# Guru Nanak Dev University

## Amritsar



**Department Of Electronics Technology**

A project based Report

**Multiple Disease Detection**

**Submitted By-**                                              **Submitted To –**

Aman Kumar                                                    Dr. Sukhdeep kaur

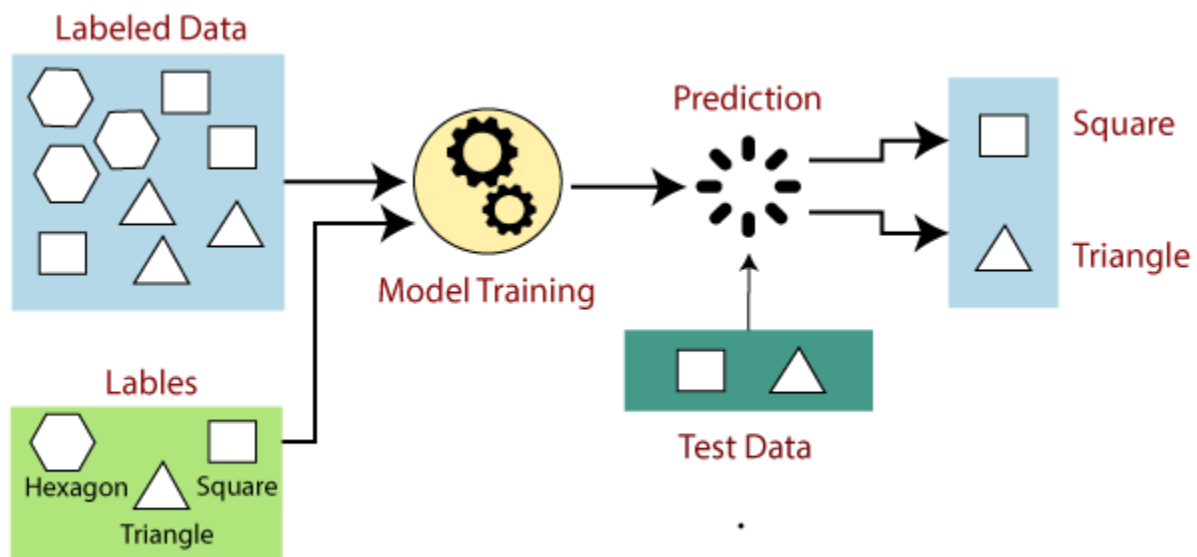Semester – VI  (ECE)

Roll number -17042011210

# Table Of Contents

# INTRODUCTION

Machine Learning Approach for Identifying Multiple Disease Prediction Using Machine Learning is based on prediction modelling that predicts disease of the patients according to the Laboratory test measures provided by the users as an i/p to the system. This report gives an idea of predicting multiple diseases using Machine Learning algorithms. Here we will use the concept of supervised Machine Learning in which implementation will be done by applying Support Vector Machine (SVM) , Logistic Regression algorithms which will help in early prediction of diseases accurately and better patients care. The results ensured that the system would be functional and user oriented for patients for timely diagnoses of diseases in a patient.

Supervised learning is the types of machine learning in which machines are trained using well "labelled" training data, and on basis of that data, machines predict the output. The labelled data means some input data is already tagged with the correct output.

In supervised learning, the training data provided to the machines work as the supervisor that teaches the machines to predict the output correctly. It applies the same concept as a student learns in the supervision of the teacher.

Supervised learning is a process of providing input data as well as correct output data to the machine learning model. The aim of a supervised learning algorithm is to **find a mapping function to map the input variable(x) with the output variable(y)**.

