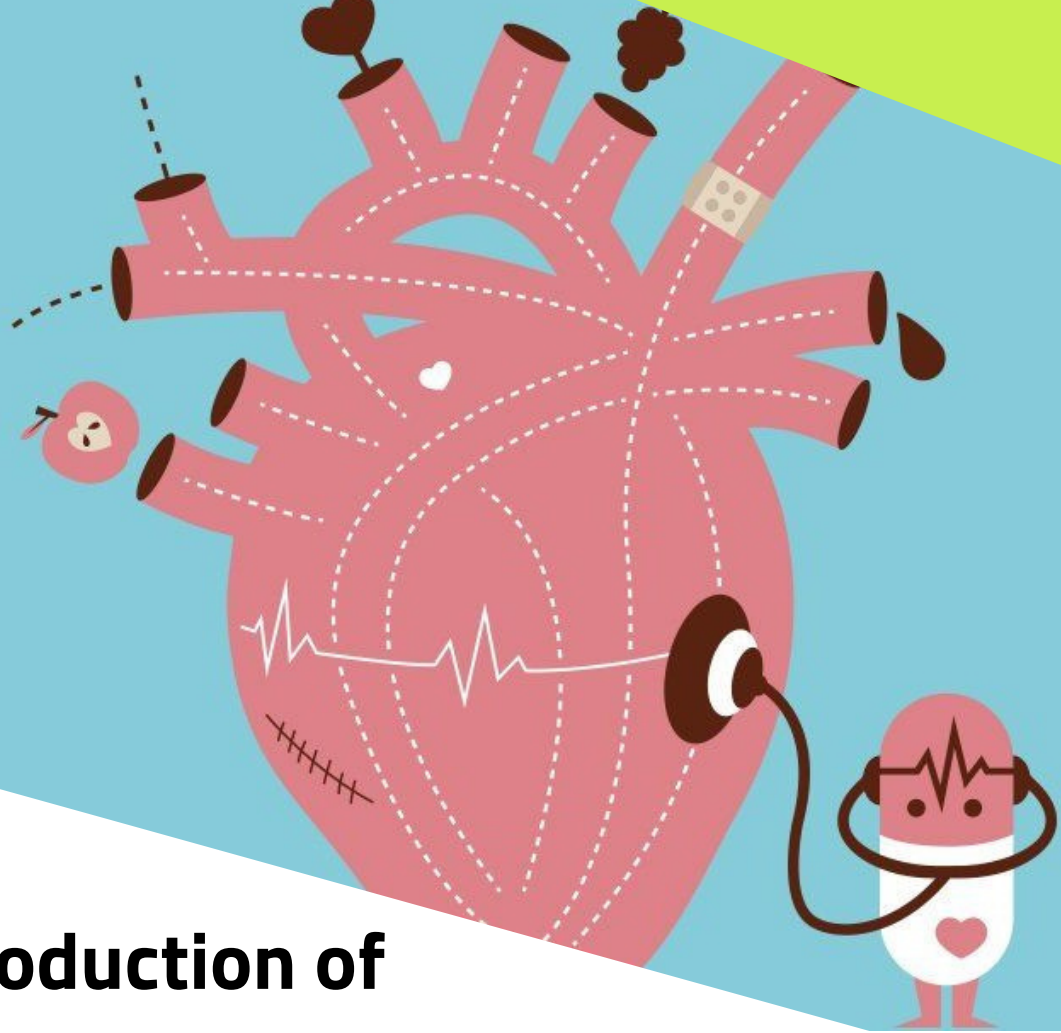**DATA SCIENCE**

# CARDIOVASCULAR DISEASE REPORT

May 2021

Prepared by
**Sarthak Chandel**

# Introduction of problem

Heart diseases or CVDs are the number one cause of death globally with more than 18 million deaths every year . CVDs are concertedly contributed by hypertension , diabetes , overweight and unhealthy lifestyles .

