

BANK MARKETING ANALYSIS

A report submitted in partial fulfillment of the requirements for the Award of Degree of

BACHELOR OF TECHNOLOGY

in

INFORMATION TECHNOLOGY

By

PENMETSA TEJA VISHNU VARMA

Regd No.: 20B91A12F4

Under Supervision of

Mr. Gundala Nagaraju

Henotic Technology Pvt Ltd, Hyderabad

**(Duration: 7th July, 2022 to 6th September,
2022)**



**DEPARTMENT OF INFORMATION TECHNOLOGY
SAGI RAMA KRISHNAM RAJUENGINEERING COLLEGE**

(An Autonomous Institution)

Approved by AICTE, NEW DELHI and Affiliated to JNTUK, Kakinada

CHINNA AMIRAM, BHIMAVARAM,

ANDHRA PRADESH

SAGI RAMA KRISHNAM RAJU ENGINEERING COLLEGE
(Autonomous)

Chinna Amiram, Bhimavaram

**DEPARTMENT OF INFORMATION
TECHNOLOGY**



CERTIFICATE

This is to certify that the Summer Internship Report “BANK MARKETING ANALYSIS” submitted by PENMETSA TEJA VISHNU VARMAVARMA bearing Register Number 20B91A12F4 at the end of second year second semester at Henotic Technology Pvt Ltd., Hyderabad from 07.07.2022 to 06.09.2022 in partial fulfillment of the requirements for the award of the degree of Bachelor of Technology in Information Technology

Department Internship Coordinator

Dean-T & P Cell

Head of the Department

Table of Contents

S. No	Name	Page.No
1.	Introduction	2
1.1	What are the different types of Machine Learning?	3
1.2	Benefits of Using Machine Learning in Bank Marketing	8
1.3	About Industry (Bank Marketing)	9
1.4	AI / ML Role in Bank Marketing	10
2.	Bank Marketing Dataset content	10
2.1	Main Drivers for AI Bank Marketing	10
2.2	Internship Project - Data Link	11
3.	AI / ML Modelling and Results	12
3.1	Your Problem of Statement	12
3.2	Data Science Project Life Cycle	13
3.2.1	Data Exploratory Analysis	13
3.2.2	Data Pre-processing	14
3.2.2.1	Check the Duplicate and low variation data	14
3.2.2.2	Identify and address the missing variables	14
3.2.2.3	Handling of Outliers	16
3.2.2.4	Categorical data and Encoding Techniques	17
3.2.2.5	Feature Scaling	18
3.2.3	Selection of Dependent and Independent variables	19
3.2.4	Data Sampling Methods	19
3.2.4.1	Stratified sampling	19
3.2.4.2	Simple random sampling	19
3.2.5	Models Used for Development	20
3.2.5.1	Model 01	20
3.2.5.2	Model 02	21
3.2.5.3	Model 03	21
3.2.5.4	Model 04	22
3.2.5.5	Model 05	
3.2.5.6	Model 06	
3.2.5.7	Model 07	
3.2.5.8	Model 08	
3.2.5.9	Model 09	
3.2.5.10	Model 10	22
3.3	AI / ML Models Analysis and Final Results	22
3.3.1	Different Model codes	23
3.3.2	Random Forest Python Code	23
3.3.3	Extra Trees Python code	24
4.	Conclusions and Future work	24
5.	References	24
6.	Appendices	25

6.1 Python code Results	30
6.2 List of Charts	
6.2.1 Chart 01: Education	
6.2.2 Chart 02: Duration	
6.2.3 Chart 03: Marital status	
6.2.4 Chart 04: Heat Map	32

Abstract

The data is related with direct marketing campaigns of a Portuguese banking institution. The marketing campaigns were based on phone calls. Often, more than one contact to the same client was required, in order to access if the product (bank term deposit) would be ('yes') or not ('no') subscribed.

A key challenge for the Bank Marketing is to make Deposits in the Bank. The goal of this program is to see how well various statistical methods perform in predicting Bank Deposit Prediction based on the different characteristics .I am trying to increase the accuracy of Bank Marketing prediction.

In a traditional bank marketing campaign, a bank will reduce a lot of cost if they know what kind of customers are likely to subscribe to their products.As a result, it's highly recommended that banks try their best to increase the duration of a call. Additionally, it will be a good choice to take action when the basic interest rate is low. Lastly, knowing whether a customer has a record of similar subscription is of great importance.

INTRODUCTION

With the increasing power of computer technology, companies and institutions can nowadays store large amounts of data at reduced cost. The amount of available data is increasing exponentially and cheap disk storage makes it easy to store data that previously was thrown away. There is a huge amount of information locked up in databases that is potentially important but has not yet been explored. The growing size and complexity of the databases makes it hard to analyse the data manually, so it is important to have automated systems to support the process. Hence there is the need of computational tools able to treat these large amounts of data and extract valuable information.

In this context, Data Mining provides automated systems capable of processing large amounts of data that are already present in databases. Data Mining is used to automatically extract important patterns and trends from databases seeking regularities or patterns that can reveal the structure of the data and answer business problems. Data mining includes learning techniques that fall into the field of Machine learning. The growth of databases in recent years brings data mining at the forefront of new business technologies.

A key challenge for the insurance industry is to charge each customer an appropriate price for the risk they represent. Risk varies widely from customer to customer and a deep understanding of different risk factors helps predict the likelihood and cost of insurance claims. The goal of this program is to see how well various statistical methods perform in predicting auto Insurance claims based on the characteristics of the driver, vehicle and driver / vehicle coverage details.

A number of factors will determine BI claims prediction among them a driver's age, past accident history, and domicile, etc. However, this contest focused on the relationship between claims and vehicle characteristics well as other characteristics associated with the auto insurance policies.

1.1. What are the different types of Machine Learning?

How Machine Learning Works?

Machine learning uses two types of techniques:

Supervised learning, which trains a model on known input and output data so that it can predict future outputs, and Unsupervised learning, which finds hidden patterns or intrinsic structures in input data.

Supervised Learning:

Supervised machine learning builds a model that makes predictions based on evidence in the presence of uncertainty. A supervised learning algorithm takes a known set of input data and known responses to the data (output) and trains a model to generate reasonable predictions for the response to new data. Use supervised learning if you have known data for the output you are trying to predict. Supervised learning uses Regression and Classification techniques to develop predictive models.

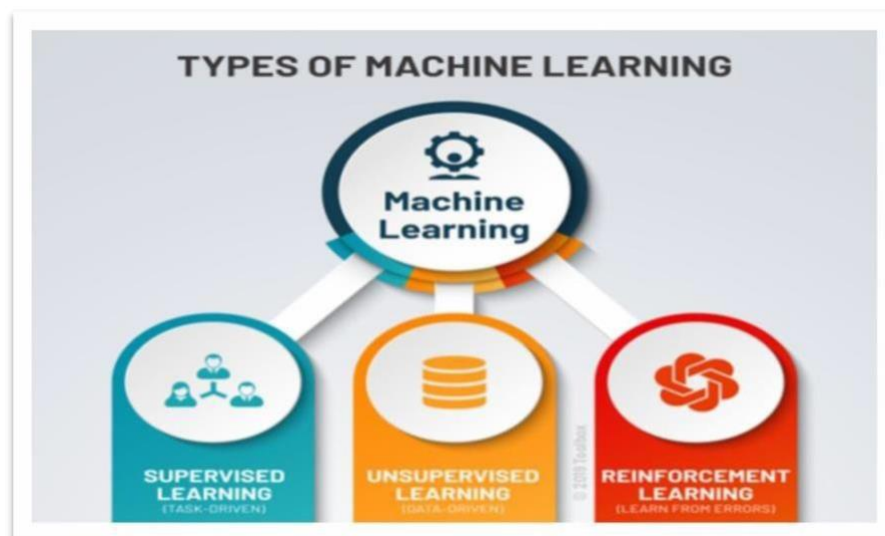


Fig 1.1 Types of machine learning

Regression techniques predict continuous responses - for example, changes in temperature or fluctuations in power demand. Typical applications include electricity load forecasting and algorithmic trading.

