



**Group - 1**

**CAPSTONE PROJECT-FINAL REPORT**

**TARGETING THE RIGHT CUSTOMERS USING CLASSIFICATION  
ALGORITHMS TO INCREASE REVENUE**

**MENTOR:**

- Mr. Subramanian P V

**SUBMITTED BY:**

- Dinesh V
- Keerthivasan V
- Muralimanohar V
- Poorani R
- Raghul kanna S
- Rukmani V

## **Table of Contents**

<b>Sl.No</b>	<b>Topic</b>	<b>Page</b>
1	Industry Review	3
2	Literature Review	4
3	Dataset and Domain	4
4	Business Problem	4
5	Objective	5
6	Dataset Information	5
7	Data Dictionary	5
8	Findings and Implications in Outset Data	6
9	Data Visualizations	7
10	Methodology Followed	20
11	Pre-Processing Steps	20
12	Algorithms Used	21
13	Assumptions for different models	21
14	Step-by-Step Walkthrough of the solution	23
15	Data Modelling and Evaluation	33
16	Feature Importance	36
17	Recommendations	37
18	Limitations	37
19	Conclusion	38
20	References	

## 1. Industry Review:

- Marketing is technique of exposing the target clients to a product via suitable systems and channels. It ultimately facilitates the way to buy the product or service and even helps in determining the need of the product and persuade customers to buy it.
- The overall aim is to increase sales of products and services for enterprise, business and financial institutions. It also helps to maintain the reputation of the company.
- Telemarketing is form of direct marketing in which salesperson approaches the customer either face to face or phone call and persuade him to buy the product. Telemarketing attains most popularity in 20th century and still gaining it.
- Nowadays, telephone (fixed-line or mobile) has been broadly used. It is cost effective and keeps the customers up to date.
- In Banking sector, marketing is the backbone to sell its product or service. Banking advertising and marketing is mostly based on an intensive knowledge of objective information about the market and the actual client needs for the bank profitable manner.
- Making right decisions in organizational operations are sometimes proved a great challenge where the quality of decision really matters.
- Decision Support Systems (DSS) are classified as a particular class of computerized facts and figures that helps the organization or administration into their decision making actions. The concept of DSS originates from a balance which lies between the data generated by computer and the judgment of human. According to Rupnik & Kukar (2007) the objective of decision support systems is to enhance the effectiveness of the decisions. This is a great tool which can analyse the sales data and provide further predictions. The purposes which can be established from the DSS are such as, analysis, optimization, forecasting and simulation. A study by Power (2008) found that research subjects who use DSS for the decision making, come-up with more effective decisions than those who did not use it. Nowadays, DSS is contributing a meaningful role in many fields such as for medical diagnosis, business and management, investment portfolios, command and control of military units, and statistics. DSS uses statistical data to overcome the deficiencies and helps the decision makers to take the right decision.
- Data mining (DM) plays vital role to support the Decision support systems which are based on the data obtained from the data mining models: rules, patterns and relationship. Data mining is the process of selecting, discovering, and modelling high volume of data to find and clarify unknown patterns. The objective of data mining in decision support systems is to suggest a tool which is easily accessible for the business users to analyse the data mining models. A specific technology used within the DSS is Machine learning (ML) that combines data and computer applications to accurately predicting the results. The fundamental principle of machine learning is to construct the algorithms that can obtain input data and then predict the results or outputs by using the statistical analysis within satisfactory interval.
- ML allows the DSS to obtain the new knowledge which helps it to make right decisions. Machine Learning can be mainly classified in 2 categories i.e. supervised learning and unsupervised learning. In supervised learning, the output of algorithm is already known and we use the input data to predict the output. The examples of supervised learning are regression and classification. In contrast, unsupervised learning we only have input data whereas no corresponding output variables are selected. The example of unsupervised learning is clustering. 2 Feature selection is the process of selecting the subset of relevant variables from the model. It identifies the most important attributes which help to predict the output. By using this techniques, we can reduce the curse of dimensionality, prevent model from overfitting and shorter the training time. In this way parsimonious model can be achieved with minimum number of parameter and good explanatory predictive power.

## **2. Literature Review:**

- Few studies examined the predicting the success of bank telemarketing for selling long-term deposits through the application of various machine learning techniques. The prior studies were shown below:
- Asare-Frempong and Jayabalan (2017) anticipated a model using four machine learning algorithms: Multilayer perceptron neural network, Decision tree (C4.5), Logistic regression, and random forest. The dataset is taken from (UCI) Machine Learning Repository, containing 45147 instances with 17 attributes. The random forest gives better accuracy is 86.8%. Palaniappan et al. (2017) suggested a model using data mining approaches. The dataset was taken from the UCI Machine Learning repository, with 41,188 instances 21 attributes. Three algorithms had been applied, which are Naïve Bayes, Random Forest, and Decision Tree. The experiments measured the accuracy percentage, precision, and recall rates. They found that the decision tree algorithm gives the best accuracy.
- Another study by Jiang (2018) explored predicting the success of bank telemarketing using data mining approaches. The dataset was taken from the (UCI) Machine Learning Repository. There are 4119 instances and 21 attributes in this dataset. They used the support vector machine, logistic regression, Naïve Bayes, Neural Network and Decision Tree. Among these five algorithms, they got the best accuracy from the logistic regression. The accuracy of the logistic regression model was 92.03%. According to Ilham et al. (2019) proposed a model using machine learning approaches. They used different techniques: Logistic Regression, Naïve Bayes, Random Forest, K-Nearest Neighbor, Support Vector Machine, Neural Network, and Decision Tree. Preprocessing was not done to the dataset features; it directly uses a ready dataset from
- St. Theresa Journal of Humanities and Social Sciences Vol.7, No.1 January-June 2021 94 the UCI repository. These models' evaluation metrics verify that the most accurate they are founded that 91.07% by using the SVM.

## **3. Dataset Domain:**

- A dataset is a collection of data and it can be structured or unstructured.
- A structured data is represented in a tabular format, where every column of the table represents a particular variable, and each row corresponds to a given record of the dataset in question.
- Unsupervised data is not represented in a tabular form, data that we fetch from Facebook, Twitter, and Netflix etc. with the help of recommendation systems are all our unsupervised data. This dataset belongs to the domain of Marketing.

## **4. Business Problem:**

- There has been a revenue decline for the Portuguese bank and they would like to know what actions to take. After investigation, they found out that the root cause is that their clients are not depositing as frequently as before. Knowing that term deposits allow banks to hold onto a deposit for a specific amount of time, so banks can invest in higher gain financial products to make a profit. In addition, banks also hold better chance to persuade term deposit clients into buying other products such as funds or insurance to further increase their revenues. As a result, the Portuguese bank would like to identify existing clients that have higher chance to subscribe for a term deposit and focus marketing effort on such clients.
- To resolve the problem, we suggest a classification approach to predict which clients are more likely to subscribe for term deposits.

## **5. Objective:**

- The basic objective of the project is to analyze Tele-Marketing data collected from a Bank and predict whether the customer will subscribe the term deposit or not, also identify the driving factors behind this. This would inform the Bank's decisions on which Customers to target for their Marketing Campaign, which would ultimately increase their Product Sales. Understanding their customers is critical for effectively executing their Marketing campaign.

## **6. Dataset Information:**

- This dataset contains 45000 rows and 23 columns, where each observation corresponds to a customer response. Among the total 45000 observations, 5289 observations (11.7%) are those who actually bought the product.

## **7. Data Dictionary:**

- We are provided with Customer details such as Age, Salary, Job, Marital Status, etc. also the campaign details such as Number of calls performed, No of days passed after and before the campaign to better understand the Problem.
- Data can be seen from the following table:

Sl. No	Column Name	Description
1	age	Age of the targeted Customers
2	age group	Age group of the targeted Customers
3	eligible	if the customer is eligible for the talk or not
4	job	what does the customer do? (Type of jobs)
5	salary	salary of the customer
6	marital	married or not?
7	education	level of education of Customers
8	marital-education	Combination of marital & education
9	targeted	Is the Customer being targeted or not ?
10	default	if the customer in default list or not (has credit in default?)
11	balance	remaining balance in their accounts
12	housing	has housing loan?
13	loan	has personal loan?
14	contact	contact communication type
15	day	last contact day of the week
16	month	last contact month of year

17	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.
18	campaign	number of contacts performed during this campaign and for this client
19	pdays	number of days that passed by after the client was last contacted from a previous campaign
20	previous	number of contacts performed before this campaign and for this client
21	poutcome	outcome of the previous marketing campaign
22	y	Response - Target column
23	response	Response - Target column (Numeric)

## 8. Findings and Implications in Outset Data:

	age	age group	salary	balance	day	duration	campaign	pdays	previous
count	45211.000000	45211.000000	45211.000000	45211.000000	45211.000000	45211.000000	45211.000000	45211.000000	45211.000000
mean	40.936210	3.645861	57006.171065	1362.272058	15.806419	258.163080	2.763841	40.197828	0.580323
std	10.618762	1.083271	32085.718415	3044.765829	8.322476	257.527812	3.098021	100.128746	2.303441
min	18.000000	1.000000	0.000000	-8019.000000	1.000000	0.000000	1.000000	-1.000000	0.000000
25%	33.000000	3.000000	20000.000000	72.000000	8.000000	103.000000	1.000000	-1.000000	0.000000
50%	39.000000	3.000000	60000.000000	448.000000	16.000000	180.000000	2.000000	-1.000000	0.000000
75%	48.000000	4.000000	70000.000000	1428.000000	21.000000	319.000000	3.000000	-1.000000	0.000000
max	95.000000	9.000000	120000.000000	102127.000000	31.000000	4918.000000	63.000000	871.000000	275.000000

	count	unique	top	freq
eligible	45211	2	Y	43380
job	45211	12	blue-collar	9732
marital	45211	3	married	27214
education	45211	4	secondary	23202
targeted	45211	2	yes	37091
default	45211	2	no	44396
housing	45211	2	yes	25130
loan	45211	2	no	37967
contact	45211	3	cellular	29285
month	45211	12	may	13766
poutcome	45211	4	unknown	36959
y	45211	2	no	39922

- There are 45000 distinct id's in the dataset which means 45000 Account holders and the average Bank Balance amount is 57006.17 with std. 32085.72 where the highest amount is 120000.
- The avg. Bank balance amount is 57006.17 with std. of 32085.72 where the minimum amount is -8019 and maximum amount is 102127.
- Account holders age ranges from 18 to 95 and the average age is 40.94 with standard deviation 10.62
- Target column which includes two distinct values 0 and 1 where 0 for response yes and 1 means response no
- There are total of 23 columns and 45000 records in the dataset and non of the columns has null values.
- There are 10 continuous columns and 13 categorical columns in the dataset.
- Though there are no null values in the dataset there were some unknown classified records which were are being pre-processed before building the base model.
- Class 0 represents Customers who did not buy the product and Class 1 represents Customers who bought the product
- The Target in the data has a large imbalance with respect to class 1 (Minority Class)
- Resampling methods needs to be applied to increase the percent of defaulters in order to achieve better model training results.
- There is no any null values in the dataset.