This report trains the model on SVM and KNN algorithms with the dataset to predict the probability of a person suffering from a heart condition.

# Dataset Description

The dataset consists of 70,000 records of patients data,
11 features + 1 target.
There are 3 types of input features:

**4**

OBJECTIVE

**4**

EXAMINATION

**3**

SUBJECTIVE

1. Age | Objective Feature | int (days)
2. Height | Objective Feature | int (cm) |
3. Weight | Objective Feature | float (kg) |
4. Gender | Objective Feature | categorical code |
5. Systolic blood pressure | Examination Feature | int |
6. Diastolic blood pressure | Examination Feature | int |
7. Cholesterol | Examination Feature | 1: normal, 2: above normal, 3: well above normal |
8. Glucose | Examination Feature | 1: normal, 2: above normal, 3: well above normal |
9. Smoking | Subjective Feature | binary |
10. Alcohol intake | Subjective Feature | binary |
11. Physical activity | Subjective Feature | binary |
12. Presence or absence of cardiovascular disease | Target Variable| binary |
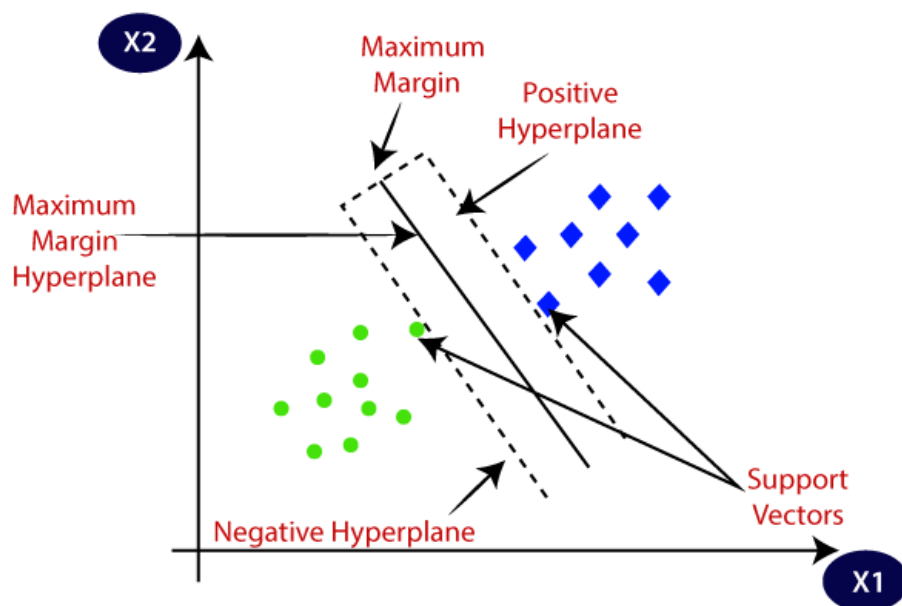
# Algorithms Used

## Support Vector Machine

Support Vector Machine or SVM is one of the most popular Supervised Learning algorithms, which is used for Classification.

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.

SVM chooses the extreme points/vectors that help in creating the hyperplane. These extreme cases are called as support vectors, and hence algorithm is termed as Support Vector Machine. Consider the below diagram in which there are two different categories that are classified using a decision boundary or hyperplane:



The distance between the hyperplane and the nearest data point from either set is known as the margin. The goal is to choose a hyperplane with the greatest possible margin between the hyperplane and any point within the training set, giving a greater chance of new data being classified correctly.
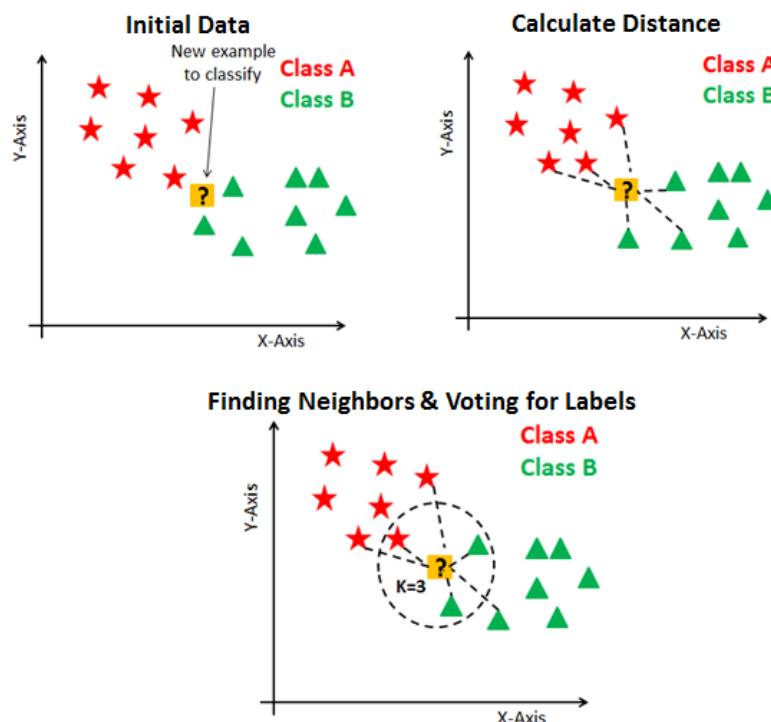
# Algorithms Used

## K- Nearest Neighbors

The K-nearest neighbors (KNN) algorithm is a supervised learning algorithm used to solve both regression and classification problems.
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.
The input consists of the k closest training example in the data set.
To make a prediction, the KNN algorithm will use the entire dataset. For an observation that is not part of the dataset and is the actual value we want to predict, the algorithm will look for the K instances of the dataset closest to our observation.
Then for these K neighbors, the algorithm will use their output to calculate the variable y of the observation that we want to predict.



A common formula used in this algorithm to find the distance.

