*Development of Statistical and Machine Learning App*

# Manual Description

Ankita Narvekar
Arvind Parameswaran
Jyoti Narnolia
Nidhi Shivgan
Varsha GanapathyRao

# Table of Contents

# Description of the App

The app is suitable for Statistical and Machine Learning analytics. It allows the user to do descriptive and predictive analytics on a real dataset.

This app contains four stages:

1. Descriptive Analysis
2. Discrete Models
3. Continuous Models
4. Classification Models

## 1. Descriptive Analysis

Descriptive statistics are used to summarize data in a way that provides insight into the information contained in the data. This might include examining the mean or median of numeric data or the frequency of observations for nominal data.

In the app, the user has the benefit of choosing from the existing datasets (from - Attitude, Cars, Islands, Rock, Seatbelts) as well as user can choose his own'.CSV' file from his system.

After uploading the file from the system, the user needs to wait until the upload is completed and a message is shown just below the uploaded file.

The user can change the number of observations to view as per his requirements. After changing or uploading the settings, press 'Update View' to update the existing summary. (Refer images)
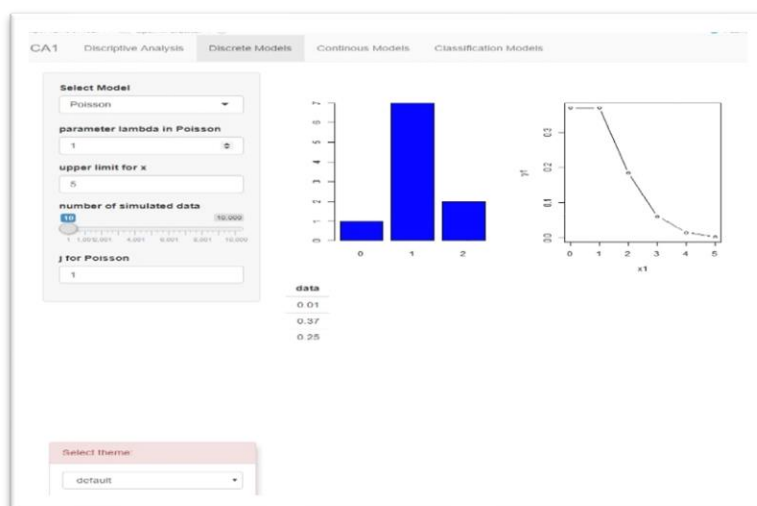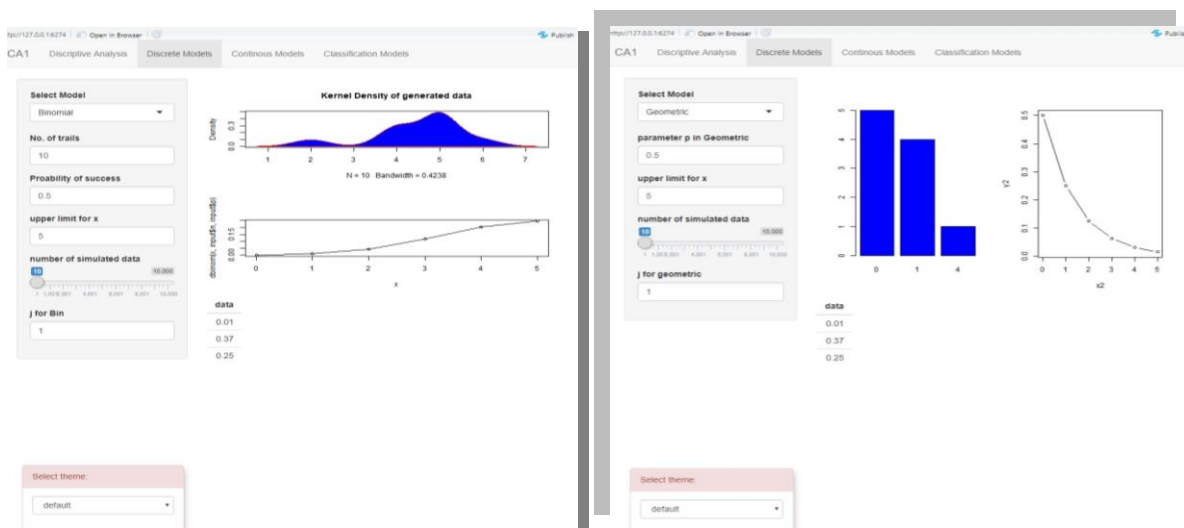
## 2. Discrete Models

