# Department of Computer Science & Engineering

## Final Year B. Tech. (CSE) – I: 2022-23

## 5CS462: PE5 - Data Mining Lab

## Assignment No. 5

## Name: Sanket Mote

## PRN: 2019BTECS00113

## Batch: B8

**Title:** Design and implement the following classifiers:

a) Regression classifier.

b) Naive Bayesian Classifier.

c) k-NN classifier

d) Three-layer Artificial Neural Network (ANN) classifier (use back propagation). Plot error graph (iteration vs error).

**Objective/Aim:**

1. To implement data analysis tool using python programming language.
2. To design the Regression, Naïve Bayes, k-NN, Three-layer Artificial Neural Network (ANN) classifier.

**Theory:**

**A) Regression classifier**:

Regression algorithms predict a continuous value based on the input variables. The main goal of regression problems is to estimate a mapping function based on the input and output variables. If your target variable is a quantity like income, scores, height or weight, or the probability of a binary category (like the probability of rain in particular regions), then you should use the regression model. However, there are various types of regressions used by data scientists and ML engineers based on different scenarios. The different types of regression algorithms include:

## 1. Simple linear regression

With simple linear regression, you can estimate the relationship between one independent variable and another dependent variable using a straight line, given both variables are quantitative.

## 2. Multiple linear regression

An extension of simple linear regression, multiple regression can predict the values of a dependent variable based on the values of two or more independent variables.

## 3. Polynomial regression

The main aim of polynomial regression is to model or find a nonlinear relationship between dependent and independent variables.

## B) Naïve Bayesian Classifier:

Naive Bayes classifiers are a collection of classification algorithms based on Bayes' Theorem. It is not a single algorithm but a family of algorithms where all of them share a common principle, i.e. every pair of features being classified is independent of each other

**Bayes' Theorem** finds the probability of an event occurring given the probability of another event that has already occurred. Bayes' theorem is stated mathematically as the following equation:

$$P(A|B) = \frac{P(B|A)P(A)}{P(B)}$$

where A and B are events and $P(B) \neq 0$.

Basically, we are trying to find probability of event A, given the event B is true. Event B is also termed as evidence.

P(A) is the priori of A (the prior probability, i.e. Probability of event before evidence is seen). The evidence is an attribute value of an unknown instance(here, it is event B).

P(A|B) is a posteriori probability of B, i.e. probability of event after evidence is seen.

## C) k-NN classifier:

K-Nearest Neighbour 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 similar to the available categories.

K-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 similar to the new data.

Example: Suppose, we have an image of a creature that looks similar to cat and dog, but we want to know either it is a cat or dog. So for this identification, we can use the KNN algorithm, as it works on a similarity measure. Our KNN model will find the similar features of the new data set to the cats and dogs images and based on the most similar features it will put it in either cat or dog category.
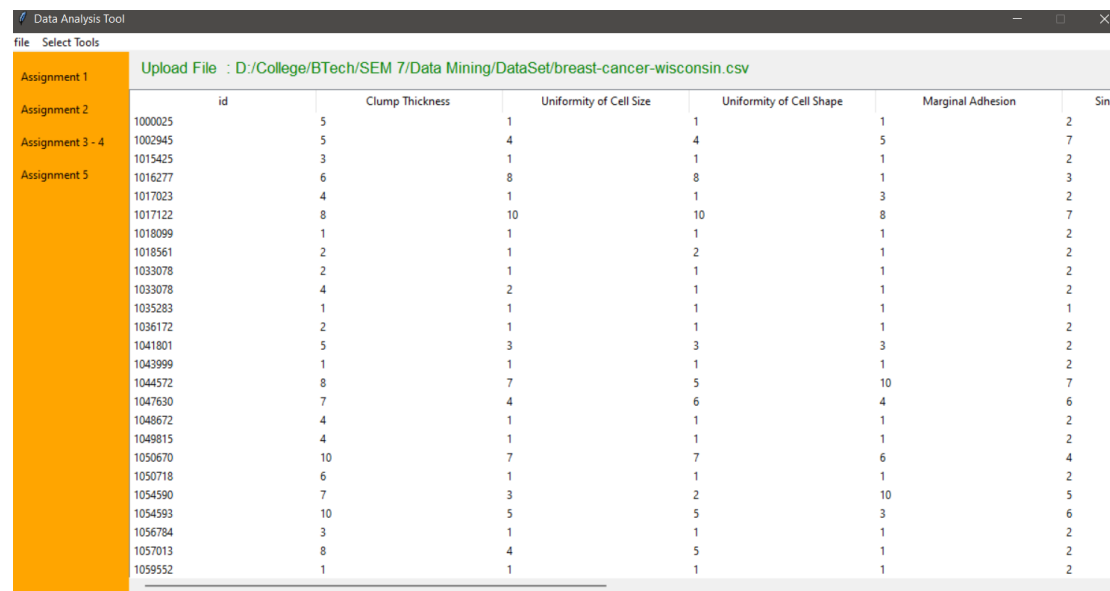
## D) Artificial Neural Network (ANN) classifier:

Classification ANNs seek to classify an observation as belonging to some discrete class as a function of the inputs. The input features (independent variables) can be categorical or numeric types, however, we require a categorical feature as the dependent variable.
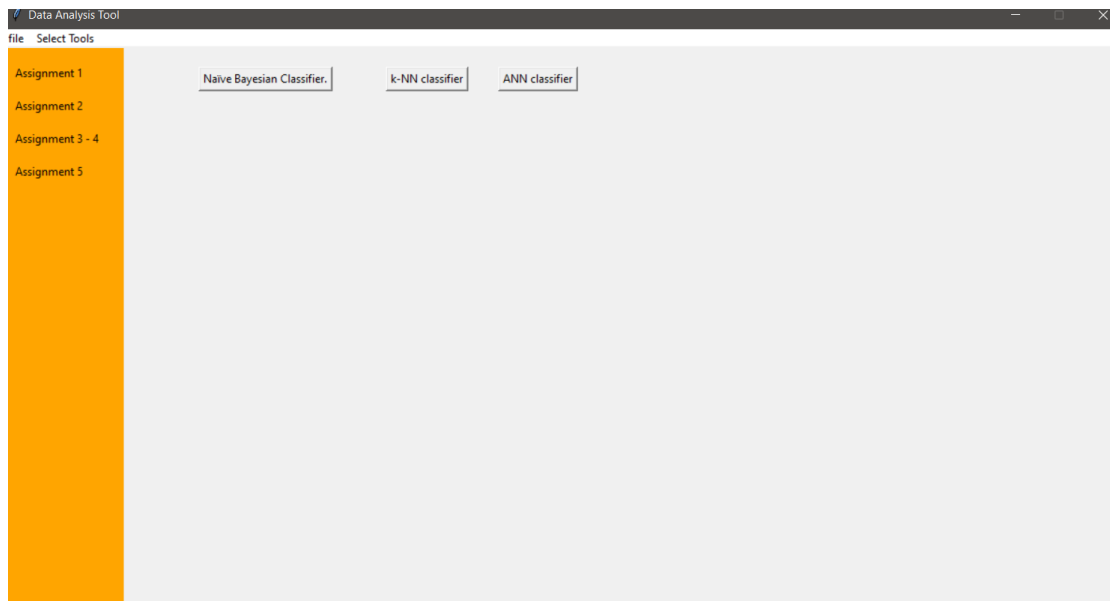
## Procedure:

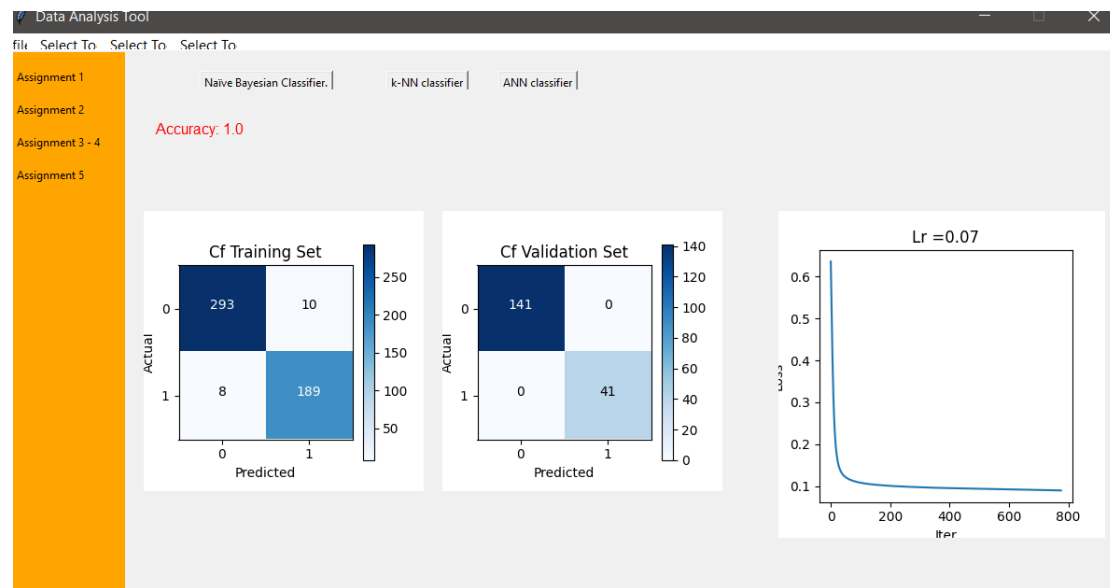Given problem statement is solved using python programming language and specifically used tkinter module to implement GUI application and pandas module to load .csv file as dataset.

## Results:

## Conclusion:

Successfully implemented data analysis tool (GUI) for extraction of rules from decision tree build based on information gain, gain ratio and gini index for selected attribute of given dataset.