## Objective

1. Our main aim is to provide a quick medical diagnosis to the patients living in rural area.
2. In present days it is very useful for post covid contactless system in rural health service .
3. The goal is to provide access to medical specialists e-Doctor.
4. This system enhance quality of health care
5. This project helps in saving the money which we give for appointments to Doctors.
6. It help a lot of people else one monitor the persons condition and take the necessary precaution thus increasing the life expectancy

## Software required
1. Google Colab Notebook
2. Anaconda.
3. Pycharm.

## Supervised Learning Algorithm

1. Support Vector Machine (SVM).

   It can classifies both linear and non linear data . it first map each data item into and n-dimentional features space where n is the number of features. It then identify the hyperplane that separates the data items into two classes while maximising the marginal distance for both classes and minimizing the classification errors.

2. Logistic Regression

   Logistic regression is a powerful and well establish method for supervised classification. It can be considers as an extension of ordinary regression and can model only a dichotomous of an event . LR helps in finding the probability that a new instance belongs to a certain class. Since it is a probability ,the outcome lies between 0 and 1.
   This type of statistical model (also known as *logit model*) is often used for classification and predictive analytics. Logistic regression estimates the probability of an event occurring, such as voted or didn't vote, based on a given dataset of independent variables.

# Block diagram



# Working

So, before execution we have some pre-requisites that we need to download or install i.e., anaconda environment, Pycharm and a code editor (Google colab Notebook) .

**Anaconda**: Anaconda is like a package of libraries and offers a great deal of information which allows a data engineer to create multiple environments and install required libraries easy and neat.

**PyCharm :** PyCharm is a dedicated Python Integrated Development Environment (IDE) providing a wide range of essential tools for Python developers, tightly integrated to create a convenient environment for productive Python, web, and data science development.

**Google colab notebook**: Google colab allows anybody to write and execute arbitrary python code through the browser, and is especially well suited to machine learning, data analysis and education.

Install the prerequisites mentioned above.

**Step1 – Set Environment in Anaconda**
Open anaconda prompt and create a new environment. To create an environment use the commands given below. Replace MachineLearning by the name of environment you want to give.

- conda create -n "MachineLearning"
- conda activate "MachineLearning"

## Step2

open the terminal of our above environment and install the package Streamlit using below command

- pip install streamlit



After that we will be able to see the default browser of streamlit

**Step3**

To run the code, type "streamlit run "path/filename.extension" in our case *streamlit run "D:\Exicuted Projects\Multiple Disease Prediction System\multiple_disease_pred.py"*

Open the environment terminal and run the code by above command .

Here we need to give the address of python code file which in our case is pycharm file.

# Data Description

In this project we will be doing multiple disease prediction the diseases are Diabetes, Heart disease and Parkinsons and Stroke and all the data sets used are tabular data .

- **Diabetes dataset**



- **Heart Disease Dataset**

- **Parkinson disease dataset**



# Preprocessing the data

Check missing values in all the disease datasets . if dataset consist of missing values .we must handle those missing values

- dataset.isnull().sum()

- Standardize the data



# Splitting the data (Training and Testing)

what we are doing is we are taking four variables means x split into x train and x test once the model is trained we are try to evaluate our model with the test data .

- ```
  X_train, X_test, Y_train, Y_test = train_test_split(X,Y, test_size =0.2 , strat
  ify= Y,random_state=2)
  ```



Test size means how much data that u want for test the data here this is 0.2 means 20 percent of the data we test here .

stratify is Y actually Y basically as the values as either one and zero so we want our dataset to be splited into same proportion .