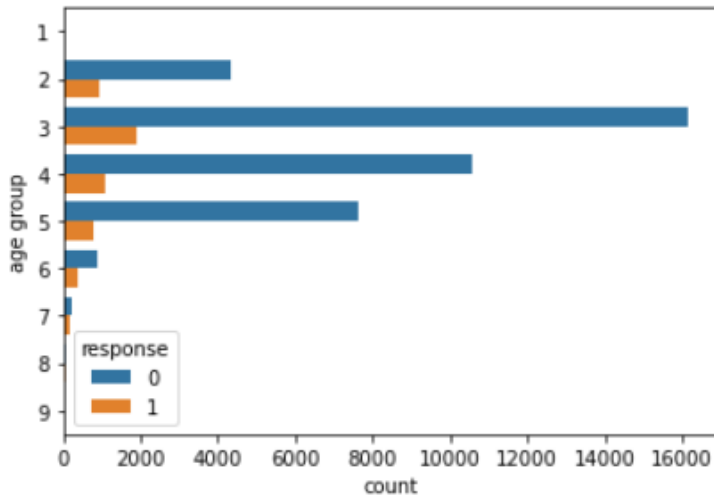
```
: 1 df.isnull().sum()
```

```
: age                0
  age group          0
  eligible            0
  job                0
  salary              0
  marital             0
  education           0
  marital-education   0
  targeted            0
  default             0
  balance             0
  housing             0
  loan                0
  contact             0
  day                 0
  month              0
  duration            0
  campaign            0
  pdays              0
  previous            0
  poutcome            0
  y                  0
  response            0
  dtype: int64
```

## 9. Data Visualizations:

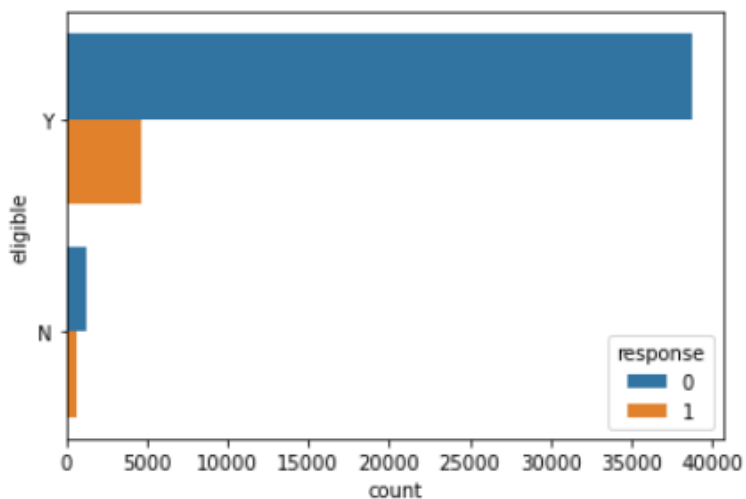
- **Categorical Variables & Target:**

**Figure 1:** Age\_group Vs Response



**Inferences:** Customers in age\_group 3 are more likely to subscribe the Term deposit. Because they age range from 31-40, who can be professional with high paying jobs and this possibilities are high because of their experience.

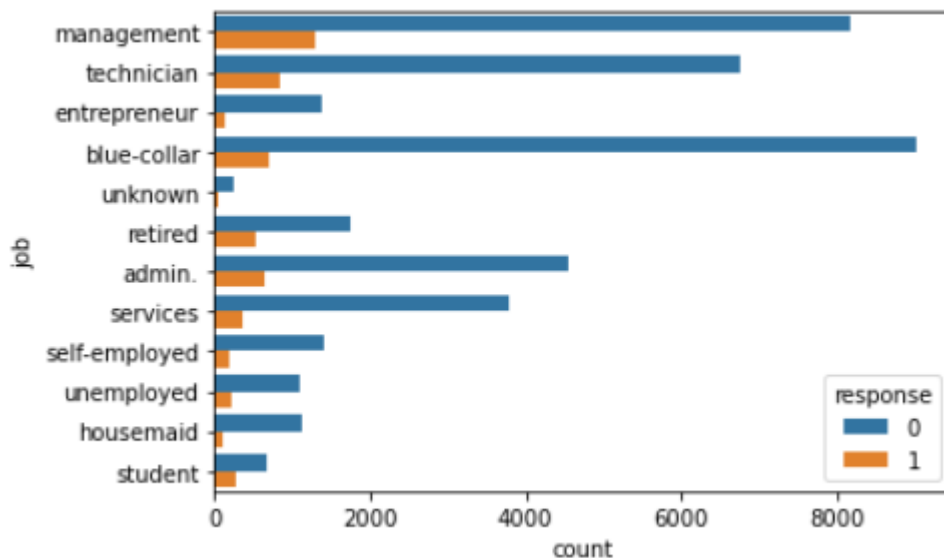
**Figure 2:** Eligible Vs Response



**Inferences:** Eligible customers are more likely to subscribe the Term deposit, Because they saves money for their children education, marriage etc.. it is a long term investment.



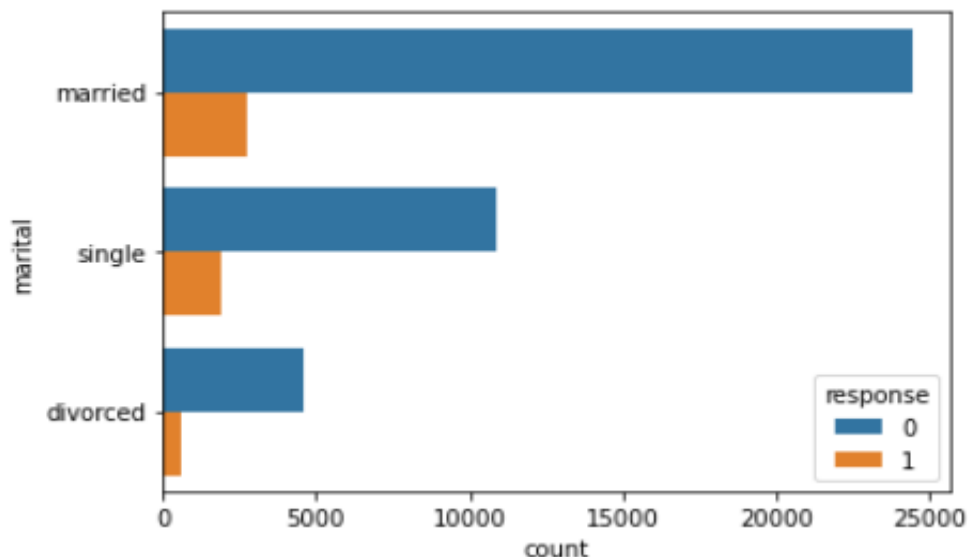
**Figure 3: Job Vs Response**



**Inferences:** Customers with “Management”, “Technician”, and “Blue-collar” job categories are more likely to subscribe the Term deposit. It is because these professionals have regular income.

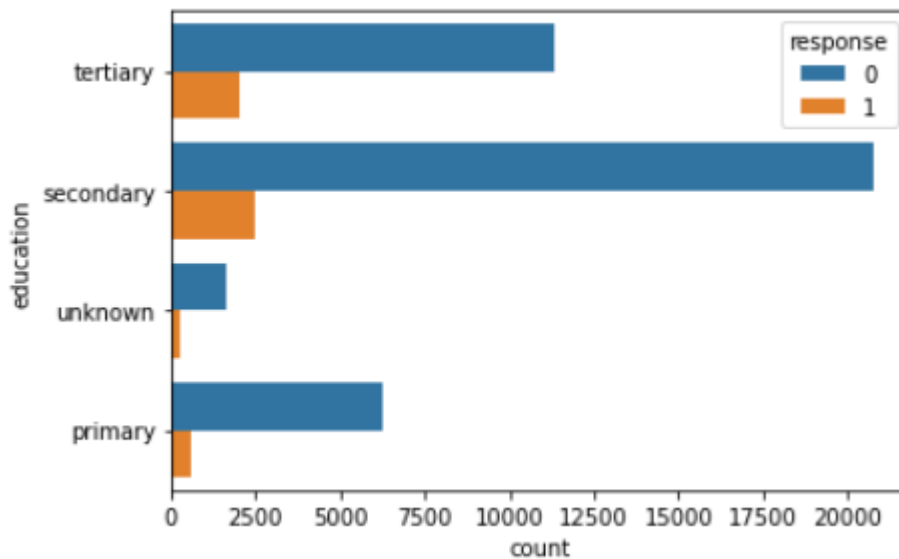
Other category workers may or may not have the regular income, which includes student (only few students work for part time jobs), housemaid (their earnings very less used to lead the normal life) etc..

**Figure 4: Marital Vs Response**



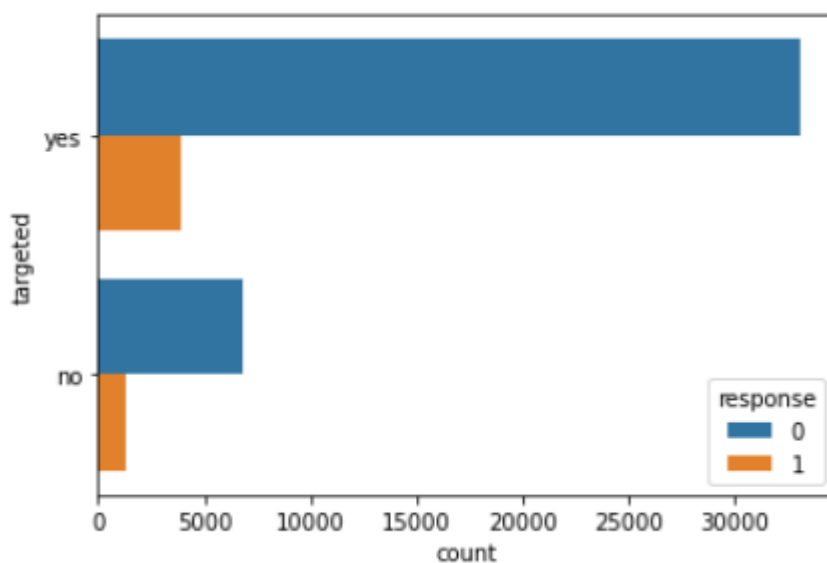
**Inferences:** Married Customers are more likely to subscribe the Term deposit, they saves money for their children education, marriage etc.. it is a long term investment.

**Figure 5: Education Vs Response**



**Inferences:** Customers with “Secondary” & Tertiary” educations are more likely to subscribe the Term deposit, , it is mainly due to they have clarity about financial education and knows the value of long term investments.