Use regression techniques if you are working with a data range or if the nature of your response is a real number, such as temperature or the time until failure for a piece of equipment.

Common regression algorithms include linear model, nonlinear model, stepwise regression, Gradient Descent Regression, Support Vector Regression, Ridge and Lasso Regressions.

Classification techniques predict discrete responses - for example, whether an email is genuine or spam, or whether a tumour is cancerous or benign. Classification models classify input data into categories. Typical applications include medical imaging, speech recognition, and credit scoring.

Use classification if your data can be tagged, categorized, or separated into specific groups or classes. For example, applications for hand-writing recognition use classification to recognize letters and numbers. In image processing and computer vision, unsupervised pattern recognition techniques are used for object detection and image segmentation.

Common algorithms for performing classification include support vector machine (SVM), boosted and bagged decision trees, k-nearest neighbour, Naïve Bayes, discriminant analysis, logistic regression, and neural networks.

Using Supervised Learning to Predict Heart Attacks: Suppose clinicians want to predict whether someone will have a heart attack within a year. They have data on previous patients, including age, weight, height, and blood pressure. They know whether the previous patients had heart attacks within a year. So, the problem is combining the existing data into a model that can predict whether a new person will have a heart attack within a year.

Unsupervised Learning:

Unsupervised learning finds hidden patterns or intrinsic structures in data. It is used to draw inferences from datasets consisting of input data without labelled responses.

Clustering is the most common unsupervised learning technique. It is used for exploratory data analysis to find hidden patterns or groupings in data. Applications for cluster analysis include gene sequence analysis and market research.

For example, if a cell phone company wants to optimize the locations where they build cell phone towers, they can use machine learning to estimate the number of clusters of people relying on their towers. A phone can only talk to one tower at a time, so the team uses clustering algorithms to design the best placement of cell towers to optimize signal reception for groups, or clusters, of their customers.

Common algorithms for performing clustering include k-means and k-medoids, Apriori algorithms, hierarchical clustering, Gaussian mixture models and hidden Markov models.

1.2. Benefits of Using ML in BANK MARKETING ANALYSIS

Exploring the Advantages and Disadvantages of Machine Learning:

- When it comes to learning technology, we should be aware of the pros and cons of that technology. The reason is so that we can understand the capabilities of that subject.
- That is exactly what we are doing here. Understanding the advantages and disadvantages of Machine Learning will help us to unlock many doors.
- The advantages of Machine Learning are vast. It helps us to create ways of modernizing technology. The disadvantages of Machine Learning tell us its limits and side effects. This helps us to find different innovative ways to reduce these problems.

Advantages of Machine Learning:

- **Automation of Everything:** Machine Learning is responsible for cutting the workload and time. By automating things, we let the algorithm do the hard work for us. Automation is now being done almost everywhere. The reason is that it is very reliable. Also, it helps us to think more creatively. Due to ML, we are now designing more advanced computers. These computers can handle various Machine Learning models and algorithms efficiently. Even though automation is spreading fast, we still don't completely rely on it. ML is slowly transforming the industry with its automation.
- **Wide Range of Applications:** ML has a wide variety of applications. This means that we can apply ML on any of the major fields. ML has its role everywhere from medical, business, banking to science and tech. This helps to create more opportunities. It plays a major role in customer interactions. Machine Learning can help in the detection of diseases more quickly. It is helping to lift up businesses. That is why investing in ML technology is worth it.
- **Scope of Improvement:** Machine Learning is the type of technology that keeps on evolving. There is a lot of scope in ML to become the top technology in the future. The reason is it has a lot of research areas in it. This helps us to improve both hardware and software. In hardware, we have various laptops and GPU. These have various ML and Deep Learning networks in them. These help in the faster processing power of the system. When

it comes to software, we have various UI and libraries in use. These help in designing more efficient algorithms.

- **Efficient Handling of Data:** Machine Learning has many factors that make it reliable. One of them is data handling. ML plays the biggest role when it comes to data currently. It can handle any type of data. Machine Learning can be multidimensional or different types of data. It can process and analyse these data those normal systems can't. Data is the most important part of any Machine Learning model. Also, studying and handling of data is a field.

- **Best for Education and Online Shopping:**

ML would be the best tool for education in the future. It provides very creative techniques to help students study.

Recently in China, a school has started to use ML to improve student focus. In online shopping, the ML model studies your searches. Based on your search history, it would provide advertisements. These will be about your search preferences in previous searches. In this, the search history is the data for the model. This is a great way to improve e-commerce with ML.

1.3. About Industry (Bank Marketing Analysis)

The bank marketing industry has been going through massive transformations in the past couple of years. With more emphasis on customized loan plans and the increasing level of market competition. Bank marketing is known for its nature of developing a unique brand image, which is treated as the capital reputation of the financial academy. It is very important for a bank to develop good relationship with valued customers accompanied by innovative ideas which can be used as measures to meet their requirements. Basically, banks engage in transaction of products and services through their retail outlets known as branches to different customers at the grassroots level. This is referred as the 'top to bottom' approach. Relationship banking can be defined as a process that includes proactively predicting the demands of individual bank customers and taking steps to meet these demands before the client shows them.

1.3.1 AI / ML Role in Bank Marketing:

Machine Learning is a sub-set of artificial intelligence where computer algorithms are used to autonomously learn from data. Machine learning (ML) is getting more and more attention and is becoming increasingly popular in many other industries. Within the insurance industry, there is more application of ML regarding the Bank Marketing.

2.0 Bank Marketing Dataset content:

The data is related with direct marketing campaigns of a Portuguese banking instructions.

The marketing campaigns were based on phone calls. Often, more than one contact to the same client was required, in order to access if the product(bank term deposit) was subscribed or not. Data set has 16 predictor variables and 45K rows . Customers who received phone calls not be unique and a same customer might receive multiple phone calls. The last column y will be the target variable.

The main factors for term deposit claims are Age, Job, Education, Marital , Default , Housing, Loan, balance, contact, month, day, duration, campaign, previous.poutcomes,pdays...

2.1 Main Drivers for AI Bank Marketing Analysis

Predictive modelling allows for simultaneous consideration of many variables and quantification of their overall effect. When a large number of claims are analysed, patterns regarding the characteristics of the claims that drive loss development begin to emerge.

2.2 Internship Project - Data Link

The data is related with direct marketing campaigns of a Portuguese banking instructions.

The marketing campaigns were based on phone calls. Often, more than one contact to the same client was required, in order to access if the product(bank term deposit) was subscribed or not. Data set has 17 predictor variables and 45K rows . Customers who received phone calls not be unique and a same customer might receive multiple phone calls.

The last column y will be the target variable.