$$\text{Euclidean distance} = \sqrt{\sum_{i=1}^{n}\left(x_i - y_i\right)^2}$$

# Data Visualization

## Box Plot



From this box plot we can spot the outliers. With extreme abnormalities in the ap_hi and ap_lo or the blood pressure columns.
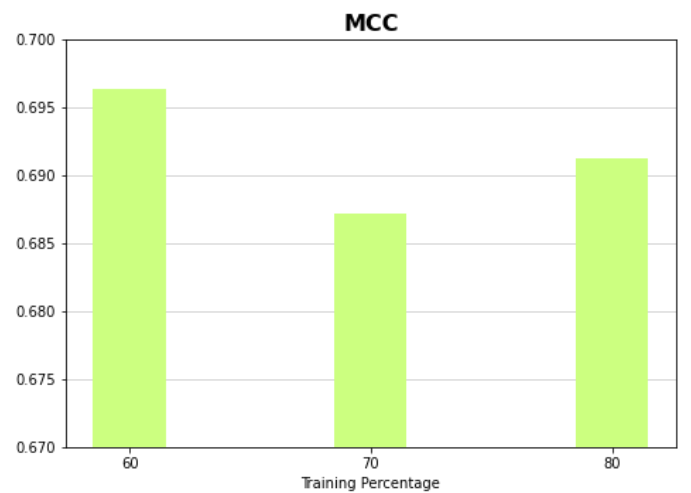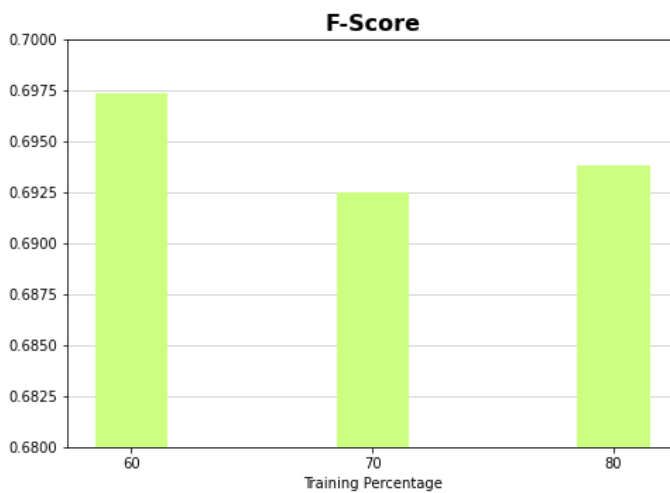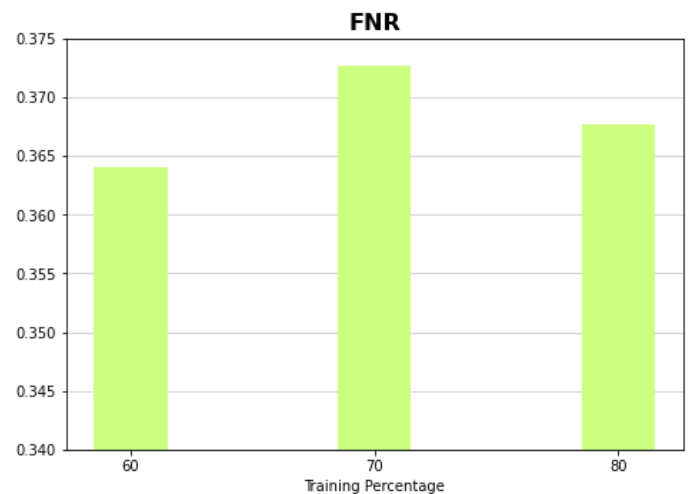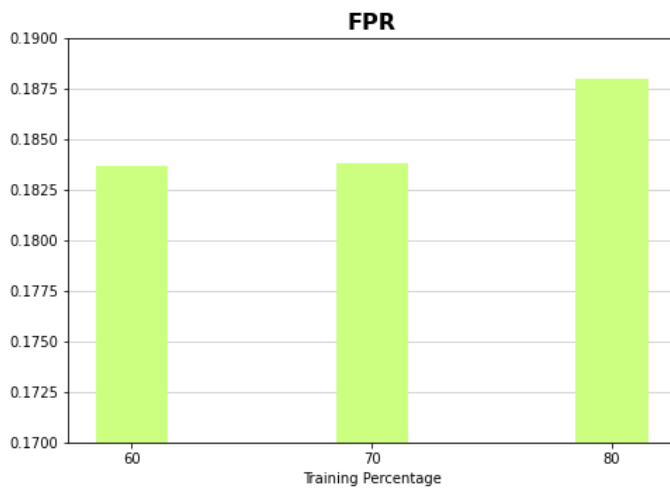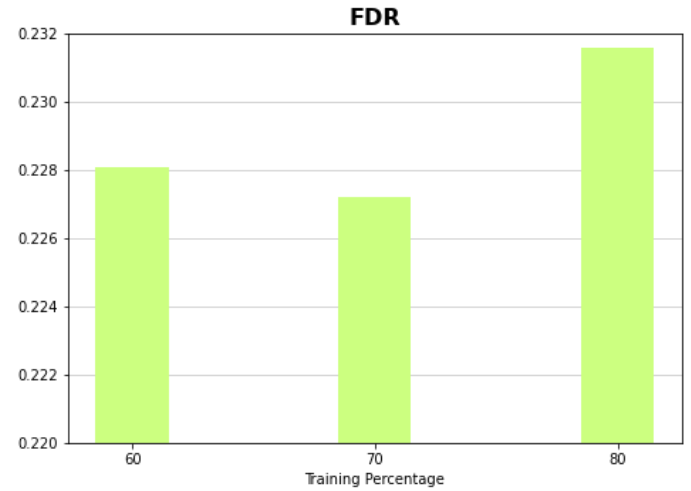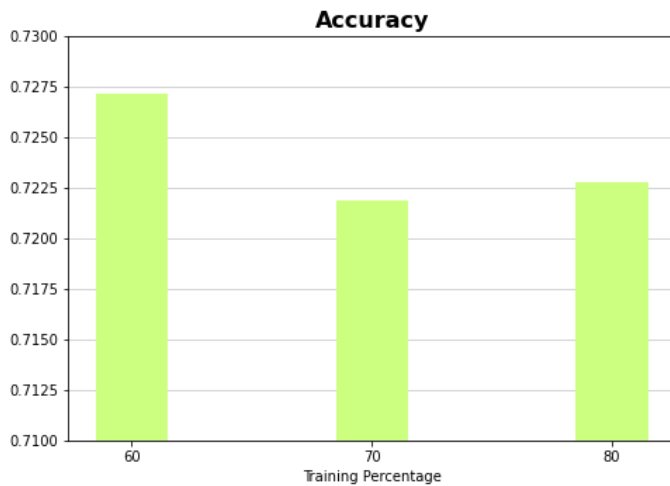
# Data Visualization