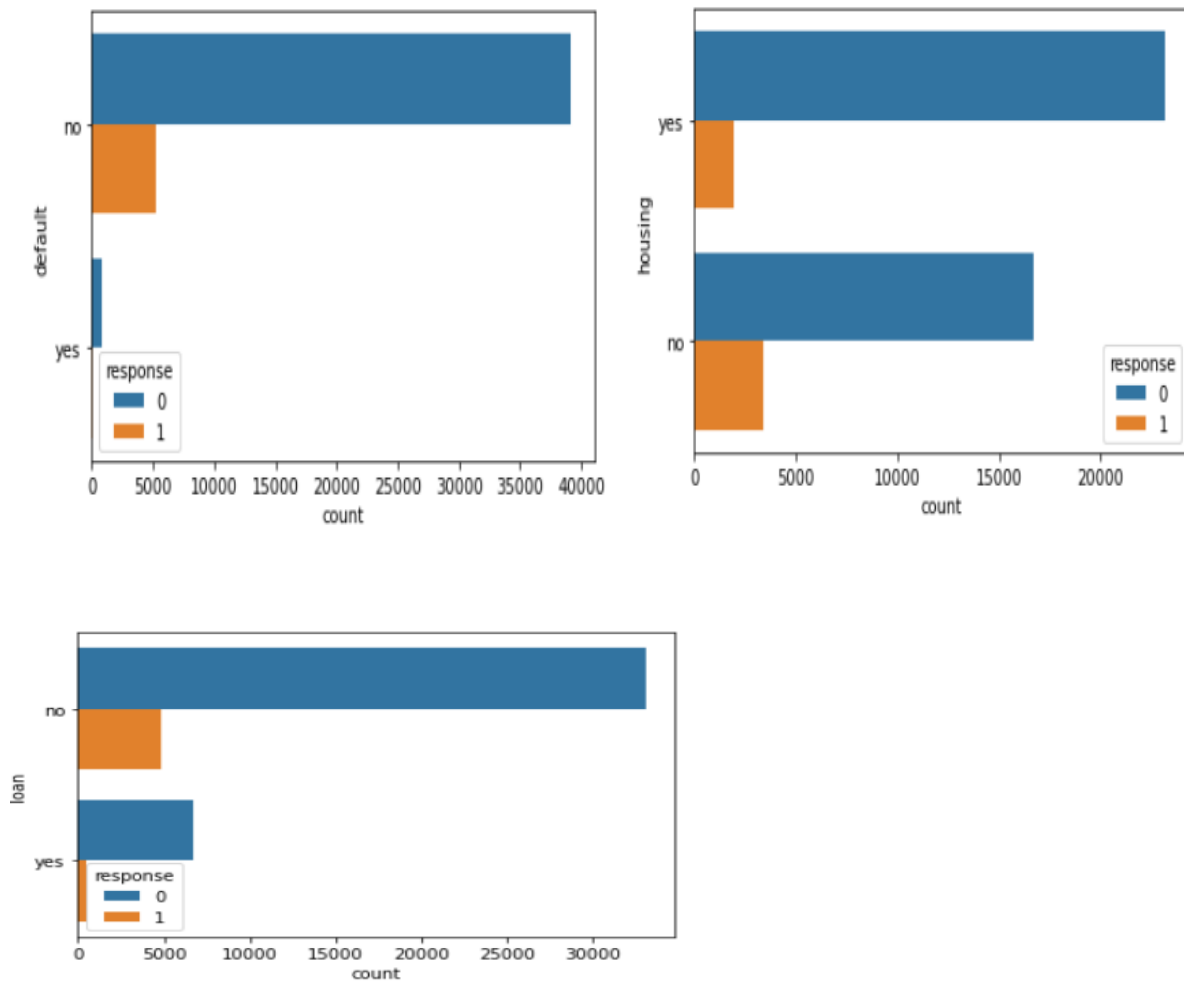
**Figure 6: Targeted Vs Response**



**Inferences:**

- Targeted Customers are more likely to subscribe the Term deposit.

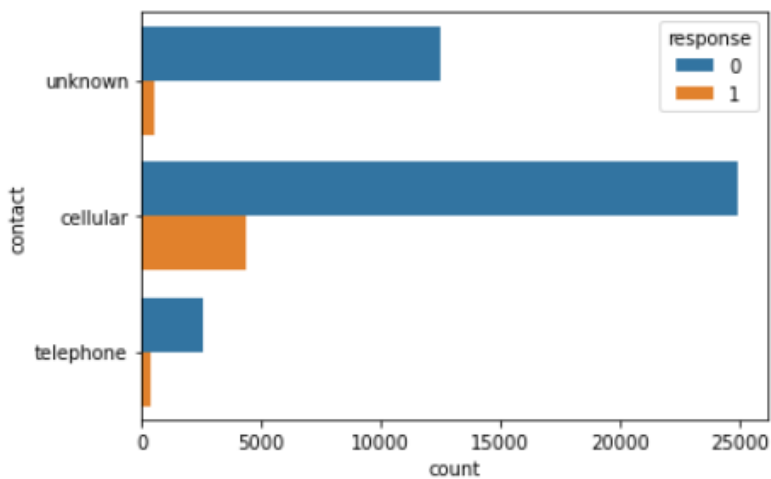
**Figure 7: Loan Vs Response**



**Inferences:**

- Customers with no loan defaults, no loans are more likely to subscribe the Term deposit.

**Figure 8: Targeted Vs Response**

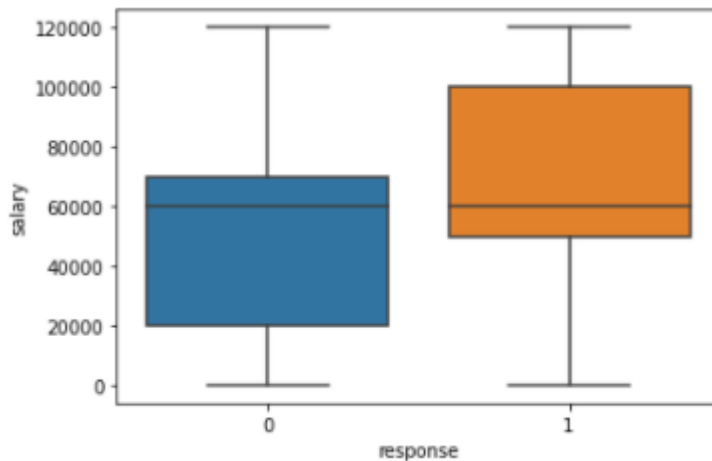


### Inferences:

- Customers contacted with Cellular phones are more likely to subscribe the Term deposit.

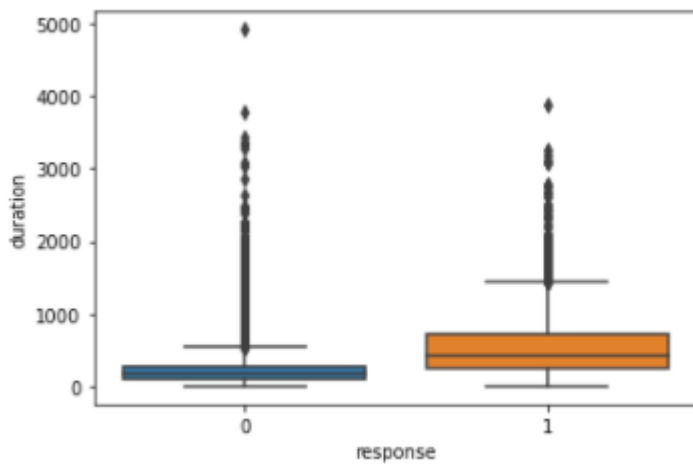
### Numerical Variables & Target:

**Figure 9:** Salary Vs Response



**Inferences:** Customers with higher salary are more likely to subscribe the Term deposit.

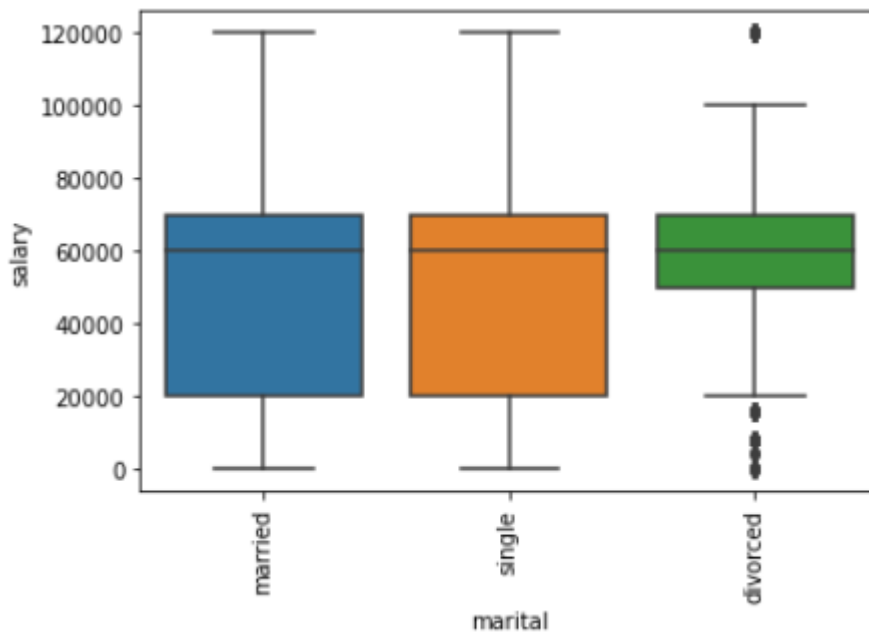
**Figure 10:** Duration Vs Response



**Inferences:** Customers with high Call durations are more likely to subscribe the Term deposit, from the call duration we can understand they are more likely to take the term deposit, they are interested to know about the plans.

## Insights from the Data:

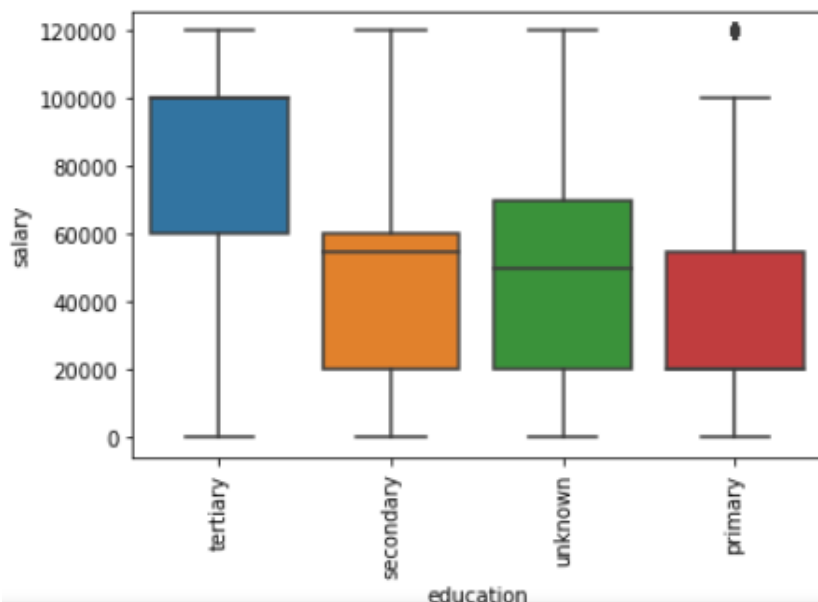
**Figure 11:** Does Salary vary with marital status??



### Inferences:

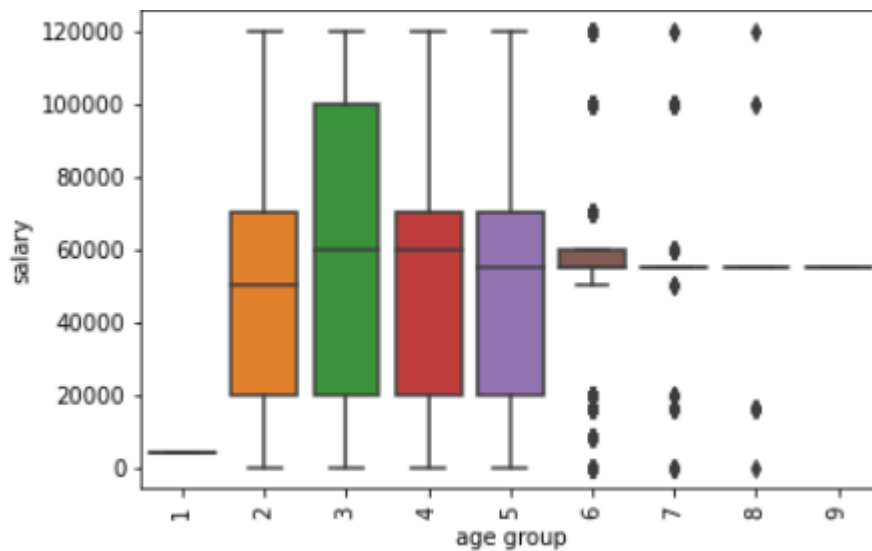
- Married and Divorced Customers tend to have equal salary but higher than that of Divorcees.

