

# Classification and Regression with Bank Marketing Campaign dataset

## 3. Initial Results and the Code

By: Marina Golberg 501072689 CIND820

Code and documentation for this project on GitHub repository as following: <https://github.com/marinagolberg/CIND820-MarGolb.git> (<https://github.com/marinagolberg/CIND820-MarGolb.git>)

```
In [884]: #!pip install SMOTE
```

```
In [885]: #!pip install imblearn
```

```
In [886]: #!pip install mlxtend
```

```
In [887]: #!pip install matplotlib
```

```
In [888]: #!pip install seaborn
```

```
In [889]: #pip install cufflinks
```

```
In [890]: import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
from pandas import DataFrame
from pandas.plotting import scatter_matrix
from sklearn.preprocessing import LabelEncoder
import plotly.express as px
from collections import Counter
from imblearn.over_sampling import SMOTE
from sklearn.tree import DecisionTreeClassifier
from sklearn.linear_model import LogisticRegression
from sklearn.feature_selection import mutual_info_classif, SelectPercentile
from sklearn.metrics import classification_report, confusion_matrix
from sklearn.ensemble import RandomForestRegressor, RandomForestClassifier
from sklearn.metrics import roc_auc_score
from mlxtend.feature_selection import SequentialFeatureSelector
from sklearn import metrics
from sklearn.model_selection import train_test_split
from sklearn.metrics import brier_score_loss
from sklearn.feature_selection import SelectFromModel
from mlxtend.plotting import plot_sequential_feature_selection as plot_sfs
from sklearn.metrics import average_precision_score
from sklearn.model_selection import StratifiedKFold
from sklearn.model_selection import cross_val_score
from sklearn.model_selection import KFold
from sklearn.model_selection import RepeatedStratifiedKFold
```

```
In [891]: bank = pd.read_csv("bank-additional-full.csv", sep=';')
```

```
In [892]: cross_validation_data = pd.read_csv("bank-additional-full.csv", sep=';')
```

Attribute Information:

Input variables:

**Bank client data:**

- 1 - age (numeric)
- 2 - job : type of job (categorical: 'admin.', 'blue-collar', 'entrepreneur', 'housemaid', 'management', 'retired', 'self-employed', 'services', 'student', 'technician', 'unemployed', 'unknown')
- 3 - marital : marital status (categorical: 'divorced', 'married', 'single', 'unknown'; note: 'divorced' means divorced or widowed)
- 4 - education (categorical: 'basic.4y', 'basic.6y', 'basic.9y', 'high.school', 'illiterate', 'professional.course', 'university.degree', 'unknown')
- 5 - default: has credit in default? (categorical: 'no', 'yes', 'unknown')
- 6 - housing: has housing loan? (categorical: 'no', 'yes', 'unknown')
- 7 - loan: has personal loan? (categorical: 'no', 'yes', 'unknown')

**Related with the last contact of the current campaign:**

- 8 - contact: contact communication type (categorical: 'cellular', 'telephone')
- 9 - month: last contact month of year (categorical: 'jan', 'feb', 'mar', ..., 'nov', 'dec')
- 10 - day\_of\_week: last contact day of the week (categorical: 'mon', 'tue', 'wed', 'thu', 'fri')
- 11 - duration: last contact duration, in seconds (numeric). Important note: this attribute highly affects the output target (e.g., if duration=0 then y='no'). Yet, 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 if the intention is to have a realistic predictive model.

**Other attributes:**

- 12 - campaign: number of contacts performed during this campaign and for this client (numeric, includes last contact)
- 13 - pdays: number of days that passed by after the client was last contacted from a previous campaign (numeric; 999 means client was not previously contacted)

14 - previous: number of contacts performed before this campaign and for this client (numeric)

15 - poutcome: outcome of the previous marketing campaign (categorical: 'failure','nonexistent','success')

### **Social and economic context attributes**

16 - emp.var.rate: employment variation rate - quarterly indicator (numeric)

17 - cons.price.idx: consumer price index - monthly indicator (numeric)

18 - cons.conf.idx: consumer confidence index - monthly indicator (numeric)

19 - euribor3m: euribor 3 month rate - daily indicator (numeric)

20 - nr.employed: number of employees - quarterly indicator (numeric)

Output variable (desired target): 21 - y - has the client subscribed a term deposit? (binary: 'yes','no')

<https://archive.ics.uci.edu/ml/datasets/bank+marketing>.

### **Exploratory Data Analysis and Cleaning**

```
In [893]: bank.head(100)
```

Out[893]:

	age	job	marital	education	default	housing	loan	contact	month	day_of_week	...	campaign
0	56	housemaid	married	basic.4y	no	no	no	telephone	may	mon	...	1
1	57	services	married	high.school	unknown	no	no	telephone	may	mon	...	1
2	37	services	married	high.school	no	yes	no	telephone	may	mon	...	1
3	40	admin.	married	basic.6y	no	no	no	telephone	may	mon	...	1
4	56	services	married	high.school	no	no	yes	telephone	may	mon	...	1
...	...	...	...	...	...	...	...	...	...	...	...	...
95	45	services	married	professional.course	no	yes	no	telephone	may	mon	...	1
96	42	management	married	university.degree	no	no	no	telephone	may	mon	...	1
97	53	admin.	divorced	university.degree	unknown	no	no	telephone	may	mon	...	1
98	37	technician	single	professional.course	no	no	no	telephone	may	mon	...	1
99	44	blue-collar	married	basic.6y	no	no	no	telephone	may	mon	...	1

100 rows × 21 columns

```
In [894]: #basic descriptive statistics
# higth sd in "duration", "campaign", "previous", emp.var.rate, cons.conf.idx which indicates a fairl.
bank.describe()
```

Out[894]:

	age	duration	campaign	pdays	previous	emp.var.rate	cons.price.idx	cons.conf.idx	
<b>count</b>	41188.00000	41188.000000	41188.000000	41188.000000	41188.000000	41188.000000	41188.000000	41188.000000	41
<b>mean</b>	40.02406	258.285010	2.567593	962.475454	0.172963	0.081886	93.575664	-40.502600	
<b>std</b>	10.42125	259.279249	2.770014	186.910907	0.494901	1.570960	0.578840	4.628198	
<b>min</b>	17.00000	0.000000	1.000000	0.000000	0.000000	-3.400000	92.201000	-50.800000	
<b>25%</b>	32.00000	102.000000	1.000000	999.000000	0.000000	-1.800000	93.075000	-42.700000	
<b>50%</b>	38.00000	180.000000	2.000000	999.000000	0.000000	1.100000	93.749000	-41.800000	
<b>75%</b>	47.00000	319.000000	3.000000	999.000000	0.000000	1.400000	93.994000	-36.400000	
<b>max</b>	98.00000	4918.000000	56.000000	999.000000	7.000000	1.400000	94.767000	-26.900000	



```
In [895]: bank.groupby('y').mean()
```

Out[895]:

	age	duration	campaign	pdays	previous	emp.var.rate	cons.price.idx	cons.conf.idx	euribor3m	nr.emp
<b>y</b>										
<b>no</b>	39.911185	220.844807	2.633085	984.113878	0.132374	0.248875	93.603757	-40.593097	3.811491	5176.16
<b>yes</b>	40.913147	553.191164	2.051724	792.035560	0.492672	-1.233448	93.354386	-39.789784	2.123135	5095.11



std bigger then mean(duration,campaign,previous,emp.var.rate,cons.conf.idx)- high variation between values, and abnormal distribution for data. A smaller standard deviation indicates that more of the data is clustered about the mean while, a larger once indicates the data are more spread out.

```
In [896]: bank.shape
```

```
Out[896]: (41188, 21)
```

```
In [897]: bank['y'].value_counts()
```

```
Out[897]: no      36548  
yes     4640  
Name: y, dtype: int64
```

```
In [898]: #Some times, we want to know what percentage of the whole is  
#for each value that appears in the column.  
#To calculate this in pandas with the value_counts()  
#method, set the argument normalize to True.  
bank['y'].value_counts(normalize=True)
```

```
Out[898]: no      0.887346  
yes     0.112654  
Name: y, dtype: float64
```

That makes it highly unbalanced, the positive class of target variable for 11.26%

In [899]: bank.info()

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 41188 entries, 0 to 41187
Data columns (total 21 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  duration           41188 non-null    int64  
 11  campaign           41188 non-null    int64  
 12  pdays              41188 non-null    int64  
 13  previous            41188 non-null    int64  
 14  poutcome           41188 non-null    object  
 15  emp.var.rate       41188 non-null    float64 
 16  cons.price.idx     41188 non-null    float64 
 17  cons.conf.idx      41188 non-null    float64 
 18  euribor3m          41188 non-null    float64 
 19  nr.employed        41188 non-null    float64 
 20  y                  41188 non-null    object  
dtypes: float64(5), int64(5), object(11)
memory usage: 6.6+ MB
```

```
In [900]: #Check the datatypes of the attributes.  
bank.dtypes
```

```
Out[900]: age           int64  
job            object  
marital        object  
education      object  
default         object  
housing         object  
loan            object  
contact         object  
month           object  
day_of_week    object  
duration        int64  
campaign        int64  
pdays           int64  
previous        int64  
poutcome        object  
emp.var.rate   float64  
cons.price.idx float64  
cons.conf.idx  float64  
euribor3m      float64  
nr.employed    float64  
y               object  
dtype: object
```

```
In [901]: #Are there any missing values in the dataset?  
bank.isnull().values.any()
```

```
Out[901]: False
```

```
In [902]: bank.isna().sum()
```

```
Out[902]: age          0  
job           0  
marital       0  
education     0  
default        0  
housing        0  
loan           0  
contact        0  
month          0  
day_of_week    0  
duration        0  
campaign        0  
pdays           0  
previous        0  
poutcome         0  
emp.var.rate    0  
cons.price.idx  0  
cons.conf.idx   0  
euribor3m        0  
nr.employed      0  
y                 0  
dtype: int64
```

```
In [903]: bank.isin([0]).any().any()
```

```
Out[903]: True
```

```
In [904]: #How many 0 values in every attribute
```

```
bank.isin([0]).sum()
```

```
Out[904]: age          0  
job           0  
marital       0  
education     0  
default        0  
housing        0  
loan           0  
contact        0  
month          0  
day_of_week    0  
duration       4  
campaign       0  
pdays          15  
previous      35563  
poutcome       0  
emp.var.rate   0  
cons.price.idx 0  
cons.conf.idx  0  
euribor3m      0  
nr.employed    0  
y               0  
dtype: int64
```

```
In [905]: #previous 35563 is "0" (35563/41188 no data in this attribute ,I will drop this attribute)  
bank = bank.drop(['previous'], axis=1)
```

```
In [906]: #previous 35563 is "0" (35563/41188 no data in this attribute ,I will drop this attribute)  
cross_validation_data = cross_validation_data.drop(['previous'], axis=1)
```

```
In [907]: #Calculating the mean  
duration_mean = bank['duration']  
durationMean = duration_mean.mean()  
durationMean
```

```
Out[907]: 258.2850101971448
```

```
In [908]: #Calculating the mean  
duration_mean = cross_validation_data['duration']  
durationMean = duration_mean.mean()  
durationMean
```

```
Out[908]: 258.2850101971448
```

```
In [909]: #replacing all 0 values with mean of that column  
cross_validation_data = cross_validation_data.replace(0, durationMean)
```

```
In [910]: #replacing all 0 values with mean of that column  
bank = bank.replace(0, durationMean)
```

```
In [911]: bank.isin([0]).any().any()
```

```
Out[911]: False
```

```
In [912]: bank.isin([0]).sum()
```

```
Out[912]: age          0  
job           0  
marital       0  
education     0  
default        0  
housing        0  
loan           0  
contact        0  
month          0  
day_of_week    0  
duration        0  
campaign        0  
pdays           0  
poutcome        0  
emp.var.rate    0  
cons.price.idx  0  
cons.conf.idx   0  
euribor3m        0  
nr.employed     0  
y                 0  
dtype: int64
```

```
In [913]: #In the 'pdays' column, it is observed that 999 makes 96% of the values of the column.  
#from attribute information 999 means client was not previously contacted.  
# I suggest to drop this column as there is not enough information for further analysis.  
  
bank['pdays'].value_counts(normalize=True)
```

```
Out[913]: 999.00000    0.963217  
3.00000    0.010658  
6.00000    0.010003  
4.00000    0.002865  
9.00000    0.001554  
2.00000    0.001481  
7.00000    0.001457  
12.00000   0.001408  
10.00000   0.001263  
5.00000    0.001117  
13.00000   0.000874  
11.00000   0.000680  
1.00000    0.000631  
15.00000   0.000583  
14.00000   0.000486  
8.00000    0.000437  
258.28501  0.000364  
16.00000   0.000267  
17.00000   0.000194  
18.00000   0.000170  
22.00000   0.000073  
19.00000   0.000073  
21.00000   0.000049  
25.00000   0.000024  
26.00000   0.000024  
27.00000   0.000024  
20.00000   0.000024  
Name: pdays, dtype: float64
```

```
In [914]: bank = bank.drop(['pdays'], axis=1)
```

```
In [915]: cross_validation_data = cross_validation_data.drop(['pdays'], axis=1)
```

```
In [916]: #In the 'poutcome' column, it is observed that nonexistent +  
#failure makes 96.6% of the values of the column.  
#from attribute information 'poutcome' is outcome of  
#the previous marketing campaign  
# I will not drop this column as a success rate 3.3%  
#might be interesting for further analysis.  
bank['poutcome'].value_counts(normalize=True)
```

```
Out[916]: nonexistent    0.863431  
failure        0.103234  
success        0.033335  
Name: poutcome, dtype: float64
```

```
In [917]: #campaign: number of contacts performed during this campaign and for this client  
bank['campaign'].value_counts(normalize=True)
```

```
Out[917]: 1    0.428329  
2    0.256628  
3    0.129674  
4    0.064363  
5    0.038822  
6    0.023769  
7    0.015271  
8    0.009712  
9    0.006871  
10   0.005463  
11   0.004297  
12   0.003035  
13   0.002234  
14   0.001675  
17   0.001408  
16   0.001238  
15   0.001238  
18   0.000801  
20   0.000728  
19   0.000631  
21   0.000583  
22   0.000413  
23   0.000388  
24   0.000364  
27   0.000267  
29   0.000243  
28   0.000194  
26   0.000194  
25   0.000194  
31   0.000170  
30   0.000170  
35   0.000121  
32   0.000097  
33   0.000097  
34   0.000073
```

```
42    0.000049  
40    0.000049  
43    0.000049  
56    0.000024  
39    0.000024  
41    0.000024  
37    0.000024  
Name: campaign, dtype: float64
```

```
In [918]: bank['default'].value_counts(normalize=True)
```

```
Out[918]: no      0.791201  
unknown  0.208726  
yes     0.000073  
Name: default, dtype: float64
```

```
In [919]: #cons.conf.idx  
bank['cons.conf.idx'].value_counts()
```

```
Out[919]: -36.4    7763  
-42.7    6685  
-46.2    5794  
-36.1    5175  
-41.8    4374  
-42.0    3616  
-47.1    2458  
-31.4     770  
-40.8     715  
-26.9     447  
-30.1     357  
-40.3     311  
-37.5     303  
-50.0     282  
-29.8     267  
-34.8     264  
-38.3     233  
-39.8     229  
-40.0     212  
-49.5     204  
-33.6     178  
-34.6     174  
-33.0     172  
-50.8     128  
-40.4      67  
-45.9      10  
Name: cons.conf.idx, dtype: int64
```

To obtain a better understanding of the dataset, the distribution of key variables and the relationships among them were plotted.

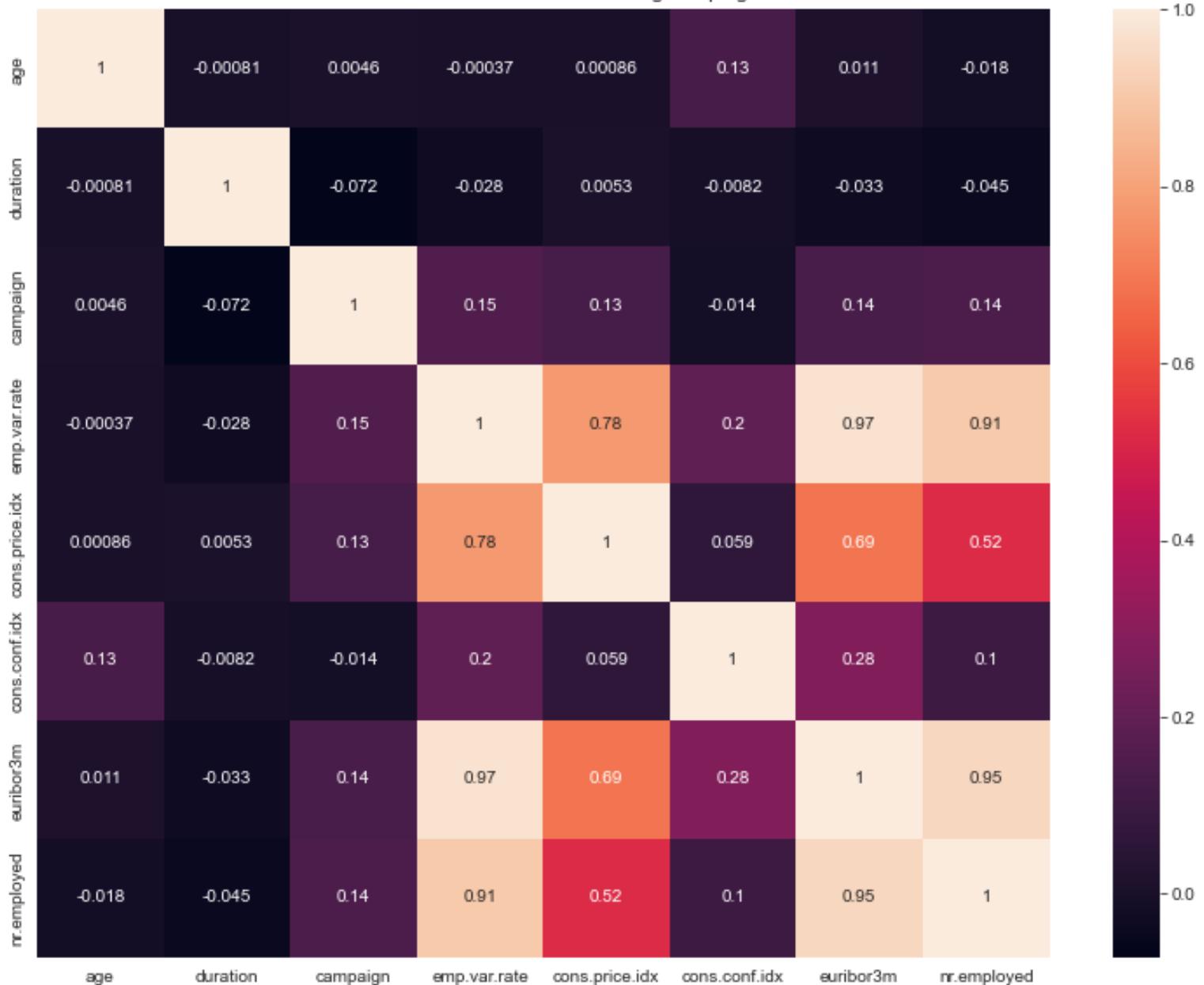
In [920]: #Creating Dataframe in Panda  
df = pd.DataFrame(bank)  
#print(df)

```
In [921]: # taking all rows and 11 columns(without y)
plt.figure(figsize=(13, 10))
df_corr = bank.iloc[:, :18]
correlation_mat = df_corr.corr()
sns.heatmap(correlation_mat, annot = True);
plt.title("Correlation matrix of bank Marketing campaign")

#plt.xlabel("attributes")
#plt.ylabel("attributes")

plt.show()
```

Correlation matrix of bank Marketing campaign



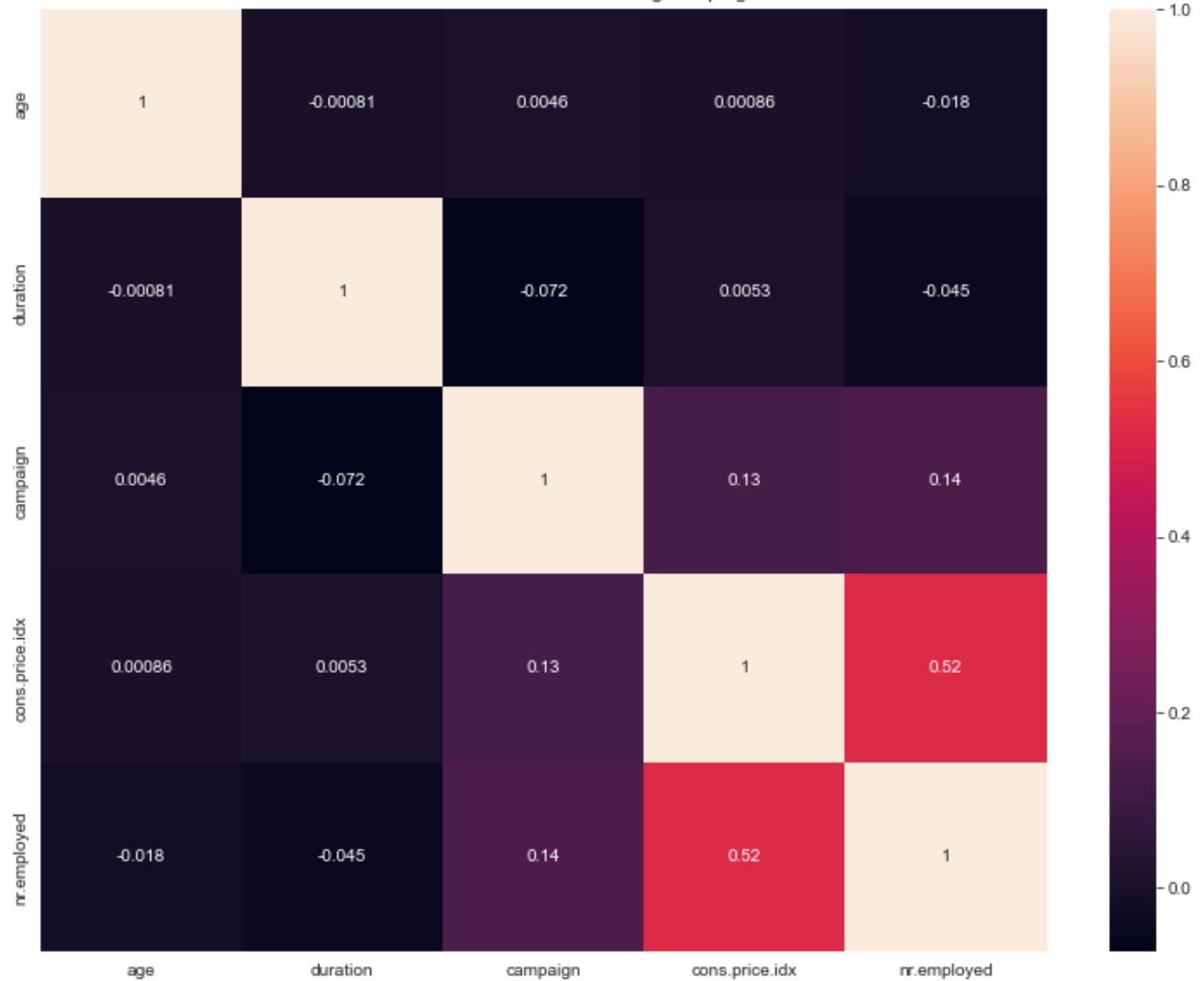
The social and economic context attributes have correlation among themselves. All columns with a high correlation will be removed to prevent Multicollinearity, it happens when predictor variable can be linearly predicted from the others with a high degree of accuracy. This can lead to skewed or misleading results. The columns are 'emp.var.rate','euribor3m', and 'cons.conf.idx'.