# Apply Machine learning Algorithm



```
CO  Copy of Diabetes.ipynb
    File  Edit  View  Insert  Runtime  Tools  Help   All changes saved

+ Code   + Text                                                          Connect ▾   ^
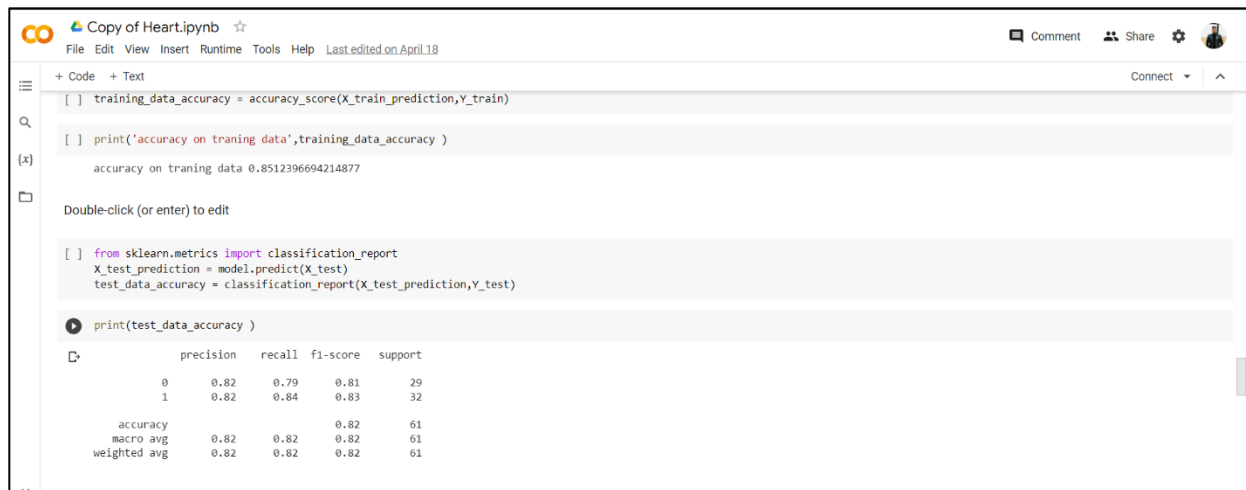
[ ]  classifier = svm.SVC(kernel='linear')

[ ]  # training the support vector machine classifier
     classifier.fit(X_train,Y_train)

     SVC(kernel='linear')
```

# Model evaluation and Result

Accuracy without SMOTE
**Accuracy is => 82 %**



```
CO  Copy of Heart.ipynb
    File  Edit  View  Insert  Runtime  Tools  Help   Last edited on April 18

+ Code   + Text                                                          Connect ▾   ^

[ ]  training_data_accuracy = accuracy_score(X_train_prediction,Y_train)

[ ]  print('accuracy on traning data',training_data_accuracy )

     accuracy on traning data 0.8512396694214877

Double-click (or enter) to edit

[ ]  from sklearn.metrics import classification_report
     X_test_prediction = model.predict(X_test)
     test_data_accuracy = classification_report(X_test_prediction,Y_test)

  ▶  print(test_data_accuracy )

              precision    recall  f1-score   support

          0       0.82      0.79      0.81        29
          1       0.82      0.84      0.83        32

   accuracy                          0.82        61
  macro avg       0.82      0.82      0.82        61
weighted avg      0.82      0.82      0.82        61
```

SMOTE :- Synthetic Minority Oversampling Technique (SMOTE) is a statistical technique for increasing the number of cases in your dataset in a balanced way. The component works by generating new instances from existing minority cases that you supply as input.

Accuracy with SMOTE
**Accuracy is => 84 %**



## Compare Accuracy with or without Smote

**Conclusion**

The main aim of this report is to predict the disease in accordance with parameters of Laboratory reports put down by the patients with proper implementation of Machine Learning algorithm. In this report we have used 2 Machine Learning algorithm for prediction and achieved the mean accuracy of more than 84% which shows remarkable rectification and high accuracy than previous work and also makes this system more reliable than the existing one for this job and hence provides better satisfaction to the user in comparison with the other one. It also stores the data entered by the user and the name of the disease the patient is suffering from in the Database which can be used as past record and will help in future for future treatment and thus contributing in easier health management .We have also created a GUI for better interaction with the system by users which is very easy to operate .This report shows that Machine Learning algorithm can be used to predict the disease easily with different parameters and models. In the end we can say that our system has no threshold of the users because everyone can use this system.

**References :-**

Datasets –
- Diabetes
- Heart
- Parkinson's

Theory-

- Machine Learning
- Supervised Learning

Software-

- Google Colab Notebook
- Anaconda.