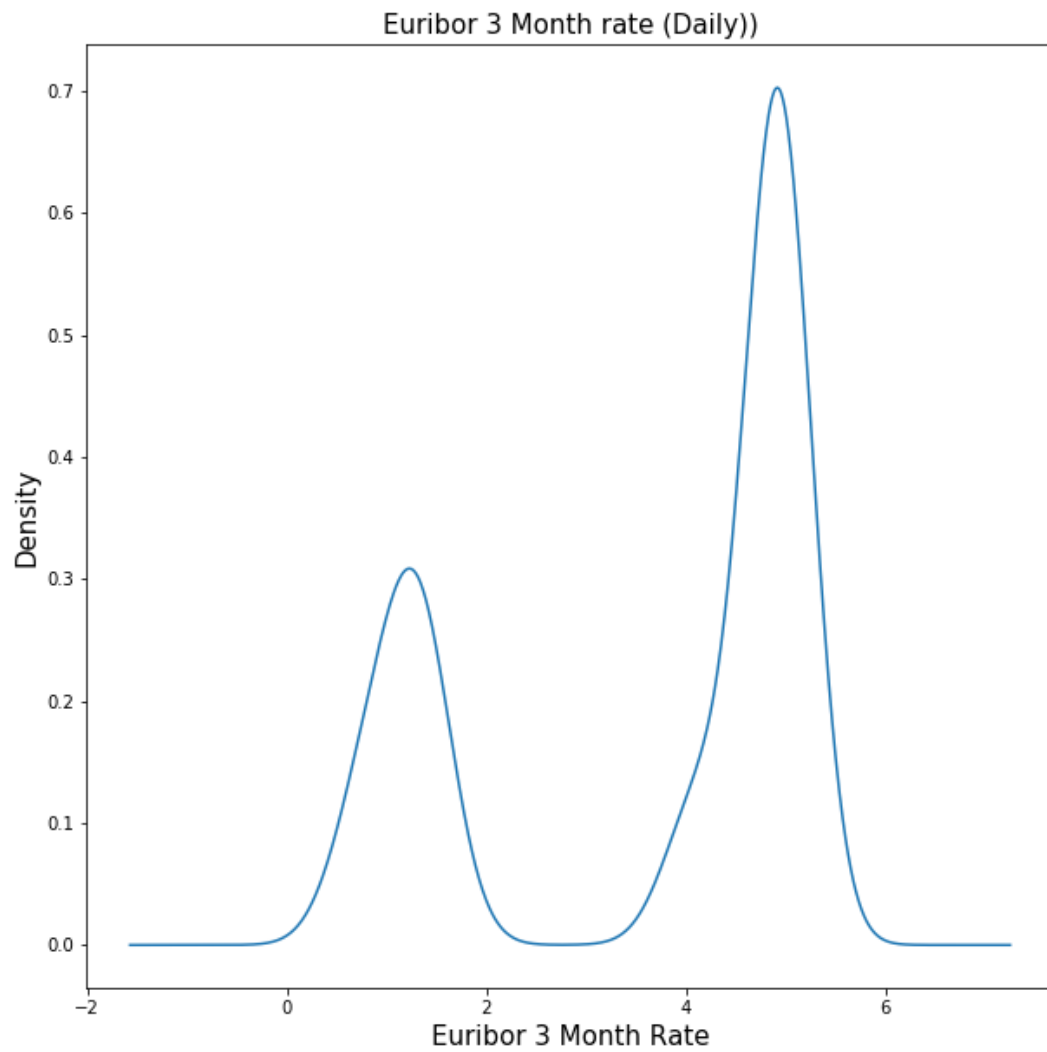Density of Number of Days After Contact

Number of Contacts Before the Campaign

# Outcome of Previous Marketing Campaign

Employment Variation Rate Distribution(Quarterly))

Consumer Price Index Distribution(Monthly))

Consumer Confidence Index Distribution(Monthly))

Euribor 3 Month rate (Daily))

Number of Employees (Quarterly))