```
In [922]: bank = bank.drop(['emp.var.rate'], axis=1)
bank = bank.drop(['euribor3m'], axis=1)
bank = bank.drop(['cons.conf.idx'], axis=1)
```

```
In [923]: cross_validation_data = cross_validation_data.drop(['emp.var.rate'], axis=1)
cross_validation_data = cross_validation_data.drop(['euribor3m'], axis=1)
cross_validation_data = cross_validation_data.drop(['cons.conf.idx'], axis=1)
```

```
In [924]: # taking all rows and 11 columns(without y)
plt.figure(figsize=(13, 10))
df_corr = bank.iloc[:, :15]
correlation_mat = df_corr.corr()
sns.heatmap(correlation_mat, annot = True);
plt.title("Correlation matrix of bank Marketing campaign")
plt.show()
```

Correlation matrix of bank Marketing campaign



```
In [925]: bank.head()
```

Out[925]:

	age	job	marital	education	default	housing	loan	contact	month	day_of_week	duration	campaign	pout
0	56	housemaid	married	basic.4y	no	no	no	telephone	may	mon	261.0	1	none
1	57	services	married	high.school	unknown	no	no	telephone	may	mon	149.0	1	none
2	37	services	married	high.school	no	yes	no	telephone	may	mon	226.0	1	none
3	40	admin.	married	basic.6y	no	no	no	telephone	may	mon	151.0	1	none
4	56	services	married	high.school	no	no	yes	telephone	may	mon	307.0	1	none

### Creating different data samples for training and testing.

In this way, we can use the training set for training our model and testing set help evaluate whether the model can generalise well to new, unseen data. In this way I will prevent overfitting.

I would divide the data set into 2 portions in the ratio of 70:30 My target variable is 'y' included in training and test data samples, next steps I will divide the data set into more 2 portions

```
In [926]: training, testing = train_test_split(bank, test_size=0.3, random_state=25)
```

```
In [927]: print(f"No. of training examples: {training.shape[0]}")  
print(f"No. of testing examples: {testing.shape[0]}")
```

No. of training examples: 28831

No. of testing examples: 12357

```
In [928]: #Let's check for duplicates  
training.duplicated().any()
```

Out[928]: True

```
In [929]: training.head()
```

Out[929]:

	age	job	marital	education	default	housing	loan	contact	month	day_of_week	duration	campaign
8726	39	admin.	single	high.school	no	no	no	telephone	jun	wed	172.0	2
30619	50	entrepreneur	married	basic.9y	no	yes	no	telephone	may	mon	331.0	5
31121	30	blue-collar	divorced	high.school	unknown	no	no	cellular	may	wed	848.0	1
37287	33	admin.	married	high.school	no	yes	no	cellular	aug	mon	252.0	1
38307	44	admin.	divorced	high.school	no	no	no	cellular	oct	thu	634.0	1

```
In [930]: #this doesn't seem like the case of some customers randomly having similar details.  
#It looks like the data duplication happened while entering the data.  
training[training.duplicated(keep = False)]
```

Out[930]:

	age	job	marital	education	default	housing	loan	contact	month	day_of_week	duration	camp
1266	39	blue-collar	married	basic.6y	no	no	no	telephone	may	thu	124.0	
20072	55	services	married	high.school	unknown	no	no	cellular	aug	mon	33.0	
18464	32	technician	single	professional.course	no	yes	no	cellular	jul	thu	128.0	
32505	35	admin.	married	university.degree	no	yes	no	cellular	may	fri	348.0	
18465	32	technician	single	professional.course	no	yes	no	cellular	jul	thu	128.0	
32516	35	admin.	married	university.degree	no	yes	no	cellular	may	fri	348.0	
20216	55	services	married	high.school	unknown	no	no	cellular	aug	mon	33.0	
1265	39	blue-collar	married	basic.6y	no	no	no	telephone	may	thu	124.0	
28476	24	services	single	high.school	no	yes	no	cellular	apr	tue	114.0	
28477	24	services	single	high.school	no	yes	no	cellular	apr	tue	114.0	

```
In [931]: #Let's remove these duplicate rows.  
training.drop_duplicates(inplace = True)
```

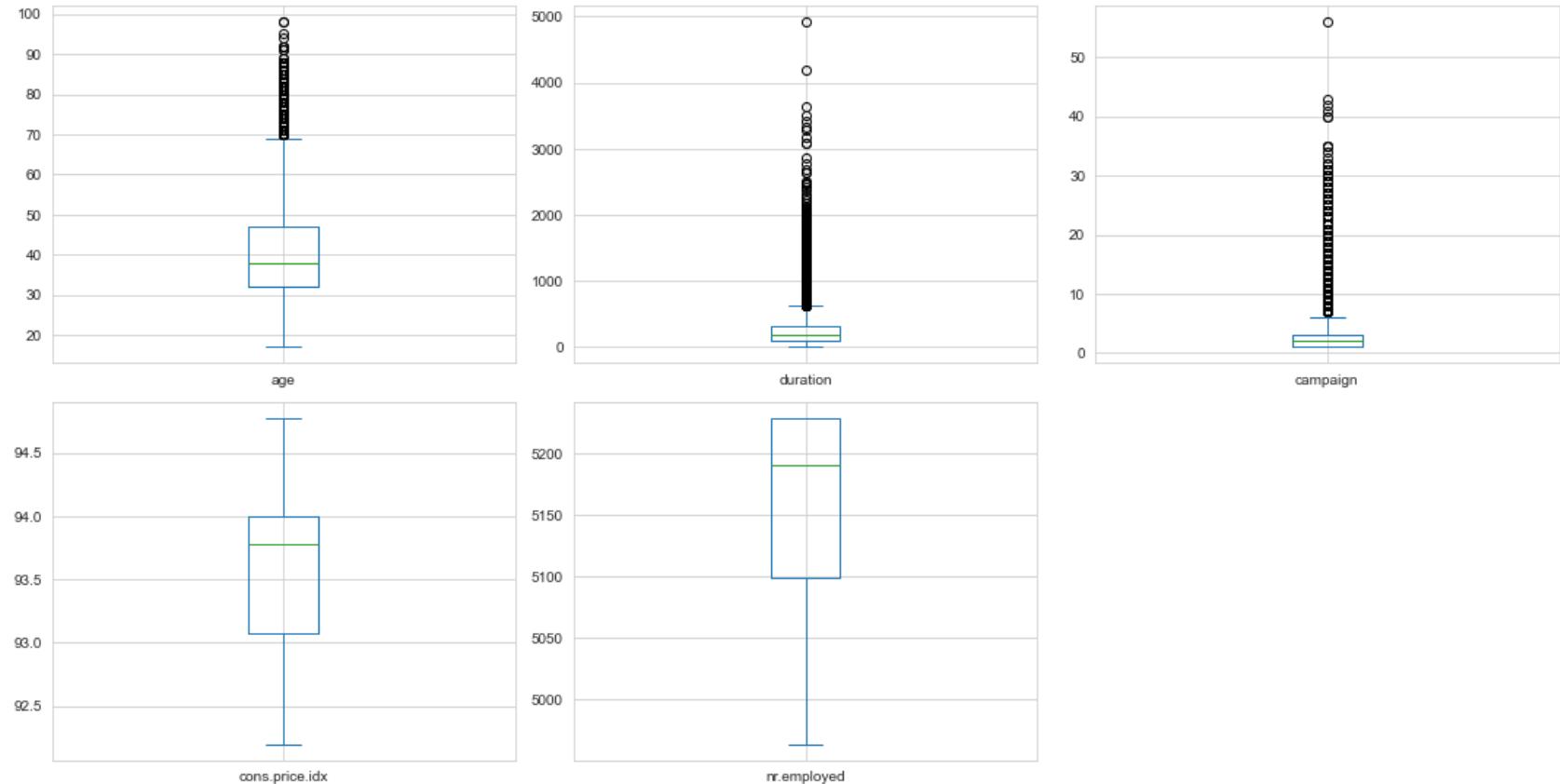
```
In [932]: training.shape
```

Out[932]: (28826, 16)

## Outliers Treatment

```
In [933]: #We need to install cufflinks to link plotly to pandas and add the iplot method:  
import cufflinks as cf  
cf.go_offline()  
cf.set_config_file(offline=False, world_readable=True)
```

```
In [934]: training.plot(kind='box', subplots=True, layout=(4,3), figsize=(15,15))  
plt.tight_layout()
```



On the boxplot above looks like there are outliers. Age-appropriate for the context of the attribute (min 17, max 98), Duration(is the last contact to the client in seconds max 4918 is 82 minutes for call it's too long but can be real), and maximum of Campaign looks very high 56 calls to the same customer very high but real, std bigger than mean(duration,

campaign)- high variation between values, and abnormal distribution for data. Probably the minimum and maximum values are the mistakes and other values in my opinion are appropriate. I will be removing only percentile 10 and percentile 90 because different deletion percentages will cause algorithms to perform worse.

In [935]: `training.describe()`

Out[935]:

	age	duration	campaign	cons.price.idx	nr.employed
<b>count</b>	28826.000000	28826.000000	28826.000000	28826.000000	28826.000000
<b>mean</b>	40.035107	256.775163	2.569208	93.577212	5166.675914
<b>std</b>	10.431859	256.988442	2.722317	0.580438	72.576247
<b>min</b>	17.000000	1.000000	1.000000	92.201000	4963.600000
<b>25%</b>	32.000000	102.000000	1.000000	93.075000	5099.100000
<b>50%</b>	38.000000	179.000000	2.000000	93.773500	5191.000000
<b>75%</b>	47.000000	317.000000	3.000000	93.994000	5228.100000
<b>max</b>	98.000000	4918.000000	56.000000	94.767000	5228.100000

In [936]: `min_duration, max_duration = training.duration.quantile([0.10, 0.90])`  
`min_duration, max_duration`

Out[936]: (59.0, 547.0)

```
In [937]: training[training.duration < min_duration]
```

Out[937]:

	age	job	marital	education	default	housing	loan	contact	month	day_of_week	duration	campa
33806	39	blue-collar	single	university.degree	no	no	no	cellular	may	wed	9.0	
9626	49	blue-collar	single	basic.4y	no	no	no	telephone	jun	mon	47.0	
19469	30	technician	married	high.school	no	no	no	cellular	aug	thu	51.0	
34378	39	blue-collar	married	basic.9y	unknown	no	no	cellular	may	thu	37.0	
17969	46	admin.	married	university.degree	no	no	no	telephone	jul	tue	53.0	
...	...	...	...	...	...	...	...	...	...	...	...	...
11803	23	blue-collar	single	basic.4y	no	yes	no	telephone	jun	fri	31.0	
27639	52	services	divorced	basic.4y	unknown	yes	yes	cellular	nov	fri	23.0	
13213	30	blue-collar	married	high.school	unknown	yes	no	cellular	jul	wed	44.0	
32394	35	blue-collar	married	basic.9y	no	no	no	cellular	may	fri	25.0	
35702	42	services	single	high.school	no	yes	no	cellular	may	mon	36.0	

2833 rows × 16 columns

```
In [938]: training[training.duration > max_duration]
```

Out[938]:

	age	job	marital	education	default	housing	loan	contact	month	day_of_week	duration	...
31121	30	blue-collar	divorced	high.school	unknown	no	no	cellular	may	wed	848.0	
38307	44	admin.	divorced	high.school	no	no	no	cellular	oct	thu	634.0	
20483	53	management	married	university.degree	no	yes	no	cellular	aug	tue	1133.0	
40652	42	admin.	single	high.school	no	no	no	cellular	sep	wed	638.0	
1086	53	retired	married	basic.4y	unknown	yes	no	telephone	may	wed	591.0	
...	...	...	...	...	...	...	...	...	...	...	...	...
39553	36	unemployed	married	professional.course	no	yes	no	cellular	apr	thu	835.0	
26767	36	blue-collar	married	basic.9y	no	no	no	cellular	nov	thu	768.0	
6618	49	management	divorced	university.degree	no	no	no	telephone	may	wed	675.0	
24894	48	services	married	high.school	no	yes	no	cellular	nov	tue	835.0	
29828	31	admin.	single	university.degree	no	yes	yes	cellular	apr	mon	885.0	

2879 rows × 16 columns

```
In [939]: training = training[(training.duration<max_duration)&(training.duration>min_duration)]  
training.shape
```

Out[939]: (23009, 16)

```
In [940]: training.describe()
```

Out[940]:

	age	duration	campaign	cons.price.idx	nr.employed
<b>count</b>	23009.000000	23009.000000	23009.000000	23009.000000	23009.000000
<b>mean</b>	40.125516	209.681466	2.349689	93.568078	5165.763549
<b>std</b>	10.550189	116.774299	2.217667	0.582286	72.986774
<b>min</b>	17.000000	60.000000	1.000000	92.201000	4963.600000
<b>25%</b>	32.000000	117.000000	1.000000	93.075000	5099.100000
<b>50%</b>	38.000000	180.000000	2.000000	93.444000	5191.000000
<b>75%</b>	47.000000	277.000000	3.000000	93.994000	5228.100000
<b>max</b>	98.000000	546.000000	56.000000	94.767000	5228.100000

```
In [941]: min_campaign, max_campaign = training.campaign.quantile([0.10, 0.90])
min_campaign, max_campaign
```

Out[941]: (1.0, 4.0)

```
In [942]: training[training.campaign < min_campaign]
```

Out[942]:

```
age job marital education default housing loan contact month day_of_week duration campaign poutcome cons.
```



```
In [943]: training[training.campaign > max_campaign]
```

Out[943]:

	age	job	marital	education	default	housing	loan	contact	month	day_of_week	duration	ca
30619	50	entrepreneur	married	basic.9y	no	yes	no	telephone	may	mon	331.0	
11493	49	admin.	divorced	high.school	no	yes	no	telephone	jun	fri	115.0	
8492	35	blue-collar	married	basic.4y	unknown	yes	no	telephone	jun	wed	180.0	
37111	31	admin.	single	high.school	no	no	yes	telephone	jul	tue	258.0	
10258	39	services	divorced	high.school	no	yes	no	telephone	jun	mon	81.0	
...	...	...	...	...	...	...	...	...	...	...	...	
30055	61	admin.	married	university.degree	no	yes	yes	cellular	apr	thu	266.0	
4148	24	admin.	single	high.school	no	yes	no	telephone	may	mon	243.0	
17224	26	services	married	high.school	no	yes	no	cellular	jul	fri	442.0	
17937	33	technician	married	professional.course	no	yes	yes	cellular	jul	tue	110.0	
16690	24	services	single	high.school	no	no	no	cellular	jul	wed	74.0	

2289 rows × 16 columns

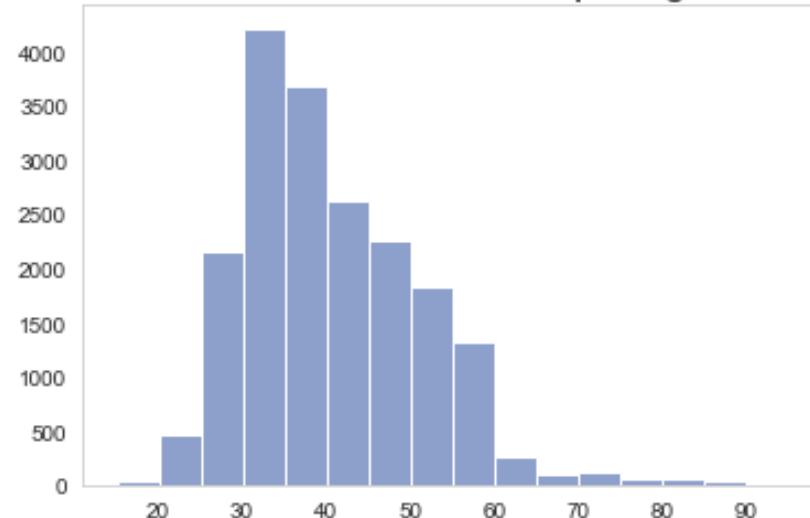


```
In [944]: training = training[(training.campaign<max_campaign)]
training.shape
```

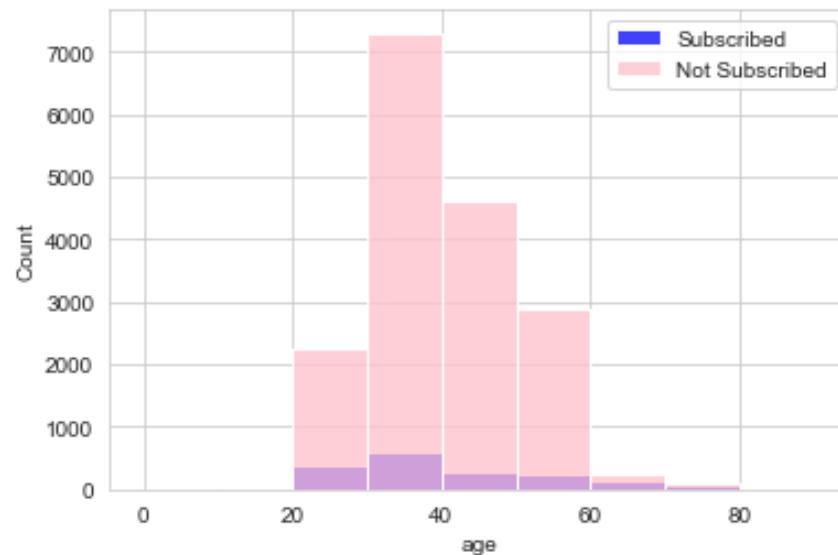
Out[944]: (19251, 16)

```
In [945]: #Let's see the distribution of clients per age.  
base_color = sns.color_palette('Set2')[2]  
  
age_bins = np.arange(15, 100, 5)  
plt.hist(data = training, x = 'age', bins = age_bins, color = base_color);  
plt.title("Distribution of clients per age", fontsize=18)  
plt.grid();
```

Distribution of clients per age



```
In [946]: #Visualize the relationship between feature category vs dependent variable y
#type of age and deposit
bins = range(0, 100, 10)
ax = sns.histplot(training.age[training.y=='yes'],
                  color='blue', kde=False, bins=bins, label='Subscribed')
sns.histplot(training.age[training.y=='no'],
              ax=ax, # Overplots on first plot
              color='pink', kde=False, bins=bins, label="Not Subscribed")
plt.legend()
plt.show()
```



Customers who in (30-40) followed by (20-30) and (40-50) had higher percentage of subscription to deposit account

```
In [947]: #Crosstab to display job stats with respect to y class variable  
pd.crosstab(index=training["job"], columns=training["y"])
```

Out[947]:

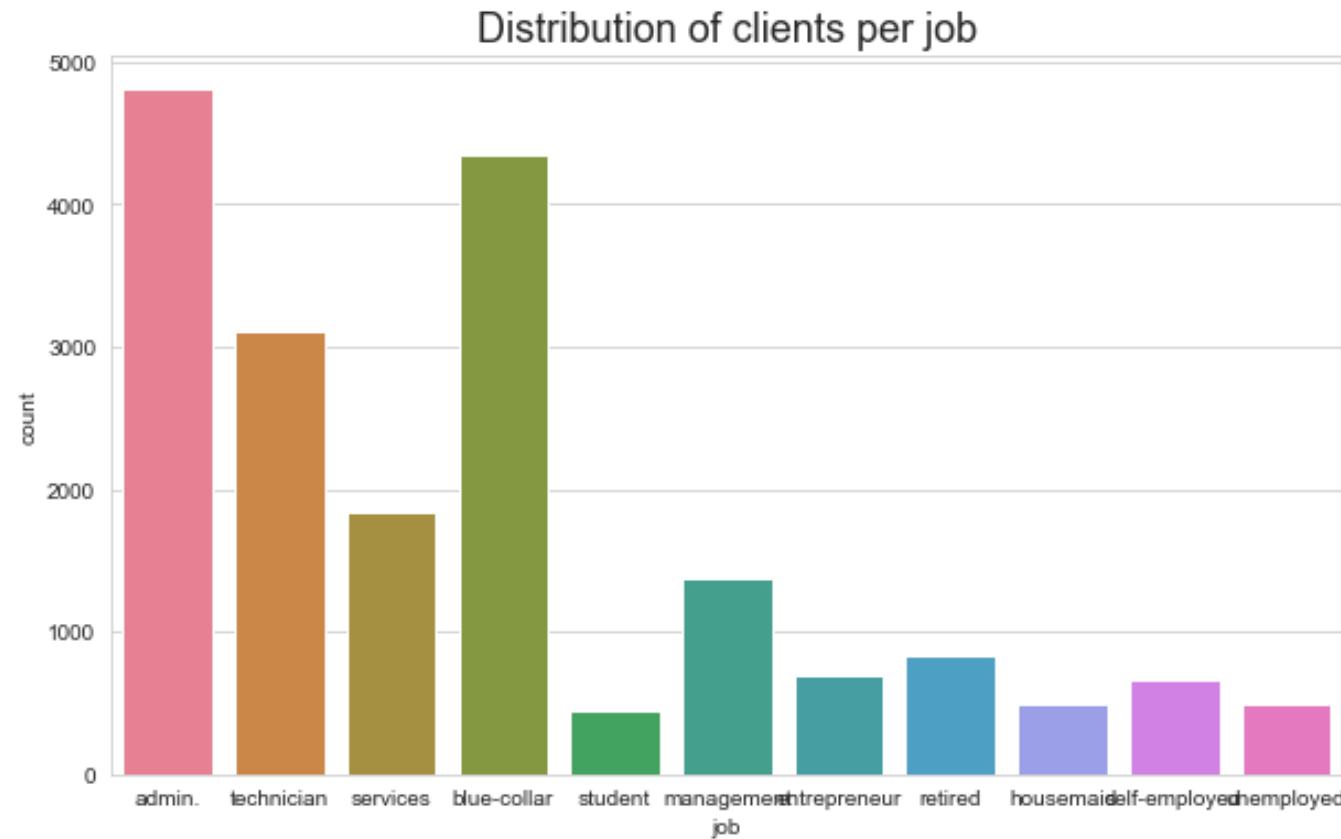
	y	no	yes
job			
admin.	4252	558	
blue-collar	4209	138	
entrepreneur	649	38	
housemaid	453	41	
management	1255	124	
retired	615	219	
self-employed	606	57	
services	1747	93	
student	305	135	
technician	2853	264	
unemployed	423	69	
unknown	132	16	

```
In [948]: # Get names of indexes for which column job has value unknown  
indexNames = training[ training['job'] == "unknown"].index
```

```
In [949]: # Delete these row indexes from DataFrame  
training.drop(indexNames , inplace=True)
```

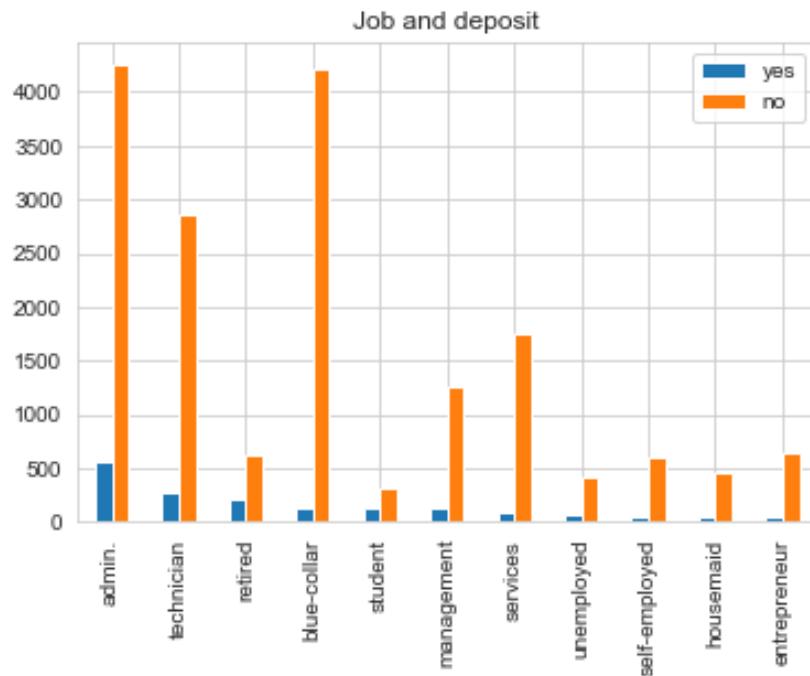
```
In [950]: #Let's see the distribution of clients per job.
```

```
sns.set_style('whitegrid')
plt.figure(figsize=(10, 6))
plt.title("Distribution of clients per job", fontsize=18)
sns.countplot(x="job", data=training, palette='husl');
```



```
In [951]: #Visualization of relationship between feature category vs dependent variable y  
#job and deposit  
j_bank = pd.DataFrame()  
  
j_bank['yes'] = training[training['y'] == 'yes']['job'].value_counts()  
j_bank['no'] = training[training['y'] == 'no']['job'].value_counts()  
  
j_bank.plot.bar(title = 'Job and deposit')
```

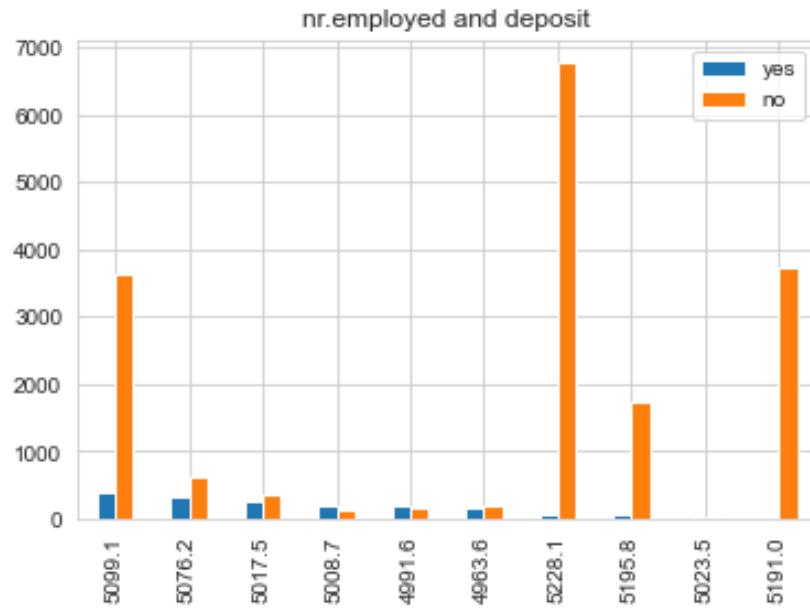
```
Out[951]: <AxesSubplot:title={'center':'Job and deposit'}>
```



Customers who worked in administrative position followed by technicians and blue collar made deposits

