



# **MSc Big Data Analytics**

## **CMM704 – Data Mining**

**IIT Student ID: 20200361**

**RGU Student ID: 2016629**

**Don Manula Ransika Udugahapattuwa**

## Table of Contents

(a).....	2
I.    Handle missing values and outliers.....	2
II.   Produce Q-Q Plots and Histograms of the features, and apply the transformations if required	10
III.  Feature Coding.....	16
IV.  Scaling data.....	18
V.   Feature discretization.....	18
(b).....	19
I.    SVD (Singular Value Decomposition) for feature reduction.....	19
II.   Significant and independent features.....	21
(c). Applying Logistic Regression and Support Vector Machine techniques.....	22
(d) applicability of SVM and LR.....	25
(e) How significant your findings are.....	28

(a)

We can start the project by setting up the prerequisites and introducing the data set.

### ▼ Setup the pre-requisites and introduce the data set

```
[3] 1 #Set Prerequisites
2 import pandas as pd
3 import numpy as np
4 from sklearn import preprocessing
5 import matplotlib.pyplot as plt
6 plt.rc("font", size=14)
7 from sklearn.linear_model import LogisticRegression
8 from sklearn.model_selection import train_test_split
9 import seaborn as sns
10 sns.set(style="white")
11 sns.set(style="whitegrid", color_codes=True)
```

```
[4] 1 #Mount gdrive
2 from google.colab import drive
3 drive.mount("/content/gdrive")
```

Mounted at /content/gdrive

```
1 #Import the dataset as a Pandas data frame
2 df = pd.read_csv('/content/gdrive/My Drive/DataMining/assignment/dataset/banking.csv')
3 df.head(10)
```

## I. Handle missing values and outliers

Drop "Duration" column first: the duration is not known before a call is performed. Also, after the end of the call y is obviously known. Thus, this Input should only be included for benchmark purposes and should be discarded.

### ▼ (a)I. Handle Missing Values and Outliers if any

Drop "Duration" column first: the duration is not known before a call is performed. Also, after the end of the call y is obviously known. Thus, this Input should only be included for benchmark purposes and should be discarded

```
[6] 1 df = df.drop(['duration'],axis=1)
2 df.head(5)
```

	age	job	marital	education	default	housing	loan	contact	month	day_of_week	campaign	pdays	previous	poutcom
0	44	blue-collar	married	basic.4y	unknown	yes	no	cellular	aug	thu	1	999	0	nonexister
1	53	technician	married	unknown	no	no	no	cellular	nov	fri	1	999	0	nonexister
2	28	management	single	university.degree	no	yes	no	cellular	jun	thu	3	6	2	success
3	39	services	married	high.school	no	no	no	cellular	apr	fri	2	999	0	nonexister
4	55	retired	married	basic.4y	no	yes	no	cellular	aug	fri	1	3	1	success

Then we can proceed with checking for missing values:

```

1 #Check for missing values
2 print(df.info())

```

```

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 41188 entries, 0 to 41187
Data columns (total 20 columns):
#   Column                Non-Null Count  Dtype
---  -
0   age                    41188 non-null  int64
1   job                    41188 non-null  object
2   marital                41188 non-null  object
3   education              41188 non-null  object
4   default                41188 non-null  object
5   housing                41188 non-null  object
6   loan                   41188 non-null  object
7   contact                41188 non-null  object
8   month                  41188 non-null  object
9   day_of_week            41188 non-null  object
10  campaign                41188 non-null  int64
11  pdays                  41188 non-null  int64
12  previous                41188 non-null  int64
13  poutcome                41188 non-null  object
14  emp_var_rate            41188 non-null  float64
15  cons_price_idx          41188 non-null  float64
16  cons_conf_idx           41188 non-null  float64
17  euribor3m               41188 non-null  float64
18  nr_employed             41188 non-null  float64
19  y                       41188 non-null  int64
dtypes: float64(5), int64(5), object(10)
memory usage: 6.3+ MB
None

```

Even though there are no NULL values, there might be other custom values like "unknown" that suggest a missing values.

- To check this, we can get the unique values of the columns.

Checking for unique values in each of the categorical features to check for other cases of missing values such as unknown and nonexistent

```

[5] 1 df['job'].unique()

array(['blue-collar', 'technician', 'management', 'services', 'retired',
      'admin.', 'housemaid', 'unemployed', 'entrepreneur',
      'self-employed', 'unknown', 'student'], dtype=object)

[6] 1 df['marital'].unique()

array(['married', 'single', 'divorced', 'unknown'], dtype=object)

[7] 1 df['education'].unique()

array(['basic.4y', 'unknown', 'university.degree', 'high.school',
      'basic.9y', 'professional.course', 'basic.6y', 'illiterate'],
      dtype=object)

[8] 1 df['default'].unique()

array(['unknown', 'no', 'yes'], dtype=object)

[9] 1 df['housing'].unique()

array(['yes', 'no', 'unknown'], dtype=object)

```

```
[10] 1 df['loan'].unique()

array(['no', 'yes', 'unknown'], dtype=object)

[11] 1 df['contact'].unique()

array(['cellular', 'telephone'], dtype=object)

[12] 1 df['month'].unique()

array(['aug', 'nov', 'jun', 'apr', 'jul', 'may', 'oct', 'mar', 'sep',
      'dec'], dtype=object)

[13] 1 df['day_of_week'].unique()

array(['thu', 'fri', 'tue', 'mon', 'wed'], dtype=object)

[14] 1 df['poutcome'].unique()

array(['nonexistent', 'success', 'failure'], dtype=object)
```

- Then we can check the missing value percentage in order to get an idea on handling the data.

```
1 #Check for missing values in "job"
2 print('Unique Values in job column: ',df['job'].unique())
3
4 #Get the unknownvalue percentage out of the total number of
5 #data in the column in order to decide how to handle the missing values.
6 print('unkown entries percentage: ', len(df[df['job'] == 'unknown']/len(df['job'])*100)
```

Unique Values in job column: ['blue-collar' 'technician' 'management' 'services' 'retired' 'admin.' 'housemaid' 'unemployed' 'entrepreneur' 'self-employed' 'unknown' 'student']

unkown entries percentage: 0.8012042342429834

As the missing value percentage is low in the feature “job” (0.8%), we can choose to remove the rows with the "unknown" entry.

```
1 #There are "unknown" fields which we might have to handle
2 #As the missing value percentage is low, we can choose to
3 #remove the rows with the "unknown" entry
4 housing_dropping_index = df[ df['job'] == 'unknown' ].index
5 df.drop(housing_dropping_index, inplace = True)
6 df.shape
```

(40858, 20)

It can be done for all the columns which show a lower percentage of missing value count.

Marital:0.1%

```
[10] 1 #Check for missing values in "marital"
2 print('Unique Values in job column: ',df['marital'].unique())
3
4 #Get the unkown value percentage out of the total number of
5 #data in the column in order to decide how to handle the missing values.
6 print('unkown entries percentage: ', len(df[df['marital'] == 'unknown']/len(df['marital'])*100)
```

```
Unique Values in job column: ['married' 'single' 'divorced' 'unknown']
unkown entries percentage: 0.17377257819766018
```

```
1 #As the missing value percentage is low, we can choose to remove the rows with the "unknown" entry
2 housing_dropping_index = df[ df['marital'] == 'unknown' ].index
3 df.drop(housing_dropping_index, inplace = True)
4 df.shape
```

```
(40787, 20)
```

## Education: 3.9%

```
1 #Check for missing values in "education"
2 print('Unique Values in job column: ',df['education'].unique())
3
4 #Get the unkown value percentage out of the total number of data in the column in order to decide how
5 print('unkown entries percentage: ', len(df[df['education'] == 'unknown']/len(df['education'])*100)
```

```
Unique Values in job column: ['basic.4y' 'unknown' 'university.degree' 'high.school' 'basic.9y'
'professional.course' 'basic.6y' 'illiterate']
unkown entries percentage: 3.9130114987618607
```

```
[13] 1 #As the missing value percentage is low, we can choose to remove the rows with the "unknown" entry
2 housing_dropping_index = df[ df['education'] == 'unknown' ].index
3 df.drop(housing_dropping_index, inplace = True)
4 df.shape
```

```
(39191, 20)
```

## Loan: 20%

20% is a considerable "unknown" percentage. Still, the percentage is not high enough to fully drop the column. However, removing the rows would result in a considerable effect on the final predictions. Thus, decided to keep the "default" column as it is.

```
[14] 1 #Check for missing values in "default"
2 print('Unique Values in job column: ',df['default'].unique())
3
4 #Get the unkown value percentage out of the total number of data in the column in order to decide how
5 print('unkown entries percentage: ', len(df[df['default'] == 'unknown']/len(df['default'])*100)
```

```
Unique Values in job column: ['unknown' 'no' 'yes']
unkown entries percentage: 20.320992064504605
```

```
[15] 1 #As the missing value percentage is low, we can choose to remove the rows with the "unknown" entry
2 housing_dropping_index = df[ df['housing'] == 'unknown' ].index
3 df.drop(housing_dropping_index, inplace = True)
4 df.shape

(38245, 20)

[16] 1 #Check for missing values in "loan"
2 print('Unique Values in job column: ',df['loan'].unique())
3
4 #Get the unknown value percentage out of the total number of data in the column in order to decide how to handle the missing values.
5 print('unknown entries percentage: ', len(df[df['loan'] == 'unknown'])/len(df['loan'])*100)

Unique Values in job column:  ['no' 'yes']
unknown entries percentage:  0.0
```

Poutcome: 86%

A very high \*86% of missing values \* "nonexistent" are there in the poutcome column. As the missing data is extremely high, it will affect the final outcome if we use this column as a feature. Therefore, we can decide to **drop** this column.