Discrete modelling is the discrete analogue of continuous modelling. In discrete modelling, formulae are fit to discrete data—data that could potentially take on only a countable set of values, such as the integers, and which are not infinitely divisible.

In the app, Discrete Models consists of three models, namely, Binomial, Geometric and Poisson.

In Binomial Model, the user can select the number of trails, the probability of success, upper limit for $x$, number of simulated data (default value - 10) and $j$ value. After choosing all the data, the final graph is to be considered for final output. The Kernel Density will be updated as soon as any of the above specified data is altered.

In Geometric Model, the user can alter the value of $p$ parameter, upper limit of $x$, number of simulated data (default value - 10) and $j$ value. After choosing all the data, the final graph is to be considered for final output. The Kernel Density will be updated as soon as any of the above specified data is altered.

In Poisson Model, the user is given an option to alter the values of the parameter *lambda*, upper limit of $x$, number of simulated data (default value - 10) and $j$ value. After choosing all the data, the final graph is to be considered for final output. The Kernel Density will be updated as soon as any of the above specified data is altered.

## 3. Continuous Models

If a random variable is a continuous variable, its probability distribution is called a **continuous probability distribution.**
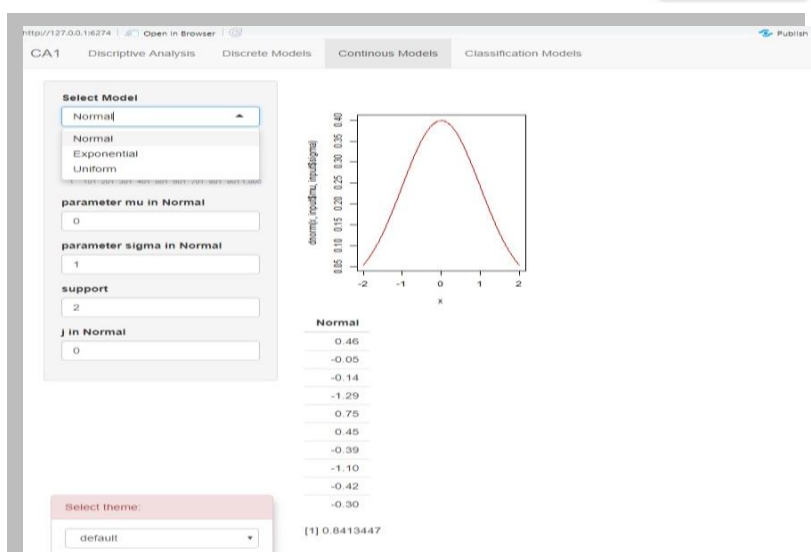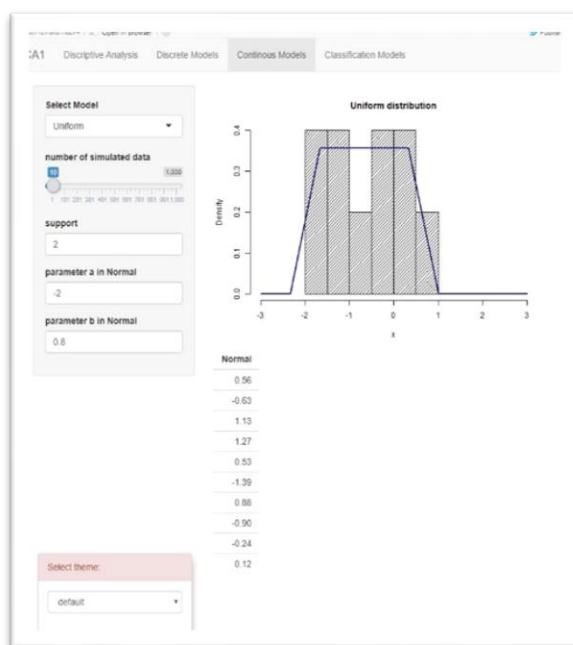
In the app, Continuous Models consists of three models, namely, Normal, Exponential and Uniform.