## Heatmap



This heatmap shows us the strong correlation between cholesterol, blood pressure, age, glycogen, cholesterol levels and cardiovascular diseases
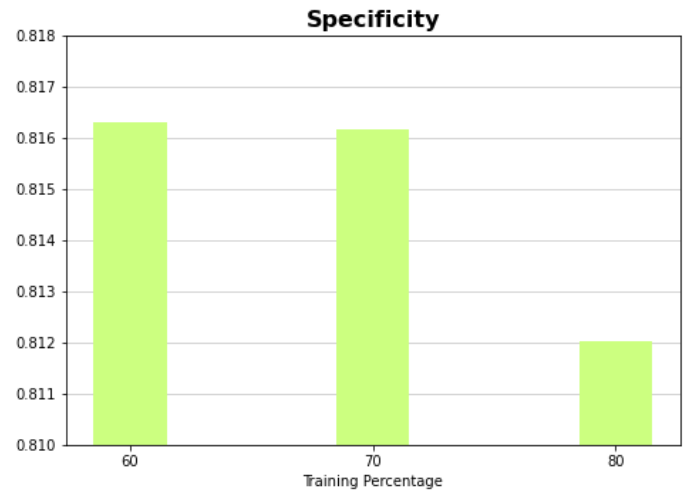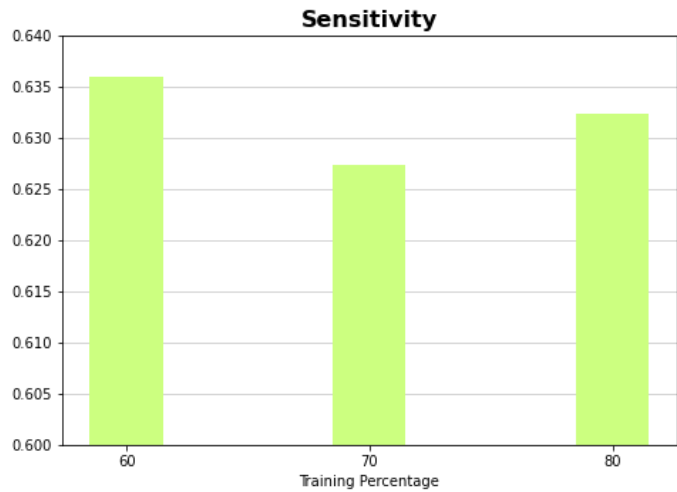