**Figure 12:** Does Salary vary with Education??



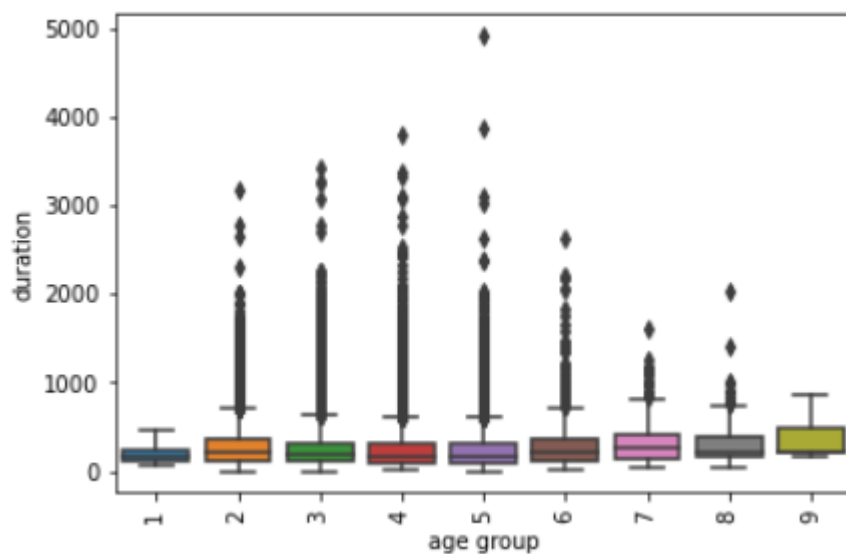
**Inferences:** Higher the Education level, higher the Salary. Clearly, Customers with Tertiary level education earns more, with the tertiary education they work as Lead in the team irrespective of domain, because they have more experience and exposure than others.

**Figure 13:** Do all age groups earn the same Salary??



**Inferences:** Customers with Age group 3 (30~39 years) earn more salary compared to others.

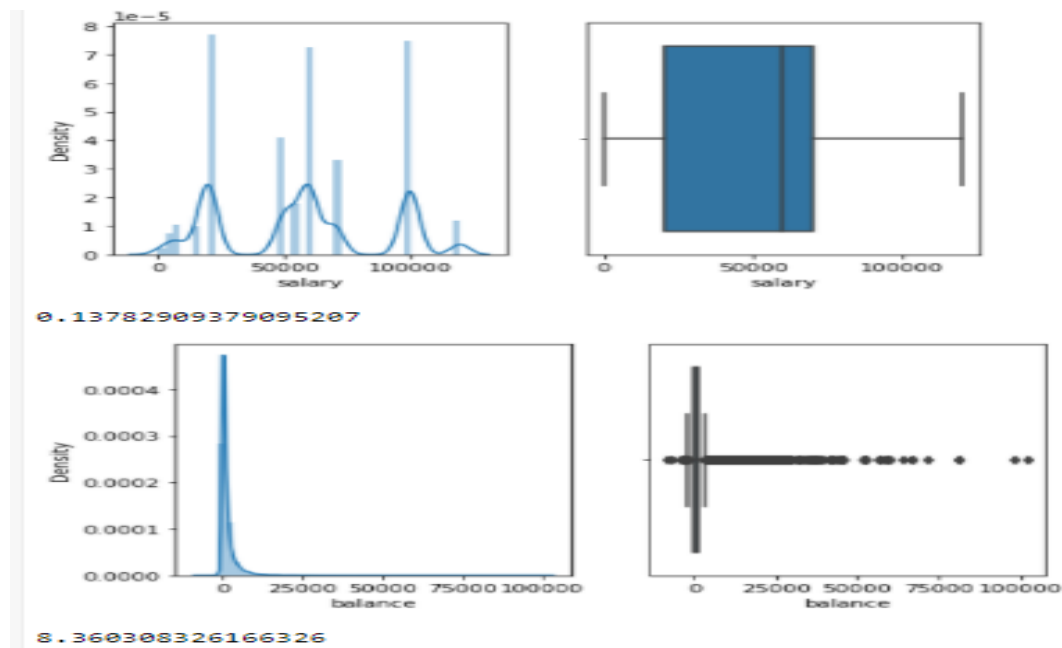
**Figure 14:** Particular Age groups are contacted more??



**Inferences:** Customers with Age group 7 (70~79 years) are contacted more and have higher call durations compared to other age groups.

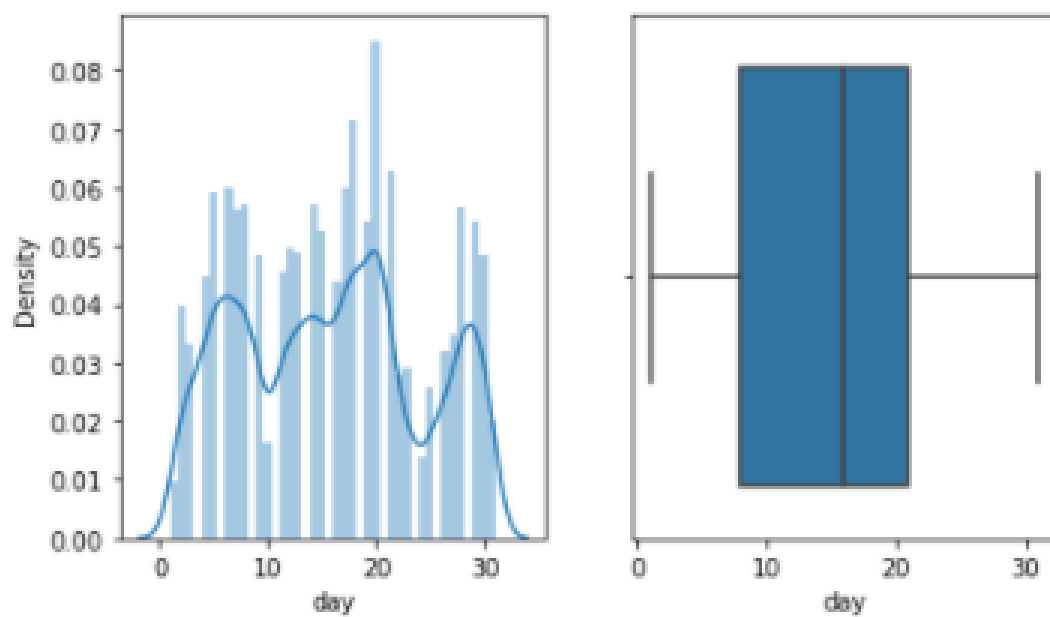
## Distribution of Numerical variables:

**Figure 15:** Distribution of “Salary” & “Balance”

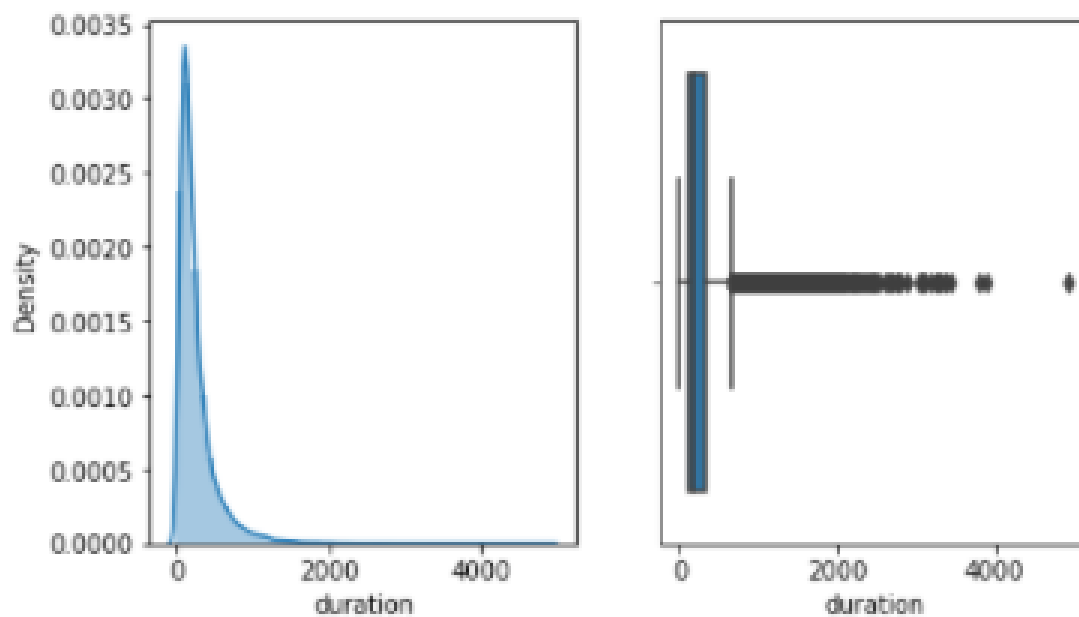


**Inferences:** Salary column is multi-modal with no outliers (skew=0.13) & Salary column is right skewed with outliers (skew=8.36)