The main factors for term deposit claims are Age, Job, Education, Marital , Default , Housing, Loan, balance, contact, month, day, duration, campaign, previous.poutcomes,pdays...

The internship project data has taken from Kaggle and the link is:

<https://www.kaggle.com/code/hariharanpavan/bank-marketing-analysis/data>

3.0 AI / ML Modelling and Results

3.1 Your Problem of Statement

Predictive models are most effective when they are constructed using a company's own historical claims data since this allows the model to recognize the specific nature of a company's exposure as well as its claims practices. The construction of the model also involves input from the company throughout the process, as well as consideration of industry leading claims practices and benchmarks.

Predictive modelling can be used to quantify the impact to the claims department resulting from the failure to meet or exceed claim service leading practices. It can also be used to identify the root cause of claim leakage. Proper use of predictive modelling will allow for potential savings across two dimensions:

Early identification of claims with the potential for high leakage, thereby allowing for the proactive management of the claim

Recognition of practices that are unnecessarily increasing claims settlement payments.

3.2 Data Science Project Life Cycle

Data Science is a multidisciplinary field of study that combines programming skills, domain expertise and knowledge of statistics and mathematics to extract useful insights and knowledge from data.

In simple terms, a data science life cycle is nothing but a repetitive set of steps that you need to take to complete and deliver a project/product to your client.

Although the data science projects and the teams involved in deploying and developing the model will be different, every data science life cycle will be slightly different in every other company.

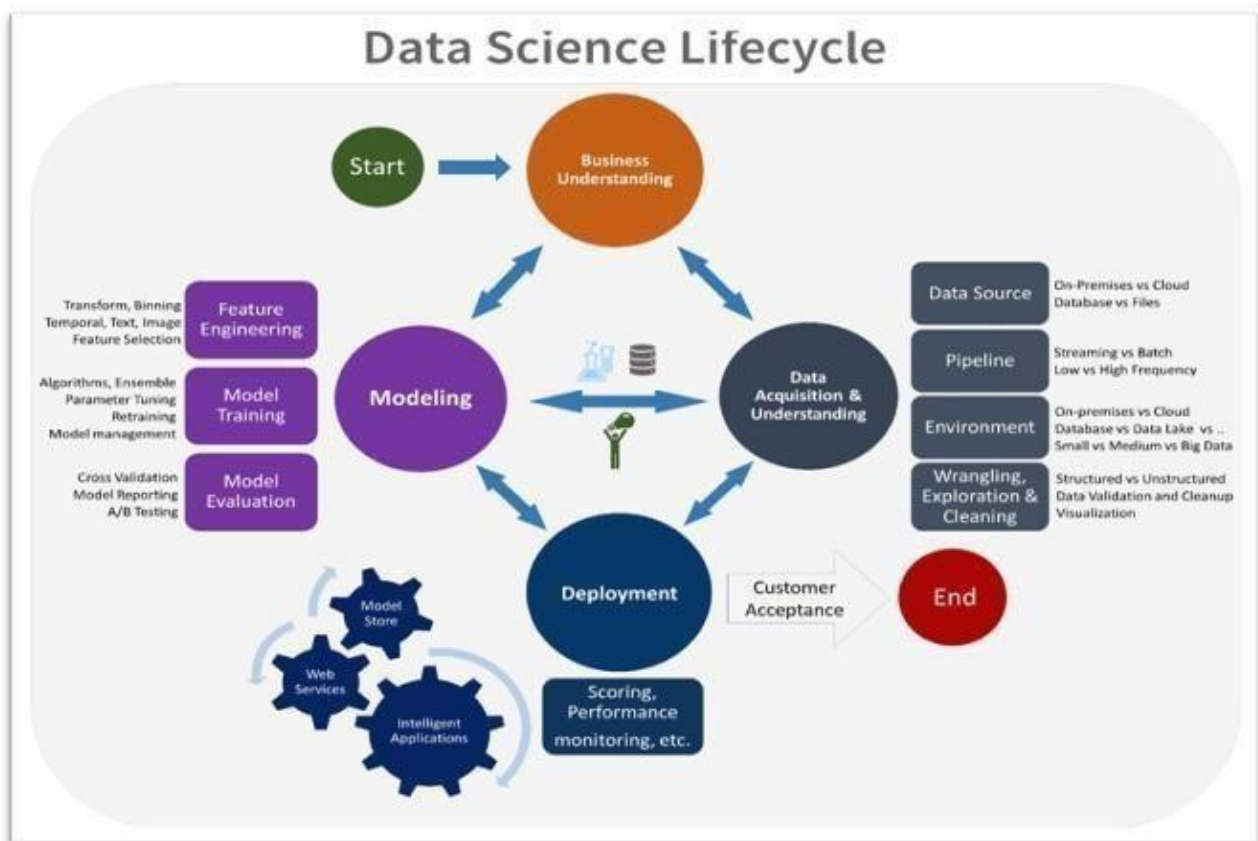
However, most of the data science projects happen to follow a somewhat similar process.

In order to start and complete a data science-based project, we need to understand the various roles and responsibilities of the people involved in building, developing the project.

Let us look at those employees who are involved in a typical data science project:

Who Are Involved in The Projects:

- Business Analyst
- Data Analyst
- Data Scientists
- Data Engineer
- Data Architect
- Machine Learning Engineer



3.2.1 Data Exploratory Analysis

Exploratory Data Analysis refers to the critical process of performing initial investigations on data so as to discover patterns, to spot anomalies, to test hypothesis and to check assumptions with the help of summary statistics and graphical representations.

3.2.2 Data Pre-processing

Data pre-processing, a component of data preparation, describes any type of processing performed on raw data to prepare it for another data processing procedure. It has traditionally been an important preliminary step for the data mining process. More recently, data preprocessing techniques have been adapted for training machine learning models and AI models and for running inferences against them. Data pre-processing transforms the data into a format that is more easily and effectively processed in data mining, machine learning and other data science tasks. The techniques are generally used at the earliest stages of the machine learning and AI development pipeline to ensure accurate results.

3.2.1.1 Check the Duplicate and low variation data

Types of Duplicate Features in Machine Learning:

Two things distinguish top data scientists from others in most cases: Feature Creation and Feature Selection. i.e., creating features that capture deeper/hidden insights about the business or customer and then making the right choices about which features to choose for your model.

1. Duplicate Values (Same value for each record)
2. Duplicate Index (value of two features are different but they occur at the same index)

Keeping duplicate features in your dataset introduces the problem of multicollinearity.

- In the case of linear models, weights distribution between the two features will be problematic.
- If you are using tree-based models, it won't matter unless you are looking at feature importance.
- In the case of distance-based models, it will make that feature count more in the distance.

3.2.2.2 Identify and address the missing variables

What are missing data?

Missing data are values that are not recorded in a dataset. They can be a single value missing in a single cell or missing of an entire observation (row). Missing data can occur both in a continuous variable (e.g., height of students) or a categorical variable (e.g., gender of a population).

Missing data are common in any field of natural or behavioral science, but it is particularly commonplace in social sciences research data.

In programming languages, missing values are represented as NA or Nan or simply an empty cell in an object.

• **The origins of missing data**

So where do the missing values come from, and why do they even exist?

Let's give an example. You are administering a questionnaire survey among a sample of respondents; and in the questionnaire, you are asking a question about household income. Now, what if a respondent refuses to answer that question? Would you make that up or rather leave the field empty? You'd probably leave that cell empty — creating an instance of missing value

• **Problems caused**

Missing values are undesirable, but it is difficult to quantify the magnitude of effects in statistical and machine learning projects. If it's a large dataset and a very small percentage of data is missing the effect may not be detectable at all.