When the housing column dropped the unknown rows, loan column has dropped its rows as well

```
[17] 1 #Check for missing values in "poutcome"
2 print('Unique Values in job column: ',df['poutcome'].unique())
3
4 #Get the unknown value percentage out of the total number of data in the column in order to decide how to handle the m
5 print('unknown entries percentage: ', len(df[df['poutcome'] == 'nonexistent'])/len(df['poutcome'])*100)

Unique Values in job column:  ['nonexistent' 'success' 'failure']
unknown entries percentage:  86.45836057000915
```

A very high **86% of missing values** "nonexistent" are there in the poutcome column. As the missing data is extremely high, it will affect the final outcome if we use this column as a feature. Therefore, we can decide to drop this column.

```
[24] 1 df = df.drop(['poutcome'],axis=1)
2 df.head(5)
```

	age	job	marital	education	default	housing	loan	contact	month	day_of_week	campaign	pdays	previous	emp
0	44	blue-collar	married	basic.4y	unknown	yes	no	cellular	aug	thu	1	999	0	
2	28	management	single	university.degree	no	yes	no	cellular	jun	thu	3	6	2	
3	39	services	married	high.school	no	no	no	cellular	apr	fri	2	999	0	
4	55	retired	married	basic.4y	no	yes	no	cellular	aug	fri	1	3	1	
5	30	management	divorced	basic.4y	no	yes	no	cellular	jul	tue	8	999	0	

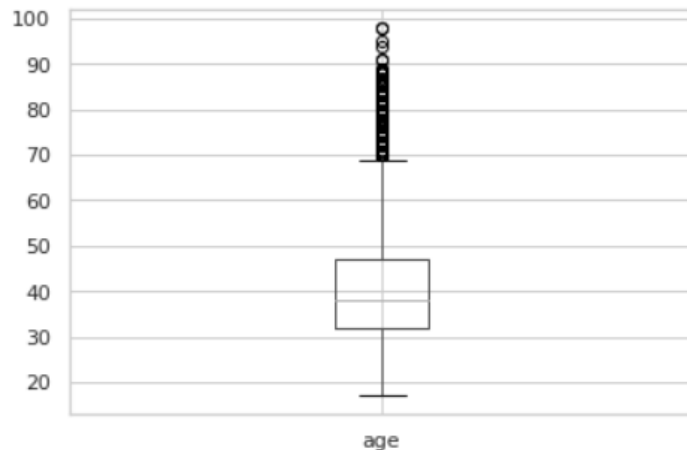
## Checking and handling Outliers

Boxplots can identify the outliers efficiently.

## Checking for Outliers and handling them.

```
[ ] 1 #Checking for outliers in "age"
    2 df.boxplot(column=['age'])
```

<matplotlib.axes.\_subplots.AxesSubplot at 0x7f3c5eb91110>

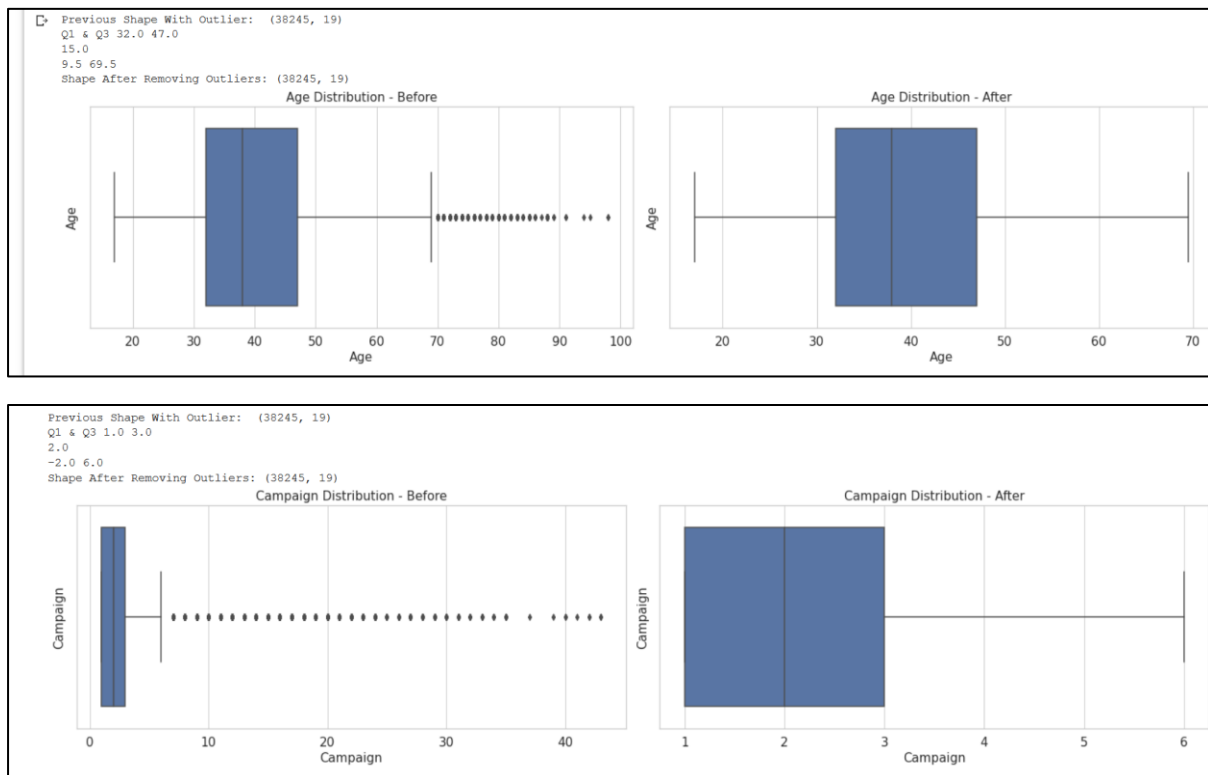


Removing the outliers in age would affect the overall result as the outlier count is almost half the total dataset. Thus, we can decide to replace the missing values using the Inter Quartile Range.

```
1 import warnings
2 warnings.filterwarnings("ignore")
3
4 numeric_col = df.select_dtypes(include=['int', 'float']).columns
5 print("Previous Shape With Outlier: ", df.shape)
6 Q1 = df.age.quantile(0.25)
7 Q3 = df.age.quantile(0.75)
8 print('Q1 & Q3', Q1, Q3)
9 IQR = Q3 - Q1
10 print(IQR)
11
12 lower_limit = Q1 - 1.5 * IQR
13 upper_limit = Q3 + 1.5 * IQR
14 print(lower_limit, upper_limit)
15
16 df_copy = df.copy()
17 df_copy['age'] = np.where(df_copy['age'] > upper_limit, upper_limit, df_copy['age'])
18 df_copy['age'] = np.where(df_copy['age'] < lower_limit, lower_limit, df_copy['age'])
19 print("Shape After Removing Outliers:", df_copy.shape)
20 fig, (ax1, ax2) = plt.subplots(nrows=1, ncols=2, figsize=(20, 5))
21
22 #Manipulating the boxplot views
23 sns.boxplot(x='age', data=df[numeric_col], orient='v', ax=ax1)
24 ax1.set_xlabel('Age', fontsize=15)
25 ax1.set_ylabel('Age', fontsize=15)
26 ax1.set_title('Age Distribution - Before', fontsize=15)
27 ax1.tick_params(labelsize=15)
28
29 sns.boxplot(x='age', data=df_copy[numeric_col], orient='v', ax=ax2)
30 ax2.set_xlabel('Age', fontsize=15)
31 ax2.set_ylabel('Age', fontsize=15)
32 ax2.set_title('Age Distribution - After', fontsize=15)
33 ax2.tick_params(labelsize=15)
34
35
36 plt.subplots_adjust(wspace=0.5)
37 plt.tight_layout()
38
```