**Figure 16:** Distribution of “Day” & “Duration”



0.09307901402122411

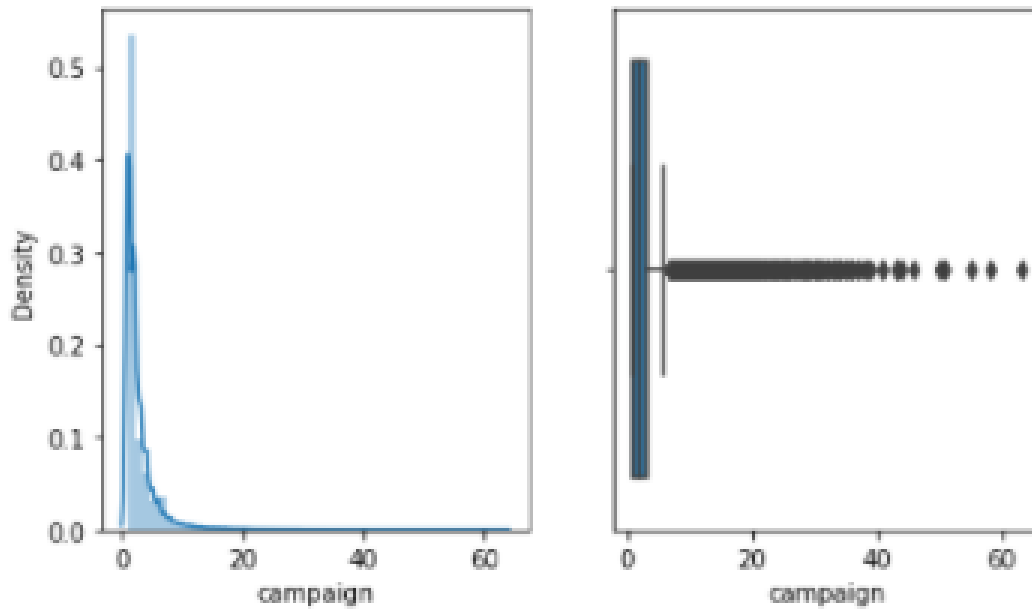


3.144318099423456

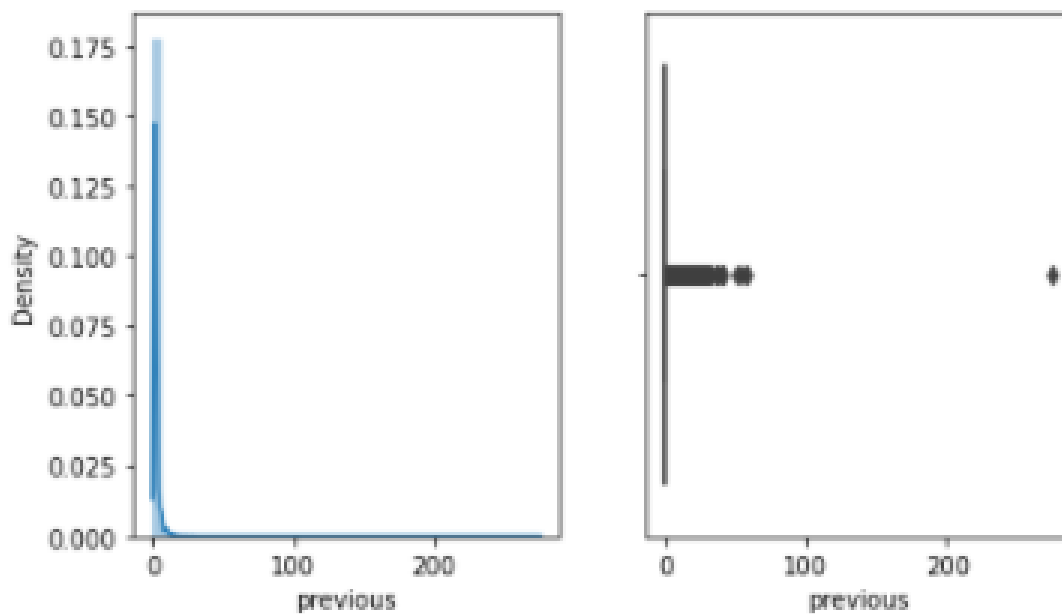
**Inferences:** Day column is multi-modal with no outliers (skew=0.09) & Duration column is right skewed with outliers (skew=3.14)



**Figure 17:** Distribution of “Campaign” & “Previous”



4.898650166179674



41.84645447266292

**Inferences:** Campaign column is right skewed with outliers (skew=4.89) & Previous column is right skewed with outliers (skew=41.84)

## Relationship between variables & target:

**Figure 18:** Correlation between variables

```
1 plt.figure(figsize=(18,7))
2 sns.heatmap(df_n.corr(),annot=True,fmt='.2f');
```

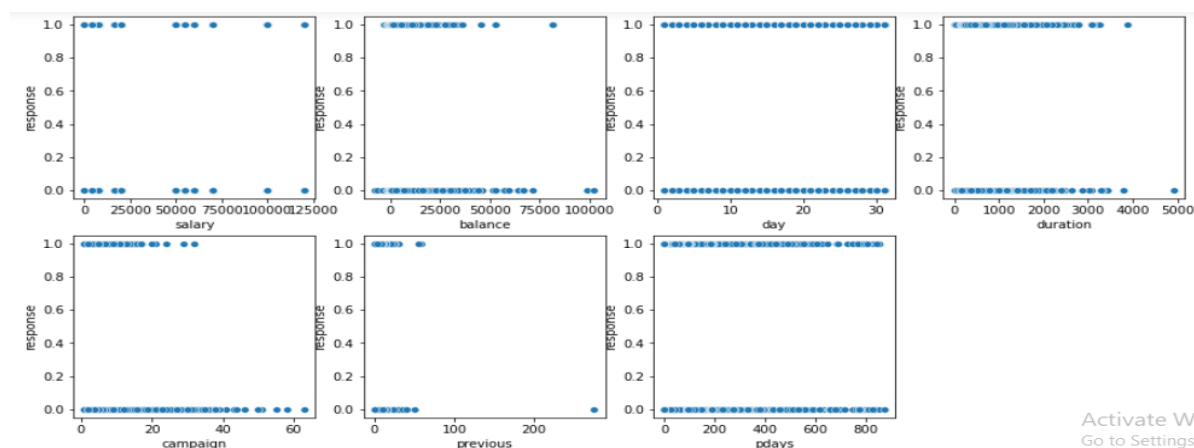


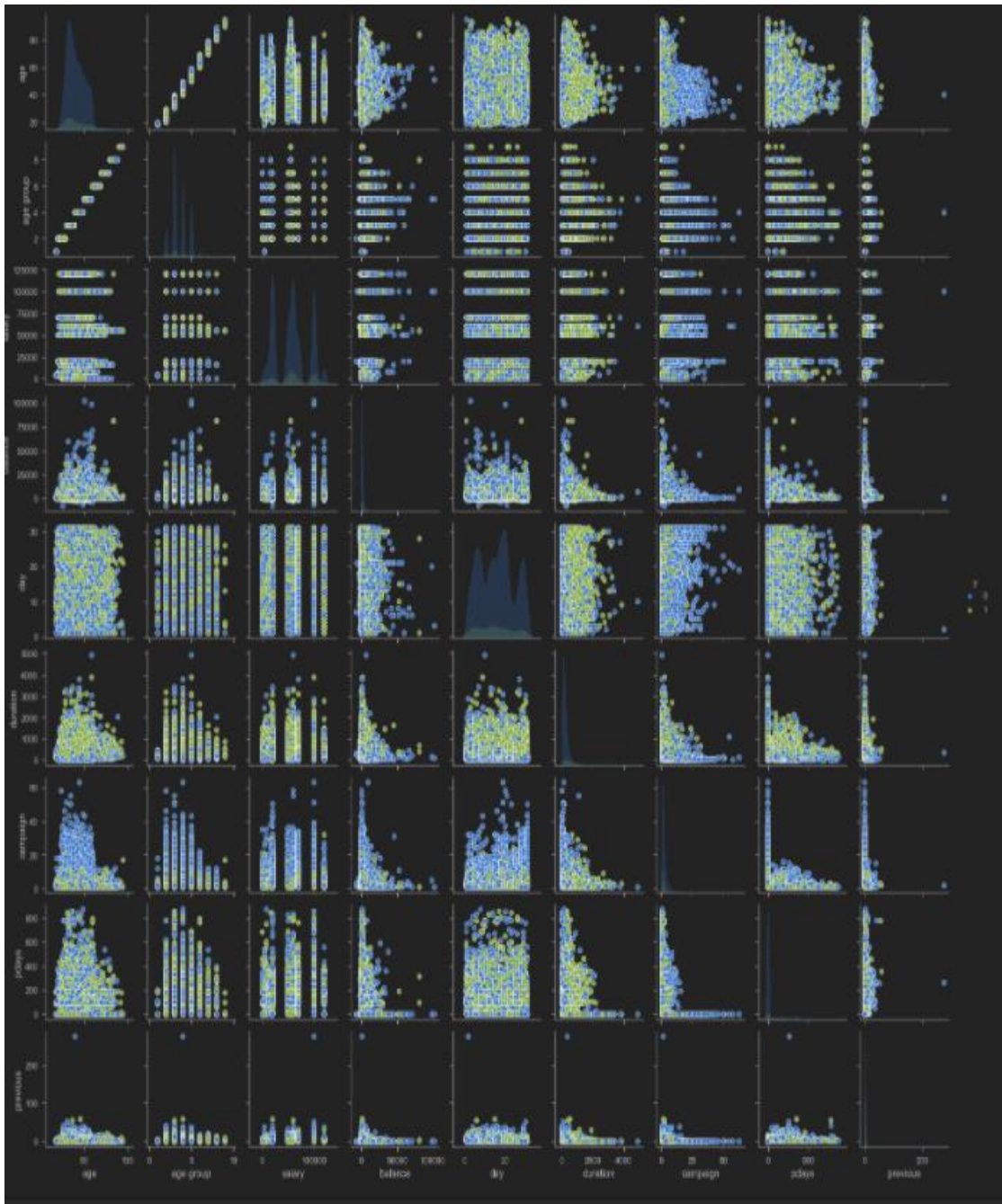
### Inferences:

- Age & Age group are related with each other (corr = 0.96) – High Multicollinearity
- Pdays & Previous are related with each other (corr = 0.45)
- Duration has positive relation with Response (corr=0.40)
- 

## Relationship between Numerical variables & target:

**Figure 19:** Relationship with Target

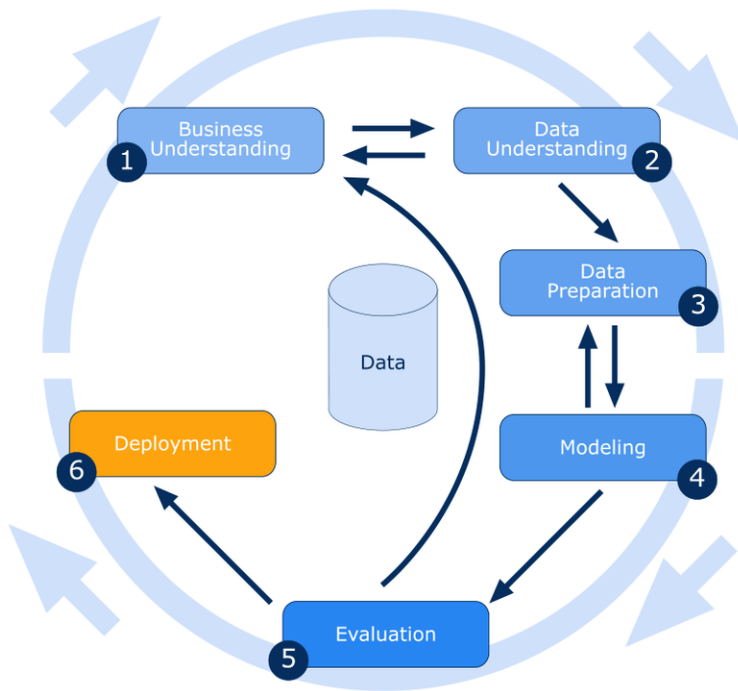




**Inferences:** There is no pattern observed in the Scatterplot between variables and Target. No relation between variables and Target.

## 10. Methodology Followed:

- In this project, we use the CRISP-DM Framework. CRISP-DM has been the most favoured methodology in data mining domain. Therefore, we have chosen it as our reference model. By using CRISP-DM, we can find interesting patterns from within the data that we want to run. Where the data will be processed through the phases of the existing phase of the business understanding phase, understanding data, data preparation, modelling, evaluation and finally deployment. With this phase, it is expected that the results of this study will get the most appropriate modelling in the data mining process, so that the information generated is more efficient. The six phases of CRISP-DM and described briefly as follows:



**Figure 1: Methodology CRISPDM**

## 11. Pre-Processing Steps:

- Data pre-processing is done to enhance the quality of data, to promote extraction of meaningful insights from the data. In simpler words, cleaning and organizing the raw data to make it suitable for building and training machine learning models. It is a data mining technique that transforms raw data into an understandable and readable format.
- Treated 'Unknown' Values in the education column in a logical perspective with job column taken as reference. For each category of the job column, based on the higher frequency in education categories, the unknown values have been replaced. And other null values in the education column have been removed.
- Treated 'Unknown' Values in the contact column with mode replacement.
- Encoded the month column with the 1-12(monthwise) order.