```
In [952]: #Visualization of relationship between feature category vs dependent variable y  
#nr.employed and deposit  
j_bank = pd.DataFrame()  
  
j_bank['yes'] = training[training['y'] == 'yes']['nr.employed'].value_counts()  
j_bank['no'] = training[training['y'] == 'no']['nr.employed'].value_counts()  
  
j_bank.plot.bar(title = 'nr.employed and deposit')
```

```
Out[952]: <AxesSubplot:title={'center':'nr.employed and deposit'}>
```



```
In [953]: #Crosstab to display job stats with respect to y class variable  
pd.crosstab(index=training["cons.price.idx"], columns=training["y"])
```

Out[953]:

	y	no	yes
cons.price.idx			
<b>92.201</b>	303	150	
<b>92.379</b>	96	58	
<b>92.431</b>	164	111	
<b>92.469</b>	55	32	
<b>92.649</b>	104	93	
<b>92.713</b>	37	45	
<b>92.756</b>	4	0	
<b>92.843</b>	82	78	
<b>92.893</b>	2509	117	
<b>92.963</b>	278	149	
<b>93.075</b>	1033	189	
<b>93.200</b>	1719	39	
<b>93.369</b>	60	91	
<b>93.444</b>	2294	29	
<b>93.749</b>	38	53	
<b>93.798</b>	6	24	
<b>93.876</b>	44	64	
<b>93.918</b>	2709	29	
<b>93.994</b>	3722	12	
<b>94.027</b>	51	64	

	y	no	yes
cons.price.idx			
94.055	53	54	
94.199	71	84	
94.215	51	79	
94.465	1772	11	
94.601	68	50	
94.767	44	31	

```
In [954]: #Crosstab to display marital stats with respect to y class variable  
pd.crosstab(index=training["marital"], columns=training["y"])
```

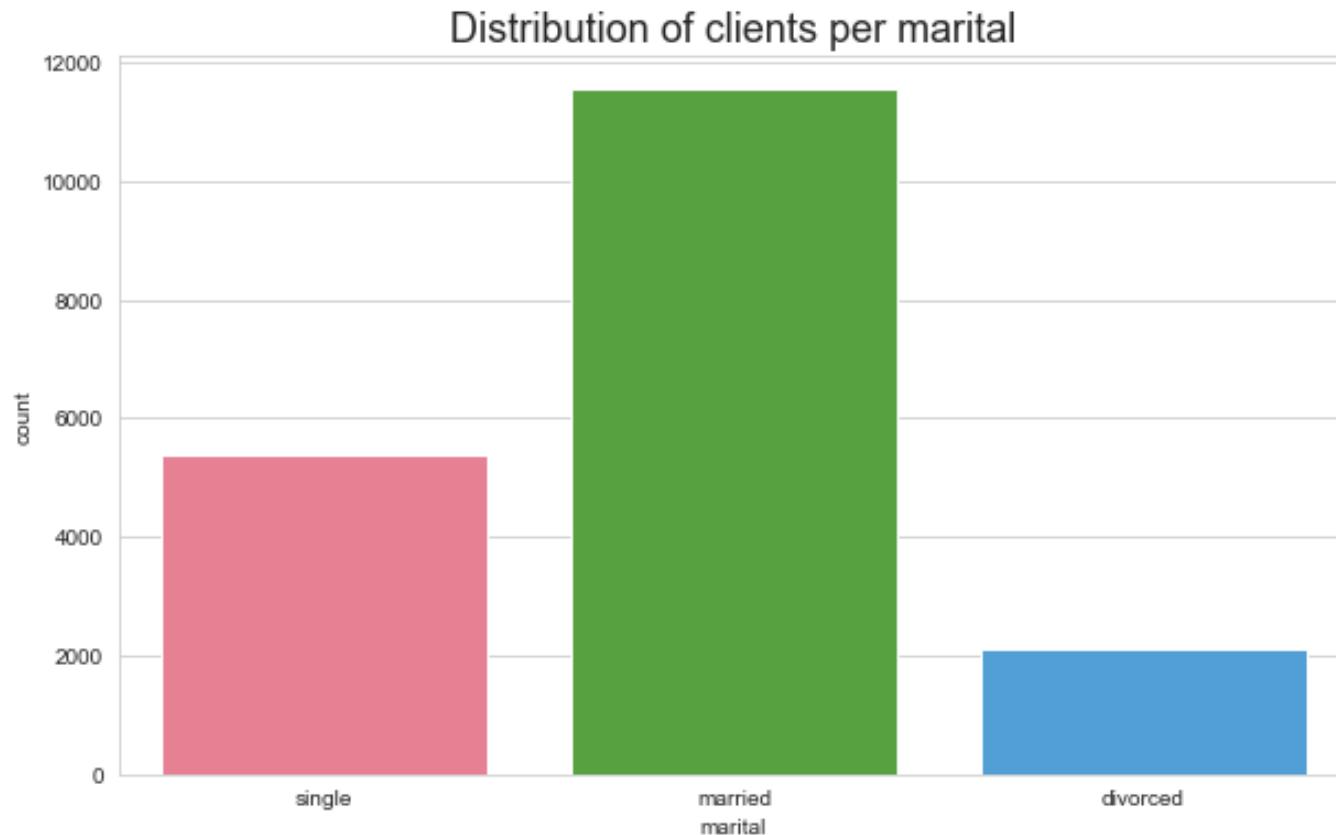
Out[954]:

	y	no	yes
marital			
divorced	1952	177	
married	10663	884	
single	4724	672	
unknown	28	3	

```
In [955]: # Get names of indexes for which column job has value unknown  
indexMarital = training[ training['marital'] == "unknown"].index
```

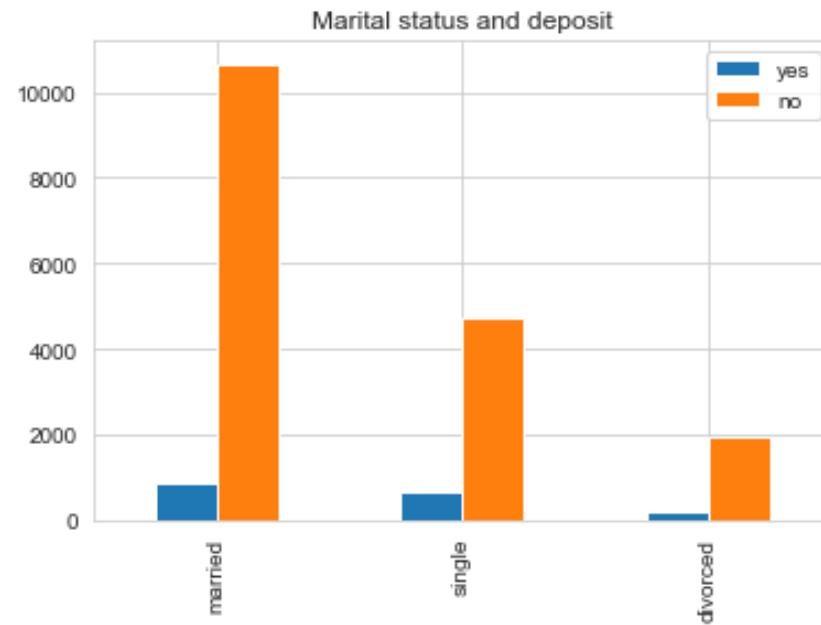
```
In [956]: # Delete these row indexes from DataFrame  
training.drop(indexMarital , inplace=True)
```

```
In [957]: #Let's see the distribution of clients per job.  
sns.set_style('whitegrid')  
plt.figure(figsize=(10, 6))  
plt.title("Distribution of clients per marital", fontsize=18)  
sns.countplot(x="marital", data=training, palette='husl');
```



```
In [958]: #Visualize the relationship between feature category vs dependent variable y  
#marital status and deposit  
j_bank = pd.DataFrame()  
  
j_bank['yes'] = training[training['y'] == 'yes']['marital'].value_counts()  
j_bank['no'] = training[training['y'] == 'no']['marital'].value_counts()  
  
j_bank.plot.bar(title = 'Marital status and deposit')
```

```
Out[958]: <AxesSubplot:title={'center':'Marital status and deposit'}>
```



Married customers followed by single had made deposits

```
In [959]: #Crosstab to display education stats with respect to y class variable  
pd.crosstab(index=training["education"], columns=training["y"])
```

Out[959]:

	y	no	yes
education			
<b>basic.4y</b>	1790	174	
<b>basic.6y</b>	1033	47	
<b>basic.9y</b>	2694	124	
<b>high.school</b>	4027	363	
<b>illiterate</b>	5	1	
<b>professional.course</b>	2184	237	
<b>university.degree</b>	4951	690	
<b>unknown</b>	655	97	

```
In [960]: # Get names of indexes for which column education has value unknown and illiterate  
indexEducation = training[ training['education'] == "illiterate"].index  
indexEducation2 = training[ training['education'] == "unknown"].index
```

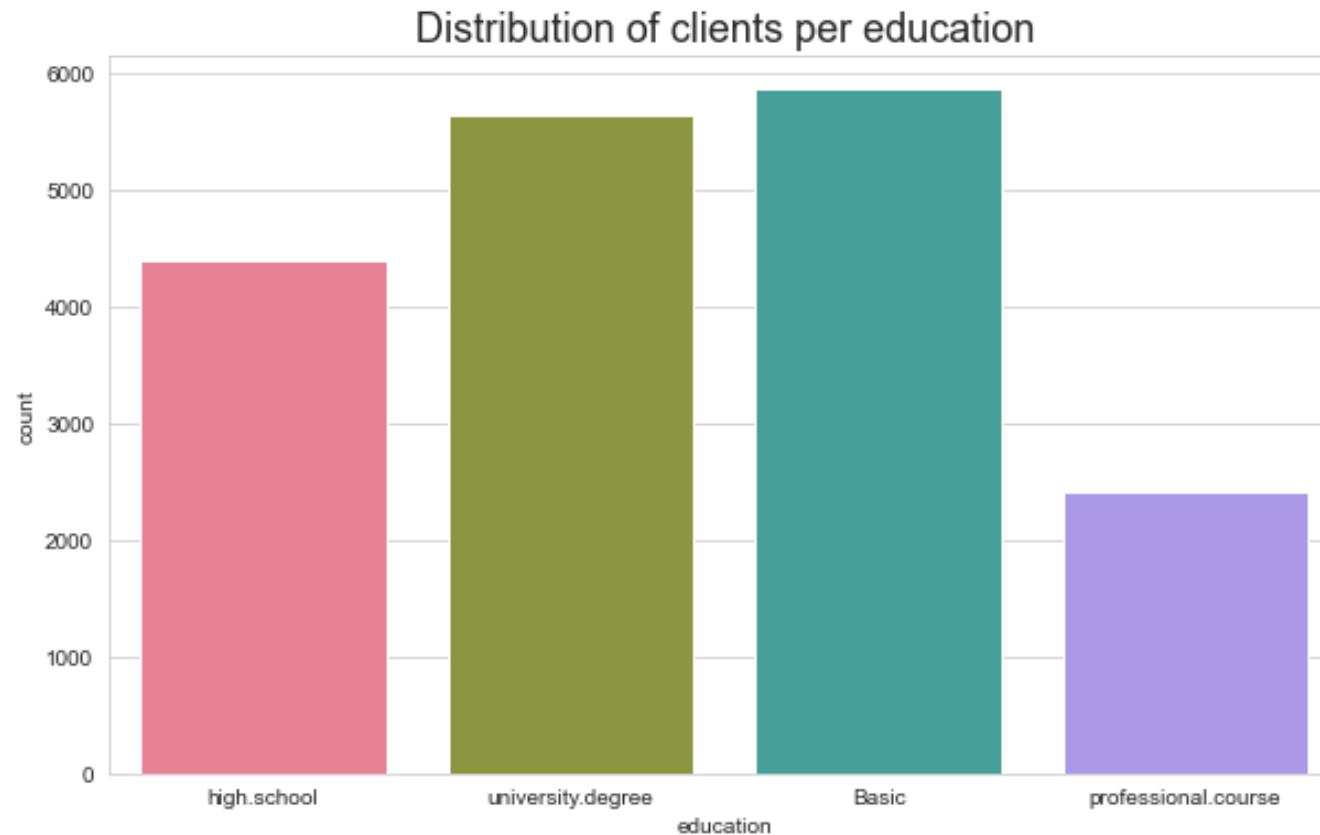
```
In [961]: # Delete these row indexes from DataFrame  
training.drop(indexEducation , inplace=True)  
training.drop(indexEducation2 , inplace=True)
```

```
In [962]: #Lets group "basic.4y", "basic.9y" and "basic.6y" together and call them "basic"  
training['education']=np.where(training['education'] =='basic.9y', 'Basic', training['education'])  
training['education']=np.where(training['education'] =='basic.6y', 'Basic', training['education'])  
training['education']=np.where(training['education'] =='basic.4y', 'Basic', training['education'])
```

```
In [963]: training['education'].unique()
```

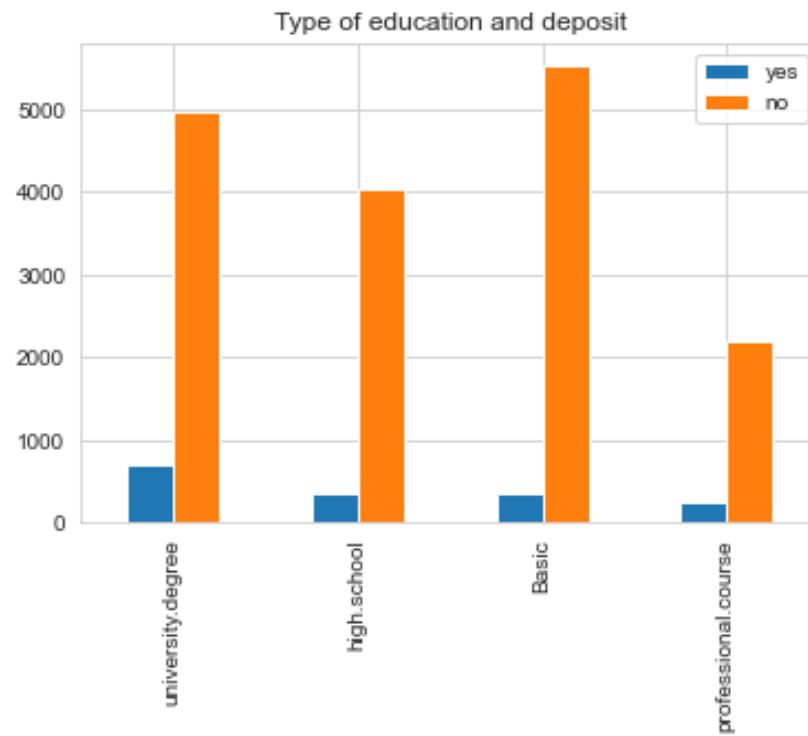
```
Out[963]: array(['high.school', 'university.degree', 'Basic', 'professional.course'],
                 dtype=object)
```

```
In [964]: sns.set_style('whitegrid')
plt.figure(figsize=(10, 6))
plt.title("Distribution of clients per education", fontsize=18)
sns.countplot(x="education", data=training, palette='husl');
```



```
In [965]: #Visualize the relationship between feature category vs dependent variable y  
#type of education and deposit  
j_bank = pd.DataFrame()  
  
j_bank['yes'] = training[training['y'] == 'yes']['education'].value_counts()  
j_bank['no'] = training[training['y'] == 'no']['education'].value_counts()  
  
j_bank.plot.bar(title = 'Type of education and deposit')
```

Out[965]: <AxesSubplot:title={'center':'Type of education and deposit'}>



Customers who had university degree followed by highschool had higher chance of making a deposit.

```
In [966]: #Crosstab to display housing stats with respect to y class variable  
pd.crosstab(index=training["housing"], columns=training["y"])
```

Out[966]:

	y	no	yes
housing			
no	7475	673	
unknown	409	40	
yes	8795	922	

```
In [967]: # Get names of indexes for which column housing has value unknown  
indexhousing = training[ training['housing'] == "unknown"].index
```

```
In [968]: # Delete these row indexes from DataFrame  
training.drop(indexhousing , inplace=True)
```

```
In [969]: sns.set_style('whitegrid')
plt.figure(figsize=(10, 6))
plt.title("Distribution of clients per housing", fontsize=18)
sns.countplot(x="housing", data=training, palette='husl');
```

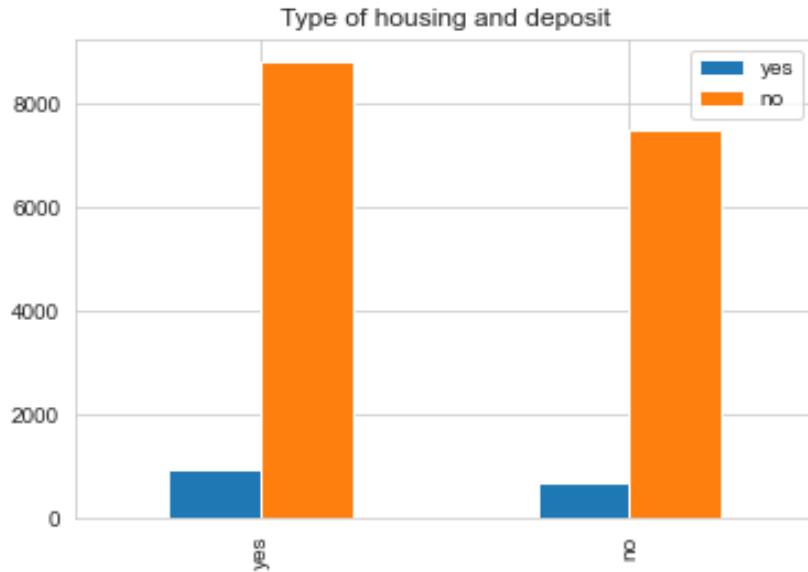


```
In [970]: #Visualize the relationship between feature category vs dependent variable y
#type of housing and deposit
j_bank = pd.DataFrame()

j_bank['yes'] = training[training['y'] == 'yes']['housing'].value_counts()
j_bank['no'] = training[training['y'] == 'no']['housing'].value_counts()

j_bank.plot.bar(title = 'Type of housing and deposit')
```

```
Out[970]: <AxesSubplot:title={'center':'Type of housing and deposit'}>
```



Customers who had house had higher chance in making a deposit.

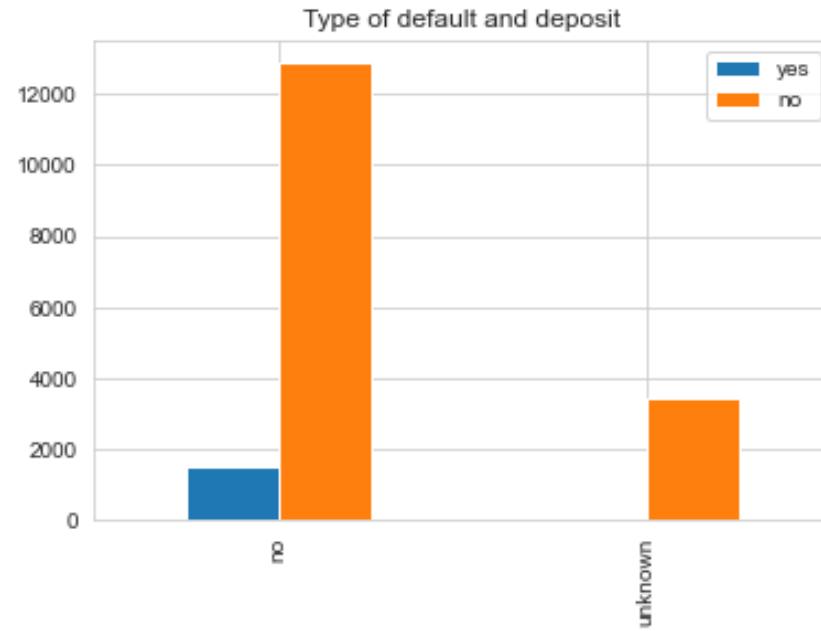
```
In [971]: #Crosstab to display default stats with respect to y class variable  
pd.crosstab(index=training["default"], columns=training["y"])
```

Out[971]:

	y	no	yes
default			
no	12854	1539	
unknown	3413	56	
yes	3	0	

```
In [972]: #Visualize the relationship between feature category vs dependent variable y  
#type of default and deposit  
j_bank = pd.DataFrame()  
  
j_bank['yes'] = training[training['y'] == 'yes']['default'].value_counts()  
j_bank['no'] = training[training['y'] == 'no']['default'].value_counts()  
  
j_bank.plot.bar(title = 'Type of default and deposit')
```

```
Out[972]: <AxesSubplot:title={'center':'Type of default and deposit'}>
```



Customers who had no default had higher chance in making a deposit.

```
In [973]: #Crosstab to display with y class variable  
pd.crosstab(index=training["loan"], columns=training["y"])
```

Out[973]:

	y	no	yes
loan			
no	13724	1375	
yes	2546	220	

```
In [974]: # Get names of indexes for which column Loan has value unknown  
indexloan = training[ training['loan'] == "unknown"].index
```

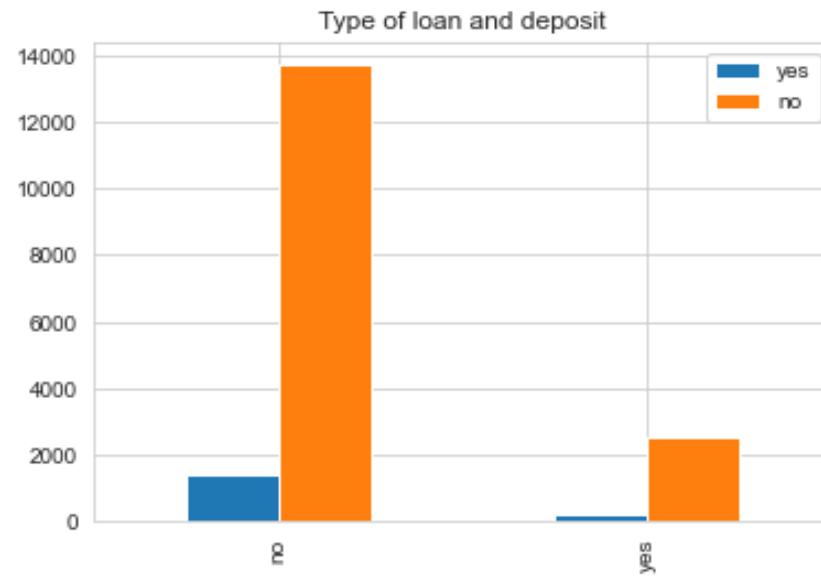
```
In [975]: # Delete these row indexes from DataFrame  
training.drop(indexloan , inplace=True)
```

```
In [976]: #Visualize the relationship between feature category vs dependent variable y
#type of loan and deposit
j_bank = pd.DataFrame()

j_bank['yes'] = training[training['y'] == 'yes']['loan'].value_counts()
j_bank['no'] = training[training['y'] == 'no']['loan'].value_counts()

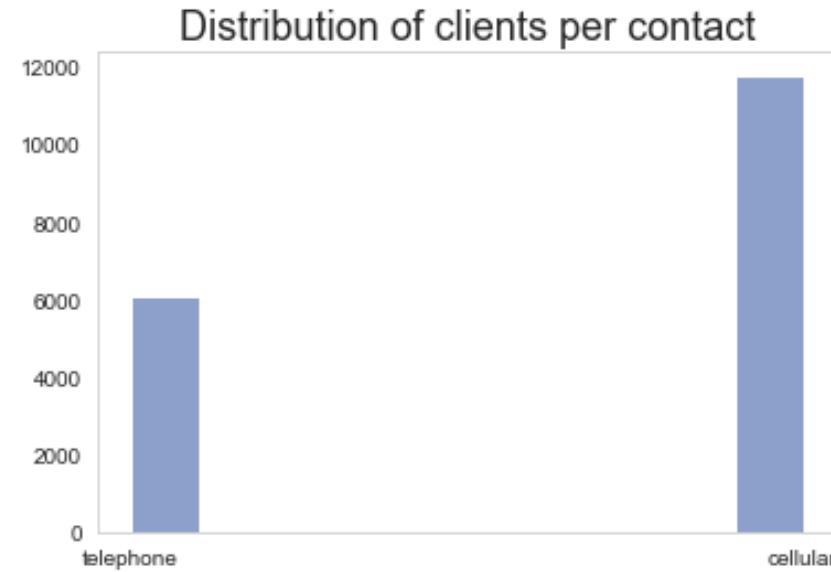
j_bank.plot.bar(title = 'Type of loan and deposit')
```

```
Out[976]: <AxesSubplot:title={'center':'Type of loan and deposit'}>
```



Clients that had no loan had a higher chance to subscribe to term deposits