## Subscribed Status

Task 2.2

Relationship between Age and Consumer Confidence Index
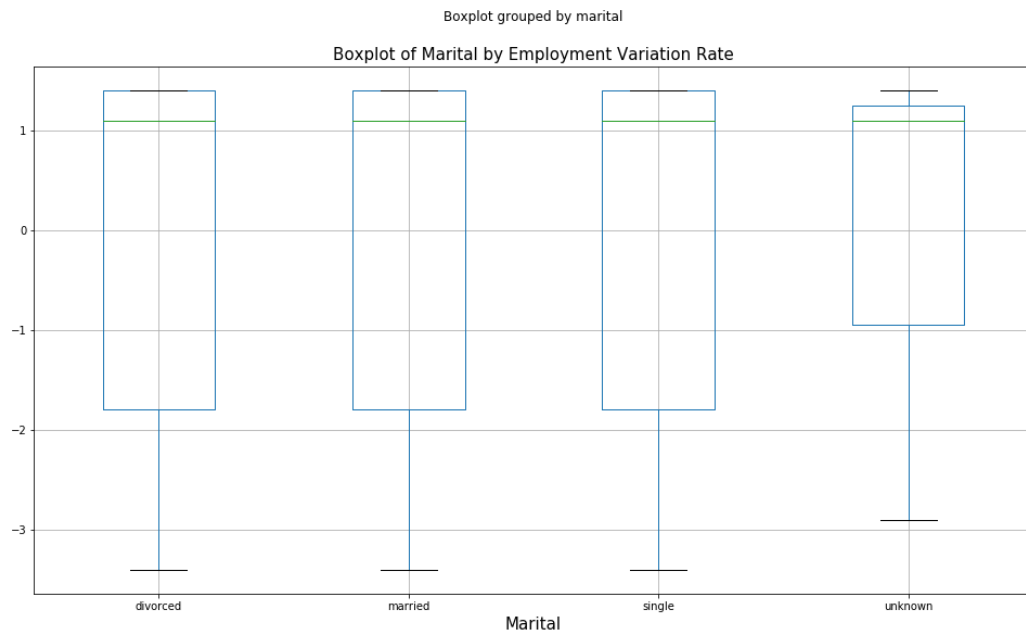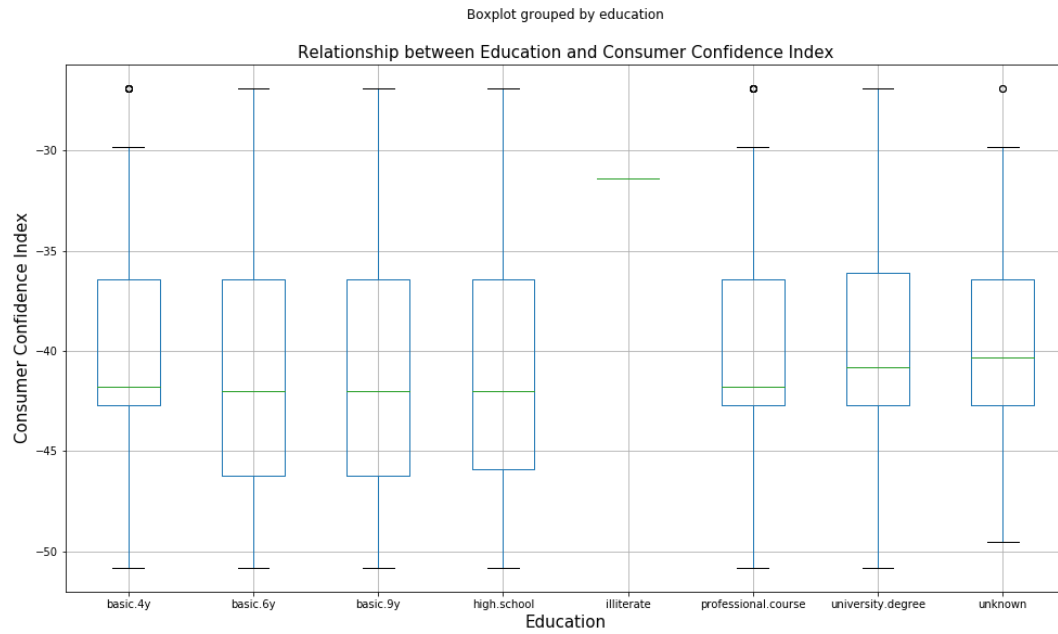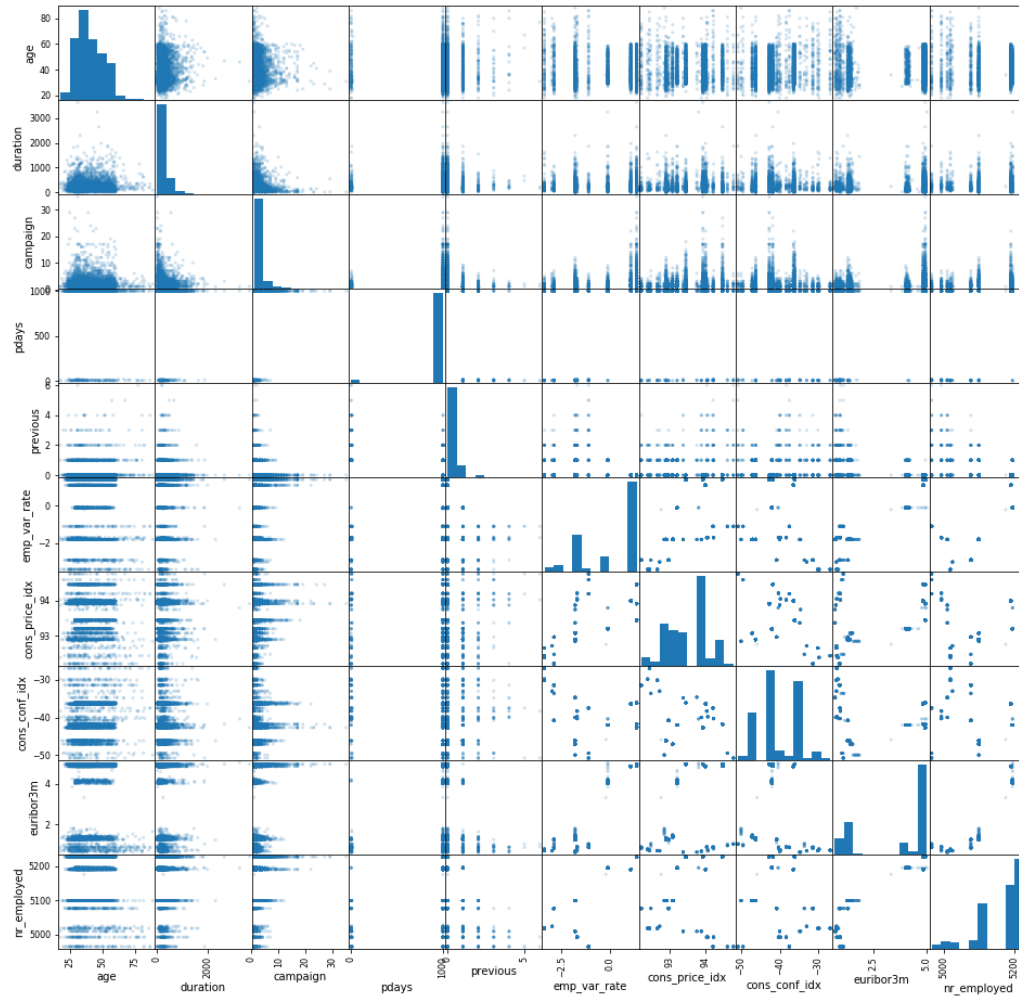


To explore the relationship, matplotlib packages is used (Hunter, 2007). The scatter above shows the relationship between consumer confidence index and ages. It shows that between age 30 and 60 age group, there is highest confidence index.

## Relationship between Education and Consumer Confidence Index

## Boxplot of Marital by Employment Variation Rate

Task 2.3

# References

Boschetti, A., & Massaron, L. (2015). *Python data science essentials: Become an efficient data science practitioner by thoroughly understanding the key concepts of Python*. Birmingham, England: Packt Publishing.

Hunter, J. D. (2007). Matplotlib: A 2D graphics environment. *Computing in science & engineering*, *9*(3), 90-95.

McKinney, W. (2010). Data structures for statistical computing in python. In *Proceedings of the 9th Python in Science Conference* (Vol. 445, pp. 51-56).

Python Software Foundation. Python Language Reference, version 2.7. Available at http://www.python.org