In Uniform Model, the user needs to feed the value of simulated data, support variable, parameter '*a*' and parameter '*b*'.

In Exponential Model, the user needs to feed the value of simulated data, support variable, parameter lambda and *j* value.

In Normal Model, the user needs to feed the value of simulated data, support variable, parameter '*mu*', parameter '*sigma*', *and* j value.

After updating the values of the input, the graph is plotted on the right hand corner space of the screen.
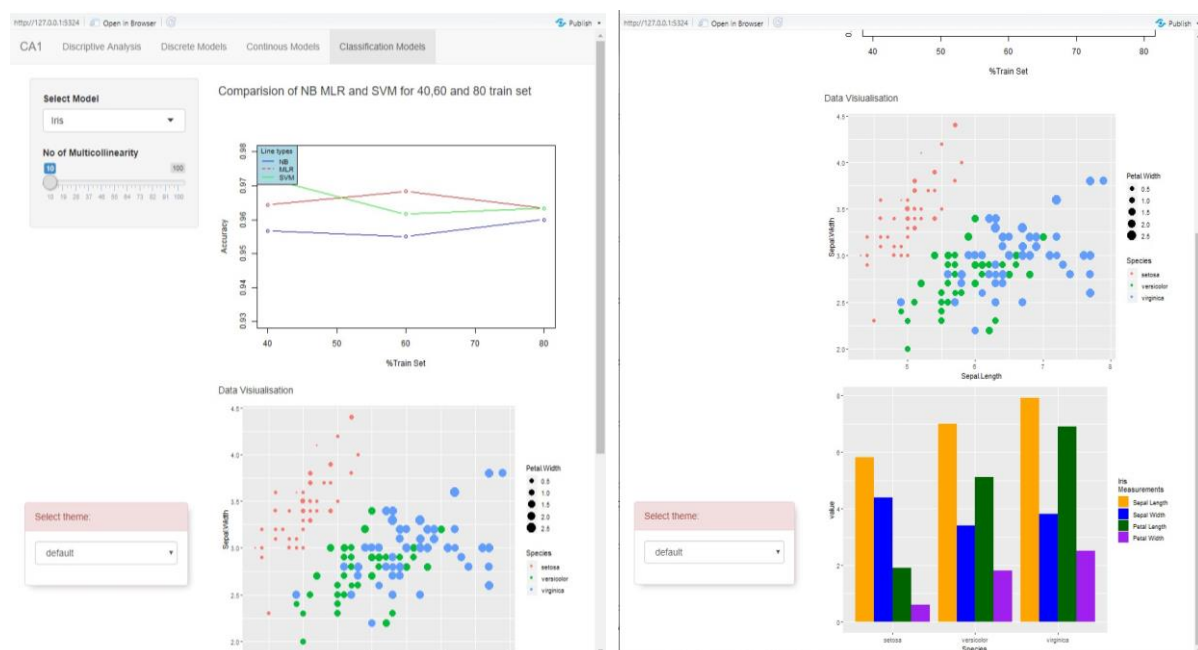
## 4.  Classification Models

A classification model tries to draw some conclusion from the input values given for training. It will predict the class labels/categories for the new data. In short Classification either predicts categorical class labels or classifies data (construct a model) based on the training set and the values (class labels) in classifying attributes and uses it in classifying new data. There are a number of classification models. Here in our app we are using classification on the basis of Naïve Bayes, Logistic Regression and Support Vector Method.