```
In [977]: base_color = sns.color_palette('Set2')[2]
plt.hist(data = training, x = 'contact', color = base_color);
plt.title("Distribution of clients per contact", fontsize=18)
plt.grid();
```

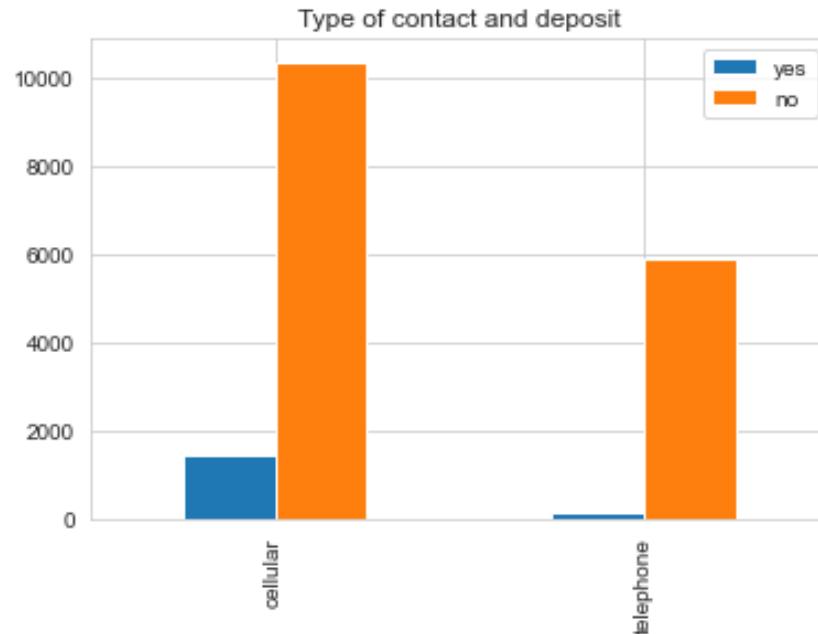


```
In [978]: #Visualize the relationship between feature category vs dependent variable y
#type of contact and deposit
j_bank = pd.DataFrame()

j_bank['yes'] = training[training['y'] == 'yes']['contact'].value_counts()
j_bank['no'] = training[training['y'] == 'no']['contact'].value_counts()

j_bank.plot.bar(title = 'Type of contact and deposit')
```

```
Out[978]: <AxesSubplot:title={'center':'Type of contact and deposit'}>
```



Clients that was contacted by cellular had a higher chance to subscribe for term deposits

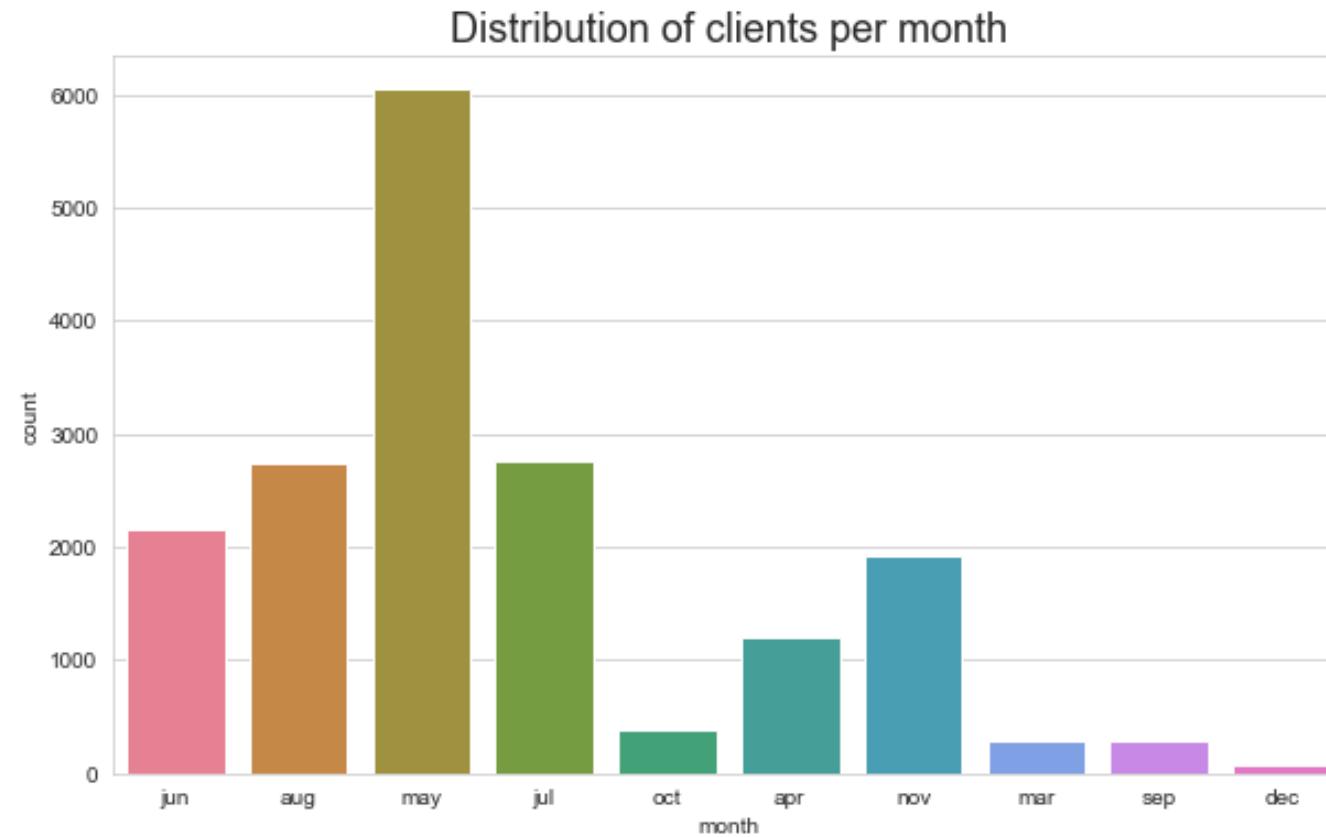
choice that was contacted by cellular has a higher chance to subscribe to term deposit.

In [979]: *#Crosstab to display contact stats with respect to y class variable*  
pd.crosstab(index=training["contact"], columns=training["y"])

Out[979]:

	y	no	yes
contact			
cellular	10360	1426	
telephone	5910	169	

```
In [980]: sns.set_style('whitegrid')
plt.figure(figsize=(10, 6))
plt.title("Distribution of clients per month", fontsize=18)
sns.countplot(x="month", data=training, palette='husl');
```



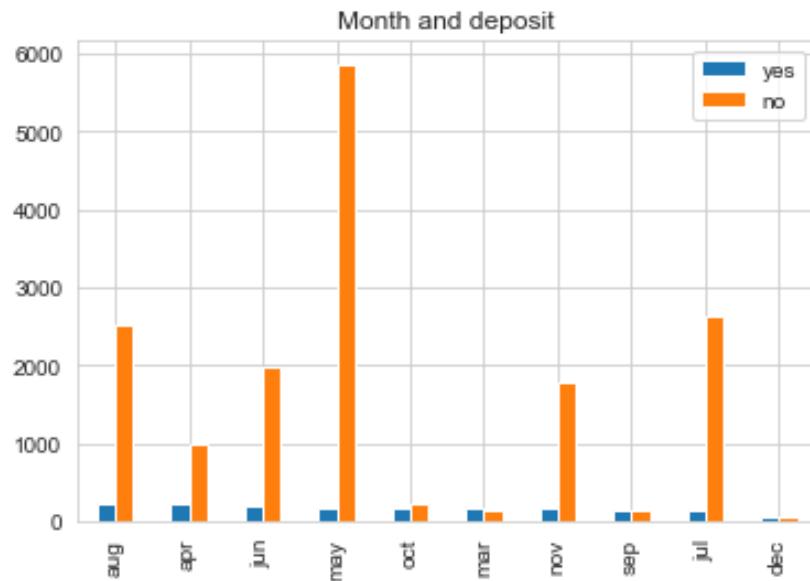
```
In [981]: #Crosstab to display contact stats with respect to y class variable  
pd.crosstab(index=training["month"], columns=training["y"])
```

Out[981]:

	y	no	yes
month			
apr	987	217	
aug	2522	224	
dec	38	37	
jul	2624	128	
jun	1965	199	
mar	128	155	
may	5868	178	
nov	1774	153	
oct	214	174	
sep	150	130	

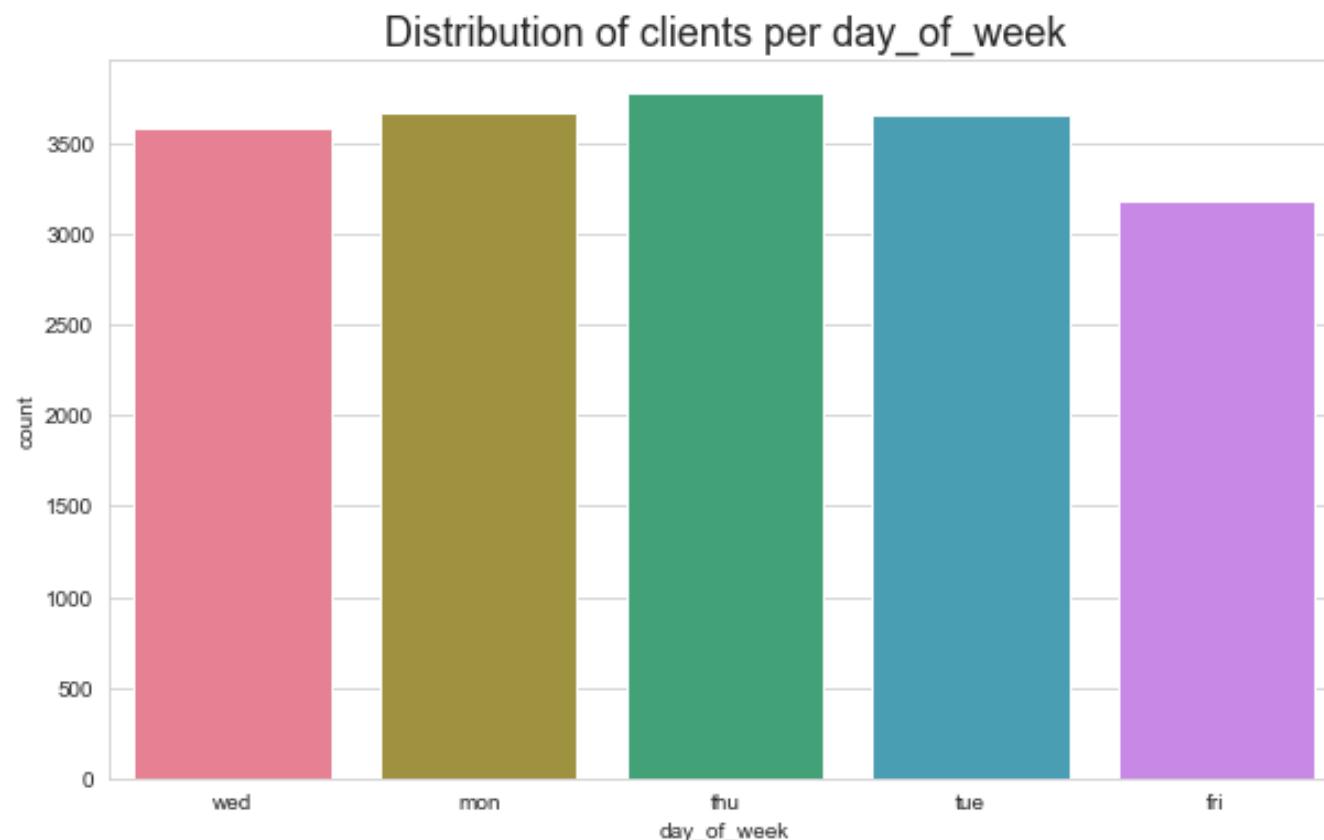
```
In [982]: #Visualize the relationship between feature category vs dependent variable y  
#Month and deposit  
j_bank = pd.DataFrame()  
  
j_bank['yes'] = training[training['y'] == 'yes']['month'].value_counts()  
j_bank['no'] = training[training['y'] == 'no']['month'].value_counts()  
  
j_bank.plot.bar(title = 'Month and deposit')
```

```
Out[982]: <AxesSubplot:title={'center':'Month and deposit'}>
```



Most of the deposits made during August followed by April and June.

```
In [983]: sns.set_style('whitegrid')
plt.figure(figsize=(10, 6))
plt.title("Distribution of clients per day_of_week", fontsize=18)
sns.countplot(x="day_of_week", data=training, palette='husl');
```



```
In [984]: #Crosstab to display day_of_week stats with respect to y class variable  
pd.crosstab(index=training["day_of_week"], columns=training["y"])
```

Out[984]:

	y	no	yes
day_of_week			
fri	2911	270	
mon	3380	286	
thu	3393	382	
tue	3314	345	
wed	3272	312	

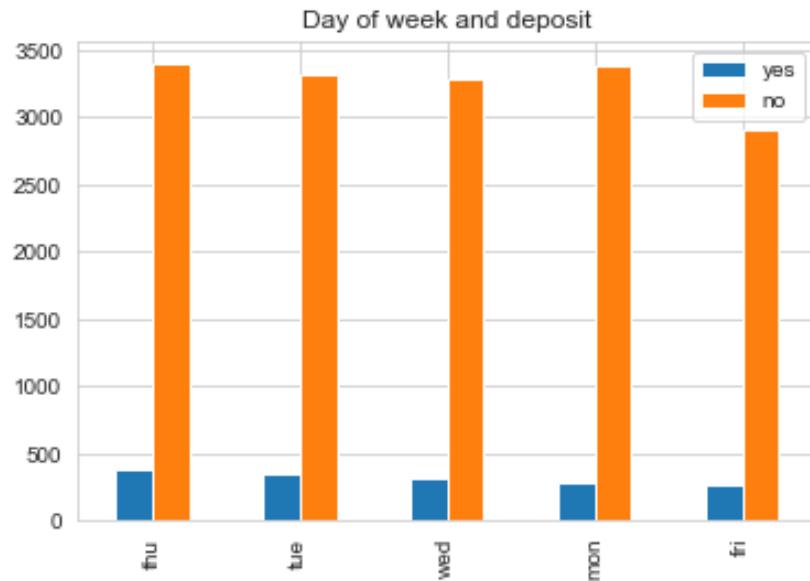
Less deposit made on Friday

```
In [985]: #Visualize the relationship between feature category vs dependent variable y
j_bank = pd.DataFrame()

j_bank['yes'] = training[training['y'] == 'yes']['day_of_week'].value_counts()
j_bank['no'] = training[training['y'] == 'no']['day_of_week'].value_counts()

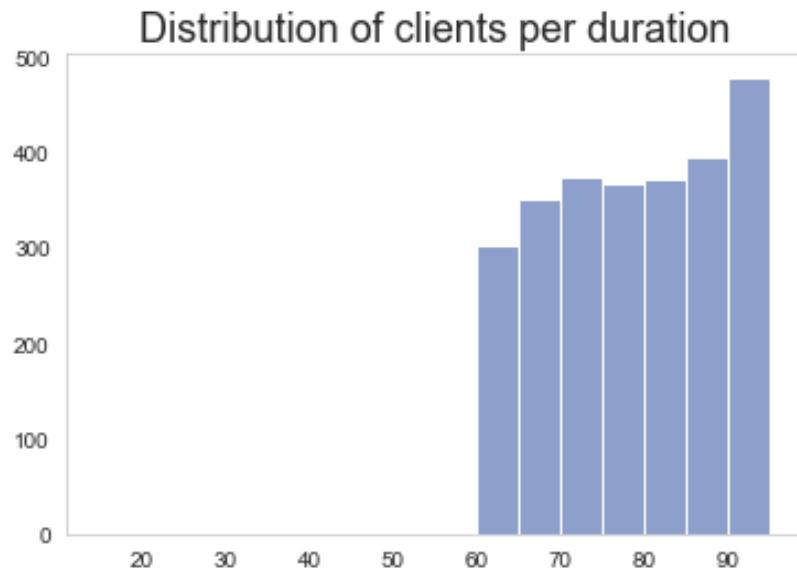
j_bank.plot.bar(title = 'Day of week and deposit')
```

```
Out[985]: <AxesSubplot:title={'center':'Day of week and deposit'}>
```



```
In [986]: base_color = sns.color_palette('Set2')[2]

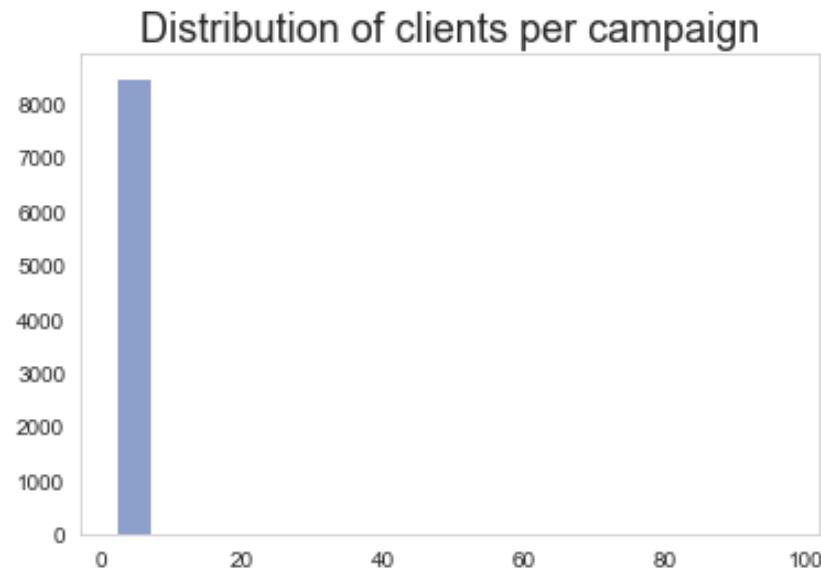
duration_bins = np.arange(15, 100, 5)
plt.hist(data = training, x = 'duration', bins = duration_bins, color = base_color);
plt.title("Distribution of clients per duration", fontsize=18)
plt.grid();
```



Duration(is the last contact to the client in seconds)

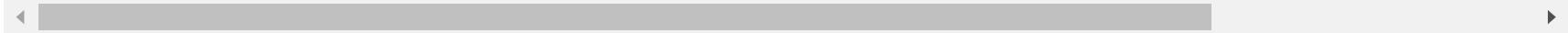
```
In [987]: base_color = sns.color_palette('Set2')[2]

campaign_bins = np.arange(2, 100, 5)
plt.hist(data = training, x = 'campaign', bins = campaign_bins, color = base_color);
plt.title("Distribution of clients per campaign", fontsize=18)
plt.grid();
```



In [988]:

```
fig = px.scatter(training, y="campaign", x="duration", color="y")
fig.show()
```



```
In [989]: #Crosstab to display default stats with respect to y class variable  
pd.crosstab(index=training[\"campaign\"], columns=training[\"y\"])
```

Out[989]:

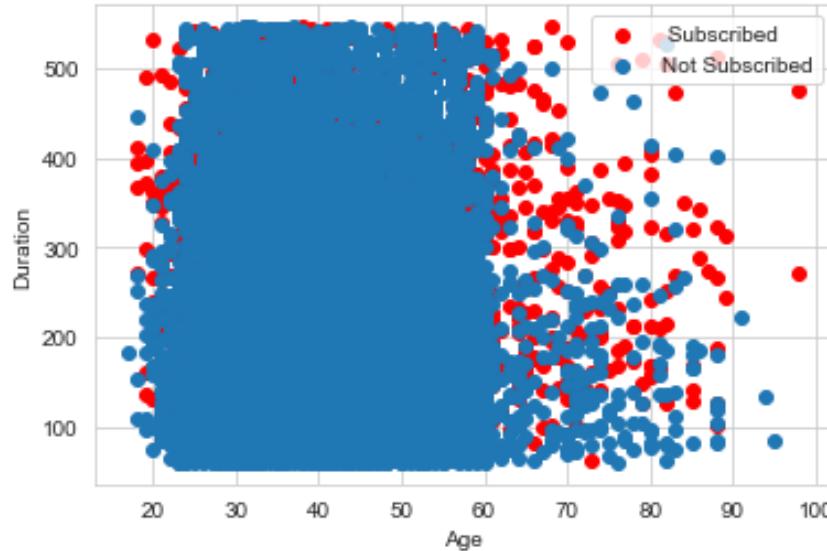
	y	no	yes
campaign			
<b>1</b>	8422	953	
<b>2</b>	5274	446	
<b>3</b>	2574	196	

campaign: is a number of contacts to client

Duration: is last contact duration, in seconds

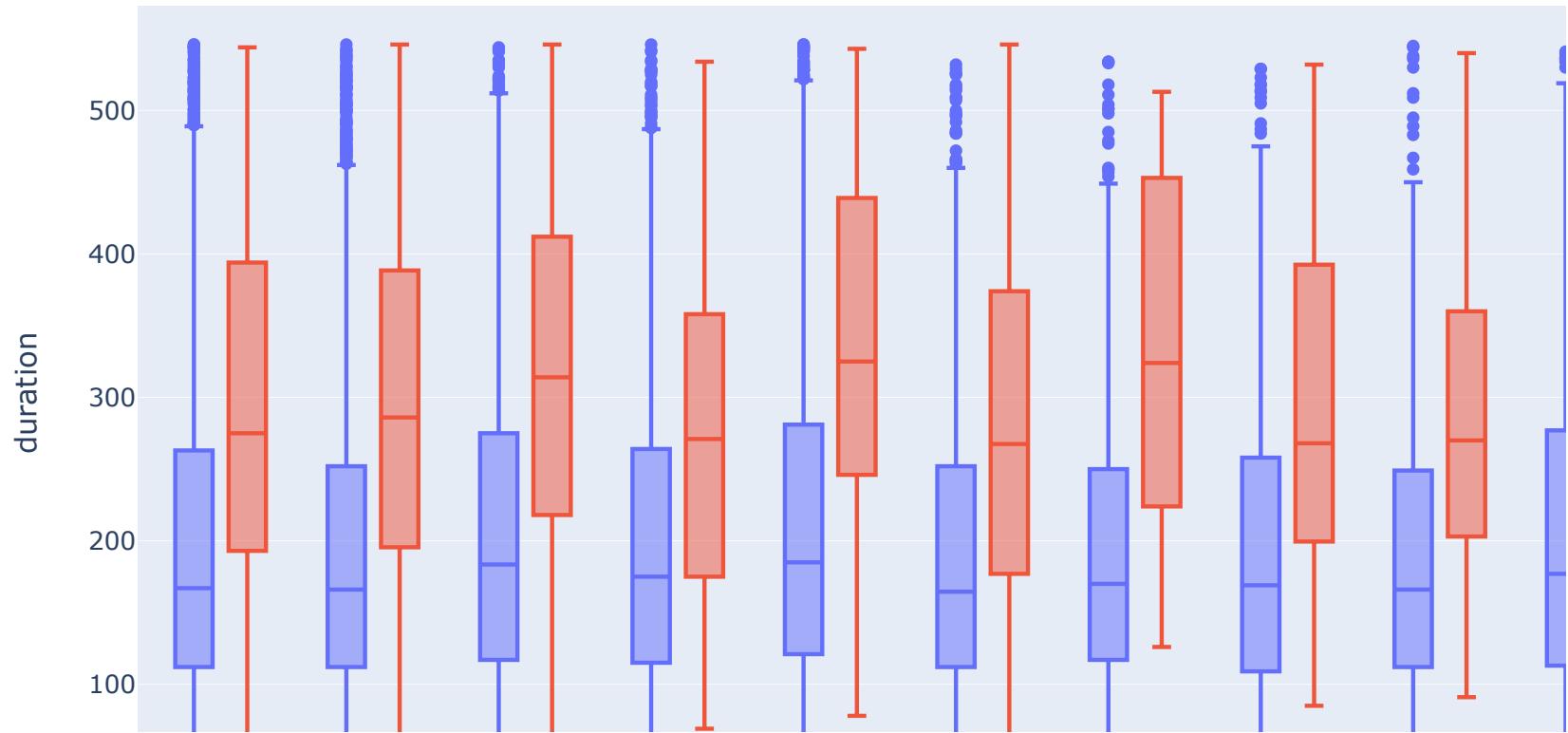
As more as employees contacted customers the less likely they made deposit

```
In [990]: plt.scatter(x=training.age[training.y=='yes'], y=training.duration[(training.y=='yes')], c="red")
plt.scatter(x=training.age[training.y=='no'], y=training.duration[(training.y=='no'))])
plt.legend([" Subscribed", "Not Subscribed"])
plt.xlabel("Age")
plt.ylabel("Duration")
plt.show()
```



In [991]:

```
fig = px.box(training, x="job", y="duration", color="y")
fig.update_traces(quartilemethod="exclusive")
fig.show()
```



The longer conversation with clients , the more likly they made deposit.

Comparing the median, the blue collar, entrepreneur and services had high duration of calls

## Categorical Treatment

The dataset contains object type variables using sklearn's preprocessing tool I will encode all variables to numerical labels.