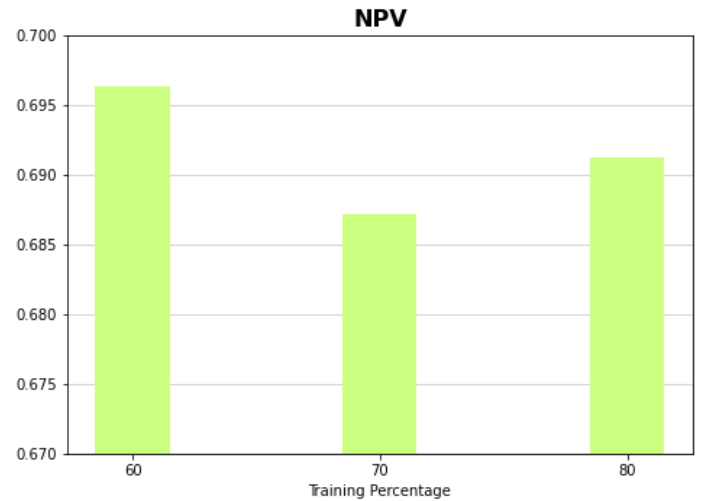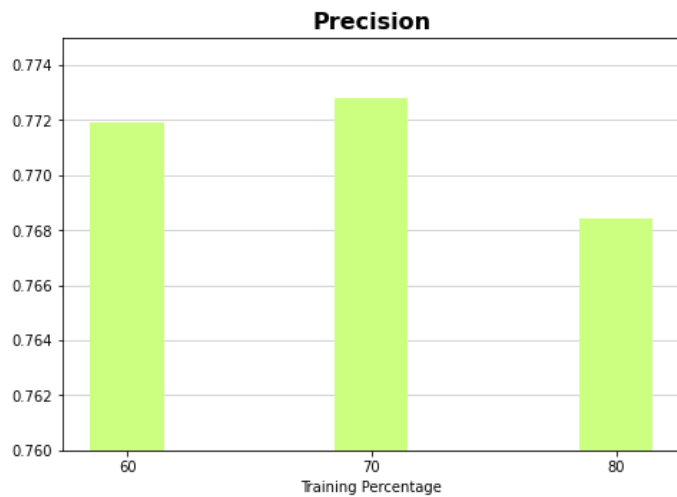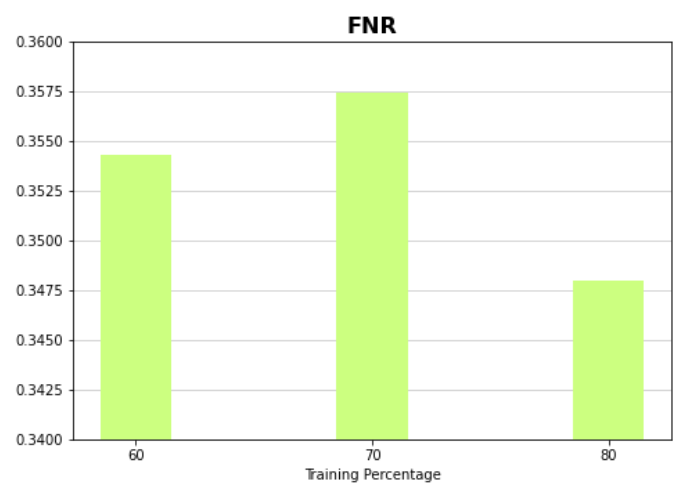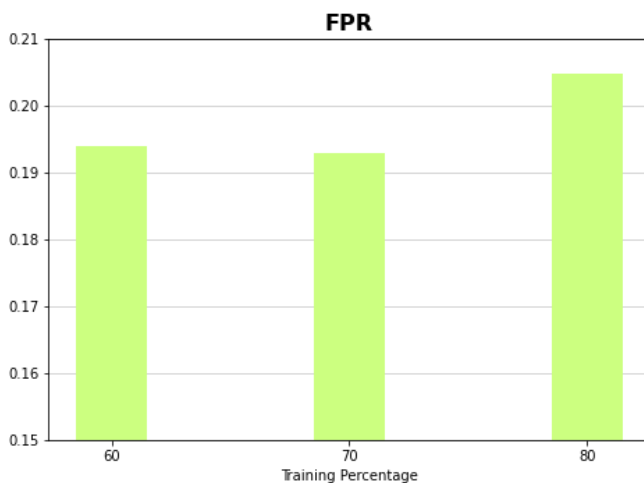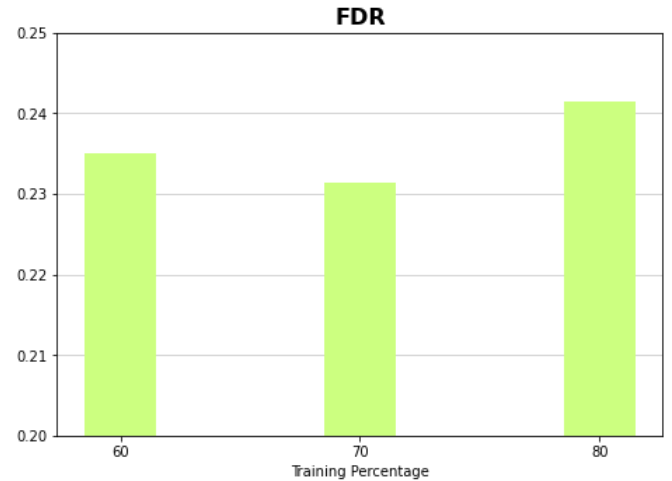
# Results