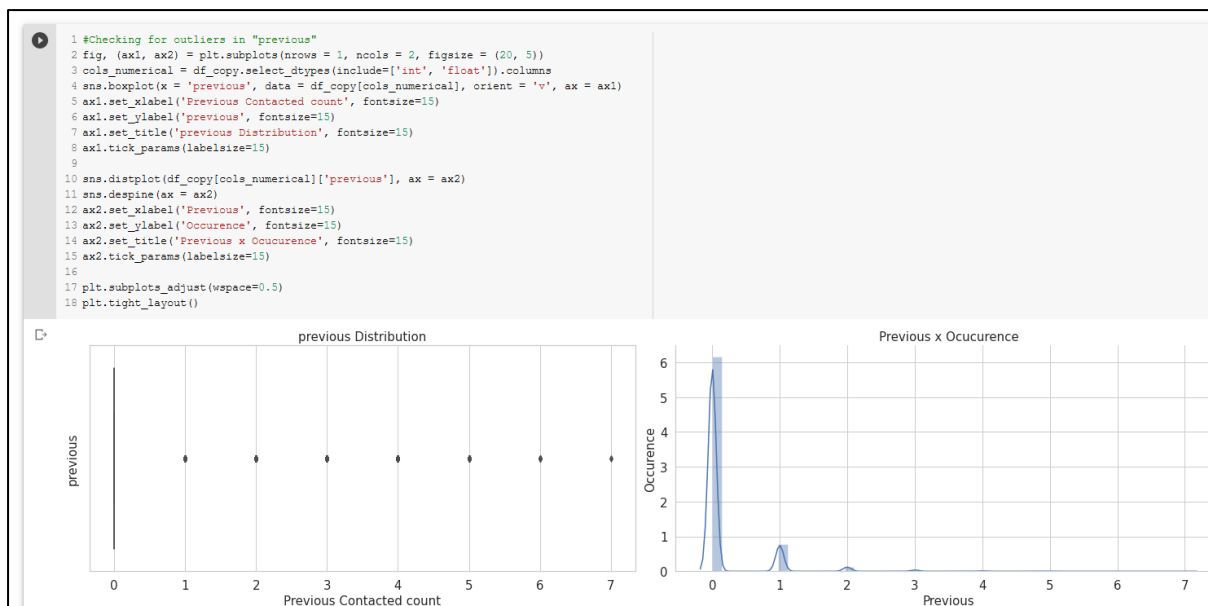


This is how the Before and after boxplots look like.

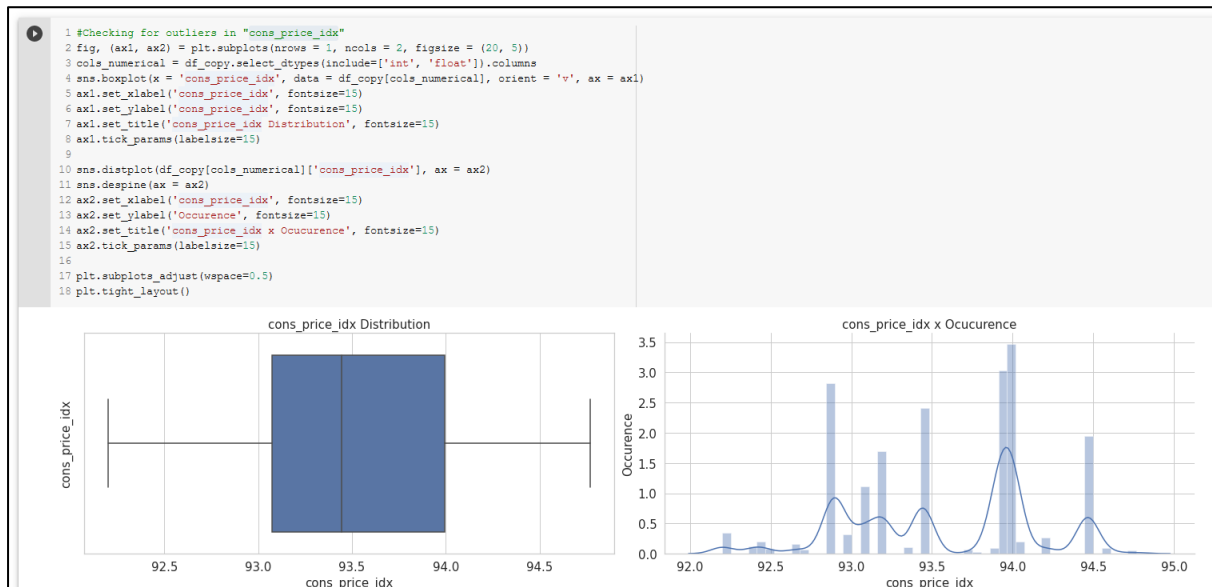
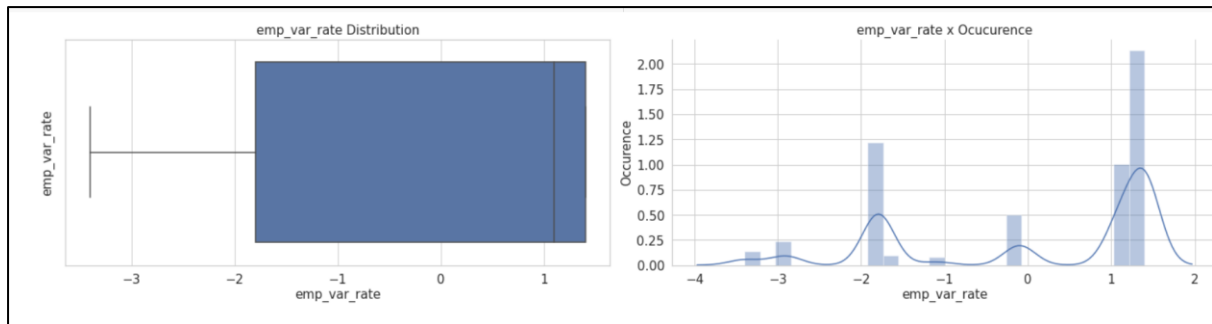


Number of previous contacts are distributed between 0 times and 7 times. Number of previous contacts are distributed between 0 times and 7 times.

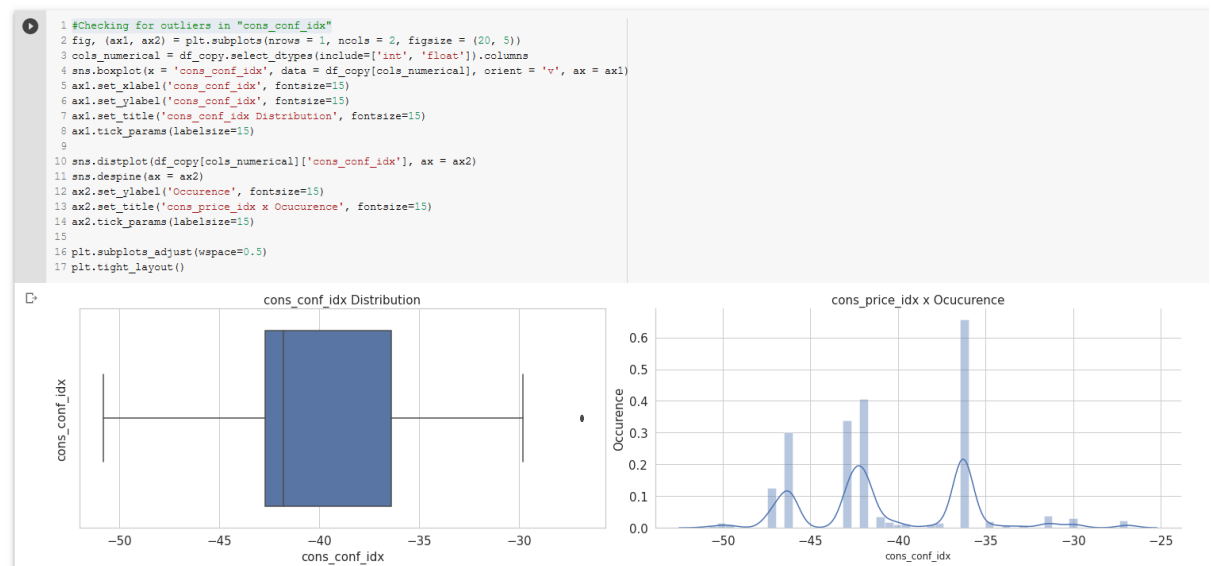
This would not be manipulated as that might make it loose almost all the data in previous column.



emp\_var\_rate and cons\_price\_idx has no visible outliers to be removed:



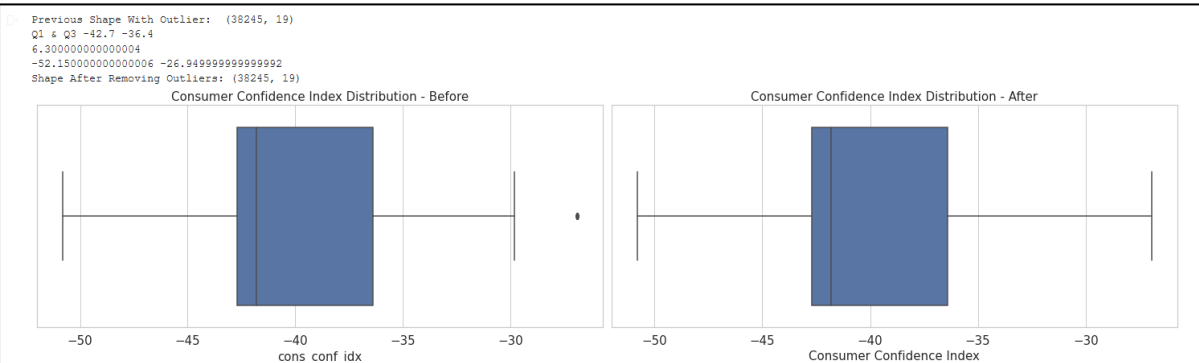
Checking for outliers in "cons\_conf\_idx" which will then be replaced using the IQR



```

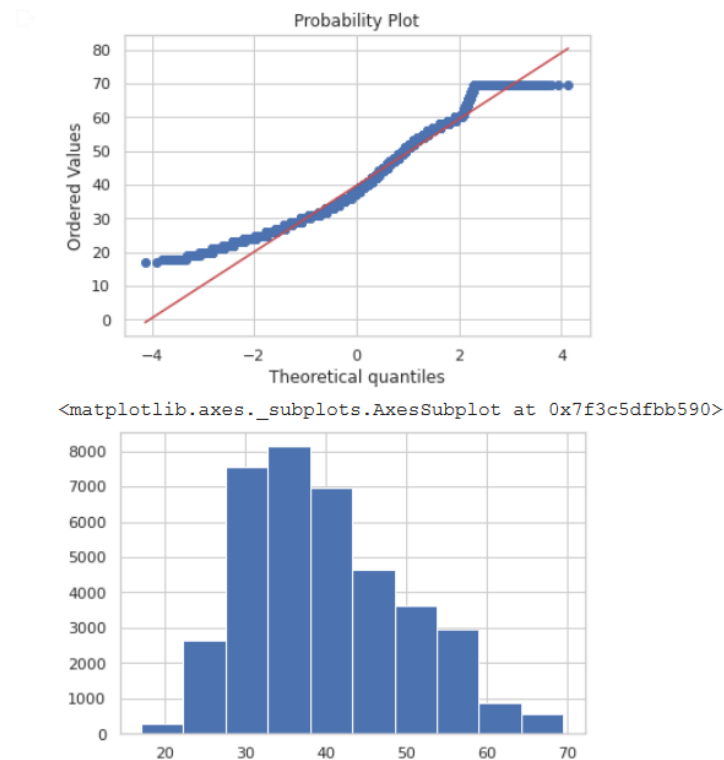
1 #Handling outliers for Consumer Confidence Index
2
3 print("Previous Shape With Outlier: ",df_copy.shape)
4 Q1 = df_copy.cons_conf_idx.quantile(0.25)
5 Q3 = df_copy.cons_conf_idx.quantile(0.75)
6 print('Q1 & Q3',Q1,Q3)
7 IQR = Q3-Q1
8 print(IQR)
9
10 lower_limit = Q1 - 1.5 * IQR
11 upper_limit = Q3 + 1.5 * IQR
12 print(lower_limit,upper_limit)
13
14
15 df_copy['cons_conf_idx'] = np.where(df_copy['cons_conf_idx']>upper_limit,upper_limit,df_copy['cons_conf_idx'])
16 df_copy['cons_conf_idx'] = np.where(df_copy['cons_conf_idx']<lower_limit,lower_limit,df_copy['cons_conf_idx'])
17 print("Shape After Removing Outliers:", df_copy.shape)
18 fig, (ax1, ax2) = plt.subplots(nrows = 1, ncols = 2, figsize = (20, 5))
19
20 sns.boxplot(x = 'cons_conf_idx', data = df[numeric_col], orient = 'v', ax = ax1)
21 ax1.set_xlabel('cons_conf_idx', fontsize=15)
22 ax1.set_title('Consumer Confidence Index Distribution - Before', fontsize=15)
23 ax1.tick_params(labelsize=15)
24
25 sns.boxplot(x = 'cons_conf_idx', data = df_copy[numeric_col], orient = 'v', ax = ax2)
26 ax2.set_xlabel('Consumer Confidence Index', fontsize=15)
27 ax2.set_title('Consumer Confidence Index Distribution - After', fontsize=15)
28 ax2.tick_params(labelsize=15)
29
30
31 plt.subplots_adjust(wspace=0.5)

```



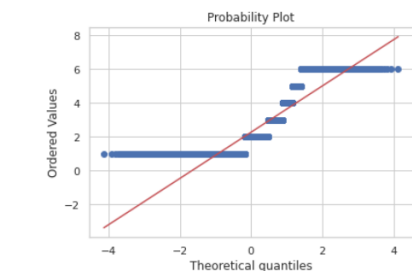
- II. Produce **Q-Q Plots** and Histograms of the features, and apply the transformations if required

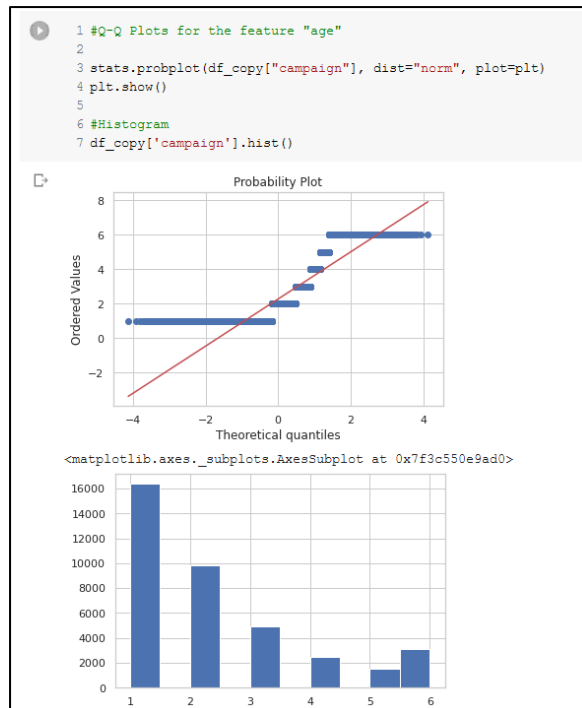
```
[ ] 1 #Q-Q Plots for the feature "age"
2 # import the libraries
3 import matplotlib.pyplot as plt
4 import scipy.stats as stats
5 import pandas as pd
6
7 #Q-Q Plot
8 stats.probplot(df_copy["age"], dist="norm", plot=plt)
9 plt.show()
10
11 #Histogram
12 df_copy['age'].hist()
```



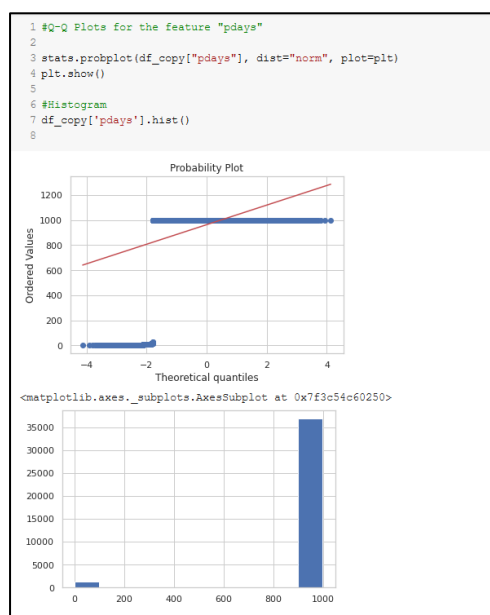
The feature "age" is right skewed when we take a look at the histogram. Thus, we can apply the technique of Transforming **Right Skewed Data using logarithm**

```
[ ] 1 #Q-Q Plots for the feature "age"
2
3 stats.probplot(df_copy["campaign"], dist="norm", plot=plt)
4 plt.show()
5
6 #Histogram
7 df_copy['campaign'].hist()
```