```
In [992]: #build a new dataframe containing only the object columns.  
obj_bank = training.select_dtypes(include=['object']).copy()  
obj_bank.head(5)
```

Out[992]:

	job	marital	education	default	housing	loan	contact	month	day_of_week	poutcome	y
8726	admin.	single	high.school	no	no	no	telephone	jun	wed	nonexistent	no
37287	admin.	married	high.school	no	yes	no	cellular	aug	mon	success	yes
20981	technician	single	university.degree	no	yes	no	cellular	aug	thu	nonexistent	no
31623	services	divorced	Basic	no	yes	no	cellular	may	thu	failure	no
36959	admin.	single	university.degree	no	yes	no	cellular	jul	thu	nonexistent	yes

```
In [993]: training["month"].value_counts()
```

```
Out[993]: may    6046  
       jul    2752  
       aug    2746  
       jun    2164  
       nov    1927  
       apr    1204  
       oct     388  
       mar     283  
       sep     280  
       dec      75  
Name: month, dtype: int64
```

```
In [994]: training.head()
```

Out[994]:

	age	job	marital	education	default	housing	loan	contact	month	day_of_week	duration	campaign
8726	39	admin.	single	high.school	no	no	no	telephone	jun	wed	172.0	
37287	33	admin.	married	high.school	no	yes	no	cellular	aug	mon	252.0	
20981	32	technician	single	university.degree	no	yes	no	cellular	aug	thu	118.0	
31623	50	services	divorced	Basic	no	yes	no	cellular	may	thu	429.0	
36959	46	admin.	single	university.degree	no	yes	no	cellular	jul	thu	309.0	

```
In [995]: training["day_of_week"].value_counts()
```

Out[995]:

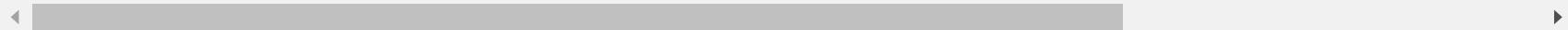
```
thu    3775  
mon   3666  
tue   3659  
wed   3584  
fri   3181  
Name: day_of_week, dtype: int64
```

```
In [996]: #I will be converting the month and day by it's corresponding number for training set
month_dict={'may':5,'jul':7,'aug':8,'jun':6,'nov':11,'apr':4,'oct':10,'sep':9,'mar':3,'dec':12}
training['month']= training['month'].map(month_dict)

day_dict={'thu':5,'mon':2,'wed':4,'tue':3,'fri':6}
training['day_of_week']= training['day_of_week'].map(day_dict)
training.head()
```

Out[996]:

	age	job	marital	education	default	housing	loan	contact	month	day_of_week	duration	campaign
8726	39	admin.	single	high.school	no	no	no	telephone	6	4	172.0	
37287	33	admin.	married	high.school	no	yes	no	cellular	8	2	252.0	
20981	32	technician	single	university.degree	no	yes	no	cellular	8	5	118.0	
31623	50	services	divorced	Basic	no	yes	no	cellular	5	5	429.0	
36959	46	admin.	single	university.degree	no	yes	no	cellular	7	5	309.0	



In [ ]:

```
In [997]: #I will be converting the month and day by it's corresponding number for testing set
month_dict={'may':5,'jul':7,'aug':8,'jun':6,'nov':11,'apr':4,'oct':10,'sep':9,'mar':3,'dec':12}
testing['month']= testing['month'].map(month_dict)

day_dict={'thu':5,'mon':2,'wed':4,'tue':3,'fri':6}
testing['day_of_week']= testing['day_of_week'].map(day_dict)
testing.head()
```

Out[997]:

	age	job	marital	education	default	housing	loan	contact	month	day_of_week	duration	campaign
37499	32	services	married	high.school	no	no	no	cellular	8	5	100.0	
23884	43	technician	married	high.school	no	yes	no	cellular	8	6	60.0	
32970	27	services	married	basic.9y	no	yes	no	cellular	5	2	220.0	
30374	28	admin.	married	university.degree	no	no	no	cellular	4	5	115.0	
12442	42	housemaid	divorced	basic.4y	unknown	no	no	cellular	7	2	461.0	



```
In [998]: #I will be converting the month and day by it's corresponding number for training set
month_dict={'may':5,'jul':7,'aug':8,'jun':6,'nov':11,'apr':4,'oct':10,'sep':9,'mar':3,'dec':12}
cross_validation_data['month']= cross_validation_data['month'].map(month_dict)

day_dict={'thu':5,'mon':2,'wed':4,'tue':3,'fri':6}
cross_validation_data['day_of_week']= cross_validation_data['day_of_week'].map(day_dict)
cross_validation_data.head()
```

Out[998]:

	age	job	marital	education	default	housing	loan	contact	month	day_of_week	duration	campaign	pout
0	56	housemaid	married	basic.4y	no	no	no	telephone	5	2	261.0	1	none
1	57	services	married	high.school	unknown	no	no	telephone	5	2	149.0	1	none
2	37	services	married	high.school	no	yes	no	telephone	5	2	226.0	1	none
3	40	admin.	married	basic.6y	no	no	no	telephone	5	2	151.0	1	none
4	56	services	married	high.school	no	no	yes	telephone	5	2	307.0	1	none

```
In [999]: #The dataset contains nine object type variables. I will use  
#a custom function by sklearn's preprocessing tool  
#to convert all nine variables to numerical labels for training set.  
LabEn=LabelEncoder()  
  
categorical_var=['job','marital', 'education','contact', 'poutcome', 'housing','default','loan','y']  
for i in categorical_var:  
    training[i]=LabEn.fit_transform(training[i])  
  
training.head()
```

Out[999]:

	age	job	marital	education	default	housing	loan	contact	month	day_of_week	duration	campaign	poutcome	
8726	39	0	2	1	0	0	0	1	6		4	172.0	2	1
37287	33	0	1	1	0	1	0	0	8		2	252.0	1	2
20981	32	9	2	3	0	1	0	0	8		5	118.0	2	1
31623	50	7	0	0	0	1	0	0	5		5	429.0	1	0
36959	46	0	2	3	0	1	0	0	7		5	309.0	1	1



```
In [1000]: #The dataset contains nine object type variables. I will use a custom function by sklearn's preproce  
#to convert all nine variables to numerical labels for testing set.  
LabEn=LabelEncoder()  
  
categorical_var=['job','marital', 'education','contact', 'poutcome', 'housing','default','loan','y']  
for i in categorical_var:  
    testing[i]=LabEn.fit_transform(testing[i])  
  
testing.head()
```

Out[1000]:

	age	job	marital	education	default	housing	loan	contact	month	day_of_week	duration	campaign	poutcome
37499	32	7	1	3	0	0	0	0	8	5	100.0	1	1
23884	43	9	1	3	0	2	0	0	8	6	60.0	1	1
32970	27	7	1	2	0	2	0	0	5	2	220.0	2	1
30374	28	0	1	6	0	0	0	0	4	5	115.0	2	0
12442	42	3	0	0	1	0	0	0	7	2	461.0	2	1



```
In [1001]: LabEn=LabelEncoder()

categorical_var=['job','marital', 'education','contact', 'poutcome', 'housing','default','loan','y']
for i in categorical_var:
    cross_validation_data[i]=LabEn.fit_transform(cross_validation_data[i])

cross_validation_data.head()
```

Out[1001]:

	age	job	marital	education	default	housing	loan	contact	month	day_of_week	duration	campaign	poutcome	cons
0	56	3	1	0	0	0	0	1	5		2	261.0	1	1
1	57	7	1	3	1	0	0	1	5		2	149.0	1	1
2	37	7	1	3	0	2	0	1	5		2	226.0	1	1
3	40	0	1	1	0	0	0	1	5		2	151.0	1	1
4	56	7	1	3	0	0	2	1	5		2	307.0	1	1



```
In [1002]: #Checking if I didn't get any NaN valueuse when new lebelS was created in training set  
training.isna().sum()
```

```
Out[1002]: age          0  
job           0  
marital       0  
education     0  
default        0  
housing        0  
loan           0  
contact        0  
month          0  
day_of_week    0  
duration        0  
campaign        0  
poutcome        0  
cons.price.idx  0  
nr.employed     0  
y               0  
dtype: int64
```

```
In [1003]: #Checking if I didn't get any NaN valueuse when new lebelS was created in testing set  
testing.isna().sum()
```

```
Out[1003]: age          0  
job           0  
marital       0  
education     0  
default        0  
housing        0  
loan           0  
contact        0  
month          0  
day_of_week    0  
duration        0  
campaign        0  
poutcome        0  
cons.price.idx  0  
nr.employed     0  
y               0  
dtype: int64
```

**Divide the dataset to training (X\_train, y\_train) and test (X\_test, y\_test )sets.**

My data set already divided into 2 portions in the ratio of 70:30, my target variable is 'y'

```
In [1004]: X_train= training.drop("y",axis=1)  
y_train= training["y"]
```

```
In [1005]: X_test= testing.drop("y",axis=1)  
y_test= testing["y"]
```

```
In [1006]: print(X_train.shape)
print(y_train.shape)
print(X_test.shape)
print(y_test.shape)
```

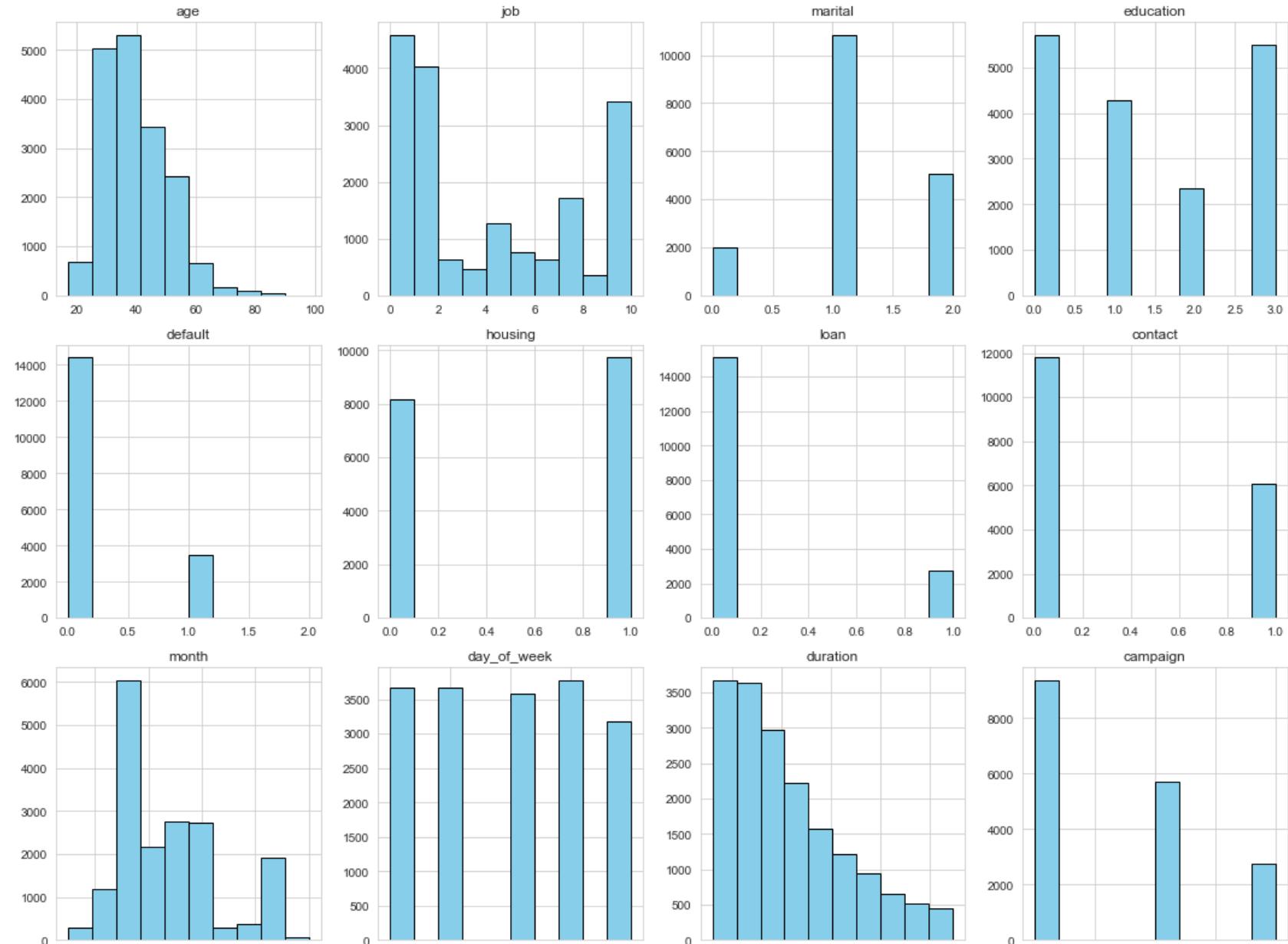
```
(17865, 15)
(17865,)
(12357, 15)
(12357,)
```

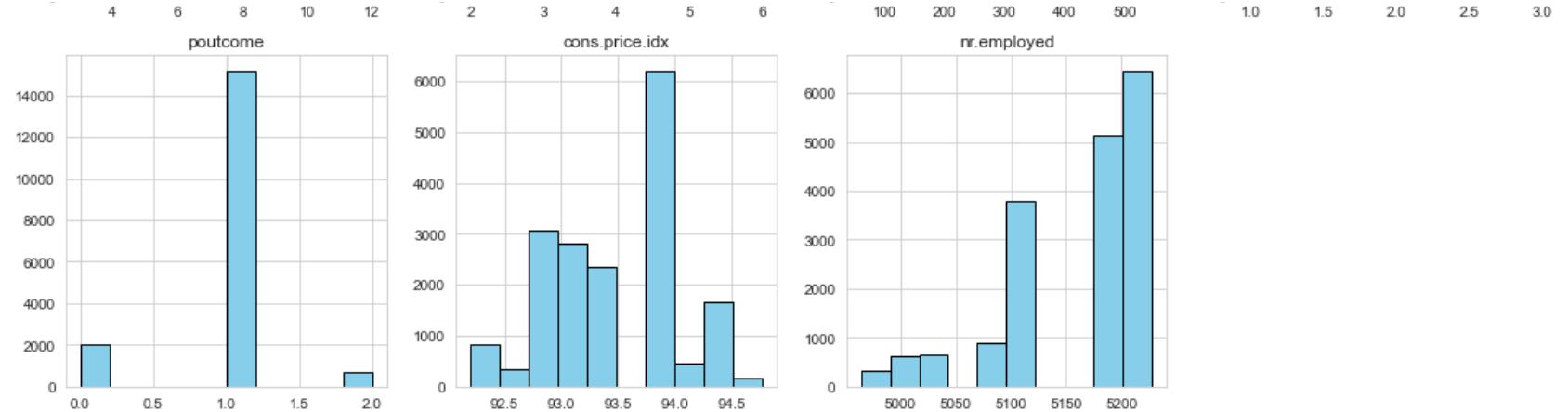
```
In [1007]: y_test.head()
```

```
Out[1007]: 37499    0
23884    0
32970    0
30374    0
12442    0
Name: y, dtype: int32
```

**Scaling** I have tried to rescale with StandardScaler(centering the variable at zero and standardizing the variance at 1) was no effect on the algorithms and I have tried PowerTransformer(method='yeo-johnson'), had no prediction of class 1 of recall and precision. Normalization (Min-Max Scalar) technic will not work as this data doesn't need to suppress outliers, I already deleted them and more than that will cause algorithms to perform worse.

```
In [1008]: # histograms of the variables  
#Histogram for the numerical attributes  
X_train.hist(figsize=(15,15),edgecolor='k',color='skyblue')  
plt.tight_layout()  
plt.show()
```





## Oversampling using SMOTE

I will over sample only on the training data, so no information bleed to test data or validation data .

```
In [1009]: counter = Counter(y_train)
```

```
In [1010]: print("Before SMOTE", counter)
```

```
Before SMOTE Counter({0: 16270, 1: 1595})
```

```
In [1011]: smt = SMOTE()
```

```
In [1012]: X_train,y_train = smt.fit_resample(X_train,y_train)
```

```
In [1013]: print("After SMOTE", Counter(y_train))
```

```
After SMOTE Counter({0: 16270, 1: 16270})
```

```
In [1014]: y_train.isna().sum()
```

```
Out[1014]: 0
```

```
In [1015]: y_test.isna().sum()
```

```
Out[1015]: 0
```

```
In [1016]: y_test.isin([0]).any().any()
```

```
Out[1016]: True
```

```
In [1017]: y_train.isin([0]).any().any()
```

```
Out[1017]: True
```

```
In [1018]: print(Counter(y_test))
```

```
Counter({0: 10949, 1: 1408})
```

```
In [1019]: y_train.shape
```

```
Out[1019]: (32540,)
```

```
In [1020]: X_train.shape
```

```
Out[1020]: (32540, 15)
```

```
In [1021]: X_test.shape
```

```
Out[1021]: (12357, 15)
```

```
In [719]: y_test.shape
```

```
Out[719]: (12357,)
```

```
In [1022]: #spliting data to X and y for cross validation  
X = cross_validation_data.drop("y",axis=1)  
y = cross_validation_data["y"]
```

# Classification with all Features

## Test 1a. RandomForest with all 15 Features

```
In [1024]: #see the classification performance of the Random Forest using all 15 features
# To improve the results of RF I tested n_estimators for 40,50,100,200,10000
#with max_depth of 2, 3 and 4

FullRandFor = RandomForestClassifier(n_estimators=50, random_state=43, max_depth=3)
FullRandFor.fit(X_train, y_train)

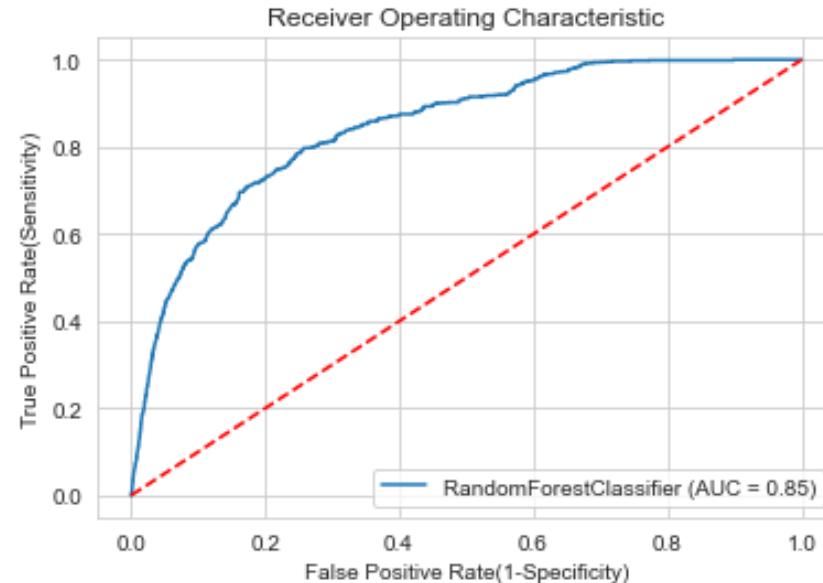
fulltrainpred = FullRandFor.predict_proba(X_train)
print('Accuracy on training set: {}'.format(roc_auc_score(y_train, fulltrainpred[:,1])))

fulltestpred = FullRandFor.predict_proba(X_test)
print('Accuracy on test set: {}'.format(roc_auc_score(y_test, fulltestpred[:,1])))

Accuracy on training set: 0.9440173410513806
Accuracy on test set: 0.8454609932206345
```

```
In [1025]: # draw the ROC-AUC chart  
metrics.plot_roc_curve(FullRandFor, X_test, y_test)  
plt.title('Receiver Operating Characteristic')  
plt.plot([0, 1], [0, 1], 'r--')  
plt.ylabel('True Positive Rate(Sensitivity)')  
plt.xlabel('False Positive Rate(1-Specificity)')
```

```
Out[1025]: Text(0.5, 0, 'False Positive Rate(1-Specificity)')
```



```
In [1026]: pred6 = FullRandFor.predict(X_test)
```

```
In [1027]: print("Random Forest with all 15 Features")
cm = confusion_matrix(y_test, pred6)
print(cm)
print('\n')
print(classification_report(y_test,pred6))
tn, fp, fn, tp=cm.ravel()
print ("TP: ", tp,", FP: ", fp,", TN: ", tn,", FN:", fn)
```

Random Forest with all 15 Features

```
[[9547 1402]
 [ 540  868]]
```

	precision	recall	f1-score	support
0	0.95	0.87	0.91	10949
1	0.38	0.62	0.47	1408
accuracy			0.84	12357
macro avg	0.66	0.74	0.69	12357
weighted avg	0.88	0.84	0.86	12357

TP: 868 , FP: 1402 , TN: 9547 , FN: 540

```
In [1028]: fulltestpred = FullRandFor.predict_proba(X_test)
prob3 = fulltestpred[:, 1]# Keeping only the values in positive label
```

```
In [1029]: #The average precision (PR AUC) is returned by passing the
#true Label & the probability estimate.
# Average precision score
PR_AUC = average_precision_score(y_test, prob3)
print(PR_AUC)
```

0.4329736397720649

```
In [1030]: #Brier skill score calculates the mean squared error between predicted  
#probabilities and the expected values(actuals).  
#compute the Brier Score-perfect skill has a score of 0.0 and the worst  
#has a score of 1.0. From this score, we can infer that our model  
#has good performance or skill.  
loss = brier_score_loss(y_test, prob3)  
loss
```

```
Out[1030]: 0.1283332997195806
```

```
In [1033]: # define the evaluation procedure  
cv = RepeatedStratifiedKFold(n_splits=10, n_repeats=5, random_state=1)  
# evaluate the model using cross-validation  
scores1 = cross_val_score(FullRandFor, X, y, scoring='accuracy', cv=cv, n_jobs=-1)
```

```
In [1034]: scores1
```

```
Out[1034]: array([0.89463462, 0.89779073, 0.89560573, 0.89681962, 0.89584851,  
0.89584851, 0.89536295, 0.89900461, 0.89460903, 0.89849441,  
0.8980335 , 0.89779073, 0.89099296, 0.89584851, 0.8980335 ,  
0.89342073, 0.89584851, 0.89827628, 0.89800874, 0.89460903,  
0.89414907, 0.89414907, 0.89876184, 0.9002185 , 0.89463462,  
0.89657684, 0.89584851, 0.89584851, 0.89509471, 0.89873725,  
0.89876184, 0.89706239, 0.89414907, 0.8948774 , 0.89681962,  
0.89706239, 0.89633406, 0.89609128, 0.89315202, 0.89679456,  
0.89851906, 0.89366351, 0.89293518, 0.89584851, 0.89512017,  
0.89706239, 0.89512017, 0.89754795, 0.89946576, 0.89558038])
```

```
In [1035]: print(scores1.mean())
```

```
0.8962173563772571
```

## Test 1b. Decision Tree with all 15 Features

```
In [1036]: %%time
dtree15 = DecisionTreeClassifier()
dtree15.fit(X_train, y_train)

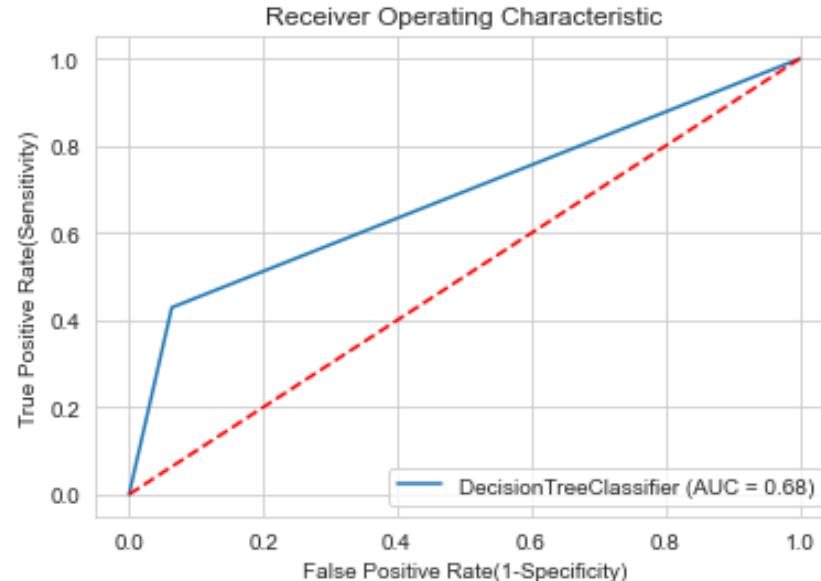
train_pred15 = dtree15.predict_proba(X_train)
print('Accuracy on training set: {}'.format(roc_auc_score(y_train, train_pred15[:,1])))

test_pred15 = dtree15.predict_proba(X_test)
print('Accuracy on test set: {}'.format(roc_auc_score(y_test, test_pred15[:,1])))
```

```
Accuracy on training set: 1.0
Accuracy on test set: 0.6826135792807978
Wall time: 180 ms
```

```
In [1037]: # draw the ROC-AUC chart  
metrics.plot_roc_curve(dtree15,X_test,y_test)  
plt.title('Receiver Operating Characteristic')  
plt.plot([0, 1], [0, 1], 'r--')  
plt.ylabel('True Positive Rate(Sensitivity)')  
plt.xlabel('False Positive Rate(1-Specificity)')
```