However, if the dataset is relatively small, every data point counts. In these situations, a missing data point means loss of valuable information.

In any case, generally missing data creates imbalanced observations, cause biased estimates, and in extreme cases, can even lead to invalid conclusion.

Case deletion: if the dataset is relatively large delete the complete record with a missing value

Substitution: substitute missing cells with (a) column mean, (b) mean of nearest neighbors, (c) moving average, or (c) filling with the last observation

Statistical imputation: a regression can be an effective way to determine the value of missing cell given other information in the dataset

Sensitivity analysis: if the sample is small or missing values are relatively large then conduct a sensitivity analysis with multiple variations of outcomes.

3.2.2.2 Identify objects and convert into numerical values:

- ✓ Defining data types when reading a CSV file
- ✓ Creating a custom function to convert data type
- ✓ `as type ()` vs. `to numeric ()`

When doing data analysis, it is important to ensure correct data types. Otherwise, you may get unexpected results or errors. In the case of Pandas, it will correctly infer data types in many cases, and you can move on with your analysis without any further thought on the topic.

Despite how well pandas works at some point in your data analysis process you will likely need to explicitly convert data from one type to another. This article will discuss how to change data to a numeric type. More specifically, you will learn how to use the Pandas built-in methods `as type ()` and `to numeric ()` to deal with the following common problems:

- ✓ Converting string/int to int/float
- ✓ Converting float to int
- ✓ Converting a column of mixed data types
- ✓ Handling missing values
- ✓ Converting a money column to float
- ✓ Converting Boolean to 0/1
- ✓ Converting multiple data columns at once

3.2.2.3 Handling of Outliers

What is an outlier?

An outlier is an observation that lies an abnormal distance from other values in a random sample from a population.

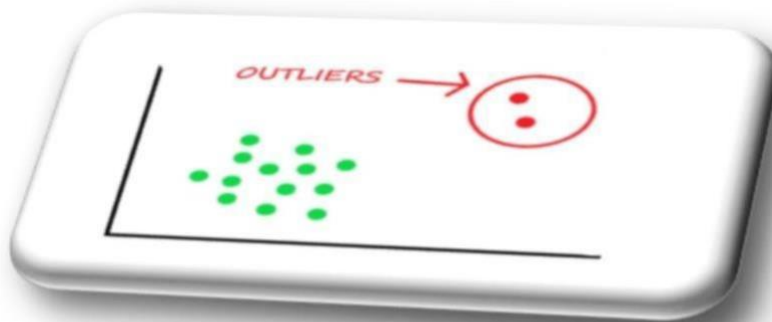
There is, of course, a degree of ambiguity. Qualifying a data point as an anomaly leaves it up to the analyst or model to determine what is abnormal—and what to do with such data points.

There are also different degrees of outliers:

- Mild outliers lie beyond an “inner fence” on either side.
- Extreme outliers are beyond an “outer fence.”

Why do outliers occur? According to Tom Barenberg, chief economist and data consultant at Unity Marketing, “It can be the result of measurement or recording errors, or the unintended and truthful outcome resulting from the set’s definition.”

Outliers may contain valuable information. Or be meaningless aberrations caused by measurement and recording errors. In any case, they can cause problems with repeatable A/B test results, so it’s important to question and analyze outliers.



3.2.2.4 Categorical data and Encoding Techniques

What is Categorical Data?

Since we are going to be working on categorical variables in this article, here is a quick refresher on the same with a couple of examples. Categorical variables are usually represented as ‘strings’ or ‘categories’ and are finite in number. Here are a few examples:

1. The city where a person lives: Delhi, Mumbai, Ahmedabad, Bangalore, etc.
2. The department a person works in: Finance, Human resources, IT, Production.
3. The highest degree a person has: High school, Diploma, Bachelors, Masters, PhD.
4. The grades of a student: A+, A, B+, B, B- etc.

In the above examples, the variables only have definite possible values. Further, we can see there are two kinds of categorical data-

- Ordinal Data: The categories have an inherent order
- Nominal Data: The categories do not have an inherent order

Label 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 representing the education qualification of a person.

Binary Encoding:

- Binary encoding is a combination of Hash encoding and one-hot encoding. In this encoding scheme, the categorical feature is first converted into numerical using an ordinal encoder. Then the numbers are transformed in the binary number. After that binary value is split into different columns.
- Binary encoding works well when there are a high number of categories. For example, the cities in a country where a company supplies its products

3.2.2.5 Feature Scaling

Why Feature Scaling?

Real Life Datasets have many features with a wide range of values like for example let's consider the house price prediction dataset. It will have many features like no. of bedrooms, square feet area of the house, etc.

As you can guess, the no. of bedrooms will vary between 1 and 5, but the square feet area will range from 500-2000. This is a huge difference in the range of both features.

Many machine learning algorithms that are using Euclidean distance as a metric to calculate the similarities will fail to give a reasonable recognition to the smaller feature, in this case, the number of bedrooms, which in the real case can turn out to be an important metric.

E.g.: Linear Regression, Logistic Regression, KNN

There are several ways to do feature scaling. I will be discussing the top 5 of the most used feature scaling techniques.

3.2.3 Selection of Dependent and Independent variables

The dependent or target variable here is satisfaction Target which tells us a

The independent variables are selected after doing exploratory data analysis. Which tells us that the Customer is satisfied, neutral or dissatisfied.

3.2.4 Data Sampling Methods

The data we have is highly unbalanced data so we used some sampling methods which are used to balance the target variable so we our model will be developed with good accuracy and precision.

We used three Sampling methods

3.2.4.1 Stratified sampling

Stratified sampling randomly selects data points from majority class so they will be equal to the data points in the minority class. So, after the sampling both the class will have same no of observations.

It can be performed using strata function from the library sampling.

3.2.4.2 Simple random sampling

Simple random sampling is a sampling technique where a set percentage of the data is selected randomly. It is generally done to reduce bias in the dataset which can occur if data is selected manually without randomizing the dataset.

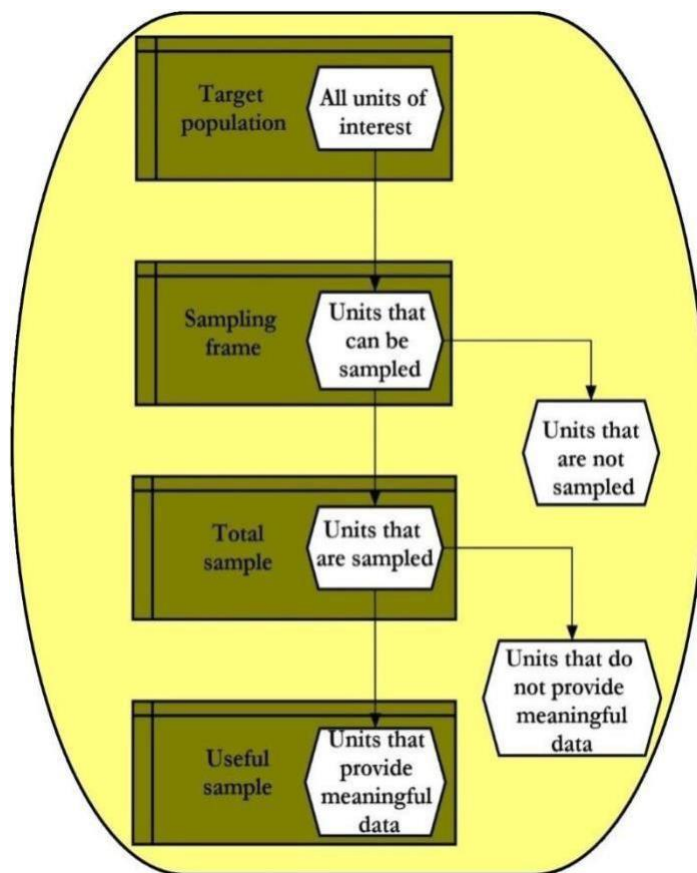
We used this method to split the dataset into train dataset which contains 70% of the total data and test dataset with the remaining 30% of the data.