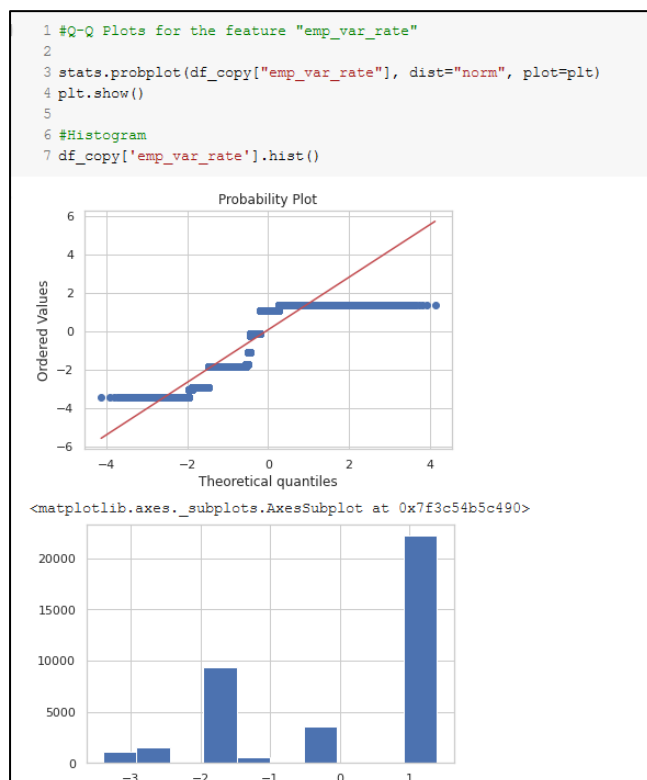
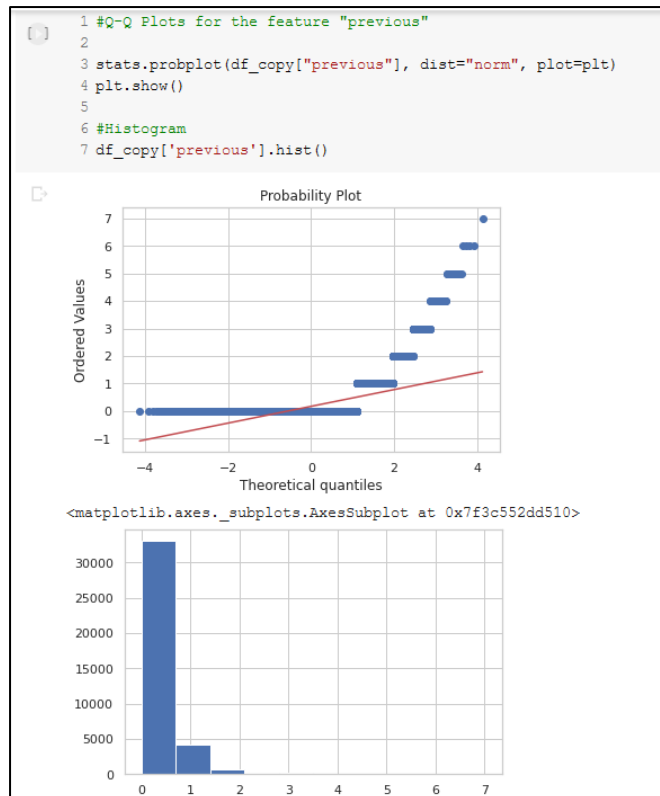




The feature "campaign" is also right skewed when we take a look at the histogram. Thus, we can apply the technique of **Transforming** Right Skewed Data using logarithm



\*Pdays\* is heavily left skewed. Exponential transformation can be used to transform for left skewed data.

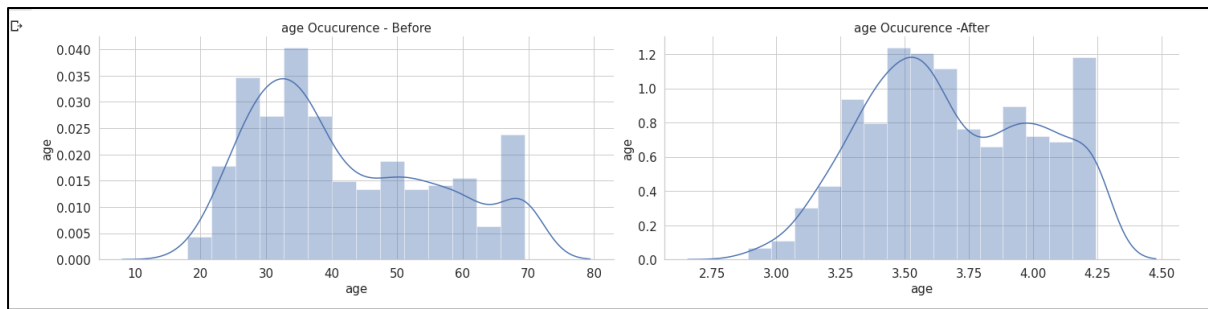


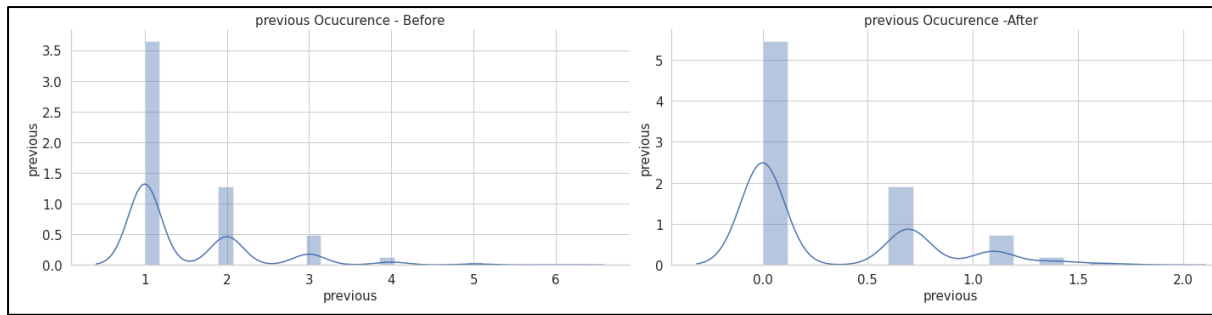
emp\_var\_rate is left skewed. Exponential transformation can be used to transform for left skewed data.

- emp\_var\_rate is left skewed. Exponential transformation can be used to transform for left skewed data.
- exponential\_transformation can be applied for left skewed data

## All Right Skewed features will be transformed using Logarithmic transformation method

```
1 from sklearn.preprocessing import FunctionTransformer
2 # load your data
3 data = df_copy
4 data = data.replace(0, np.nan) #replace 0 s with nan because ln
5 data = data.dropna() #drop all nan
6 columns = ['age', 'campaign', 'previous', 'cons_price_idx']
7
8 # create the function transformer object with logarithm transformation
9 logarithm_transformer = FunctionTransformer(np.log, validate=True)
10
11 # apply the transformation to the data
12 data_new = logarithm_transformer.transform(data[columns])
13 df_new = pd.DataFrame(data_new, columns=columns)
14
15 #histograms for 'age', 'campaign' Before & After the transformation
16 fig, (ax1, ax2) = plt.subplots(nrows = 1, ncols = 2, figsize = (20, 5))
17
18 sns.distplot(data[numeric_col]['age'], ax = ax1)
19 sns.despine(ax = ax1)
20 ax1.set_xlabel('age', fontsize=15)
21 ax1.set_ylabel('age', fontsize=15)
22 ax1.set_title('age Ocurence - Before', fontsize=15)
23 ax1.tick_params(labelsize=15)
24
25 sns.distplot(df_new[columns]['age'], ax = ax2)
26 sns.despine(ax = ax2)
27 ax2.set_xlabel('age', fontsize=15)
28 ax2.set_ylabel('age', fontsize=15)
29 ax2.set_title('age Ocurence -After', fontsize=15)
30 ax2.tick_params(labelsize=15)
31
32 plt.subplots_adjust(wspace=0.5)
33 plt.tight_layout()
```





```

1 #exponential_transformation for left skewed data
2 exponential_transformer = FunctionTransformer(np.exp)
3 columns2 = ['pdays', 'emp_var_rate', 'nr_employed']
4
5 data2 = df_copy
6
7 data_new2 = exponential_transformer.transform(data2[columns2])
8 df_new2 = pd.DataFrame(data_new2, columns=columns2)
9 df_new2.head(2)
10

```

	pdays	emp_var_rate	nr_employed
0	inf	4.055200	inf
2	403.428793	0.182684	inf

Now we can combine the transformed features with the complete dataset we have:

```

[46] 1 #replacing age & campaign column data with transformed data
2 data['age'] = df_new['age']
3 data['campaign'] = df_new['campaign']

```

```

[48] 1 #replacing 'pdays', 'emp_var_rate' column data with transformed data
2 data['pdays'] = df_new2['pdays']
3 data['emp_var_rate'] = df_new2['emp_var_rate']

```

```

1 data.head(10)

```

	age	job	marital	education	default	housing	loan	contact	month	day_of_week	campaign	pdays	previous
2	3.583519	management	single	university.degree	no	yes	no	cellular	jun	thu	0.000000	403.428793	2.0



### III. Feature Coding

Integer (Label) Encoding is used as it will not add new columns and does not expand feature space

```
[50] 1 #Analyze the data  
      2 data.dtypes
```

```
age          float64  
job           object  
marital       object  
education     object  
default       object  
housing       object  
loan          object  
contact       object  
month         object  
day_of_week   object  
campaign      float64  
pdays        float64  
previous      float64  
emp_var_rate  float64  
cons_price_idx float64  
cons_conf_idx float64  
euribor3m     float64  
nr_employed   float64  
y             float64  
dtype: object
```

```
1 data['job_code'] = data['job'].astype('category').cat.codes  
2 data['marital_code'] = data['marital'].astype('category').cat.codes  
3 data['education_code'] = data['education'].astype('category').cat.codes  
4 data['default_code'] = data['default'].astype('category').cat.codes  
5 data['housing_code'] = data['housing'].astype('category').cat.codes  
6 data['loan_code'] = data['loan'].astype('category').cat.codes  
7 data['contact_code'] = data['contact'].astype('category').cat.codes  
8 data['month_code'] = data['month'].astype('category').cat.codes  
9 data['day_of_week_code'] = data['day_of_week'].astype('category').cat.codes  
10 data.head(10)
```