```
Out[1037]: Text(0.5, 0, 'False Positive Rate(1-Specificity)')
```



```
In [1038]: pred15 = dtree15.predict(X_test)
```

```
In [1039]: print("DecisionTree with all 15 Features")
cm = confusion_matrix(y_test, pred15)
print(cm)
print('\n')
print(classification_report(y_test,pred15))
tn, fp, fn, tp=cm.ravel()
print ("TP: ", tp," , FP: ", fp," , TN: ", tn," , FN:", fn)
```

DecisionTree with all 15 Features

```
[[10251  698]
 [ 804  604]]
```

	precision	recall	f1-score	support
0	0.93	0.94	0.93	10949
1	0.46	0.43	0.45	1408
accuracy			0.88	12357
macro avg	0.70	0.68	0.69	12357
weighted avg	0.87	0.88	0.88	12357

TP: 604 , FP: 698 , TN: 10251 , FN: 804

```
In [1040]: test_pred15 = dtree15.predict_proba(X_test)
prob15 = test_pred15[:, 1]# Keeping only the values in positive label
```

```
In [1041]: #The average precision (PR AUC) is returned by passing t
#he true Label & the probability estimate.
# Average precision score
PR_AUC = average_precision_score(y_test, prob15)
print(PR_AUC)
```

0.26406761766845643

```
In [1042]: #Brier skill score calculates the mean squared error between  
#predicted probabilities and the expected values(actuals).  
#compute the Brier Score-perfect skill has a score of 0.0 and the worst has a score of 1.0  
#From this score, we can infer that our model has good performance or skill.  
loss15 = brier_score_loss(y_test, prob15)  
loss15
```

```
Out[1042]: 0.12155053815651048
```

```
In [1043]: # define the evaluation procedure  
cv2 = RepeatedStratifiedKFold(n_splits=10, n_repeats=5, random_state=1)  
# evaluate the model using cross-validation  
scores2 = cross_val_score(dtree15, X, y, scoring='accuracy', cv=cv2, n_jobs=-1)
```

```
In [1044]: scores2
```

```
Out[1044]: array([0.8851663 , 0.88540908, 0.88079631, 0.87666909, 0.89633406,  
0.88783685, 0.8776402 , 0.87958242, 0.88440991, 0.88465274,  
0.89172129, 0.88225297, 0.88128186, 0.88565186, 0.88710852,  
0.89147851, 0.88152464, 0.88079631, 0.87858184, 0.88732394,  
0.87788298, 0.88322408, 0.88662297, 0.88589463, 0.88152464,  
0.88103909, 0.88031076, 0.8873513 , 0.88683827, 0.87858184,  
0.88273853, 0.88176742, 0.88832241, 0.89123574, 0.89026463,  
0.88710852, 0.88225297, 0.88807963, 0.88513842, 0.88028169,  
0.89414907, 0.88395242, 0.88783685, 0.88103909, 0.87958242,  
0.88419519, 0.88783685, 0.8820102 , 0.88489558, 0.88271005])
```

```
In [1045]: print(scores2.mean())
```

```
0.8844177393264325
```

### Test 1c. Logistic Regression with all 15 Features

In [1046]: #Logistic Regression

```
LR15 = LogisticRegression (solver='liblinear')
LR15.fit(X_train, y_train)

LRtrain_pred15 = LR15.predict_proba(X_train)
print('Accuracy on training set: {}'.format(roc_auc_score(y_train, LRtrain_pred15[:,1])))

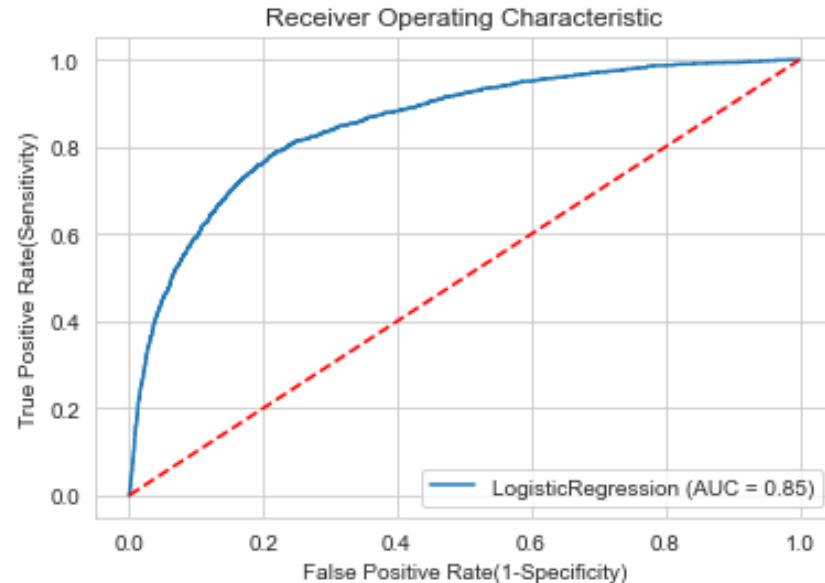
LRtest_pred15 = LR15.predict_proba(X_test)
print('Accuracy on test set: {}'.format(roc_auc_score(y_test, LRtest_pred15[:,1])))
```

Accuracy on training set: 0.9326002510644551

Accuracy on test set: 0.8533663177002465

```
In [1047]: # draw the ROC-AUC chart
metrics.plot_roc_curve(LR15, X_test, y_test)
plt.title('Receiver Operating Characteristic')
plt.plot([0, 1], [0, 1], 'r--')
plt.ylabel('True Positive Rate(Sensitivity)')
plt.xlabel('False Positive Rate(1-Specificity)')
```

```
Out[1047]: Text(0.5, 0, 'False Positive Rate(1-Specificity)')
```



```
In [1048]: pred16 = LR15.predict(X_test)
```

```
In [1049]: print("LogisticRegression with all 15 Features")
cm = confusion_matrix(y_test, pred16)
print(cm)
print('\n')
print(classification_report(y_test,pred16))
tn, fp, fn, tp=cm.ravel()
print ("TP: ", tp,", FP: ", fp,", TN: ", tn,", FN:", fn)
```

LogisticRegression with all 15 Features

```
[[9546 1403]
 [ 487  921]]
```

	precision	recall	f1-score	support
0	0.95	0.87	0.91	10949
1	0.40	0.65	0.49	1408
accuracy			0.85	12357
macro avg	0.67	0.76	0.70	12357
weighted avg	0.89	0.85	0.86	12357

TP: 921 , FP: 1403 , TN: 9546 , FN: 487

```
In [1050]: LRtest_pred15 = LR15.predict_proba(X_test)
prob16 = LRtest_pred15[:, 1]# Keeping only the values in positive Label
```

```
In [1051]: #The average precision (PR AUC) is returned by passing the true Label & the probability estimate.
# Average precision score
PR_AUC = average_precision_score(y_test, prob16)
print(PR_AUC)
```

0.46745802463781755

```
In [1052]: #Brier skill score calculates the mean squared error between predicted  
#probabilities and the expected values(actuals).  
#compute the Brier Score-perfect skill has a score of 0.0 and the  
#worst has a score of 1.0. From this score, we can infer that our model  
#has good performance or skill.  
loss16 = brier_score_loss(y_test, prob16)  
loss16
```

```
Out[1052]: 0.11360507592726816
```

```
In [1054]: # define the evaluation procedure  
cv3 = RepeatedStratifiedKFold(n_splits=10, n_repeats=5, random_state=1)  
# evaluate the model using cross-validation  
scores3 = cross_val_score(LR15, X, y, scoring='accuracy', cv=cv3, n_jobs=-1)
```

```
In [1055]: scores3
```

```
Out[1055]: array([0.91235737, 0.90725904, 0.90531682, 0.90847293, 0.90337461,  
       0.90483127, 0.89584851, 0.90483127, 0.9111219 , 0.90286547,  
       0.90968682, 0.90410294, 0.89851906, 0.90216072, 0.91138626,  
       0.89827628, 0.9024035 , 0.89730517, 0.90189412, 0.90820787,  
       0.90386016, 0.89924739, 0.90823015, 0.90725904, 0.90653071,  
       0.90143239, 0.90580238, 0.90386016, 0.90723652, 0.89946576,  
       0.90264627, 0.90653071, 0.90895849, 0.90386016, 0.90337461,  
       0.89851906, 0.9109007 , 0.90628793, 0.89728023, 0.90432249,  
       0.9024035 , 0.91017237, 0.9024035 , 0.90580238, 0.89851906,  
       0.90458849, 0.90871571, 0.89827628, 0.90553667, 0.90140845])
```

```
In [1056]: print(scores3.mean())
```

```
0.9042730727821567
```

## Classification with Filter Methods for Feature Selection-Mutual Information Gain

```
In [1057]: #I have tested deferent qty of Features and found that 10 has the best accuracy for all algorithms
```

```
MI=mutual_info_classif(X_train, y_train)
```

```
In [1058]: len(MI)
```

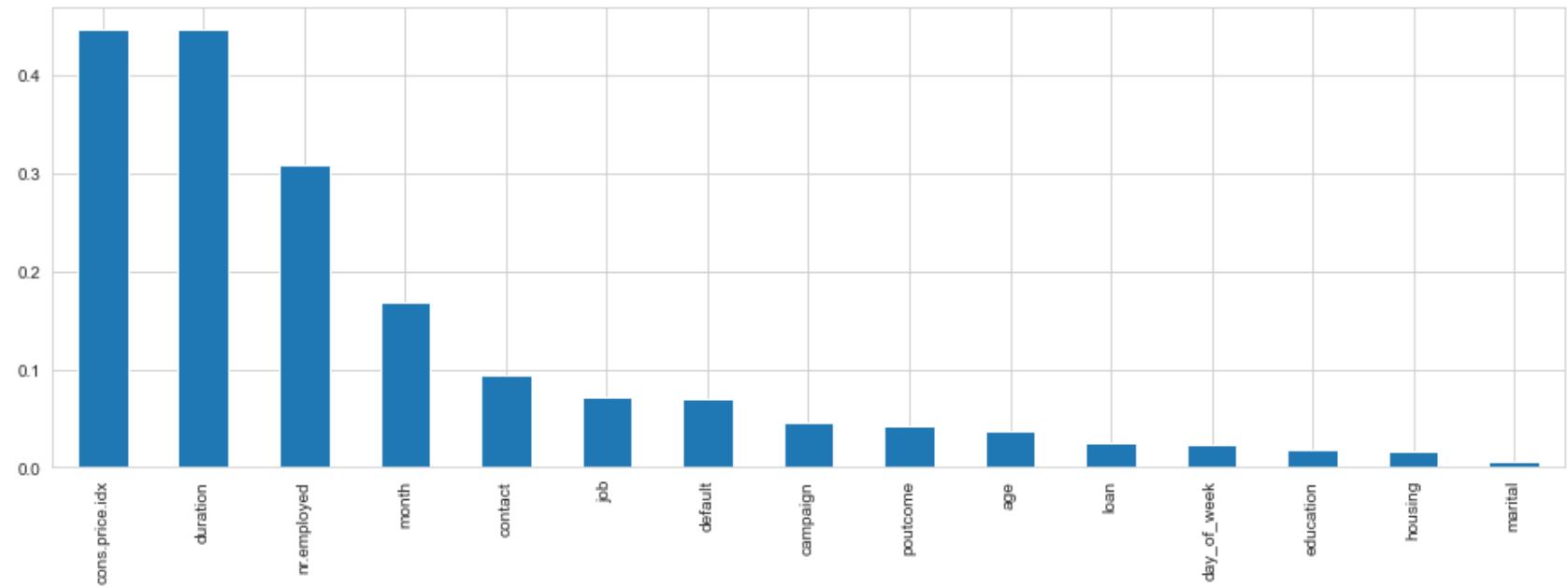
```
Out[1058]: 15
```

```
In [1059]: MI = pd.Series(MI)  
MI.index = X_train.columns
```

```
In [1060]: MI.sort_values(ascending=False, inplace = True)
```

```
In [1061]: MI.plot.bar(figsize = (16,5))
```

```
Out[1061]: <AxesSubplot:>
```



```
In [1062]: #percentile=65
sel = SelectPercentile(mutual_info_classif, percentile=65).fit(X_train, y_train)
X_train.columns[sel.get_support()]
```

```
Out[1062]: Index(['age', 'job', 'default', 'contact', 'month', 'duration', 'campaign',
       'poutcome', 'cons.price.idx', 'nr.employed'],
      dtype='object')
```

```
In [1063]: len(X_train.columns[sel.get_support()])
```

```
Out[1063]: 10
```

```
In [1064]: X_trainMI = sel.transform(X_train)
X_testMI = sel.transform(X_test)
```

```
In [1065]: X_trainMI.shape
```

```
Out[1065]: (32540, 10)
```

```
In [1066]: X_testMI.shape
```

```
Out[1066]: (12357, 10)
```

## Test 2a. Random Forests with Mutual Information Gain-Filter

```
In [1067]: ## To improve the results of RF I tested n_estimators for 40,50,100,200,10000
#with max_depth of 2, 3 and 4 and found that
#n_estimators=50 and max_depth=3 is the best combination.

RandFor1 = RandomForestClassifier(n_estimators=50, random_state=41, max_depth=3)
RandFor1.fit(X_trainMI, y_train)

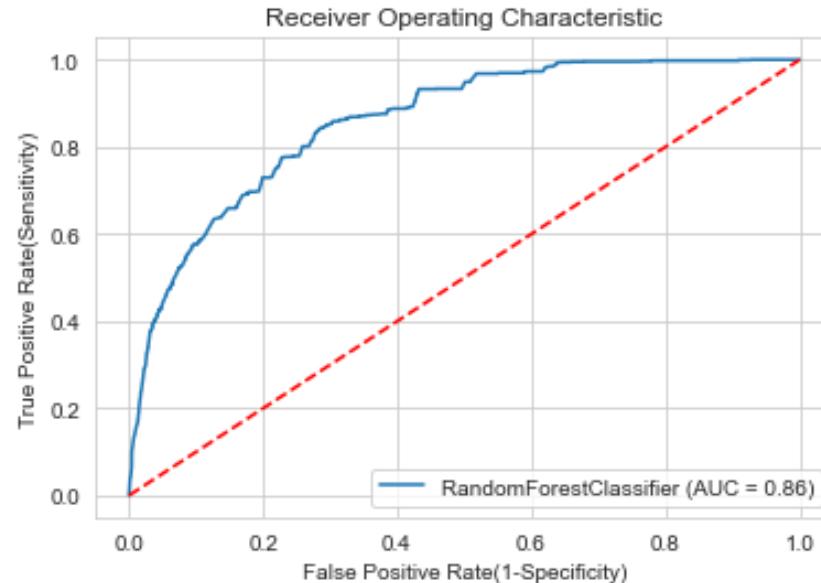
trainpred = RandFor1.predict_proba(X_trainMI)
print('Accuracy on training set: {}'.format(roc_auc_score(y_train, trainpred[:,1])))

testpred = RandFor1.predict_proba(X_testMI)
print('Accuracy on test set: {}'.format(roc_auc_score(y_test, testpred[:,1])))
```

```
Accuracy on training set: 0.9430749105918148
Accuracy on test set: 0.8585076651873563
```

```
In [1068]: # draw the ROC-AUC chart  
metrics.plot_roc_curve(RandFor1,X_testMI,y_test)  
plt.title('Receiver Operating Characteristic')  
plt.plot([0, 1], [0, 1], 'r--')  
plt.ylabel('True Positive Rate(Sensitivity)')  
plt.xlabel('False Positive Rate(1-Specificity)')
```

```
Out[1068]: Text(0.5, 0, 'False Positive Rate(1-Specificity)')
```



```
In [1069]: pred3 = RandFor1.predict(X_testMI)
```

```
In [1070]: print("RandomForest with Mutual Information Gain")
cm = confusion_matrix(y_test, pred3)
print(cm)
print('\n')
print(classification_report(y_test,pred3))
tn, fp, fn, tp=cm.ravel()
print ("TP: ", tp,", FP: ", fp,", TN: ", tn,", FN:", fn)
```

RandomForest with Mutual Information Gain

```
[[9436 1513]
 [ 507  901]]
```

	precision	recall	f1-score	support
0	0.95	0.86	0.90	10949
1	0.37	0.64	0.47	1408
accuracy			0.84	12357
macro avg	0.66	0.75	0.69	12357
weighted avg	0.88	0.84	0.85	12357

TP: 901 , FP: 1513 , TN: 9436 , FN: 507

```
In [1071]: #our aim is to find the brier score loss, so we will first
#calculate the probabilities for each data entry in
#X using the predict_proba() function.
```

```
In [1072]: testpred = RandFor1.predict_proba(X_testMI)
prob = testpred[:, 1] # Keeping only the values in positive label
```

```
In [1073]: #The average precision (PR AUC) is returned by passing the true Label & the probability estimate.  
# Average precision score  
PR_AUC = average_precision_score(y_test, prob)  
print(PR_AUC)
```

```
0.4745795331146324
```

```
In [1074]: #Brier skill score calculates the mean squared error between predicted  
#probabilities and the expected values(actuals).  
#compute the Brier Score-perfect skill has a score of 0.0 and the worst has a score of 1.0  
#From this score, we can infer that our model has good performance or skill.  
loss = brier_score_loss(y_test, prob)  
loss
```

```
Out[1074]: 0.11874581138372281
```

```
In [1075]: # define the evaluation procedure  
cv4 = RepeatedStratifiedKFold(n_splits=10, n_repeats=5, random_state=1)  
# evaluate the model using cross-validation  
scores4 = cross_val_score(RandFor1, X, y, scoring='accuracy', cv=cv4, n_jobs=-1)
```

```
In [1076]: scores4
```

```
Out[1076]: array([0.90337461, 0.89900461, 0.89584851, 0.89851906, 0.8980335 ,  
0.89949017, 0.89876184, 0.89827628, 0.89946576, 0.90189412,  
0.90191794, 0.9002185 , 0.89633406, 0.89681962, 0.90070405,  
0.89924739, 0.89876184, 0.90046128, 0.8997086 , 0.8997086 ,  
0.89609128, 0.89949017, 0.90167516, 0.90410294, 0.89924739,  
0.89876184, 0.89827628, 0.89827628, 0.89728023, 0.9023798 ,  
0.90191794, 0.8980335 , 0.90070405, 0.89827628, 0.90337461,  
0.90288905, 0.89827628, 0.89754795, 0.89558038, 0.89946576,  
0.9002185 , 0.90288905, 0.89876184, 0.89754795, 0.89657684,  
0.89924739, 0.8980335 , 0.89949017, 0.89946576, 0.89922292])
```

```
In [1077]: print(scores4.mean())
```

```
0.8993930282686483
```

## Test 2b. Decision Tree with Mutual Information Gain-Filter

```
In [1078]: %%time
```

```
dtreeMI = DecisionTreeClassifier()
dtreeMI.fit(X_trainMI, y_train)

train_predMI = dtreeMI.predict_proba(X_trainMI)
print('Accuracy on training set: {}'.format(roc_auc_score(y_train, train_predMI[:,1])))

test_predMI = dtreeMI.predict_proba(X_testMI)
print('Accuracy on test set: {}'.format(roc_auc_score(y_test, test_predMI[:,1])))
```

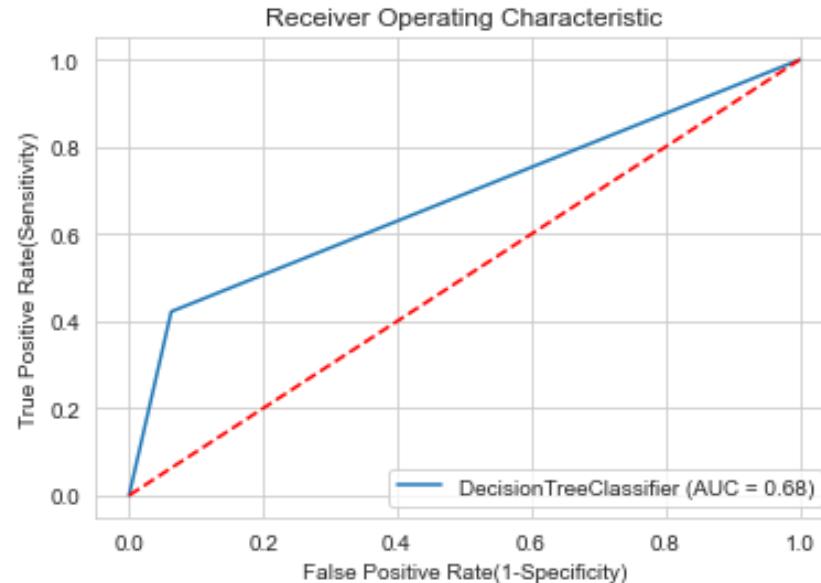
```
Accuracy on training set: 0.9999999924446448
```

```
Accuracy on test set: 0.6794734393551923
```

```
Wall time: 134 ms
```

```
In [1079]: # draw the ROC-AUC chart  
metrics.plot_roc_curve(dtreeMI,X_testMI,y_test)  
plt.title('Receiver Operating Characteristic')  
plt.plot([0, 1], [0, 1], 'r--')  
plt.ylabel('True Positive Rate(Sensitivity)')  
plt.xlabel('False Positive Rate(1-Specificity)')
```

```
Out[1079]: Text(0.5, 0, 'False Positive Rate(1-Specificity)')
```



```
In [1080]: pred4 = dtreeMI.predict(X_testMI)
```

```
In [1081]: print("DecisionTree with Mutual Information Gain")
cm = confusion_matrix(y_test, pred4)
print(cm)
print('\n')
print(classification_report(y_test,pred4))
tn, fp, fn, tp=cm.ravel()
print ("TP: ", tp,", FP: ", fp,", TN: ", tn,", FN:", fn)
```

DecisionTree with Mutual Information Gain

```
[[10260  689]
 [ 814  594]]
```

	precision	recall	f1-score	support
0	0.93	0.94	0.93	10949
1	0.46	0.42	0.44	1408
accuracy			0.88	12357
macro avg	0.69	0.68	0.69	12357
weighted avg	0.87	0.88	0.88	12357

TP: 594 , FP: 689 , TN: 10260 , FN: 814

```
In [1082]: test_predMI = dtreeMI.predict_proba(X_testMI)
prob1 = test_predMI[:, 1]# Keeping only the values in positive Label
```

```
In [1083]: #The average precision (PR AUC) is returned by passing the true
#label & the probability estimate.
# Average precision score
PR_AUC = average_precision_score(y_test, prob1)
print(PR_AUC)
```

0.26119218315834

```
In [1084]: #Brier skill score calculates the mean squared error between predicted  
#probabilities and the expected values(actuals).  
#compute the Brier Score-perfect skill has a score of 0.0 and the worst has a score of 1.0  
#From this score, we can infer that our model has good performance or skill.  
loss = brier_score_loss(y_test, prob1)  
loss
```

```
Out[1084]: 0.12163146394756008
```

```
In [1087]: # define the evaluation procedure  
cv5 = RepeatedStratifiedKFold(n_splits=10, n_repeats=5, random_state=1)  
# evaluate the model using cross-validation  
scores5 = cross_val_score(dtreeMI, X, y, scoring='accuracy', cv=cv5, n_jobs=-1)
```

```
In [1088]: scores5
```

```
Out[1088]: array([0.88686574, 0.88443797, 0.87715465, 0.87739743, 0.89463462,  
0.88759408, 0.87958242, 0.88249575, 0.88562409, 0.88489558,  
0.88565186, 0.88492353, 0.8820102 , 0.88152464, 0.88468075,  
0.88832241, 0.87861131, 0.88249575, 0.88076736, 0.8836814 ,  
0.87861131, 0.87933965, 0.88905074, 0.8895363 , 0.88419519,  
0.88298131, 0.88565186, 0.88807963, 0.8810102 , 0.88392424,  
0.88006798, 0.88103909, 0.89123574, 0.88905074, 0.88929352,  
0.88905074, 0.8873513 , 0.88710852, 0.88271005, 0.87761049,  
0.8895363 , 0.88638019, 0.88662297, 0.88395242, 0.88322408,  
0.88395242, 0.88977907, 0.88468075, 0.88586693, 0.88295289])
```

```
In [1089]: print(scores5.mean())
```

```
0.8845439623366101
```

## Test 2c. Logistic Regression with Mutual Information Gain-Filter

In [1090]: #Logistic Regression

```
LR = LogisticRegression (solver='liblinear')
LR.fit(X_trainMI, y_train)