Steps involved in Sampling:

The first stage in the sampling process is to clearly define the target population.

- So, to carry out opinion polls, polling agencies consider only the people who are above 18 years of age and are eligible to vote in the population. Sampling Frame – It is a list of items or people forming a population from which the sample is taken.
- So, the sampling frame would be the list of all the people whose names appear on the voter list of a constituency.

- Generally, probability sampling methods are used because every vote has equal value and any person can be included in the sample irrespective of his caste, community, or religion. Different samples are taken from different regions all over the country.
- Sample Size – It is the number of individuals or items to be taken in a sample that would be enough to make inferences about the population with the desired level of accuracy and precision
- Once the target population, sampling frame, sampling technique, and sample **size** have been established, the next step is to collect data from the sample.



3.2.5 Models Used for Development

3.2.5.1 Model 01(Logistic regression)

Logistic uses logistic link function to convert the likelihood values to probabilities so we can get a good estimate on the probability of a particular observation to be positive class or negative class. The also gives us p-value of the variables which tells us about significance of each independent variable.

3.2.5.2 Model 02(Decision Tree Classifier)

Decision Tree is a Supervised learning technique that can be used for both classification and Regression problems, but mostly it is preferred for solving Classification problems. It is a tree-structured classifier, where internal nodes represent the features of a dataset, branches represent the decision rules, and each leaf node represents the outcome.

In a Decision tree, there are two nodes, which are the Decision Node and Leaf Node. Decision nodes are used to make any decision and have multiple branches, whereas Leaf nodes are the output of those decisions and do not contain any further branches.

The decisions or the test are performed based on features of the given dataset.

It is a graphical representation for getting all the possible solutions to a problem/decision based on given conditions.

It is called a decision tree because, like a tree, it starts with the root node, which expands on further branches and constructs a tree-like structure.

Decision Trees (DTs) are a non-parametric supervised learning method used for classification and regression. The goal is to create a model that predicts the value of a target variable by learning simple decision rules inferred from the data features. A tree can be seen as a piece wise constant approximation

3.2.5.3 Model 03(Random Forest Classifier)

Random forest is an algorithm that consists of many decision trees. It was first developed by Leo Breiman and Adele Cutler. The idea behind it is to build several trees, to have the instance classified by each tree, and to give a "vote" at each class. The model uses a "bagging" approach and the random selection of features to build a collection of decision trees with controlled variance. The instance's class is to the class with the highest number of votes, the class that occurs the most within the leaf in which the instance is placed.

The error of the forest depends on:

- Trees correlation: the higher the correlation, the higher the forest error rate.
- The strength of each tree in the forest. A strong tree is a tree with low error. By using trees that classify the instances with low error the error rate of the forest decreases.

3.2.5.4 Model 04(Extra Trees Classifier)

This class implements a meta estimator that fits several randomized decision trees (a.k.a. extra-trees) on various sub-samples of the dataset and uses averaging to improve the predictive accuracy and control over-fitting.

3.2.5.5 Model 05(KNN Classifier)

K-Nearest Neighbor is one of the simplest Machine Learning algorithms based on Supervised Learning technique K-NN algorithm assumes the similarity between the new case/data and available cases and put the new case into the category that is most like the available categories'-NN algorithm stores all the available data and classifies a new data point based on the similarity. This means when new data appears then it can be easily classified into a well suite category by using K- NN algorithm.

K-NN algorithm can be used for Regression as well as for Classification but mostly it is used for the Classification problems K-NN is a non-parametric algorithm, which means it does not make any assumption on underlying data. It is also called a lazy learner algorithm because it does not learn from the training set immediately instead it stores the dataset and at the time of classification, it performs an action on the dataset KNN-algorithm at the training phase just stores the dataset and when it gets new data, then it classifies that data into a category that is much like the new data.

3.2.5.6 Model 06(Naïve Bayes)

Naïve Bayes is a probabilistic machine learning algorithm used for many classification functions and is based on the Bayes theorem. Gaussian Naïve Bayes is the extension of naïve Bayes. While other functions are used to estimate data distribution, Gaussian or normal distribution is the simplest to implement as you will need to calculate the mean and standard deviation for the training data.

What is the Naive Bayes Algorithm?

Naive Bayes is a probabilistic machine learning algorithm that can be used in several classification tasks. Typical applications of Naive Bayes are classification of documents,

filtering spam, prediction and so on. This algorithm is based on the discoveries of Thomas Bayes and hence its name.

The name “Naïve” is used because the algorithm incorporates features in its model that are independent of each other. Any modifications in the value of one feature do not directly impact the value of any other feature of the algorithm. The main advantage of the Naïve Bayes algorithm is that it is a simple yet powerful algorithm.

It is based on the probabilistic model where the algorithm can be coded easily, and predictions did quickly in real-time. Hence this algorithm is the typical choice to solve real- world problems as it can be tuned to respond to user requests instantly. But before we dive deep into Naïve Bayes and Gaussian Naïve Bayes, we must know what is meant by conditional probability.

3.2.5.7 Model 07(XG Boost)

XG Boost is an implementation of Gradient Boosted decision trees. This library was written in C++. It is a type of Software library that was designed basically to improve speed and model performance. It has recently been dominating in applied machine learning. XG Boost models majorly dominate in many Kaggle Competitions. In this algorithm, decision trees are created in sequential form. Weights play an important role in XG Boost. Weights are assigned to all the independent variables which are then fed into the decision tree which predicts results. The weight of variables predicted wrong by the tree is increased and the variables are then fed to the second decision tree. These individual classifiers/predictors then ensemble to give a strong and more precise model. It can work on regression, classification, ranking, and user-defined prediction problems.

3.2.5.8 Model 08(Light GBM)

Light GBM is a gradient boosting framework based on decision trees to increases the efficiency of the model and reduces memory usage.

It uses two novel techniques: Gradient-based One Side Sampling and Exclusive Feature Bundling (EFB) which fulfills the limitations of histogram-based algorithm that is primarily used in all GBDT (Gradient Boosting Decision Tree) frameworks. The two techniques of GOSS and EFB described below form the characteristics of Light GBM Algorithm. They

comprise together to make the model work efficiently and provide it a cutting edge over other GBDT frameworks.

Gradient-based One Side Sampling Technique for Light GBM:

Different data instances have varied roles in the computation of information gain. The instances with larger gradients (i.e., under-trained instances) will contribute more to the information gain. GOSS keeps those instances with large gradients (e.g., larger than a predefined threshold, or among the top percentiles), and only randomly drop those instances with small gradients to retain the accuracy of information gain estimation. This treatment can lead to a more accurate gain estimation than uniformly random sampling, with the same target sampling rate, especially when the value of information gain has a large range.

3.2.5.9 Model 09 (SVC)

Support Vector Machine or SVM is one of the most popular Supervised Learning algorithms, which is used for Classification as well as Regression problems. However, primarily, it is used for Classification problems in Machine Learning.

The goal of the SVM algorithm is to create the best line or decision boundary that can segregate n-dimensional space into classes so that we can easily put the new data point in the correct category in the future. This best decision boundary is called a hyperplane.

3.2.5.10 Model 10 (Gradient Boosting Classifier)

This algorithm builds an additive model in a forward stage-wise fashion; it allows for the optimization of arbitrary differentiable loss functions. In each stage `n_classes_` regression trees are fit on the negative gradient of the loss function, e.g. binary or multiclass log loss. Binary classification is a special case where only a single regression tree is induced.

3.3 AI / ML Model Analysis and Final Results

We used our train dataset to build the above models and used our test data to check the accuracy and performance of our models.

We used confusion matrix to check accuracy, Precision, Recall and F1 score of our models and compare and select the best model for given auto dataset of size ~ 272252 policies.

3.3.1 Different Model codes

This section in which we used different types of model code as follows :

```
# Build the Classification models and compare the results
```

```
from sklearn.linear_model import LogisticRegression
```

```
from sklearn.tree import DecisionTreeClassifier
```

```
from sklearn.ensemble import RandomForestClassifier
```

```
from sklearn.ensemble import ExtraTreesClassifier
```

```
from sklearn.neighbors import KNeighborsClassifier
```

```
from sklearn.naive_bayes import GaussianNB
```

```
from xgboost import XGBClassifier
```

```
from sklearn.svm import SVC
```

```
import lightgbm as lgb
```

```
# Create objects of classification algorithm with default hyper-parameters
```

```
ModelLR = LogisticRegression()
```

```
ModelDC = DecisionTreeClassifier()
```

```
ModelRF = RandomForestClassifier()
```

```
ModelET = ExtraTreesClassifier()
```

```
ModelKNN = KNeighborsClassifier(n_neighbors=5)
```

```
ModelGNB = GaussianNB()
```

```
ModelXGB = XGBClassifier(n_estimators=100, max_depth=3, eval_metric='mlogloss')
```

```
ModelSVM = SVC(kernel='linear', random_state=1)
```

```
ModelLGB = lgb.LGBMClassifier()
```

```
# Evaluation matrix for all the algorithms
```

```
#MM = [ModelLR, ModelDC, ModelRF, ModelET, ModelKNN, ModelGNB, ModelSVM,  
modelXGB, modelLGB]
```

```
MM = [ModelLR, ModelDC, ModelRF, ModelET, ModelKNN, ModelGNB, ModelXGB]
```

```
for models in MM:
```

```
    # Fit the model
```

```
    models.fit(x_train, y_train)
```

```
    # Prediction
```

```
    y_pred = models.predict(x_test)
```

```
    y_pred_prob = models.predict_proba(x_test)
```

```
    # Print the model name
```

```
    print('Model Name: ', models)
```

```
    # confusion matrix in sklearn
```

```
    from sklearn.metrics import confusion_matrix
```

```
    from sklearn.metrics import classification_report
```

```
    # actual values
```

```
    actual = y_test
```

```
    # predicted values
```

```
    predicted = y_pred
```

```
    # confusion matrix
```

```
    matrix = confusion_matrix(actual, predicted, labels=[1,0], sample_weight=None, normalize=None)
```

```
print('Confusion matrix : \n', matrix)
```

```
# outcome values order in sklearn
```

```
tp, fn, fp, tn = confusion_matrix(actual,predicted,labels=[1,0]).reshape(-1)
```

```
print('Outcome values : \n', tp, fn, fp, tn)
```

```
# classification report for precision, recall f1-score and accuracy
```

```
C_Report = classification_report(actual,predicted,labels=[1,0])
```

```
print('Classification report : \n', C_Report)
```

```
# calculating the metrics
```

```
sensitivity = round(tp/(tp+fn), 3);
```

```
specificity = round(tn/(tn+fp), 3);
```

```
accuracy = round((tp+tn)/(tp+fp+tn+fn), 3);
```

```
balanced_accuracy = round((sensitivity+specificity)/2, 3);
```

```
precision = round(tp/(tp+fp), 3);
```

```
f1Score = round((2*tp/(2*tp + fp + fn)), 3);
```

```
# Matthews Correlation Coefficient (MCC). Range of values of MCC lie between -1 to +1.
```

```
# A model with a score of +1 is a perfect model and -1 is a poor model
```

```
from math import sqrt
```

```
mx = (tp+fp) * (tp+fn) * (tn+fp) * (tn+fn)
```

```
MCC = round(((tp * tn) - (fp * fn)) / sqrt(mx), 3)
```

```

print('Accuracy :', round(accuracy*100, 2), '%')

print('Precision :', round(precision*100, 2), '%')

print('Recall :', round(sensitivity*100,2), '%')

print('F1 Score :', f1Score)

print('Specificity or True Negative Rate :', round(specificity*100,2), '%')

print('Balanced Accuracy :', round(balanced_accuracy*100, 2), '%')

print('MCC :', MCC)

# Area under ROC curve

from sklearn.metrics import roc_curve, roc_auc_score

print('roc_auc_score:', round(roc_auc_score(y_test, y_pred), 3))


# ROC Curve

from sklearn.metrics import roc_auc_score

from sklearn.metrics import roc_curve

logit_roc_auc = roc_auc_score(y_test, y_pred)

fpr, tpr, thresholds = roc_curve(y_test, models.predict_proba(x_test)[:,-1])

plt.figure()

# plt.plot(fpr, tpr, label='Logistic Regression (area = %0.2f)' % logit_roc_auc)

plt.plot(fpr, tpr, label= 'Classification Model' % logit_roc_auc)

plt.plot([0, 1], [0, 1], 'r--')

plt.xlim([0.0, 1.0])

plt.ylim([0.0, 1.05])

plt.xlabel('False Positive Rate')

plt.ylabel('True Positive Rate')

plt.title('Receiver operating characteristic')

plt.legend(loc="lower right")

```

```
plt.savefig('Log_ROC')

plt.show()

print('.....')

#_____

new_row = {'Model Name' : models,
           'True Positive' : tp,
           'False Negative' : fn,
           'False Positive' : fp,
           'True Negative' : tn,
           'Accuracy' : accuracy,
           'Precision' : precision,
           'Recall' : sensitivity,
           'F1 Score' : f1Score,
           'Specificity' : specificity,
           'MCC':MCC,
           'ROC_AUC_Score':roc_auc_score(y_test, y_pred),
           'Balanced Accuracy':balanced_accuracy,
           }

HTResults = HTResults.append(new_row, ignore_index=True)
```

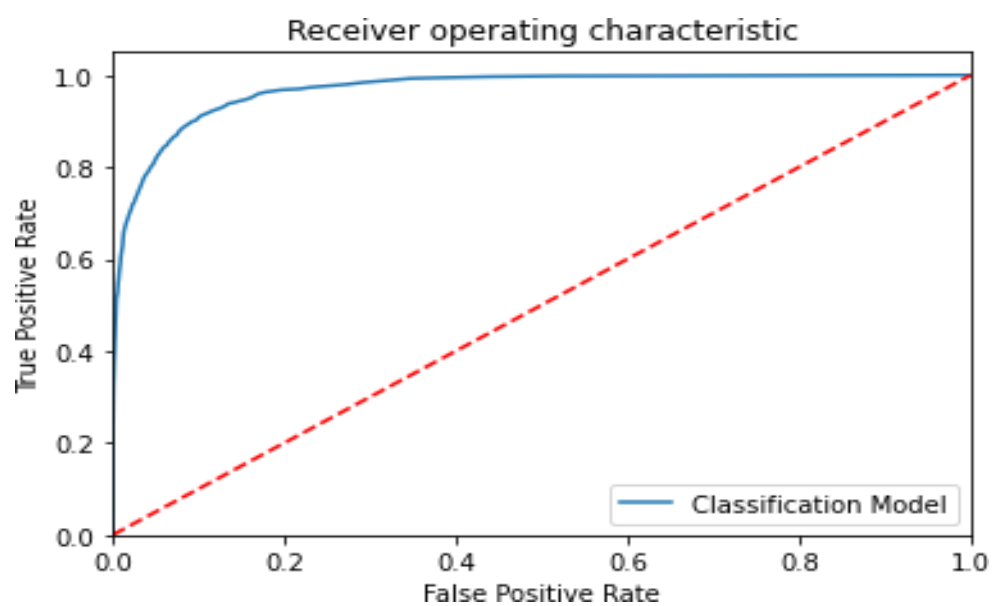

4.0 Conclusions and Future work

The model results in the following order by considering the model accuracy, F1 score and RoC AUC score.

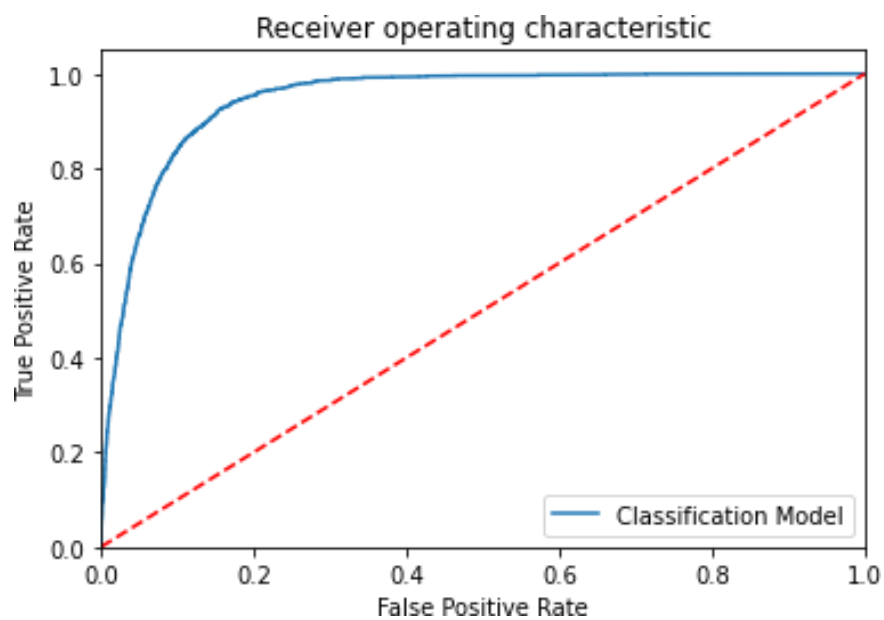
- Random tree Classifier
- LGBM Classifier
- Decision tree classifier

We recommend model – Extra Trees Classifier as a best fit for the given dataset. We considered Random Forest because it uses bootstrap aggregation which can reduce bias and variance in the data and can leads to good predictions with airline passengers satisfaction dataset.

	Model Name	True_Positive	False_Negative	False_Positive	True_Negative	Accuracy	Precision	Recall	F1 Score	Specificity	MCC	ROC_AUC_Score	Balanced Accuracy
0	LogisticRegression(C=100, class_weight='balanc...	1617	441	2165	8612	0.797	0.428	0.786	0.554	0.799	0.471	0.792412	0.792
1	DecisionTreeClassifier()	1547	511	832	9945	0.895	0.650	0.752	0.697	0.923	0.637	0.837250	0.838
2	(DecisionTreeClassifier(max_features='sqrt', r...	1592	466	385	10392	0.934	0.805	0.774	0.789	0.964	0.750	0.868921	0.869
3	(ExtraTreeClassifier(criterion='entropy', max_...	804	1254	210	10567	0.886	0.793	0.391	0.523	0.981	0.505	0.685592	0.686
4	KNeighborsClassifier()	1111	947	603	10174	0.879	0.648	0.540	0.589	0.944	0.522	0.741946	0.742
5	GaussianNB()	1077	981	1382	9395	0.816	0.438	0.523	0.477	0.872	0.368	0.697544	0.698
6	XGBClassifier(base_score=0.5, booster='gbtree'...	1207	851	485	10292	0.896	0.713	0.586	0.644	0.955	0.587	0.770744	0.770
7	LGBMClassifier()	1361	697	550	10227	0.903	0.712	0.661	0.686	0.949	0.629	0.805144	0.805



Random Forest Classifier



LGBM Classifier

5.0 References

Information gathered from :

<https://www.javatpoint.com/types-of-machine-learning>

<https://www.kaggle.com>

6.0 Appendices

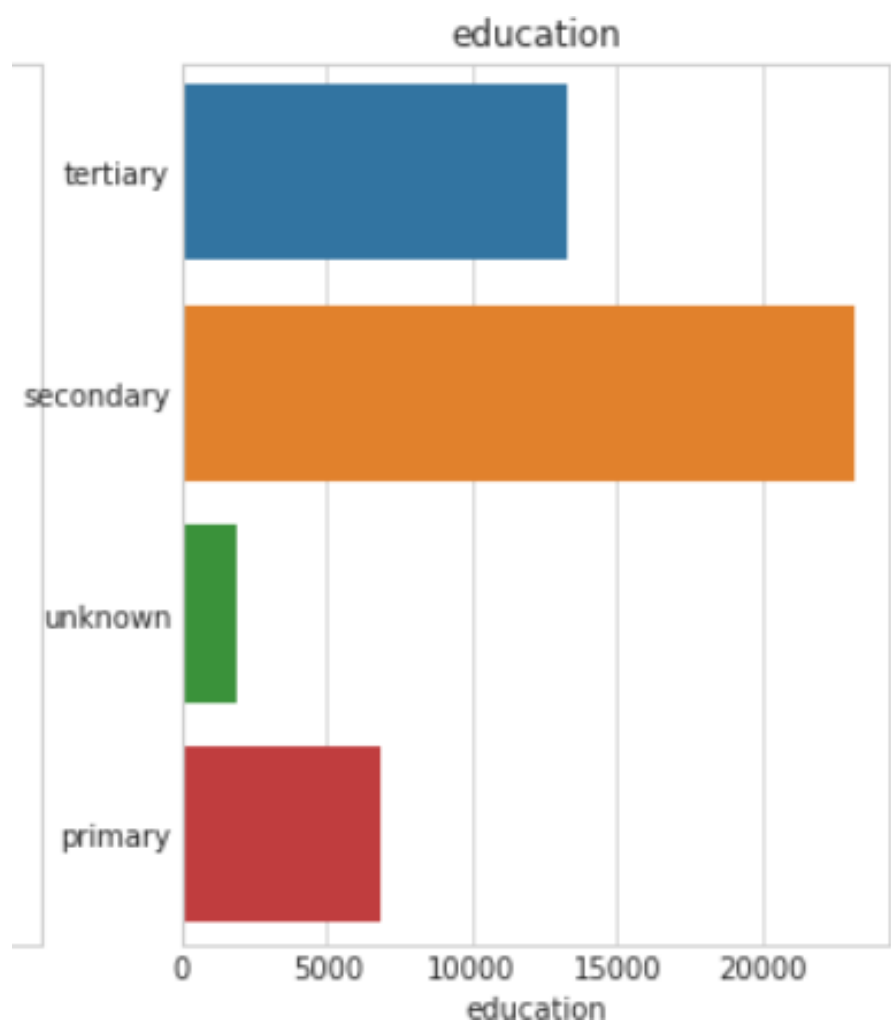
- Kernel methods
- Graphical models
- Reinforcement learning
- Convex analysis
- Optimization
- Bio informatics
- Minimal description length principle
- Topics in information theory
- Decision theory
- Network algorithms
- Computational Social Science
- Natural Language Processing Recent trends in deep learning and representation learning
- Facial recognition
- Voice Recognition
- Virtual Personal Assistants

6.1 Python code Results

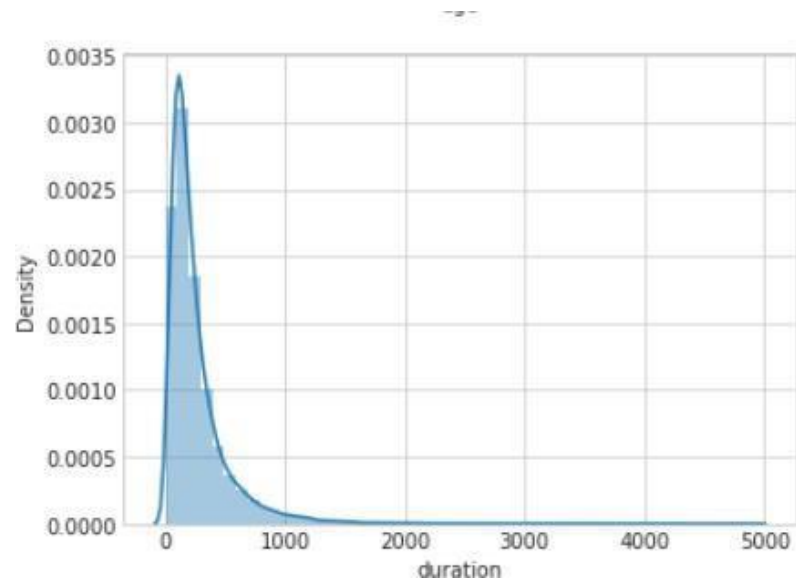
	Model Name	True_Positive	False_Negative	False_Positive	True_Negative	Accuracy	Precision	Recall	F1 Score	Specificity	MCC	ROC_AUC_Score	Balanced Accuracy
0	LogisticRegression(C=100, class_weight='balanc...	1617	441	2165	8612	0.797	0.428	0.786	0.554	0.799	0.471	0.792412	0.792
1	DecisionTreeClassifier()	1547	511	832	9945	0.895	0.650	0.752	0.697	0.923	0.637	0.837250	0.838
2	(DecisionTreeClassifier(max_features='sqrt', r...	1592	466	385	10392	0.934	0.805	0.774	0.789	0.964	0.750	0.868921	0.869
3	(ExtraTreeClassifier(criterion='entropy', max_...	804	1254	210	10567	0.886	0.793	0.391	0.523	0.981	0.505	0.685592	0.686
4	KNeighborsClassifier()	1111	947	603	10174	0.879	0.648	0.540	0.589	0.944	0.522	0.741946	0.742
5	GaussianNB()	1077	981	1382	9395	0.816	0.438	0.523	0.477	0.872	0.368	0.697544	0.698
6	XGBClassifier(base_score=0.5, booster='gbtree'...	1207	851	485	10292	0.896	0.713	0.586	0.644	0.955	0.587	0.770744	0.770
7	LGBMClassifier()	1361	697	550	10227	0.903	0.712	0.661	0.686	0.949	0.629	0.805144	0.805

6.2 List of Charts

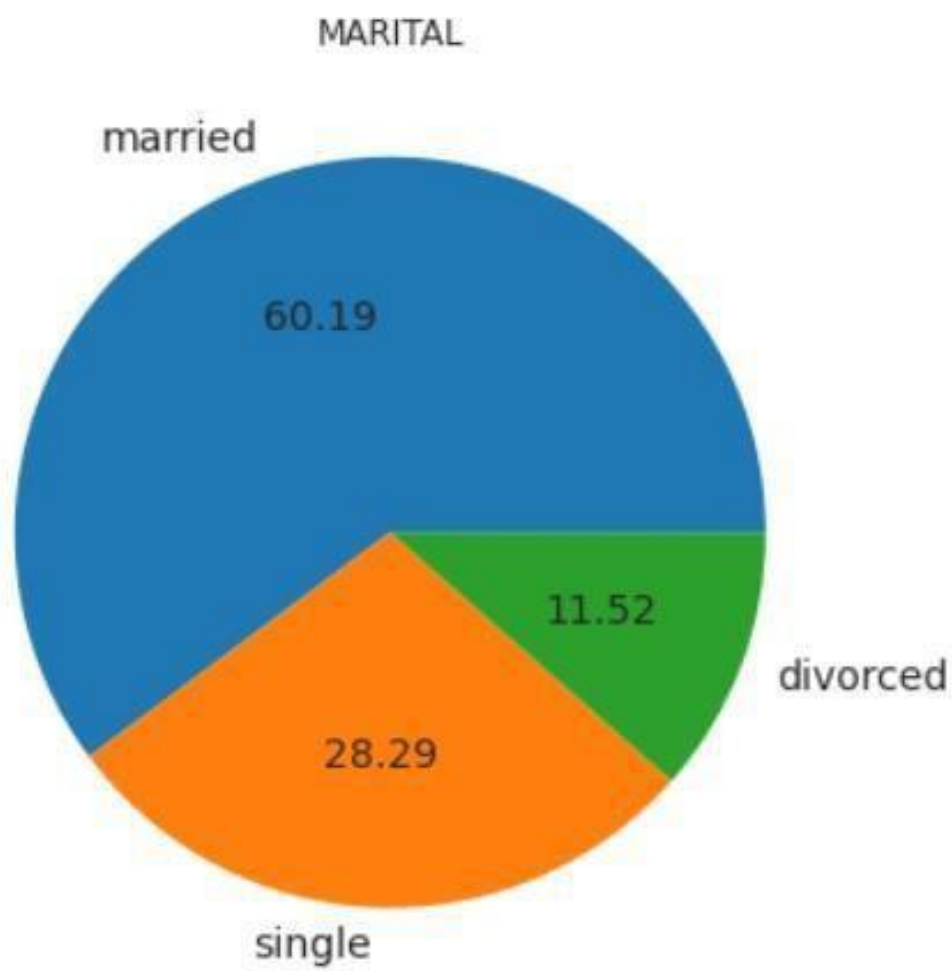
6.2.1 Chart 01: Education



6.2.2 Chart 02: Duration



6.2.3 Chart 03: Marital status



6.2.4 Chart 04: Heat Map