After data encoding:

```
1 #Verify the new categorical columns exist
2 data.dtypes
```

age	float64
job	object
marital	object
education	object
default	object
housing	object
loan	object
contact	object
month	object
day_of_week	object
campaign	float64
pdays	float64
previous	float64
emp_var_rate	float64
cons_price_idx	float64
cons_conf_idx	float64
euribor3m	float64
nr_employed	float64
y	float64
job_code	int8
marital_code	int8
education_code	int8
default_code	int8
housing_code	int8
loan_code	int8
contact_code	int8
month_code	int8
day_of_week_code	int8
dtype:	object

Now we can drop all the unencoded textual columns from the dataframe:

```
Drop the textual columns as we have already encoded them and added as new columns
```

```
1 data = data.drop(['job', 'marital', 'education', 'default', 'housing', 'loan', 'contact', 'month', 'day_of_week'], axis=1)
2 data.head(5)
```

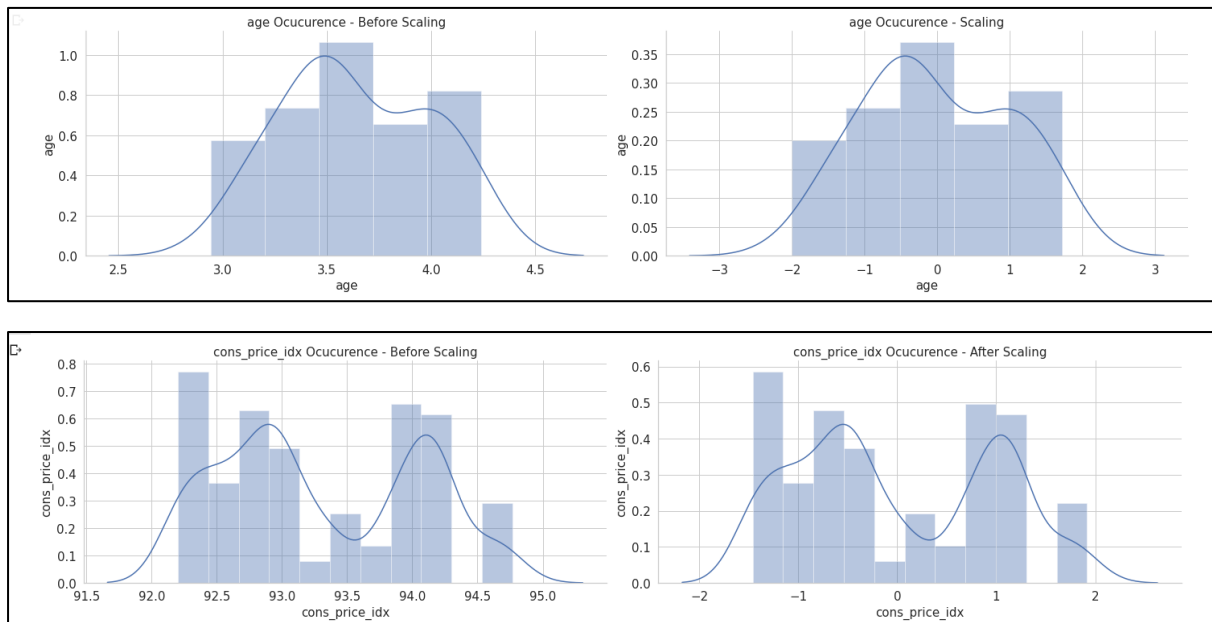
	age	campaign	pdays	previous	emp_var_rate	cons_price_idx	cons_conf_idx	euribor3m	nr_employed	y	job_code
2	3.583519	0.000000	403.428793	2.0	0.182684	94.055	-39.8	0.729	4991.6	1.0	
4	3.178054	0.000000	20.085537	1.0	0.055023	92.201	-31.4	0.869	5076.2	1.0	
8	3.713572	1.098612	20.085537	1.0	0.055023	92.963	-40.8	1.266	5076.2	1.0	
96	3.218876	0.693147	403.428793	1.0	0.165299	93.749	-34.6	0.659	5008.7	1.0	
117	3.583519	0.000000	inf	2.0	0.182684	94.027	-38.3	0.905	4991.6	1.0	

## IV. Scaling data

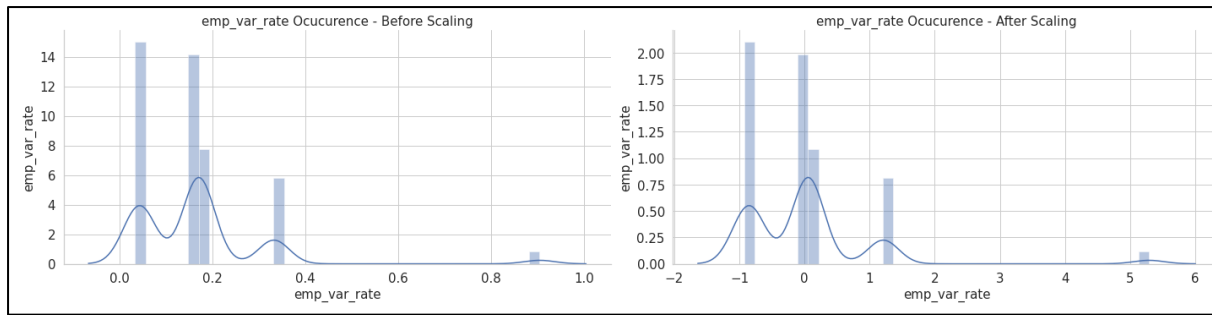
### Scaling features

```
1 # Applying scaling
2
3 from sklearn.preprocessing import StandardScaler
4
5 # create the scaler object
6 scaler = StandardScaler()
7 scaling_columns=['age', 'emp_var_rate', 'cons_price_idx', 'cons_conf_idx', 'euribor3m', 'nr_employed']
8 scaled = pd.DataFrame(data, columns = ['age', 'emp_var_rate', 'cons_price_idx', 'cons_conf_idx', 'euribor3m', 'nr_employed'])
9 # fit the scaler to the data
10 scaler.fit(scaled)
11
12 train_scaled = scaler.transform(scaled)
13 scaled_df = pd.DataFrame(train_scaled)
14 # scaled_df.info()
15
16
17 fig, (ax1, ax2) = plt.subplots(nrows = 1, ncols = 2, figsize = (20, 5))
18
19 sns.distplot(data[numeric_col]['age'], ax = ax1)
20 sns.despine(ax = ax1)
21 ax1.set_xlabel('age', fontsize=15)
22 ax1.set_ylabel('age', fontsize=15)
23 ax1.set_title('age Ocurence - Before Scaling', fontsize=15)
24 ax1.tick_params(labelsize=15)
25
26 sns.distplot(scaled_df[0], ax = ax2)
27 sns.despine(ax = ax2)
28 ax2.set_xlabel('age', fontsize=15)
29 ax2.set_ylabel('age', fontsize=15)
30 ax2.set_title('age Ocurence - Scaling', fontsize=15)
31 ax2.tick_params(labelsize=15)
32
33
34 plt.subplots_adjust(wspace=0.5)
35 plt.tight_layout()
```

This is how the scaled age column looks like:



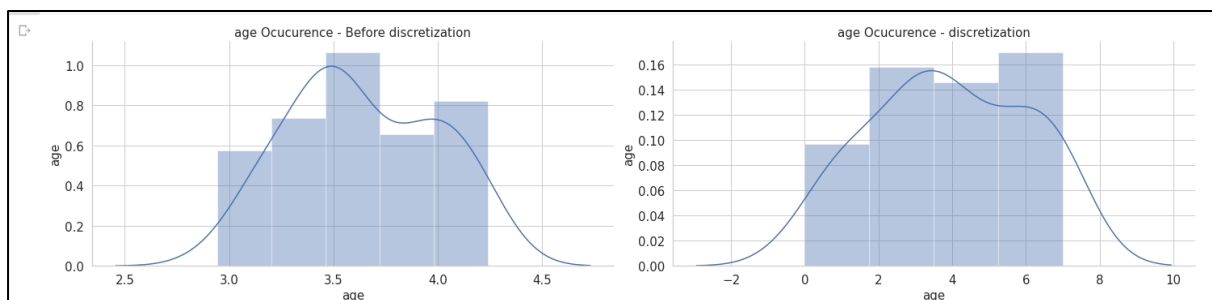
## V. Feature discretization



KBinsDiscretizer was used on age to discretize the feature:

### ▼ (a)V. Feature discretization

```
[59] 1 import pandas as pd
      2 from sklearn.preprocessing import KBinsDiscretizer
      3
      4
      5 # create the scaler object
      6 scaler = StandardScaler()
      7
      8
      9 data5 = pd.DataFrame(data, columns=['age'])
     10 data5 = data5.dropna()
     11 # fit the scaler to the data
     12 discretizer = KBinsDiscretizer(n_bins=8, encode='ordinal', strategy='kmeans')
     13 discretizer.fit(data5)
     14 _discretize = discretizer.transform(data5)
     15 df_discret = pd.DataFrame(_discretize)
     16 df_discret
     17
```



(b)

I. SVD (Singular Value Decomposition) for feature reduction.

```

1 # calculate the SVD for the data X
2 U_boston, S_boston, Vt_boston = np.linalg.svd(df_copy, full_matrices=False)
3
4 #numpy.linalg.svd actually returns a E that is not a diagonal matrix, but a list of the entries on the diagonal.
5 num_sv_boston = np.arange(1, S_boston.size+1)
6 cum_var_explained_boston = [np.sum(np.square(S_boston[0:n])) / np.sum(np.square(S_boston)) for n in num_sv_boston]
7 print(S_boston)
8 print(S_boston.size)

```

```

[1.02834811e+06 3.45687932e+04 1.96412421e+03 9.39371486e+02
 6.97669711e+02 4.83694968e+02 4.17055250e+02 3.64590642e+02
 2.99496657e+02 2.70172453e+02 1.45441507e+02 1.07597114e+02
 9.72502617e+01 7.60077278e+01 7.39901478e+01 7.07496323e+01
 6.61331454e+01 5.48248662e+01 3.10460592e+01]
19

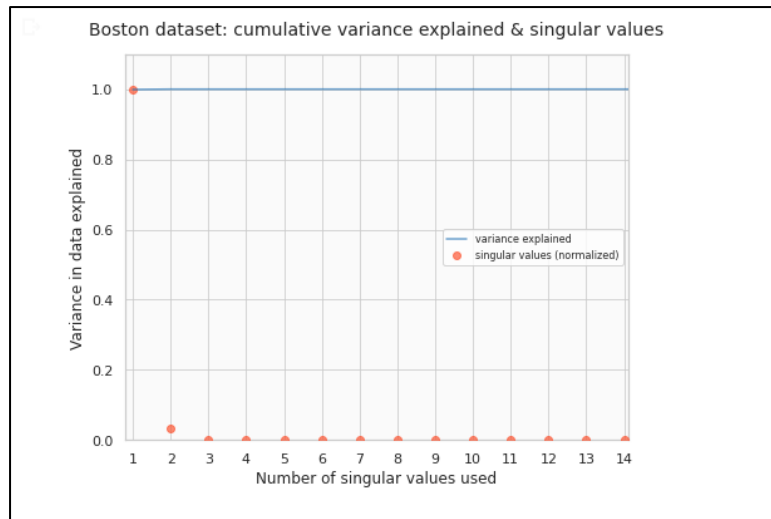
```

SVD gives a weighted value for each feature so that we can choose the most suitable ones.

```

1 import matplotlib.pyplot as plt
2 import sklearn.datasets
3 import sklearn.preprocessing
4 fig = plt.figure(figsize=(7.0,5.5))
5 ax = fig.add_subplot(111)
6
7 plt.plot(num_sv_boston,
8         cum_var_explained_boston,
9         color='#2171b5',
10        label='variance explained',
11        alpha=0.65,
12        zorder=1000)
13
14 plt.scatter(num_sv_boston,
15            sklearn.preprocessing.normalize(S_boston.reshape((1,-1))),
16            color='#fc4e2a',
17            label='singular values (normalized)',
18            alpha=0.65,
19            zorder=1000)
20
21 plt.legend(loc='center right', scatterpoints=1, fontsize=8)
22
23 ax.set_xticks(num_sv_boston)
24 ax.set_xlim(0.8, 14.1)
25 ax.set_ylim(0.0, 1.1)
26 ax.set_xlabel('Number of singular values used')
27 ax.set_ylabel('Variance in data explained')
28 ax.set_title('Boston dataset: cumulative variance explained & singular values',
29             fontsize=14,
30             y=1.03)
31
32 ax.set_facecolor('0.98')
33
34 plt.grid(alpha=0.8, zorder=1)
35 plt.tight_layout()

```



## II. Significant and independent features.

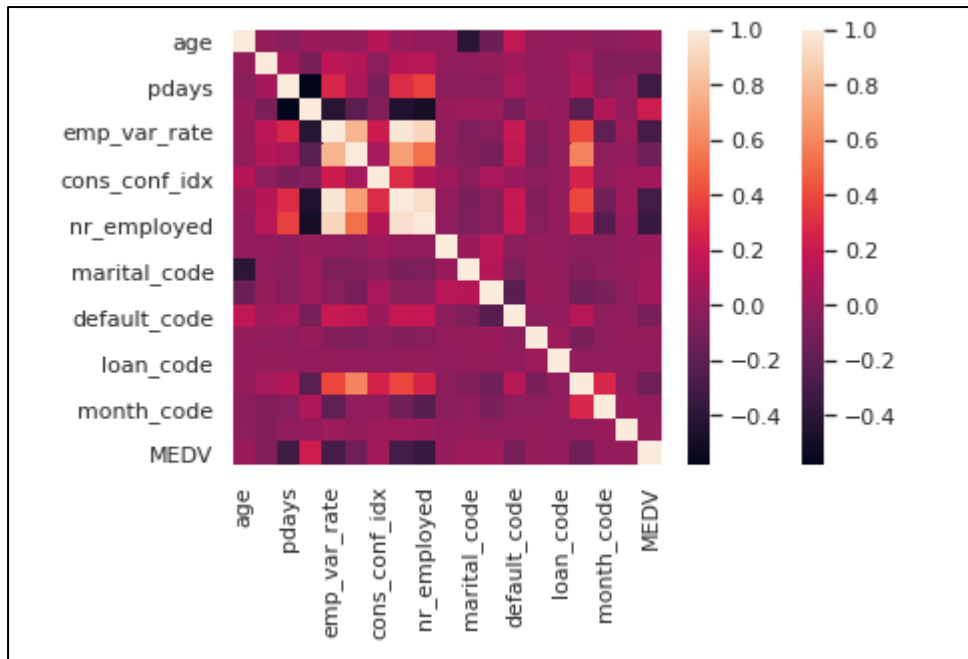
```

1 #test the signifncance of the features
2 import seaborn as sns
3 sns.heatmap(X.corr()); #Seems they can be assumed as independent
4
5 d_data = X.copy()
6 d_data['MEDV'] = y_true
7 d_data.head(10)
8 print(d_data.corr())
9 sns.heatmap(d_data.corr()) #Seems we are ok with the signifncancy

```

	age	campaign	...	day_of_week_code	MEDV
age	1.000000	0.004124	...	-0.018721	0.021775
campaign	0.004124	1.000000	...	-0.053519	-0.066827
pdays	-0.031600	0.057497	...	-0.008175	-0.319351
previous	0.020978	-0.081660	...	-0.004482	0.221159
emp_var_rate	0.008260	0.148144	...	0.032926	-0.292265
cons_price_idx	0.004577	0.115014	...	0.002583	-0.133084
cons_conf_idx	0.125476	-0.015081	...	0.040086	0.051331
euribor3m	0.020149	0.128578	...	0.038921	-0.300580
nr_employed	-0.006760	0.140767	...	0.029894	-0.347830
job_code	-0.004145	-0.005976	...	-0.002870	0.020569
marital_code	-0.397646	-0.011171	...	0.003796	0.041535
education_code	-0.146225	0.003764	...	-0.013620	0.054755
default_code	0.166061	0.038111	...	-0.009527	-0.096116
housing_code	-0.000294	-0.010134	...	0.001921	0.009996
loan_code	-0.006548	0.011173	...	-0.010177	-0.005603
contact_code	0.009825	0.073329	...	-0.011993	-0.140866
month_code	-0.025781	-0.063189	...	0.025514	-0.002884
day_of_week_code	-0.018721	-0.053519	...	1.000000	0.014957
MEDV	0.021775	-0.066827	...	0.014957	1.000000

[19 rows x 19 columns]



Darker the color the higher the correlation between the two variables.

### (c). Applying Logistic Regression and Support Vector Machine techniques

First drop the y column from the data set for testing:

```

[ ]  1 y_true = df_copy['y']
     2 X = df_copy.drop('y', axis=1)

```

Split the dataset to **80% and 20%**

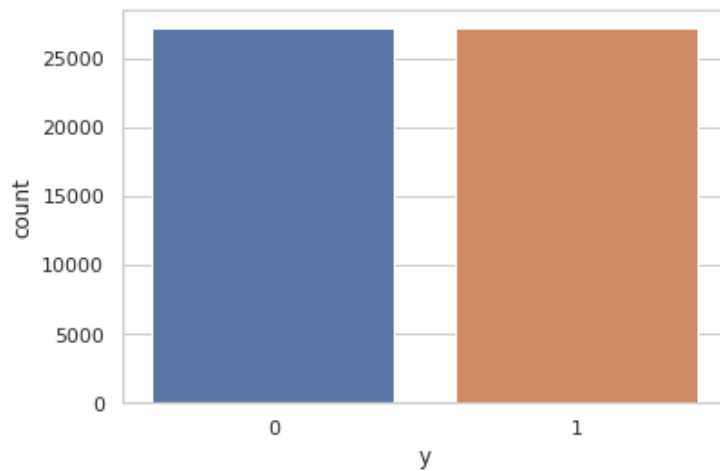
```

Splitting test and training data sets

1 from sklearn.model_selection import train_test_split
2 #Dont apply smote on test data
3 from imblearn.over_sampling import SMOTE
4
5 os = SMOTE(random_state=0)
6 X_class_train, X_test, y_class_train, y_test = train_test_split(X, y_true, test_size=0.2, random_state=0)
7 columns = X_class_train.columns
8
9 data_X, data_y = os.fit_sample(X_class_train, y_class_train)
10
11 smoted_X = pd.DataFrame(data=data_X, columns=columns )
12 smoted_y= pd.DataFrame(data=data_y, columns=['y'])

```

```
[ ] 1 sns.countplot(x='y', data=smoted_y)
    2 plt.show()
```



## Training the dataset with Logistic Regression

```
[ ] 1 from sklearn.linear_model import LogisticRegression
    2 logreg = LogisticRegression()
    3 model = logreg.fit(X_train, y_train)
    4 y_pred = logreg.predict(X_test)
    5
    6 from sklearn.metrics import classification_report, confusion_matrix
    7 print(classification_report(y_test, y_pred))
```

	precision	recall	f1-score	support
0	0.95	0.75	0.84	6794
1	0.26	0.68	0.37	855
accuracy			0.74	7649
macro avg	0.60	0.72	0.61	7649
weighted avg	0.87	0.74	0.79	7649

```
[ ] 1 print(confusion_matrix(y_test, y_pred))
```

```
[[5113 1681]
 [ 274  581]]
```

Confusion matrix is as follows:

```
[ ] 1 print(confusion_matrix(y_test, y_pred))
```

```
[[5113 1681]
 [ 274  581]]
```



ROC Curve is used to identify the diagnostic capability of a model.

```
1 #ROC Curve
2
3 from sklearn.metrics import roc_auc_score
4 from sklearn.metrics import roc_curve
5 logit_roc_auc = roc_auc_score(y_test, logreg.predict(X_test))
6 fpr, tpr, thresholds = roc_curve(y_test, logreg.predict_proba(X_test)[:,1])
7 plt.figure()
8 plt.plot(fpr, tpr, label='Logistic Regression (area = %0.2f)' % logit_roc_auc)
9 plt.plot([0, 1], [0, 1], 'r--')
10 plt.xlim([0.0, 1.0])
11 plt.ylim([0.0, 1.05])
12 plt.xlabel('False Positive Rate')
13 plt.ylabel('True Positive Rate')
14 plt.title('Receiver operating characteristic')
15 plt.legend(loc="lower right")
16 plt.savefig('Log_ROC')
17 plt.show()
```



Training the data **with SVM**

```
[157] 1 import numpy as np
      2 import matplotlib.pyplot as plt
      3 from sklearn import svm, datasets
      4 from sklearn.svm import SVC
      5
      6 model_svm = SVC(kernel='linear')
      7 model_svm.fit(X_train, y_train)
      8 preds = model_svm.predict(X_test)
      9 metrics.accuracy_score(y_test, preds)
     10
     11
```

0.8109556804811087

```
[158] 1 #confusion matrix , evaluation
      2 print(confusion_matrix(y_test, preds))
```

```
[[5701 1093]
 [ 353  502]]
```

```
1 print(classification_report(y_test, preds))
```

	precision	recall	f1-score	support
0	0.94	0.84	0.89	6794
1	0.31	0.59	0.41	855
accuracy			0.81	7649
macro avg	0.63	0.71	0.65	7649
weighted avg	0.87	0.81	0.83	7649

## (d) applicability of SVM and LR

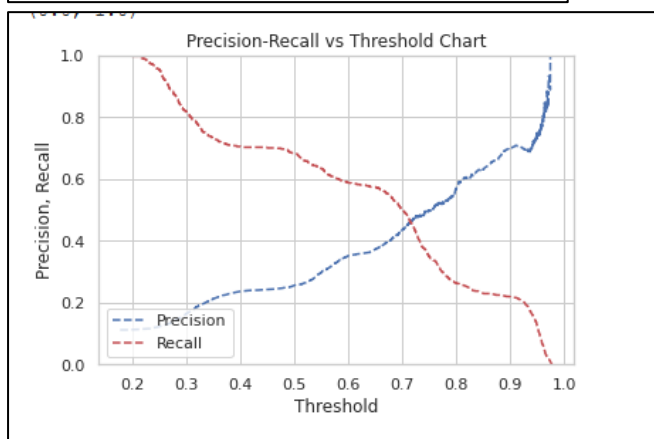
### Logistic Regression:

```
[ ] 1 from sklearn.linear_model import LogisticRegression
2 logreg = LogisticRegression()
3 model = logreg.fit(X_train, y_train)
4 y_pred = logreg.predict(X_test)
5
6 from sklearn.metrics import classification_report, confusion_matrix
7 print(classification_report(y_test, y_pred))
```

	precision	recall	f1-score	support
0	0.95	0.75	0.84	6794
1	0.26	0.68	0.37	855
accuracy			0.74	7649
macro avg	0.60	0.72	0.61	7649
weighted avg	0.87	0.74	0.79	7649

```
[ ] 1 print(confusion_matrix(y_test, y_pred))
```

```
[[5113 1681]
 [ 274  581]]
```



## Support Vector Machine

```
[158] 1 #confusion matrix , evaluation
      2 print(confusion_matrix(y_test,preds))
```

```
[[5701 1093]
 [ 353  502]]
```

```
1 print(classification_report(y_test, preds))
```

	precision	recall	f1-score	support
0	0.94	0.84	0.89	6794
1	0.31	0.59	0.41	855
accuracy			0.81	7649
macro avg	0.63	0.71	0.65	7649
weighted avg	0.87	0.81	0.83	7649

In Logistic Regression, the confusion matrix shows that negative-negative (true negative) count is 5113 while positive-positive count is 502. In SVM, negative-negative is 5701 and positive-positive (true-positive) count is 502 which is equal.

		Predicted	
		Negative	Positive
Actual	Negative	True Negative	False Positive
	Positive	False Negative	True Positive

According to the confusion matrix, SVM seems to be a better fit for our data set.

When comparing the 2 classification reports, LR shows an accuracy of 0.75 while SVM shows 0.84 which is a better score than LR.

Recall can be stated as the best metric to evaluate algorithms. This is because it takes the value of true positives over the total actual positives which gives a clear picture on the correctly predicted positive values.

In LR , recall value is 0.75 while in SVM it's 0.84, which suggests that SVM works better with the banking dataset.

$$\text{Recall} = \frac{\text{True Positive}}{\text{True Positive} + \text{False Negative}}$$

$$= \frac{\text{True Positive}}{\text{Total Actual Positive}}$$

Even in F1-score which is a good measure to check the balance between the two values, recall and precision, SVM has served our data set better.

$$\text{F1} = 2 \times \frac{\text{Precision} * \text{Recall}}{\text{Precision} + \text{Recall}}$$

Finally, we can conclude the fact that **SVM works best** with the banking dataset

### (e) How significant your findings are

Logistic Regression		precision	recall	f1-score	support
	0	0.95	0.75	0.84	6794
	1	0.26	0.68	0.37	855
	accuracy			0.74	7649
	macro avg	0.60	0.72	0.60	7649
	weighted avg	0.87	0.74	0.79	7649

Support Vector Machine		precision	recall	f1-score	support
	0	0.94	0.84	0.89	6794
	1	0.31	0.59	0.41	855
	accuracy			0.81	7649
	macro avg	0.63	0.71	0.65	7649
	weighted avg	0.87	0.81	0.83	7649

The accuracy for SVM is 81% which is a considerable finding.

ROC Curve of svm is as follows.



```
[176] 1 #root mean squared error SVM
      2 from sklearn.metrics import mean_squared_error
      3 rms = np.sqrt(mean_squared_error(y_test,preds))
      4 print(rms)

0.4347922716871717
```

RMSE shows a value of 0.43. An RMSE value lower than 0.5 reflects a good prediction rate of the model.

In Logistic Regression, recall value is 0.75 while in SVM it's 0.84. Both the values are considerably good scores which signifies the findings.

\*\*\* END OF REPORT \*\*\*