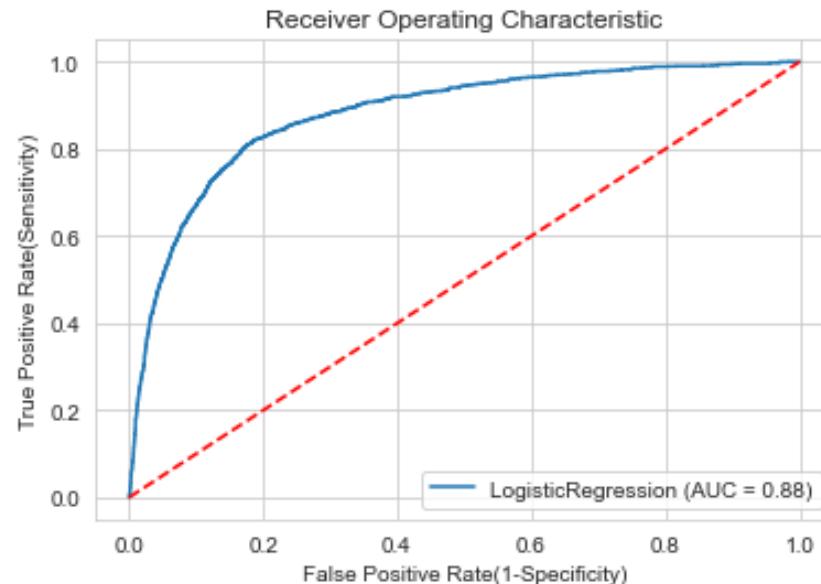
LRtrain_predMI = LR.predict_proba(X_trainMI)
print('Accuracy on training set: {}'.format(roc_auc_score(y_train, LRtrain_predMI[:,1])))

LRtest_predMI = LR.predict_proba(X_testMI)
print('Accuracy on test set: {}'.format(roc_auc_score(y_test, LRtest_predMI[:,1])))
```

Accuracy on training set: 0.922561254098308  
Accuracy on test set: 0.8806078050922044

```
In [1091]: # draw the ROC-AUC chart  
metrics.plot_roc_curve(LR, X_testMI, y_test)  
plt.title('Receiver Operating Characteristic')  
plt.plot([0, 1], [0, 1], 'r--')  
plt.ylabel('True Positive Rate(Sensitivity)')  
plt.xlabel('False Positive Rate(1-Specificity)')
```

```
Out[1091]: Text(0.5, 0, 'False Positive Rate(1-Specificity)')
```



```
In [1092]: pred5 = LR.predict(X_testMI)
```

```
In [1093]: print("LogisticRegression with Mutual Information Gain")
cm = confusion_matrix(y_test, pred5)
print(cm)
print('\n')
print(classification_report(y_test,pred5))
tn, fp, fn, tp=cm.ravel()
print ("TP: ", tp,", FP: ", fp,", TN: ", tn,", FN:", fn)
```

```
LogisticRegression with Mutual Information Gain
[[9140 1809]
 [ 298 1110]]
```

	precision	recall	f1-score	support
0	0.97	0.83	0.90	10949
1	0.38	0.79	0.51	1408
accuracy			0.83	12357
macro avg	0.67	0.81	0.70	12357
weighted avg	0.90	0.83	0.85	12357

```
TP: 1110 , FP: 1809 , TN: 9140 , FN: 298
```

```
In [1094]: LRtest_predMI = LR.predict_proba(X_testMI)
prob2 = LRtest_predMI[:, 1]# Keeping only the values in positive label
```

```
In [1095]: #The average precision (PR AUC) is returned by passing the true Label & the probability estimate.
# Average precision score
PR_AUC = average_precision_score(y_test, prob2)
print(PR_AUC)
```

```
0.5159068902810182
```

```
In [1096]: #probabilities and the expected values(actuals).
#compute the Brier Score-perfect skill has a score of 0.0 and the worst has a score of 1.0
#From this score, we can infer that our model has good performance or skill.
loss = brier_score_loss(y_test, prob2)
loss
```

```
Out[1096]: 0.12463763075043371
```

```
In [1097]: # define the evaluation procedure
cv6 = RepeatedStratifiedKFold(n_splits=10, n_repeats=5, random_state=1)
# evaluate the model using cross-validation
scores6 = cross_val_score(LR, X, y, scoring='accuracy', cv=cv6, n_jobs=-1)
```

```
In [1098]: scores6
```

```
Out[1098]: array([0.91235737, 0.90725904, 0.90531682, 0.90847293, 0.90337461,
       0.90483127, 0.89584851, 0.90483127, 0.9111219 , 0.90286547,
       0.90968682, 0.90410294, 0.89851906, 0.90216072, 0.91138626,
       0.89827628, 0.9024035 , 0.89730517, 0.90189412, 0.90820787,
       0.90386016, 0.89924739, 0.90823015, 0.90725904, 0.90653071,
       0.90143239, 0.90580238, 0.90386016, 0.90723652, 0.89946576,
       0.90264627, 0.90653071, 0.90895849, 0.90386016, 0.90337461,
       0.89851906, 0.9109007 , 0.90628793, 0.89728023, 0.90432249,
       0.9024035 , 0.91017237, 0.9024035 , 0.90580238, 0.89851906,
       0.90458849, 0.90871571, 0.89827628, 0.90553667, 0.90140845])
```

```
In [1099]: print(scores6.mean())
```

```
0.9042730727821567
```

## Classification Embedded Methods -LASSO Regularization (L1):

These methods encompass the benefits of both the wrapper and filter methods

```
In [1100]: # I have tested to find the best accuracy rate match of C = 0.002, 0.003, 0.01, 0.1, 0.5  
#with max_inter = 10000 and execution time is 39.8s  
#max_inter = 10000 was chosen because of the warning "ConvergenceWarning: Liblinear failed to converge  
#increase the number of iterations."
```

```
Sel = SelectFromModel(LogisticRegression(penalty ='l1', C = 0.001, solver ='liblinear',max_iter=1000
```

```
In [1101]:
```

```
Sel.fit(X_train, y_train)  
Sel.get_support()
```

```
Out[1101]: array([False, False, False, False, False, False, False, True, True,  
       True, True, True, False, True, True])
```

```
In [1102]: Sel.estimator_.coef_
```

```
Out[1102]: array([[ 0.          ,  0.          ,  0.          ,  0.          ,  0.          ,  
       0.          ,  0.          , -0.46490835,  0.03106993, -0.02571676,  
      0.00622425, -0.27336042,  0.          ,  0.97382011, -0.01797869]])
```

```
In [1103]: X_train_l1 = Sel.transform(X_train)  
X_test_l1 = Sel.transform(X_test)  
X_train_l1.shape
```

```
Out[1103]: (32540, 7)
```

### Test 3a. Random Forest with LASSO Regularization (L1)

```
In [1104]: # To improve the results of RF I tested n_estimators for 40,50,100,200,10000
#with max_depth of 2, 3 and 4

L1RandFor = RandomForestClassifier(n_estimators=50, random_state=43, max_depth=3)
L1RandFor.fit(X_train_l1, y_train)

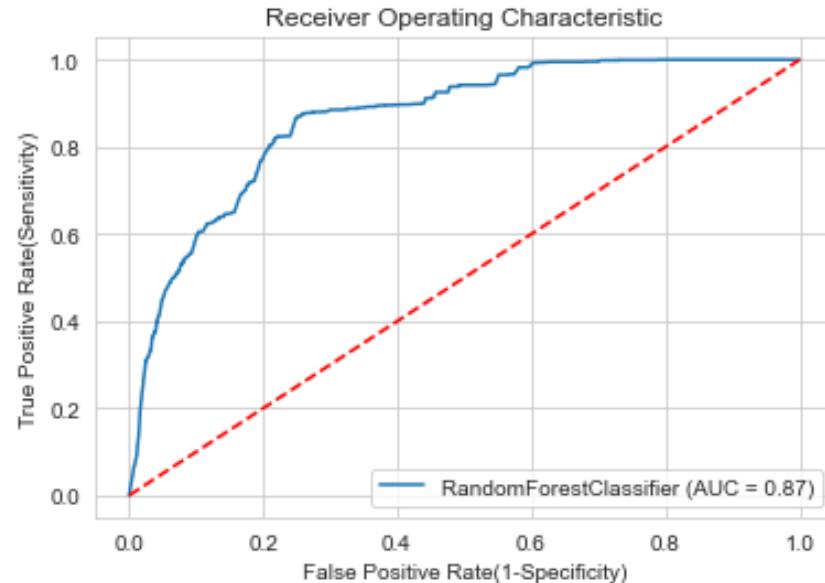
L1trainpred = L1RandFor.predict_proba(X_train_l1)
print('Accuracy on training set: {}'.format(roc_auc_score(y_train, L1trainpred[:,1])))

L1testpred = L1RandFor.predict_proba(X_test_l1)
print('Accuracy on test set: {}'.format(roc_auc_score(y_test, L1testpred[:,1])))

Accuracy on training set: 0.9455951787767048
Accuracy on test set: 0.8659482510337182
```

```
In [1105]: # draw the ROC-AUC chart  
metrics.plot_roc_curve(L1RandFor, X_test_l1, y_test)  
plt.title('Receiver Operating Characteristic')  
plt.plot([0, 1], [0, 1], 'r--')  
plt.ylabel('True Positive Rate(Sensitivity)')  
plt.xlabel('False Positive Rate(1-Specificity)')
```

```
Out[1105]: Text(0.5, 0, 'False Positive Rate(1-Specificity)')
```



```
In [1106]: pred7 = L1RandFor.predict(X_test_l1)
```

```
In [1107]: print("RandomForest with LASSO Regularization (L1)")  
cm = confusion_matrix(y_test, pred7)  
print(cm)  
print('\n')  
print(classification_report(y_test,pred7))  
tn, fp, fn, tp=cm.ravel()  
print ("TP: ", tp,", FP: ", fp,", TN: ", tn,", FN:", fn)
```

RandomForest with LASSO Regularization (L1)

```
[[9499 1450]  
 [ 515  893]]
```

	precision	recall	f1-score	support
0	0.95	0.87	0.91	10949
1	0.38	0.63	0.48	1408
accuracy			0.84	12357
macro avg	0.66	0.75	0.69	12357
weighted avg	0.88	0.84	0.86	12357

TP: 893 , FP: 1450 , TN: 9499 , FN: 515

```
In [1108]: L1testpred = L1RandFor.predict_proba(X_test_l1)  
prob4 = L1testpred[:, 1]# Keeping only the values in positive Label
```

```
In [1109]: #The average precision (PR AUC) is returned by passing the true Label  
## the probability estimate.  
# Average precision score  
PR_AUC = average_precision_score(y_test, prob4)  
print(PR_AUC)
```

0.44549859665982305

```
In [1110]: #Brier skill score calculates the mean squared error between predicted  
#probabilities and the expected values(actuals).  
#compute the Brier Score-perfect skill has a score of 0.0 and the  
#worst has a score of 1.0. From this score, we can infer that our model has good performance or skil  
loss = brier_score_loss(y_test, prob4)  
loss
```

```
Out[1110]: 0.11701306291673759
```

```
In [1112]: # define the evaluation procedure  
cv7 = RepeatedStratifiedKFold(n_splits=10, n_repeats=5, random_state=1)  
# evaluate the model using cross-validation  
scores7 = cross_val_score(L1RandFor, X, y, scoring='accuracy', cv=cv7, n_jobs=-1)
```

```
In [1113]: scores7
```

```
Out[1113]: array([0.89463462, 0.89779073, 0.89560573, 0.89681962, 0.89584851,  
0.89584851, 0.89536295, 0.89900461, 0.89460903, 0.89849441,  
0.8980335 , 0.89779073, 0.89099296, 0.89584851, 0.8980335 ,  
0.89342073, 0.89584851, 0.89827628, 0.89800874, 0.89460903,  
0.89414907, 0.89414907, 0.89876184, 0.9002185 , 0.89463462,  
0.89657684, 0.89584851, 0.89584851, 0.89509471, 0.89873725,  
0.89876184, 0.89706239, 0.89414907, 0.8948774 , 0.89681962,  
0.89706239, 0.89633406, 0.89609128, 0.89315202, 0.89679456,  
0.89851906, 0.89366351, 0.89293518, 0.89584851, 0.89512017,  
0.89706239, 0.89512017, 0.89754795, 0.89946576, 0.89558038])
```

```
In [1114]: print(scores7.mean())
```

```
0.8962173563772571
```

### Test 3b. Decision Tree with LASSO Regularization (L1)

```
In [1115]: %%time
dtreeL1 = DecisionTreeClassifier()
dtreeL1.fit(X_train_l1, y_train)

train_predL1 = dtreeL1.predict_proba(X_train_l1)
print('Accuracy on training set: {}'.format(roc_auc_score(y_train, train_predL1[:,1])))

test_predL1 = dtreeL1.predict_proba(X_test_l1)
print('Accuracy on test set: {}'.format(roc_auc_score(y_test, test_predL1[:,1])))
```

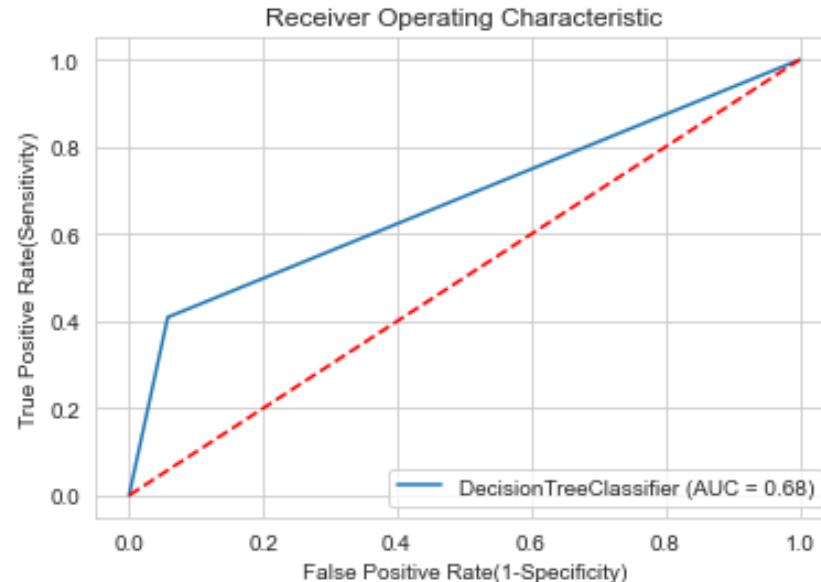
Accuracy on training set: 0.9999860377034893

Accuracy on test set: 0.6756096122829813

Wall time: 114 ms

```
In [1116]: # draw the ROC-AUC chart
metrics.plot_roc_curve(dtreet1, X_test_l1, y_test)
plt.title('Receiver Operating Characteristic')
plt.plot([0, 1], [0, 1], 'r--')
plt.ylabel('True Positive Rate(Sensitivity)')
plt.xlabel('False Positive Rate(1-Specificity)')
```

```
Out[1116]: Text(0.5, 0, 'False Positive Rate(1-Specificity)')
```



```
In [1117]: pred8= dtreet1.predict(X_test_l1)
```

```
In [1118]: print("DecisionTree with LASSO Regularization (L1)")  
cm = confusion_matrix(y_test, pred8)  
print(cm)  
print('\n')  
print(classification_report(y_test,pred8))  
tn, fp, fn, tp=cm.ravel()  
print ("TP: ", tp,", FP: ", fp,", TN: ", tn,", FN:", fn)
```

DecisionTree with LASSO Regularization (L1)

```
[[10316  633]  
 [ 832  576]]
```

	precision	recall	f1-score	support
0	0.93	0.94	0.93	10949
1	0.48	0.41	0.44	1408
accuracy			0.88	12357
macro avg	0.70	0.68	0.69	12357
weighted avg	0.87	0.88	0.88	12357

TP: 576 , FP: 633 , TN: 10316 , FN: 832

```
In [1119]: test_predL1 = dtreeL1.predict_proba(X_test_l1)  
prob5 = test_predL1[:, 1]# Keeping only the values in positive Label
```

```
In [1120]: #The average precision (PR AUC) is returned by passing the true Label & the probability estimate.  
# Average precision score  
PR_AUC = average_precision_score(y_test, prob5)  
print(PR_AUC)
```

0.2631439845778194

```
In [1121]: #Brier skill score calculates the mean squared error between predicted probabilities  
#and the expected values(actuals).  
#compute the Brier Score-perfect skill has a score of 0.0 and the worst has a score of 1.0  
#From this score, we can infer that our model has good performance or skill.  
loss = brier_score_loss(y_test, prob5)  
loss
```

```
Out[1121]: 0.11902243295958831
```

```
In [1125]: # define the evaluation procedure  
cv8 = RepeatedStratifiedKFold(n_splits=10, n_repeats=5, random_state=1)  
# evaluate the model using cross-validation  
scores8 = cross_val_score(dtreetL1, X, y, scoring='accuracy', cv=cv8, n_jobs=-1)
```

```
In [1126]: scores8
```

```
Out[1126]: array([0.88662297, 0.88468075, 0.88128186, 0.87933965, 0.89414907,  
0.89123574, 0.8776402 , 0.88103909, 0.88538125, 0.88416707,  
0.8895363 , 0.88370964, 0.88468075, 0.88273853, 0.88905074,  
0.88783685, 0.87958242, 0.88249575, 0.87931034, 0.88586693,  
0.87594076, 0.88176742, 0.88662297, 0.88249575, 0.8851663 ,  
0.88055353, 0.88176742, 0.88710852, 0.88926663, 0.8810102 ,  
0.88370964, 0.88006798, 0.88710852, 0.89026463, 0.89317796,  
0.88880796, 0.8873513 , 0.88783685, 0.88198154, 0.87736765,  
0.89293518, 0.88443797, 0.89123574, 0.8820102 , 0.87812576,  
0.88225297, 0.88565186, 0.8798252 , 0.88513842, 0.88343856])
```

```
In [1127]: print(scores8.mean())
```

```
0.8844954257276335
```

### Test 3c. Logistic Regression with LASSO Regularization (L1)

```
In [1128]: logreg = LogisticRegression(solver='liblinear')
logreg.fit(X_train_l1, y_train)

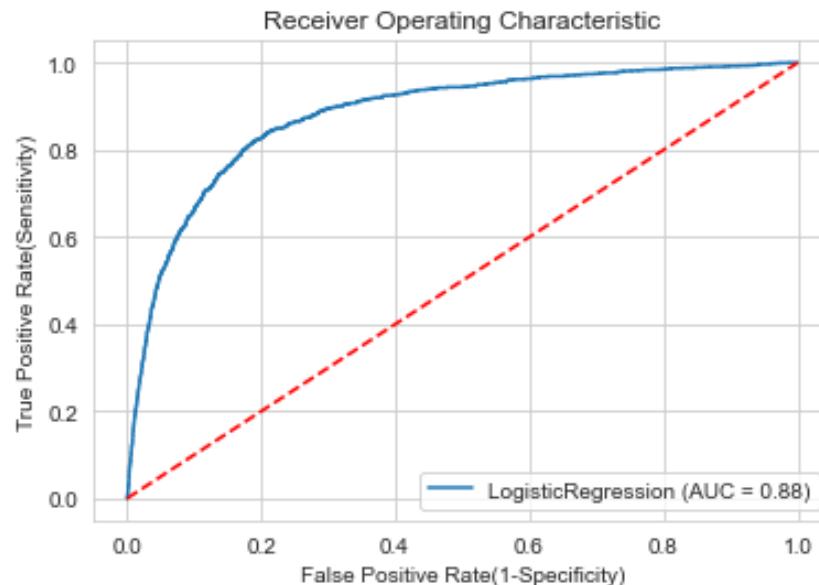
train_predLR = logreg.predict_proba(X_train_l1)
print('Accuracy on training set: {}'.format(roc_auc_score(y_train, train_predLR[:,1])))

test_predLR = logreg.predict_proba(X_test_l1)
print('Accuracy on test set: {}'.format(roc_auc_score(y_test, test_predLR[:,1])))
```

```
Accuracy on training set: 0.914914290161152
Accuracy on test set: 0.880972162256412
```

```
In [1129]: # draw the ROC-AUC chart
metrics.plot_roc_curve(logreg, X_test_l1, y_test)
plt.title('Receiver Operating Characteristic')
plt.plot([0, 1], [0, 1], 'r--')
plt.ylabel('True Positive Rate(Sensitivity)')
plt.xlabel('False Positive Rate(1-Specificity)')
```

```
Out[1129]: Text(0.5, 0, 'False Positive Rate(1-Specificity)')
```



```
In [1130]: pred9= logreg.predict(X_test_l1)
```

```
In [1131]: print("Logistic Regression with LASSO Regularization (L1)")  
cm = confusion_matrix(y_test, pred9)  
print(cm)  
print('\n')  
print(classification_report(y_test,pred9))  
tn, fp, fn, tp=cm.ravel()  
print ("TP: ", tp," , FP: ", fp," , TN: ", tn," , FN:", fn)
```

Logistic Regression with LASSO Regularization (L1)

```
[[9089 1860]  
 [ 296 1112]]
```

	precision	recall	f1-score	support
0	0.97	0.83	0.89	10949
1	0.37	0.79	0.51	1408
accuracy			0.83	12357
macro avg	0.67	0.81	0.70	12357
weighted avg	0.90	0.83	0.85	12357

TP: 1112 , FP: 1860 , TN: 9089 , FN: 296

```
In [1132]: test_predLR = logreg.predict_proba(X_test_l1)  
prob6 = test_predLR[:, 1]# Keeping only the values in positive label
```

```
In [1133]: #The average precision (PR AUC) is returned by passing the true Label  
#& the probability estimate.  
# Average precision score  
PR_AUC = average_precision_score(y_test, prob6)  
print(PR_AUC)
```

0.5133228548276421

```
In [1134]: #Brier skill score calculates the mean squared error between predicted  
#probabilities and the expected values(actuals).  
#compute the Brier Score-perfect skill has a score of 0.0 and the worst  
#has a score of 1.0. From this score, we can infer that our model has good performance or skill.  
loss = brier_score_loss(y_test, prob6)  
loss
```

```
Out[1134]: 0.1286185710853763
```

```
In [1135]: # define the evaluation procedure  
cv9 = RepeatedStratifiedKFold(n_splits=10, n_repeats=5, random_state=1)  
# evaluate the model using cross-validation  
scores9 = cross_val_score(logreg, X, y, scoring='accuracy', cv=cv9, n_jobs=-1)
```

```
In [1136]: scores9
```

```
Out[1136]: array([0.91235737, 0.90725904, 0.90531682, 0.90847293, 0.90337461,  
       0.90483127, 0.89584851, 0.90483127, 0.9111219 , 0.90286547,  
       0.90968682, 0.90410294, 0.89851906, 0.90216072, 0.91138626,  
       0.89827628, 0.9024035 , 0.89730517, 0.90189412, 0.90820787,  
       0.90386016, 0.89924739, 0.90823015, 0.90725904, 0.90653071,  
       0.90143239, 0.90580238, 0.90386016, 0.90723652, 0.89946576,  
       0.90264627, 0.90653071, 0.90895849, 0.90386016, 0.90337461,  
       0.89851906, 0.9109007 , 0.90628793, 0.89728023, 0.90432249,  
       0.9024035 , 0.91017237, 0.9024035 , 0.90580238, 0.89851906,  
       0.90458849, 0.90871571, 0.89827628, 0.90553667, 0.90140845])
```

```
In [1137]: print(scores9.mean())
```

```
0.9042730727821567
```

## Classification with Wrapper Feature Selection - Forward feature selection

Filter methods measure the relevance of features by their correlation with dependent variable while wrapper methods measure the usefulness of a subset of feature by actually training a model on it.