In the app developed, we have used an existing in-built dataset – 'Iris' and an existing dataset (which is needed to be uploaded) – 'Supermarket'. In both the datasets, the user can change the number of multi-collinearities as per requirement and can see the respective output and see a comparison between all the three attributes; i.e.; Naïve Bayes, Logistic Regression and Support Vector Method.
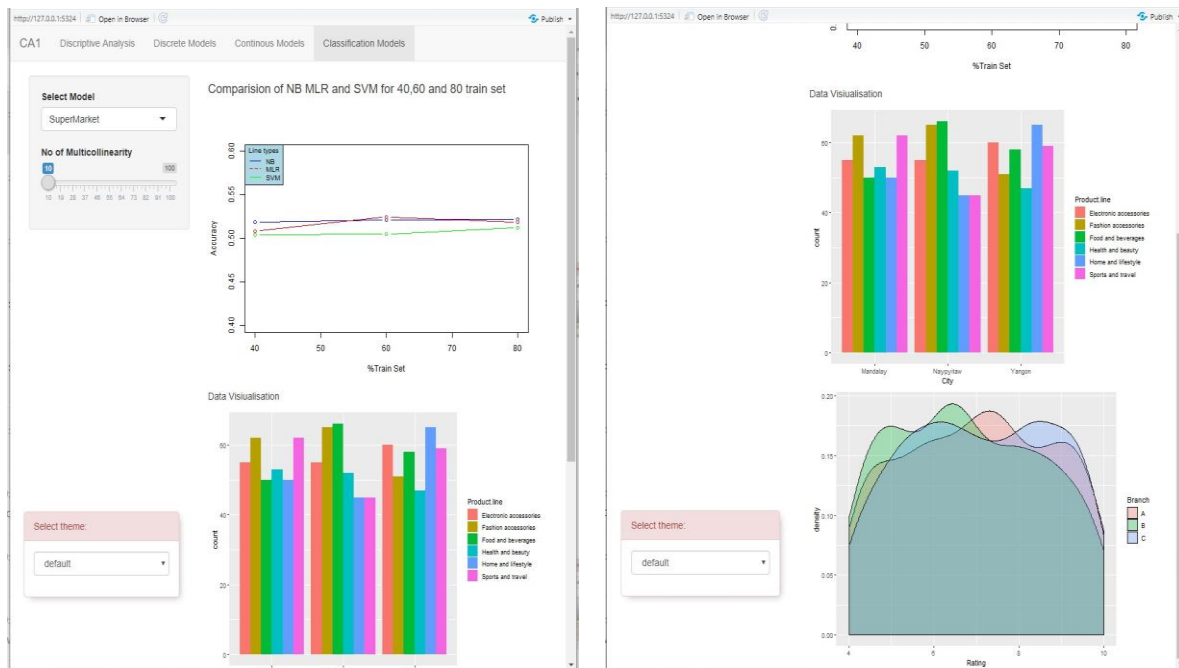
*Description of 'Iris' dataset:*

In Iris dataset, the comparison of NB, LR and SVM for 40, 60 and 80 train set is plotted in graph showing the relationship between percentage of train set and accuracy. Here, we can see as to which method have the maximum accuracy and at what percentage. Also, the Data Visualisation exemplifies the species, their size as per Petal Width and their sepal. It is distributed on the basis of Sepal length and Sepal Width, but also shows the species and the width of the petals of a particular species and hence covers four-dimensions in one graph. On the basis of this, a new graph is plotted which shows the value of all the four characteristics of the species.

*Description of* ***'Supermarket'*** *dataset:*

For Supermarket dataset, the user needs to download the dataset attached in his system and put its path in the program and he can then run the app. In this dataset the comparison of NB, LR and SVM for 40, 60 and 80 train set are plotted in graph showing the relationship between percentage of train set and accuracy (similarly, as it is done in 'Iris'). Here, we can see as to which method have the maximum accuracy and at what percentage. In this dataset, the Data Visualisation exemplifies the relation between city and count and represents six different product lines. It basically shows the product used by number of people in a particular city. The final graph visualises the density of rating verse each branch.



# Contribution towards Project:

Nidhi: Descriptive Analysis

Ankita: Discrete Models

Jyoti and Arvind: Continuous models

Classification Models and consolidating all tabs to shiny app : Varsha Ganapathy Rao