## SVM

# Results

## SVM

# Results

## KNN



**Accuracy**

**FDR**

**FPR**

**FNR**

**F-Score**

**MCC**

# Results

## KNN

# Result Analysis

## Confusion Matrix

| 90:10 Split | Predicted Negative | Predicted Positive |
|---|---|---|
| Actual Negative | 2855 | 628 |
| Actual Positive | 1265 | 2153 |

**SVM**

| 90:10 Split | Predicted Negative | Predicted Positive |
|---|---|---|
| Actual Negative | 2767 | 681 |
| Actual Positive | 1189 | 2264 |

**KNN**

## ROC



**SVM**



**KNN**

# Result Analysis

## Table for evaluation parameter

| Algorithms | SVM | | | KNN | | |
|---|---|---|---|---|---|---|
| Measures | Train test ratio 1 | Train test ratio 2 | Train test ratio 3 | Train test ratio 1 | Train test ratio 2 | Train test ratio 3 |
| Specificity | 0.8163 | 0.8162 | 0.8120 | 0.8060 | 0.8072 | 0.7953 |
| Sensitivity | 0.6359 | 0.6273 | 0.6323 | 0.6457 | 0.6426 | 0.6520 |
| Accuracy | 0.7271 | 0.7219 | 0.7228 | 0.7268 | 0.7250 | 0.7242 |
| Precision | 0.7719 | 0.7728 | 0.7684 | 0.7649 | 0.7686 | 0.7586 |
| FPR | 0.1837 | 0.1838 | 0.1880 | 0.1940 | 0.1928 | 0.2047 |
| FNR | 0.3641 | 0.3727 | 0.3677 | 0.3543 | 0.3574 | 0.3480 |
| NPV | 0.6964 | 0.6872 | 0.6913 | 0.6995 | 0.6938 | 0.6985 |
| FDR | 0.2281 | 0.2272 | 0.2316 | 0.2351 | 0.2314 | 0.2414 |
| F-Score | 0.6973 | 0.6925 | 0.6938 | 0.7003 | 0.6999 | 0.7013 |
| MCC | 0.4602 | 0.4516 | 0.4520 | 0.4580 | 0.4560 | 0.4522 |

SVM and KNN both seem to have similar values across all measures. In this case, using KNN would be beneficial as it has a lower FNR rate and false negatives in this case could be life threatening.

Having a slightly higher FPR is not a major issue in this study.