## **12. Algorithms Used:**

- Here, our objective is to predict whether a customer is capable of making the purchase or not.
- So, the classification techniques we used are as follows :
  - **Logistic Regression**
  - **K-Nearest Neighbors Classifier**
  - **Random forest Classifier**
  - **Adaboost Classifier**
  - **Gradient Boost Classifier**
  - **Extreme Gradient Boost Classifier**
  - **Stacking classifier**
  - **Voting classifier**

## **13. Assumptions for different models:**

### **1.) Logistic Regression:**

- The logistic regression assumptions are quite different from OLS regression in that:
  1. There is no need for a linear relationship between the independent and dependent variables.
  2. There is no need for residuals to be normal.
  3. There is no need to meet the homoscedasticity assumption
- So what are the assumptions that need to be met for logistic regression?
- Here are the 5 key assumptions for logistic regression.

### **Assumption 1: Appropriate dependent variable structure**

- This assumption simply states that a binary logistic regression requires your dependent variable to be dichotomous and an ordinal logistic regression requires it to be ordinal.
- In addition, the dependent variable should neither be an interval nor ratio scale.

**Inference:** Here, our Target is binary categorical type.

### **Assumption 2: No Multicollinearity**

- Multicollinearity refers to the high correlation between your independent variables.
- Multicollinearity is a problem because it creates redundant information that will cause the results of your regression model to be unreliable.
- To circumvent this issue, we can deploy two techniques:
- Run a correlation analysis across all your independent variables.
- Remove independent variables with high variance inflation factor(VIF). As a general rule of thumb a  $vif > 10$  is a strong indication of multicollinearity
- $VIF = 1/(1-R^2)$

**Inference:** Age and Age-group had multicollinearity, removed age feature to overcome the multicollinearity.

### **Assumption 3: No Influential Outliers**

- Influential outliers are extreme data points that affect the quality of the logistic regression model.
- Not all outliers are influential.
- We need to check for which points are the influential ones before removing or transforming them for analysis.
- We have checked for it while checking for outliers in EDA section and found some of the independent variables has influential outliers so removed them using IQR method.

**Inference:** Had outliers in 'salary','balance','duration','campaign' variables. Treated them using IQR method.

### **Assumption 4: Observation Independence**

- This assumption requires logistic regression observations to be independent of each other.
- That is, observations should not come from a repeated measure design.
- A repeated measure design refers to multiple measures of the same variable taken for the same person under different experimental conditions or across time. A good example of repeated measures is longitudinal studies — tracking progress of a subject over years.

**Inference:** There are no repeated measures in our feature variables

## **2.) Tree Based models:**

- For tree-based models such as *Decision Trees*, *Random Forest* & *Gradient Boosting* there are no model assumptions to validate.
- Unlike OLS regression or logistic regression, tree-based models are robust to outliers and do not require the dependent variables to meet any normality assumptions.

## **3.) KNN:**

- KNN is a non-parametric lazy learning algorithm. When you say a technique is non-parametric, it means that it does not make any assumptions on the underlying data distribution.

## 14. Step-by-Step Walkthrough of the solution:

### Step 1-Import the Dataset into Jupiter Workbook:

- Before we import our sample dataset into the notebook, we will import the panda's library. pandas is an open-source Python library that provides “high-performance, easy-to-use data structures and data analysis tools.”

```
df = pd.read_csv('bank-marketing.csv')
df.head()
```

	age	age group	eligible	job	salary	marital	education	marital-education	targeted	default	balance	housing	loan	contact	day	month	duration	campaign
0	58	5	Y	management	100000	married	tertiary	married-tertiary	yes	no	2143	yes	no	unknown	5	may	261	
1	44	4	Y	technician	60000	single	secondary	single-secondary	yes	no	29	yes	no	unknown	5	may	151	
2	33	3	Y	entrepreneur	120000	married	secondary	married-secondary	yes	no	2	yes	yes	unknown	5	may	76	
3	47	4	Y	blue-collar	20000	married	unknown	married-unknown	no	no	1506	yes	no	unknown	5	may	92	
4	33	3	Y	unknown	0	single	unknown	single-unknown	no	no	1	no	no	unknown	5	may	198	

### Step 2-Explore the data set:

#### Shape:

- shape is a tuple that gives you an indication of the number of dimensions in the array.
- We can get the shape(Total rows & Total columns) of the data using .shape.

```
df.shape
```

```
(45211, 23)
```

#### Describe:

```
df.describe()
```

	age	age group	salary	balance	day	duration	campaign	pdays	previous
count	45211.000000	45211.000000	45211.000000	45211.000000	45211.000000	45211.000000	45211.000000	45211.000000	45211.000000
mean	40.936210	3.645861	57006.171065	1362.272058	15.806419	258.163080	2.763841	40.197828	0.580323
std	10.618762	1.083271	32085.718415	3044.765829	8.322476	257.527812	3.098021	100.128746	2.303441
min	18.000000	1.000000	0.000000	-8019.000000	1.000000	0.000000	1.000000	-1.000000	0.000000
25%	33.000000	3.000000	20000.000000	72.000000	8.000000	103.000000	1.000000	-1.000000	0.000000
50%	39.000000	3.000000	60000.000000	448.000000	16.000000	180.000000	2.000000	-1.000000	0.000000
75%	48.000000	4.000000	70000.000000	1428.000000	21.000000	319.000000	3.000000	-1.000000	0.000000
max	95.000000	9.000000	120000.000000	102127.000000	31.000000	4918.000000	63.000000	871.000000	275.000000

- Age of the Customers start from 18 to 95.
- Salary of the Customers range from 0 to 1,20,000.
- Minimum balance by a Customer is -8019 and max 102127.
- There are 1 to 31 days.
- Duration of calls lasts from 0 to 4918 seconds (82 Minutes or 1.3 Hrs).
- Campaign (No of Calls performed "during this campaign") ranges from 1 to 63 calls.
- Pdays(No of days passed after the Customer was last contacted from a previous campaign) ranges from -1 to 871 days ---> -1 means the Customer is not contacted.
- Previous (number of contacts performed "before this campaign") ranges from 0 to 275 calls.
- Response is the Target Variable (1-Subscribed a Term deposit, 0 - Not subscribed a Term deposit).

### Null Values:

- `isnull().sum()` returns the number of missing values in the data set. A simple way to deal with data containing missing values is to skip rows with missing values in the dataset.

```
df.isnull().sum()
```

```
age          0
age group    0
eligible     0
job          0
salary       0
marital      0
education    0
targeted     0
default      0
balance      0
housing      0
loan         0
contact      0
day          0
month        0
duration     0
campaign     0
pdays       0
previous     0
poutcome     0
y            0
dtype: int64
```

- We haven't found any null values in the given data, to know whether we have any unknown values present in the given data we can use `.unique()` .



## Unique values:

- `unique()` Function to Get Unique Values from a Dataframe. The **`.unique()`** function returns the unique values present in a dataset. It basically uses a technique based on hash tables to return the non-redundant values from the set of values present in the data frame/series data structure.

```
for i in df.select_dtypes(include=np.object).columns:  
    print(df[i].unique())  
    print()
```

```
['Y' 'N']
```

```
['management' 'technician' 'entrepreneur' 'blue-collar' 'unknown'  
'retired' 'admin.' 'services' 'self-employed' 'unemployed' 'housemaid'  
'student']
```

```
['married' 'single' 'divorced']
```

```
['tertiary' 'secondary' 'unknown' 'primary']
```

```
['married-tertiary' 'single-secondary' 'married-secondary'  
'married-unknown' 'single-unknown' 'single-tertiary' 'divorced-tertiary'  
'married-primary' 'divorced-secondary' 'single-primary'  
'divorced-primary' 'divorced-unknown']
```

```
['yes' 'no']
```

```
['no' 'yes']
```

```
['yes' 'no']
```

```
['no' 'yes']
```

```
['unknown' 'cellular' 'telephone']
```

```
['may' 'jun' 'jul' 'aug' 'oct' 'nov' 'dec' 'jan' 'feb' 'mar' 'apr' 'sep']
```

```
['unknown' 'failure' 'other' 'success']
```

```
['no' 'yes']
```

- Using `.unique()` we came to know that we have some unknown values in the data, as we have unknown values we are changing the unknown values to NaN values.

### Step 3-Treating 'Unknown' Values:

```
ind = df[(df['job']=='unknown') & (df['education']=='primary')]['job'].index
df.iloc[ind,3] = 'blue-collar'
```

```
ind = df[(df['job']=='unknown') & (df['education']=='secondary')]['job'].index
df.iloc[ind,3] = 'blue-collar'
```

```
ind = df[(df['job']=='unknown') & (df['education']=='tertiary')]['job'].index
df.iloc[ind,3] = 'management'
```

- As we have unknown values in the education column, based on each category in education column the unknown values have been imputed.
- Contact column has unknown values and those unknown values are imputed with mode.

### Step 4- Dropping unwanted columns:

```
df.drop('marital-education',1,inplace=True)
```

```
df.drop('response',1,inplace=True)
```

```
df.drop('age',1,inplace=True)
```

```
df.drop('poutcome',1,inplace=True)
```

- As we have education and marital as a different columns we can remove marital-education column.
- As we have age group column we can remove age column.
- As we have y column with yes and no we can remove response column.
- Poutcome column has null values more than 80% we can remove that column.

### Step 5-Feature engineering:

- It is the process of using domain knowledge of the data to create new features that make the machine learning model perform better.
- Feature engineering is the essential art in machine learning, which creates a massive difference between a good model and a bad model.

```
def pdays(x):
    if (x<=0):
        return 'Not.Previously.Contacted'
    elif (x>0 and x<=150):
        return '1-150 days'
    elif (x>150 and x<=300):
        return '151-300 days'
    else:
        return '>300 days'
```

```
df['pdays'] = df['pdays'].apply(pdays)
```

- Pdays columns contains values such as -1,0 or more than 500, to decrease complexity of the column, we are converting them to categorical column with four categories.
- Pdays <=0 Not previously connected.

### Step 6-Outlier treatment:

- The difference between Q3 and Q1 is called the Inter-Quartile Range or IQR. Any data point less than the Lower Bound or more than the Upper Bound is considered as an outlier.

```
for i in num_out:
    q1 = df[i].quantile(0.25)
    q3 = df[i].quantile(0.75)
    iqr = q3 - q1
    ll = q1 - (3*iqr)
    ul = q3 + (3*iqr)
    df[i] = df[(df[i]>=ll) & (df[i]<=ul)][i]
```

### Step 7- Statistical Test:

- The Chi-square test of independence is a statistical hypothesis test used to determine whether two categorical or nominal variables are likely to be related or not.

**NOTE:** Significance value of 0.05 is considered for Statistical testing

```
# Chi-sqr Test of Independence
# Hypothesis Formation
# Ho : Variables are Independent (NO relation)
# Ha : Variables are Not independent (Relation)

def chi(obs):
    chi_stat,pval,df,exp_tab = stats.chi2_contingency(obs)
    return pval

not_sig_features = []
sig_features = []

for i in cat_col:
    obs = pd.crosstab(df[i],df['y'])
    pval = chi(obs)
    if (pval > 0.05):
        not_sig_features.append(i) # Accept H0
    else:
        sig_features.append(i) # Reject Ho
```

sig\_features

```
['loan',
 'job',
 'default',
 'day',
 'month',
 'age group',
 'housing',
 'marital',
 'pdays',
 'contact',
 'targeted',
 'education',
 'eligible']
```

- As we have target variable as categorical with 2 classes, so we can use chi-square test of independence, with that we have found the significant features.

### Step 8-Multi Collinearity Check:

- Variance inflation factor (VIF) is a measure of the amount of multicollinearity in a set of multiple regression variables.

```
v = df[num_col]
vif = [VIF(v.values,i) for i in range(v.shape[1])]
vif_df = pd.DataFrame()
vif_df['numeric_features'] = v.columns
vif_df['VIF'] = vif
vif_df.sort_values('VIF',ascending=False)
```

	numeric_features	VIF
0	age	6.379388
3	day	4.102940
1	salary	3.717447
4	duration	1.911755
5	campaign	1.831290
6	pdays	1.455194
7	previous	1.342436
2	balance	1.216352

- To check whether multi-collinearity is present or not we have used VIF and we have found that all the columns are significant features.

### Step 9- Encoding:

- We use this categorical data encoding technique when the categorical feature is ordinal. In this case, retaining the order is important. Hence encoding should reflect the sequence.
- In Label encoding, each label is converted into an integer value. We will create a variable that contains the categories representation.

```
le = LabelEncoder()
```

```
for i in cat_le:
    df[i] = le.fit_transform(df[i])
```

- For all the categorical columns encoding has been done using Label Encoding method.

## Step 10-Train Test Split:

- The train-test split is a technique for evaluating the performance of a machine learning algorithm.
- It can be used for classification or regression problems and can be used for any supervised learning algorithm.
  - **Train Dataset:** Used to fit the machine learning model.
  - **Test Dataset:** Used to evaluate the fit machine learning model.

```
x = df.drop('y',1)
y = df['y']

x_sm = df_sm.drop('y',1)
y_sm = df_sm['y']
```

```
x_train,x_test,y_train,y_test = train_test_split(x,y,test_size=0.3,random_state=10)
x_sm_train,x_sm_test,y_sm_train,y_sm_test = train_test_split(x_sm,y_sm,test_size=0.3,random_state=10)
```

- Train test split has been done with the test size of 30%.

## Step 11-Scaling:

- In many machine learning algorithms, to bring all features in the same standing, we need to do scaling so that one significant number doesn't impact the model just because of their large magnitude.

```
ss = StandardScaler()
```

```
sc = ['salary','balance','duration','campaign','previous']
```

```
# Only used for Cross Validation Score
```

```
x_scaled = x.copy(deep=True)
x_sm_scaled = x_sm.copy(deep=True)

x_scaled[sc] = ss.fit_transform(x_scaled[sc])
x_sm_scaled[sc] = ss.fit_transform(x_sm_scaled[sc])
```

```
# For Model Building
```

```
x_train[sc] = ss.fit_transform(x_train[sc])
x_test[sc] = ss.fit_transform(x_test[sc])

x_sm_train[sc] = ss.fit_transform(x_sm_train[sc])
x_sm_test[sc] = ss.fit_transform(x_sm_test[sc])
```

- We have done standard scaler for the data.

## Step 12-Class Imbalance:

- Imbalanced classification refers to a classification predictive modelling problem where the number of examples in the training dataset for each class label is not balanced.

```
x = df.drop('y',1)
y = df['y']
```

```
smote = SMOTE(sampling_strategy=0.5,random_state=10)
x_sm,y_sm = smote.fit_resample(x,y)
```

```
df_sm = pd.DataFrame(x_sm,columns=x.columns)
df_sm['y']=y_sm
df_sm.head()
```

- SMOTE is used to oversample the minority class to balance the class distribution.

## Step 13-Model Building:

- A machine learning model is built by learning and generalizing from training data, then applying that acquired knowledge to new data it has never seen before to make predictions and fulfil its purpose.
- We have built following models,

### 1. Base model - Logistic Regression:

```
print(classification_report(y_test,y_test_pred))
```

	precision	recall	f1-score	support
0	0.91	0.98	0.94	10659
1	0.56	0.21	0.30	1273
accuracy			0.90	11932
macro avg	0.73	0.59	0.62	11932
weighted avg	0.87	0.90	0.88	11932

## 2. Knn Classifier:

```
print(classification_report(y_sm_test,y_sm_test_pred))
```

	precision	recall	f1-score	support
0	0.90	0.92	0.91	10672
1	0.83	0.79	0.81	5371
accuracy			0.88	16043
macro avg	0.87	0.86	0.86	16043
weighted avg	0.88	0.88	0.88	16043

## 3. Random Forest Classifier:

```
print(classification_report(y_sm_test,y_sm_test_pred))
```

	precision	recall	f1-score	support
0	0.94	0.86	0.90	10672
1	0.77	0.90	0.83	5371
accuracy			0.88	16043
macro avg	0.86	0.88	0.86	16043
weighted avg	0.88	0.88	0.88	16043

## 4. Ada Boosting:

```
print(classification_report(y_sm_test,y_sm_test_pred))
```

	precision	recall	f1-score	support
0	0.92	0.70	0.80	10672
1	0.60	0.87	0.71	5371
accuracy			0.76	16043
macro avg	0.76	0.79	0.75	16043
weighted avg	0.81	0.76	0.77	16043

## 5. Gradient Boosting:

```
print(classification_report(y_sm_test,y_sm_test_pred))
```

	precision	recall	f1-score	support
0	0.93	0.79	0.85	10672
1	0.68	0.89	0.77	5371
accuracy			0.82	16043
macro avg	0.80	0.84	0.81	16043
weighted avg	0.85	0.82	0.83	16043

## 6. XGB (Hyperparameter tuning):

```
print(classification_report(y_sm_test,y_sm_test_pred))
```

	precision	recall	f1-score	support
0	0.96	0.77	0.86	10672
1	0.68	0.93	0.78	5371
accuracy			0.83	16043
macro avg	0.82	0.85	0.82	16043
weighted avg	0.86	0.83	0.83	16043

## 7. Stacking Classifier:

```
print(classification_report(y_sm_test,y_sm_test_pred))
```

	precision	recall	f1-score	support
0	0.94	0.82	0.88	10672
1	0.71	0.90	0.80	5371
accuracy			0.85	16043
macro avg	0.83	0.86	0.84	16043
weighted avg	0.86	0.85	0.85	16043

## 8. Voting Classifier:

```
print(classification_report(y_sm_test,y_sm_test_pred))
```

	precision	recall	f1-score	support
0	0.95	0.81	0.87	10672
1	0.71	0.92	0.80	5371
accuracy			0.85	16043
macro avg	0.83	0.86	0.84	16043
weighted avg	0.87	0.85	0.85	16043



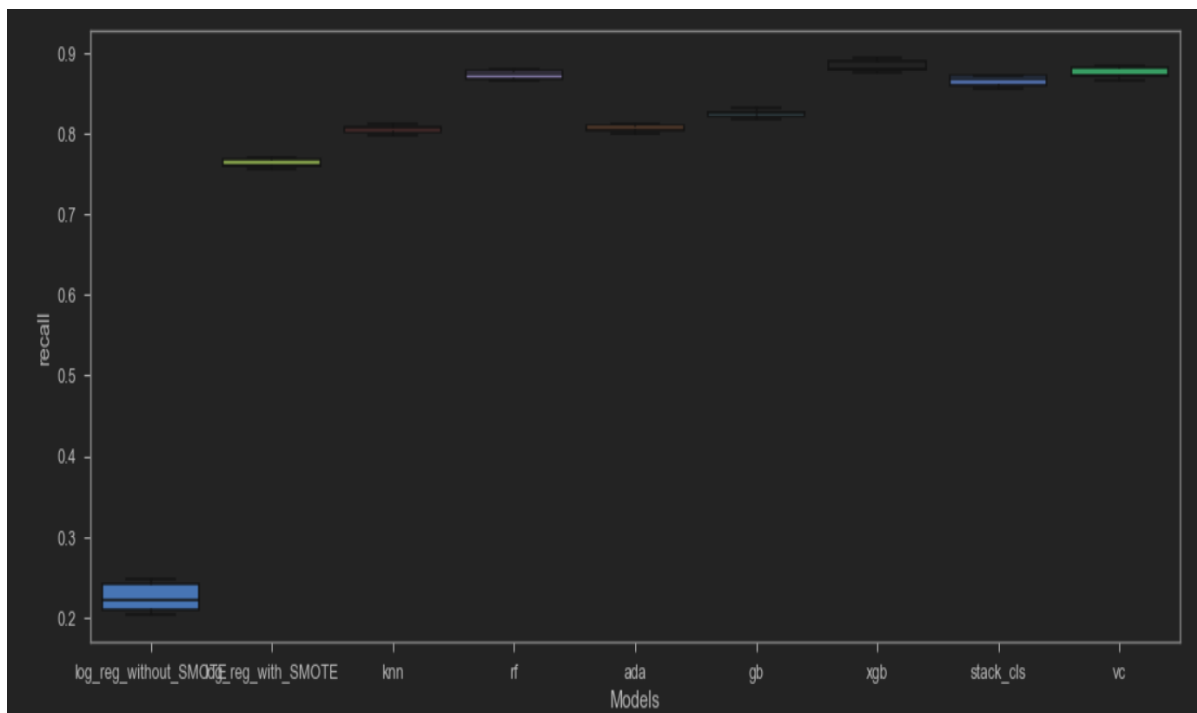
## 15. Data Modelling and Evaluation:

- The basic objective of the project is to analyse Tele-Marketing data collected from a Bank and predict whether the customer will subscribe the term deposit or not, also identify the driving factors behind this. This would inform the Bank's decisions on which Customers to target for their Marketing Campaign, which would ultimately increase their Product Sales. Understanding their customers is critical for effectively executing their Marketing campaign.
- As a part of data modelling, we are building a model which can whether the customer will subscribe the term deposit or not. As classification aims at categorizing the target data to which category it belongs, we have tried various base line models and ensemble models used for data modelling. We are interested in finding the customers who will subscribe the term deposit, so false negative is to be minimized and minority class in the dataset is to be balanced.
- Following Performance measures are to be used for our models:
  - Recall  $[TP/(TP+FN)]$  – Sensitivity of the model
  - F1 Weighted
  - ROC\_AUC
- Initially the dataset is split into predictors and dependent variable followed by dividing data into two groups, one for training the model and one for validating the model performance. The data is divided in different proportion of training and test data (that is, 70% data for training and 30% data for testing), In order to achieve higher accuracy, sampling techniques (over sampling - SMOTE) has been used (Since, our data set is imbalanced). The splitted training data is further fed into various classification algorithms. Using Cross-Validated score of recall, f1\_weighted and roc\_auc along with bias\_error and variance\_error of several models , we can choose the best model which gives higher Sensitivity (recall\_score) with low bias and variance error.

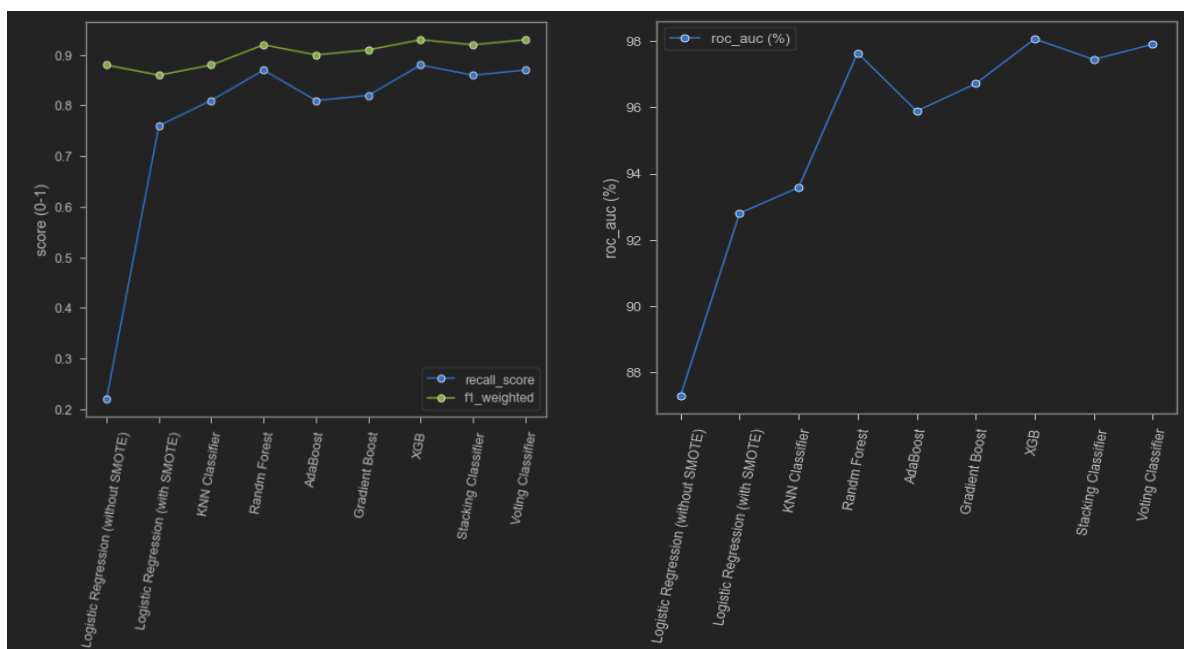
The following table shows the cross-validated results of all the models, we have used.

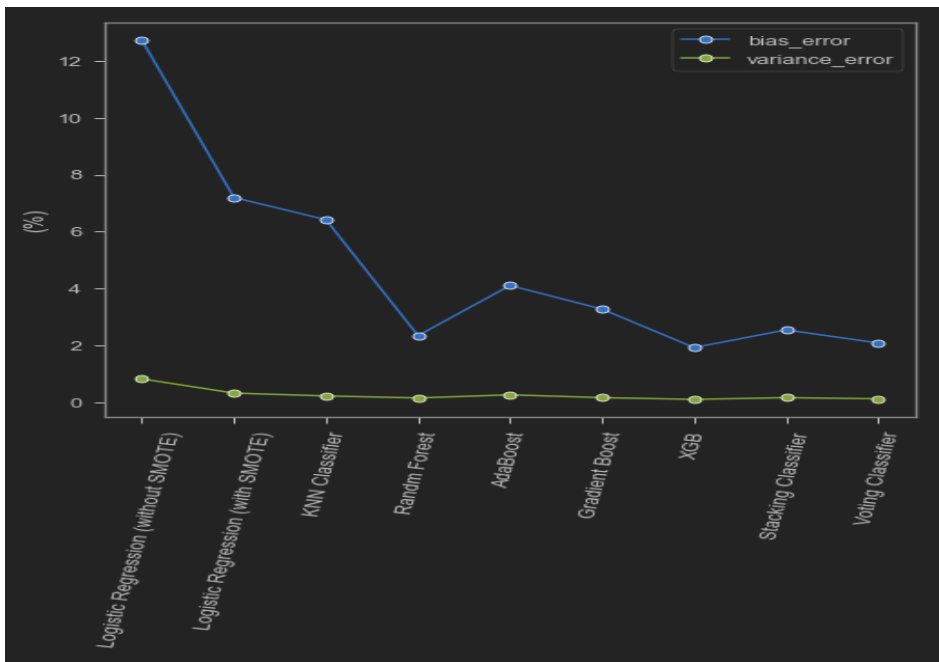
	recall_score	f1_weighted	roc_auc (%)	bias_error %(roc_auc)	variance_error %(roc_auc)
Logistic Regression (without SMOTE)	0.22	0.88	87.28	12.72	0.83
Logistic Regression (with SMOTE)	0.76	0.86	92.80	7.20	0.33
KNN Classifier	0.81	0.88	93.58	6.42	0.23
Randm Forest	0.87	0.92	97.65	2.35	0.16
AdaBoost	0.81	0.90	95.89	4.11	0.27
Gradient Boost	0.82	0.91	96.72	3.28	0.17
XGB	0.88	0.93	98.07	1.93	0.11
Stacking Classifier	0.86	0.92	97.45	2.55	0.17
Voting Classifier	0.87	0.93	97.91	2.09	0.13

Now, we can visualize these results as follows:



- From this plot (Boxplot of 5-fold cross-validated recall\_score for each models), we can clearly tell that without SMOTE the model performance was very poor.
- Now compare other metrics for each models and then conclude the final model as follows.





- From these plots, we can find that XGB model's performance is better when compared to other models.

For XGB model we attain the following metrics as:

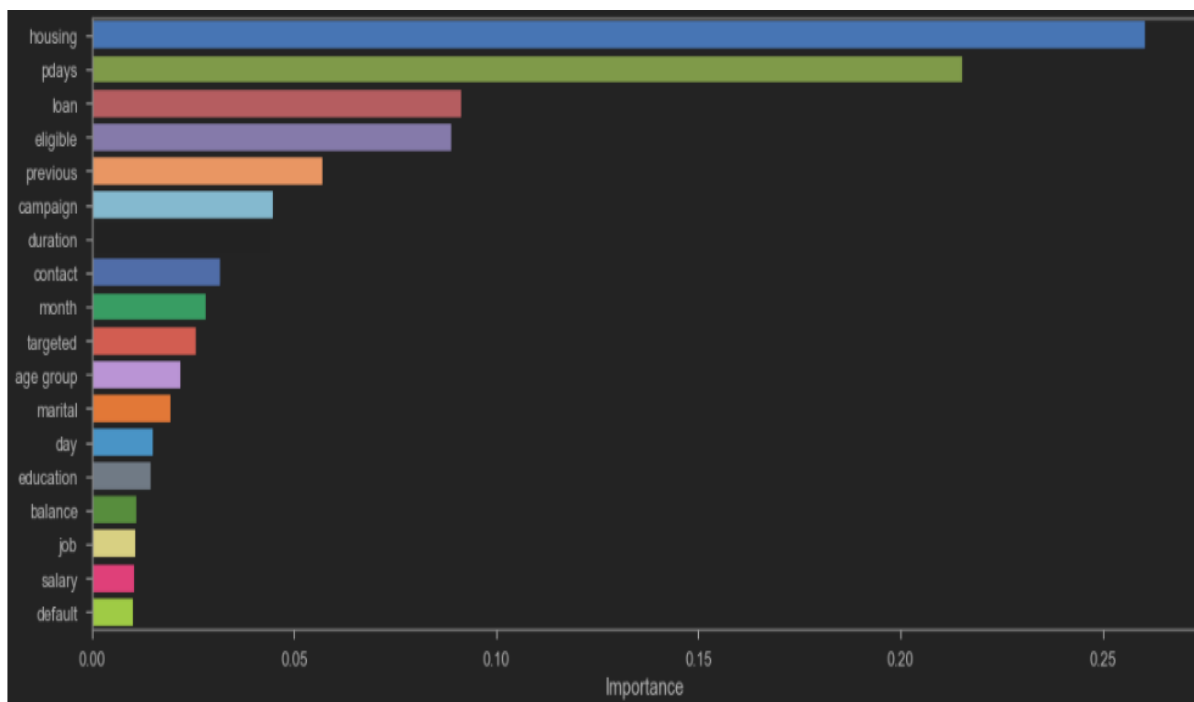
- recall\_score = 0.88
  - f1\_weighted = 0.93
  - roc\_auc (%) = 98.07 %
  - bias\_error (%) = 1.93 %
  - variance\_error (%) = 0.11%
- Our Final Model is XGB Classifier, which gives us a better results when compared to others. XGBoost uses decision trees as base learners, combining many weak learners to make a strong learner. As a result it is referred to as an ensemble learning method since it uses the output of many models in the final prediction. XGBoost or Extreme Gradient Boosting! It can be put into various use cases such as ranking, classification, regression and user-defined prediction problems. It can be referred to as an "ALL in One" algorithm. It is an ideal blend of software and hardware optimization techniques to yield prevalent outcomes by using fewer computing resources in shortest amount of time.
  - There are many advantages of XGBoost, some of them are mentioned below:
    - It is Highly Flexible
    - It uses the power of parallel processing
    - It is faster than Gradient Boosting
    - It supports regularization
    - It is designed to handle missing data with its in-build features.
    - The user can run a cross-validation after each iteration.
    - It Works well in small to medium dataset

## 16. Feature Importance:

From the final model, we have got the salient features are as follows:

Importance	
housing	0.260568
pdays	0.215173
loan	0.091399
eligible	0.088875
previous	0.056874
campaign	0.044604
duration	0.043998
contact	0.031698
month	0.028142
targeted	0.025629
age group	0.021936
marital	0.019317
day	0.015114
education	0.014506
balance	0.010980
job	0.010689
salary	0.010309
default	0.010190

- The below graph shows the import features from our best model (XG-Boost). We can see that the top five key drivers of whether someone will “subscribe a Term deposit” are influenced by **Housing**, **Pdays**, **Loan**, **Eligible** & **Previous**.



## **17. Recommendations:**

- Based on the key influencers, following recommendations are suggested to enhance our Business motive **“Identifying the right Customers and Increase the odds of subscribing our Term Deposit”**.
  1. Customers who do not prefer **“housing loans”** are more likely to subscribe our Term deposit. Thus, we can offer benefitting schemes coupled with housing loans to make them subscribe two of our products.
  2. Customers with less **“pdays”** (Frequently contacted from the previous campaign) will definitely know about the latest Term deposit schemes as they are being contacted more frequently. Minimising the pdays is essential to make customers subscribe our Term deposit. This can be achieved by frequent reminders and calls to the customers to make them aware of our latest products.
  3. Customers who do not prefer any **“loans”** are more likely to subscribe our Term deposit. We can use customer details, identify their background and then recommend them suitable loans coupled with Term deposit to make them subscribe two of our products.
  4. **“Eligible”** Customers must be focussed and contacted more frequently. More Schemes and products can be developed targeting Eligible customers for maximum benefit.
  5. Customers with more **“previous”** (No of contacts performed) will definitely know about the latest Term deposit schemes as they are being contacted more frequently. Maximising the “previous” is essential to make customers subscribe our Term deposit. This can be achieved by frequent reminders and calls to the customers to make them aware of our latest products.

## **18. Limitations:**

- The Final Model was chosen to be Extreme Gradient Boosting Model by considering all the model performance metrics. Our final Model was tuned with only ‘max-depth’.
- This final model was not a generalized model but this was business oriented model that is specifically designed to be used for targeting the customers who would subscribe to the term deposit. This Model will not be used for all the market analysis problems.
- We have used ‘recall’ as the important metric to determine the model performance which decreases the false negatives in the data specifically to correctly classify the customers who are possible to subscribe the product but were wrongly classified as not possible for subscription.
- This may not be applied in all the Marketing Analysis Problems. So, it only applies in targeting the customers.

## **19. Conclusion:**

- The best model chosen from the previous section is XGBOOST. Hence it can be taken into production where a different set of test data will be used and the classification can be performed with the obtained confidence as told by the model's metrics.
- Hence we can expect our model to classify the customers who are about to take the term deposit with the confidence of 88%.
- From the analysis, it allows the bank to better anticipate and address the potential customers, while improving their strategic marketing campaigns.

## **20. References:**

- Predicting the Success of Bank Telemarketing using various Classification Algorithms  
<https://www.diva-portal.org/smash/get/diva2:1233529/FULLTEXT01.pdf>
- Customer Profiling using classification approach for bank telemarketing  
<http://www.joiv.org/index.php/joiv/article/view/68>
- Analysis of bank customers and prediction of bank marketing strategies  
<http://serisc.org/journals/index.php/IJAST/article/view/15497/7810>