```
In [1138]: # Build step forward feature selection
# estimator is the RandomForestClassifier as passes to the SequentialFeatureSelector function.
#The k_features specifies the number of features to select. The forward parameter, set to True, per-
#step forward feature selection. The verbose parameter is used for logging the progress of
#the feature selector, the scoring parameter defines the performance evaluation criteria
#cv refers to cross-validation folds.
sfs = SequentialFeatureSelector(RandomForestClassifier(n_jobs=-1),
                                 k_features=(1, 15),
                                 forward=True,
                                 verbose=2,
                                 scoring='roc_auc',
                                 cv=4)
```

```
In [1139]: # call the fit method on our feature selector
#480 seconds execution time

sfs1 = sfs.fit(np.array(X_train.fillna(0)), y_train)

[Parallel(n_jobs=1)]: Using backend SequentialBackend with 1 concurrent workers.
[Parallel(n_jobs=1)]: Done  1 out of  1 | elapsed:    1.0s remaining:  0.0s
[Parallel(n_jobs=1)]: Done 15 out of 15 | elapsed: 16.7s finished

[2021-12-02 22:02:04] Features: 1/15 -- score: 0.959475151524163[Parallel(n_jobs=1)]: Using back
end SequentialBackend with 1 concurrent workers.
[Parallel(n_jobs=1)]: Done  1 out of  1 | elapsed:    1.6s remaining:  0.0s
[Parallel(n_jobs=1)]: Done 14 out of 14 | elapsed: 19.7s finished

[2021-12-02 22:02:23] Features: 2/15 -- score: 0.9717682058073968[Parallel(n_jobs=1)]: Using back
end SequentialBackend with 1 concurrent workers.
[Parallel(n_jobs=1)]: Done  1 out of  1 | elapsed:    2.0s remaining:  0.0s
[Parallel(n_jobs=1)]: Done 13 out of 13 | elapsed: 26.7s finished

[2021-12-02 22:02:50] Features: 3/15 -- score: 0.9798326636870762[Parallel(n_jobs=1)]: Using back
end SequentialBackend with 1 concurrent workers.
[Parallel(n_jobs=1)]: Done  1 out of  1 | elapsed:    2.7s remaining:  0.0s
[Parallel(n_jobs=1)]: Done 12 out of 12 | elapsed: 31.5s finished
```

```
In [1140]: pd.DataFrame.from_dict(sfs1 .get_metric_dict()).T
```

Out[1140]:

	feature_idx	cv_scores	avg_score	feature_names	ci_bound	std_dev	std_err
1	(13,)	[0.941880640374997, 0.9658901997732668, 0.9644...	0.959475	(13,)	0.016305	0.010172	0.005873
2	(10, 13)	[0.9594332419679319, 0.9774158339456195, 0.973...	0.971768	(10, 13)	0.011631	0.007256	0.004189
3	(0, 10, 13)	[0.973614039566852, 0.9834963235036347, 0.9809...	0.979833	(0, 10, 13)	0.005966	0.003722	0.002149
4	(0, 1, 10, 13)	[0.9792225611856855, 0.9896708016824386, 0.988...	0.986596	(0, 1, 10, 13)	0.006851	0.004274	0.002468
5	(0, 1, 10, 13, 14)	[0.9795754265028327, 0.9905193889760475, 0.991...	0.98807	(0, 1, 10, 13, 14)	0.007871	0.00491	0.002835
6	(0, 1, 9, 10, 13, 14)	[0.9815377336206546, 0.9928309952832824, 0.992...	0.989734	(0, 1, 9, 10, 13, 14)	0.007607	0.004746	0.00274
7	(0, 1, 3, 9, 10, 13, 14)	[0.9822345187141921, 0.9940380992998543, 0.993...	0.99069	(0, 1, 3, 9, 10, 13, 14)	0.007869	0.004909	0.002834
8	(0, 1, 3, 9, 10, 12, 13, 14)	[0.9829835566454609, 0.9944994292986769, 0.993...	0.99115	(0, 1, 3, 9, 10, 12, 13, 14)	0.007591	0.004736	0.002734
9	(0, 1, 2, 3, 9, 10, 12, 13, 14)	[0.984683813817669, 0.9953595309538679, 0.9943...	0.992301	(0, 1, 2, 3, 9, 10, 12, 13, 14)	0.007071	0.004411	0.002547
10	(0, 1, 2, 3, 8, 9, 10, 12, 13, 14)	[0.9835884081748703, 0.9954761251979201, 0.995...	0.992391	(0, 1, 2, 3, 8, 9, 10, 12, 13, 14)	0.008149	0.005083	0.002935
11	(0, 1, 2, 3, 8, 9, 10, 11, 12, 13, 14)	[0.98387844315677, 0.9960811882772798, 0.99608...	0.992939	(0, 1, 2, 3, 8, 9, 10, 11, 12, 13, 14)	0.008389	0.005234	0.003022
12	(0, 1, 2, 3, 7, 8, 9, 10, 11, 12, 13, 14)	[0.9837038540049066, 0.9962089946686994, 0.996...	0.992986	(0, 1, 2, 3, 7, 8, 9, 10, 11, 12, 13, 14)	0.008591	0.005359	0.003094
13	(0, 1, 2, 3, 5, 7, 8, 9, 10, 11, 12, 13, 14)	[0.9829001153007673, 0.9961940652864905, 0.996...	0.992786	(0, 1, 2, 3, 5, 7, 8, 9, 10, 11, 12, 13, 14)	0.009162	0.005716	0.0033
14	(0, 1, 2, 3, 4, 5, 7, 8, 9, 10, 11, 12, 13, 14)	[0.9824365186953341, 0.9964640332445306, 0.996...	0.992712	(0, 1, 2, 3, 4, 5, 7, 8, 9, 10, 11, 12, 13, 14)	0.009518	0.005937	0.003428

	feature_idx	cv_scores	avg_score	feature_names	ci_bound	std_dev	std_err
15	(0, 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13,...)	[0.9813930334546301, 0.9961901667231202, 0.996...]	0.992439	(0, 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13,...)	0.010235	0.006385	0.003686

```
In [1141]: sfs1.k_feature_names_
```

```
Out[1141]: ('0', '1', '2', '3', '7', '8', '9', '10', '11', '12', '13', '14')
```

```
In [1142]: #The best combination of features  
sfs1.k_score_
```

```
Out[1142]: 0.9929855023005756
```

```
In [1143]: filtered_features= X_train.columns[list(sfs1.k_feature_idx_)]  
filtered_features
```

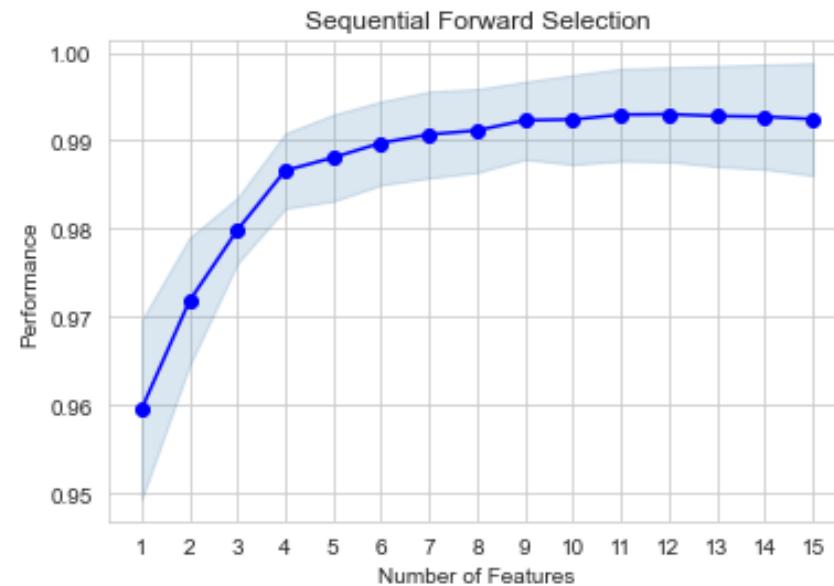
```
Out[1143]: Index(['age', 'job', 'marital', 'education', 'contact', 'month', 'day_of_week',  
       'duration', 'campaign', 'poutcome', 'cons.price.idx', 'nr.employed'],  
       dtype='object')
```

```
In [1144]: len(filtered_features)
```

```
Out[1144]: 12
```

In [1145]:

```
fig1 = plot_sfs(sfs1.get_metric_dict(), kind='std_dev')
plt.title('Sequential Forward Selection')
#plt.grid()
plt.show()
```



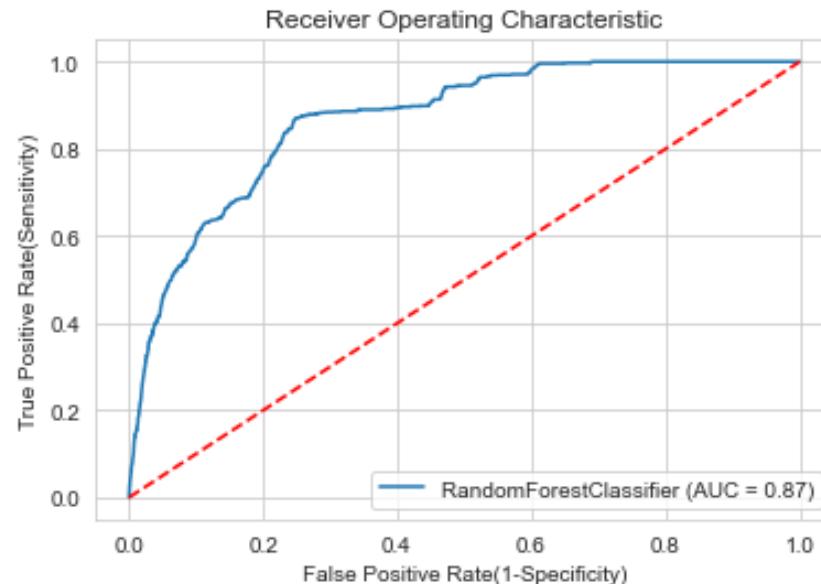
**Test 4a. Random Forest with Wrapper Forward feature selection**

```
In [1146]: #see the classification performance of the Random Forest using optimized  
#amount of features that was chosen by forward feature selection.  
  
# To improve the results of RF I tested n_estimators for 40,50,100,200,10000  
#with max_depth of 2, 3 and 4  
  
clf = RandomForestClassifier(n_estimators=50, random_state=41, max_depth=3)  
clf.fit(X_train[filtered_features].fillna(0), y_train)  
  
train_pred = clf.predict_proba(X_train[filtered_features].fillna(0))  
print('Accuracy on training set: {}'.format(roc_auc_score(y_train, train_pred[:,1])))  
  
test_pred = clf.predict_proba(X_test[filtered_features].fillna(0))  
print('Accuracy on test set: {}'.format(roc_auc_score(y_test, test_pred [:,1])))
```

Accuracy on training set: 0.9454497060777922  
Accuracy on test set: 0.8681434429462218

```
In [1147]: # draw the ROC-AUC chart
metrics.plot_roc_curve(clf, X_test[filtered_features], y_test)
plt.title('Receiver Operating Characteristic')
plt.plot([0, 1], [0, 1], 'r--')
plt.ylabel('True Positive Rate(Sensitivity)')
plt.xlabel('False Positive Rate(1-Specificity)')
```

```
Out[1147]: Text(0.5, 0, 'False Positive Rate(1-Specificity)')
```



```
In [1148]: pred = clf.predict(X_test[filtered_features])
```

```
In [1149]: print("Random Forest with Forward feature selection")
cm = confusion_matrix(y_test, pred)
print(cm)
print('\n')
print(classification_report(y_test,pred))
tn, fp, fn, tp=cm.ravel()
print ("TP: ", tp,", FP: ", fp,", TN: ", tn,", FN:", fn)
```

Random Forest with Forward feature selection

```
[[9517 1432]
 [ 509  899]]
```

	precision	recall	f1-score	support
0	0.95	0.87	0.91	10949
1	0.39	0.64	0.48	1408
accuracy			0.84	12357
macro avg	0.67	0.75	0.69	12357
weighted avg	0.89	0.84	0.86	12357

TP: 899 , FP: 1432 , TN: 9517 , FN: 509

```
In [1150]: test_pred = clf.predict_proba(X_test[filtered_features].fillna(0))
prob7 = test_pred [:, 1]# Keeping only the values in positive label
```

```
In [1151]: #The average precision (PR AUC) is returned by passing the true
#label & the probability estimate.
# Average precision score
PR_AUC = average_precision_score(y_test, prob7)
print(PR_AUC)
```

0.4808319916694562

```
In [1152]: #Brier skill score calculates the mean squared error between  
#predicted probabilities and the expected values(actuals).  
#compute the Brier Score-perfect skill has a score of 0.0  
#and the worst has a score of 1.0  
#From this score, we can infer that our model has good performance or skill.  
loss = brier_score_loss(y_test, prob7)  
loss
```

```
Out[1152]: 0.11870399234325582
```

```
In [1153]: # define the evaluation procedure  
cv10 = RepeatedStratifiedKFold(n_splits=10, n_repeats=5, random_state=1)  
# evaluate the model using cross-validation  
scores10 = cross_val_score(clf, X, y, scoring='accuracy', cv=cv10, n_jobs=-1)
```

```
In [1154]: scores10
```

```
Out[1154]: array([0.90337461, 0.89900461, 0.89584851, 0.89851906, 0.8980335 ,  
0.89949017, 0.89876184, 0.89827628, 0.89946576, 0.90189412,  
0.90191794, 0.9002185 , 0.89633406, 0.89681962, 0.90070405,  
0.89924739, 0.89876184, 0.90046128, 0.8997086 , 0.8997086 ,  
0.89609128, 0.89949017, 0.90167516, 0.90410294, 0.89924739,  
0.89876184, 0.89827628, 0.89827628, 0.89728023, 0.9023798 ,  
0.90191794, 0.8980335 , 0.90070405, 0.89827628, 0.90337461,  
0.90288905, 0.89827628, 0.89754795, 0.89558038, 0.89946576,  
0.9002185 , 0.90288905, 0.89876184, 0.89754795, 0.89657684,  
0.89924739, 0.8980335 , 0.89949017, 0.89946576, 0.89922292])
```

```
In [1155]: print(scores10.mean())
```

```
0.8993930282686483
```

#### Test 4b. Decision Tree with Wrapper Forward feature selection

```
In [1156]: %%time
dtree = DecisionTreeClassifier()
dtree.fit(X_train[filtered_features], y_train)

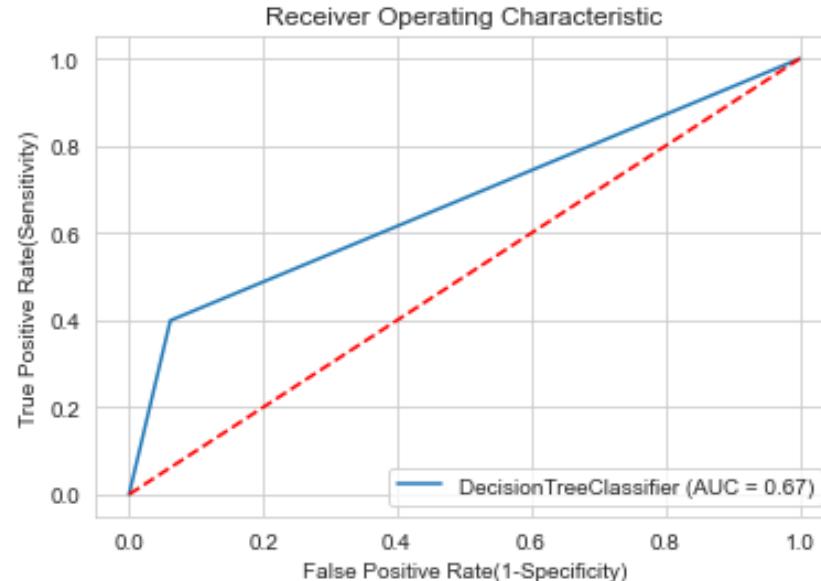
train_predTree = dtree.predict_proba(X_train[filtered_features].fillna(0))
print('Accuracy on training set: {}'.format(roc_auc_score(y_train, train_predTree[:,1])))

test_predTree = dtree.predict_proba(X_test[filtered_features].fillna(0))
print('Accuracy on test set: {}'.format(roc_auc_score(y_test, test_predTree[:,1])))
```

```
Accuracy on training set: 1.0
Accuracy on test set: 0.6687947970549406
Wall time: 153 ms
```

```
In [1157]: # draw the ROC-AUC chart
metrics.plot_roc_curve(dtree, X_test[filtered_features], y_test)
plt.title('Receiver Operating Characteristic')
plt.plot([0, 1], [0, 1], 'r--')
plt.ylabel('True Positive Rate(Sensitivity)')
plt.xlabel('False Positive Rate(1-Specificity)')
```

```
Out[1157]: Text(0.5, 0, 'False Positive Rate(1-Specificity)')
```



```
In [1158]: pred1 = dtree.predict(X_test[filtered_features])
```

```
In [1159]: print("DecisionTree with Forward feature selection")
cm = confusion_matrix(y_test, pred1)
print(cm)
print('\n')
print(classification_report(y_test,pred1))
tn, fp, fn, tp=cm.ravel()
print ("TP: ", tp,", FP: ", fp,", TN: ", tn,", FN:", fn)
```

DecisionTree with Forward feature selection

```
[[10275  674]
 [ 846  562]]
```

	precision	recall	f1-score	support
0	0.92	0.94	0.93	10949
1	0.45	0.40	0.43	1408
accuracy			0.88	12357
macro avg	0.69	0.67	0.68	12357
weighted avg	0.87	0.88	0.87	12357

TP: 562 , FP: 674 , TN: 10275 , FN: 846

```
In [1160]: test_predTree = dtree.predict_proba(X_test[filtered_features].fillna(0))
prob8 = test_predTree[:, 1]# Keeping only the values in positive label
```

```
In [1161]: #The average precision (PR AUC) is returned by passing the true
#label & the probability estimate.
# Average precision score
PR_AUC = average_precision_score(y_test, prob8)
print(PR_AUC)
```

0.2499527198163763

```
In [1162]: #Brier skill score calculates the mean squared error  
#between predicted probabilities and the expected values(actuals).  
#compute the Brier Score-perfect skill has a score of 0.0 and the worst has a score of 1.0  
#From this score, we can infer that our model has good performance or skill.  
loss = brier_score_loss(y_test, prob8)  
loss
```

```
Out[1162]: 0.12300720239540341
```

```
In [1164]: # define the evaluation procedure  
cv11 = RepeatedStratifiedKFold(n_splits=10, n_repeats=5, random_state=1)  
# evaluate the model using cross-validation  
scores11 = cross_val_score(dtree, X, y, scoring='accuracy', cv=cv11, n_jobs=-1)
```

```
In [1165]: scores11
```

```
Out[1165]: array([0.88710852, 0.88346686, 0.87909687, 0.87642632, 0.89657684,  
0.88905074, 0.87909687, 0.88176742, 0.88246722, 0.88392424,  
0.88710852, 0.88103909, 0.87909687, 0.88273853, 0.88710852,  
0.88686574, 0.87909687, 0.88298131, 0.88052453, 0.8890238 ,  
0.87569798, 0.88443797, 0.88710852, 0.88492353, 0.88176742,  
0.88249575, 0.88540908, 0.88880796, 0.88659543, 0.8836814 ,  
0.88443797, 0.8798252 , 0.88419519, 0.88807963, 0.89172129,  
0.88929352, 0.88419519, 0.88783685, 0.88125304, 0.88052453,  
0.89244962, 0.88540908, 0.8895363 , 0.88759408, 0.87691187,  
0.88370964, 0.89099296, 0.8820102 , 0.88465274, 0.88052453])
```

```
In [1166]: print(scores11.mean())
```

```
0.8844128825998662
```

#### Test 4c. Logistic Regression with Wrapper Forward feature selection

In [1167]: #Logistic Regression

```
LogitReg = LogisticRegression(solver='liblinear')
LogitReg.fit(X_train[filtered_features], y_train)

train_predReg = LogitReg .predict_proba(X_train[filtered_features].fillna(0))
print('Accuracy on training set: {}'.format(roc_auc_score(y_train, train_predReg[:,1])))

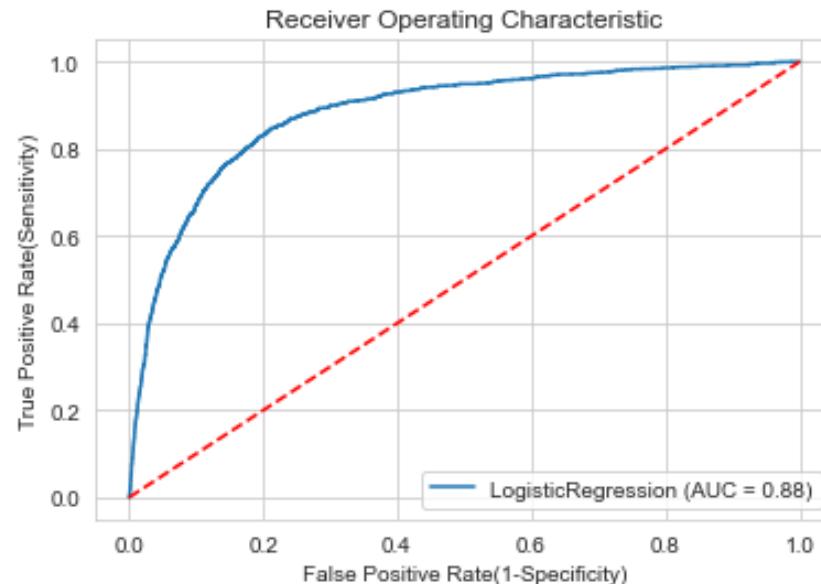
test_predReg = LogitReg .predict_proba(X_test[filtered_features].fillna(0))
print('Accuracy on test set: {}'.format(roc_auc_score(y_test, test_predReg [:,1])))
```

Accuracy on training set: 0.9183572579953602

Accuracy on test set: 0.8836393578907165

```
In [1168]: # draw the ROC-AUC chart
metrics.plot_roc_curve(LogitReg, X_test[filtered_features], y_test)
plt.title('Receiver Operating Characteristic')
plt.plot([0, 1], [0, 1], 'r--')
plt.ylabel('True Positive Rate(Sensitivity)')
plt.xlabel('False Positive Rate(1-Specificity)')
```

```
Out[1168]: Text(0.5, 0, 'False Positive Rate(1-Specificity)')
```



```
In [1169]: pred2 = LogitReg.predict(X_test[filtered_features])
```

```
In [1170]: print("Logistic Regression with Forward feature selection")
cm = confusion_matrix(y_test, pred2)
print(cm)
print('\n')
print(classification_report(y_test,pred2))
tn, fp, fn, tp=cm.ravel()
print ("TP: ", tp,", FP: ", fp,", TN: ", tn,", FN:", fn)
```

```
Logistic Regression with Forward feature selection
[[8846 2103]
 [ 247 1161]]
```

	precision	recall	f1-score	support
0	0.97	0.81	0.88	10949
1	0.36	0.82	0.50	1408
accuracy			0.81	12357
macro avg	0.66	0.82	0.69	12357
weighted avg	0.90	0.81	0.84	12357

```
TP: 1161 , FP: 2103 , TN: 8846 , FN: 247
```

```
In [1171]: test_predReg = LogitReg .predict_proba(X_test[filtered_features])
prob9 = test_predReg[:, 1]# Keeping only the values in positive label
```

```
In [1172]: #The average precision (PR AUC) is returned by passing the
#true label & the probability estimate.
# Average precision score
PR_AUC = average_precision_score(y_test, prob9)
print(PR_AUC)
```

```
0.5209186484625814
```

```
In [1173]: #Brier skill score calculates the mean squared error  
#between predicted probabilities and the expected values(actuals).  
#compute the Brier Score-perfect skill has a score of 0.0 and the worst has a score of 1.0  
#From this score, we can infer that our model has good performance or skill.  
loss = brier_score_loss(y_test, prob9)  
loss
```

```
Out[1173]: 0.13839182894594096
```

```
In [1174]: # define the evaluation procedure  
cv12 = RepeatedStratifiedKFold(n_splits=10, n_repeats=5, random_state=1)  
# evaluate the model using cross-validation  
scores12 = cross_val_score(LogitReg, X, y, scoring='accuracy', cv=cv12, n_jobs=-1)
```

```
In [1175]: scores12
```

```
Out[1175]: array([0.91235737, 0.90725904, 0.90531682, 0.90847293, 0.90337461,  
0.90483127, 0.89584851, 0.90483127, 0.9111219 , 0.90286547,  
0.90968682, 0.90410294, 0.89851906, 0.90216072, 0.91138626,  
0.89827628, 0.9024035 , 0.89730517, 0.90189412, 0.90820787,  
0.90386016, 0.89924739, 0.90823015, 0.90725904, 0.90653071,  
0.90143239, 0.90580238, 0.90386016, 0.90723652, 0.89946576,  
0.90264627, 0.90653071, 0.90895849, 0.90386016, 0.90337461,  
0.89851906, 0.9109007 , 0.90628793, 0.89728023, 0.90432249,  
0.9024035 , 0.91017237, 0.9024035 , 0.90580238, 0.89851906,  
0.90458849, 0.90871571, 0.89827628, 0.90553667, 0.90140845])
```

```
In [1176]: print(scores12.mean())
```

```
0.9042730727821567
```

TP-True positives - are when you predict an observation belongs to a class and it actually does belong to that class.

FP-False positives - occur when you predict an observation belongs to a class when in reality it does not.

TN-True negatives - are when you predict an observation does not belong to a class and it actually does not belong to that class.

FN-False negatives - occur when you predict an observation does not belong to a class when in fact it does.

The confusion matrix is in the form of the array object. The dimension of the matrixes is 2\*2 because the models are binary classification. It has two classes 0's that are "No" and 1's that are "Yes". Diagonal values represent accurate predictions, while non-diagonal elements are inaccurate predictions.

When to Use ROC vs. Precision-Recall Curves?

ROC curves should be used when there are roughly equal numbers of observations for each class. Precision-Recall curves should be used when there is a moderate to large class imbalance.

## Results

The best combination is Filter Feature selection with method of Mutual Information Gain and the classifier Logistic Regression. Filter is the fastest among Wrapper and Embedded Methods. Logistic Regression has the best scores with this imbalanced dataset.

AUC score for the best performed combination is 0.88. AUC score 1 represents perfect classifier, and 0.5 represents a worthless classifier.

Precision: how accurate your model is. In other words, when a model makes a prediction, how often it is correct. It is correct at 38%.

Recall: If there are patients who subscribed in the test set. The Logistic Regression model can identify it 79% of the time.

F1 is 0.51

PR\_AUC can be interpreted as the probability that the scores given by a classifier will rank a randomly chosen positive instance higher than a randomly chosen negative one. LR has the best score of 52%. This score might be the most commonly used for comparing classification models for imbalanced problems.

Brier score is 0.12

cross-validation score is 0.9

In [ ]: