# PR REPORT

# WINE CLASSIFICATION

Course Instructor: Dr. Anish Chand Turlapaty

Group: 6

KAKARLA V S S PAVAN TEJA        (S20190020216)

K SREENIVASULU REDDY        (S20190020217)

K LITHEESH KUMAR        (S20190020218)

**Introduction:**

Here we will classify the wine on the basis of giving features. We used the wine quality dataset downloaded from Kaggle. This dataset has the fundamental features which are responsible for affecting the quality of the wine. By the use of several Machine learning models, we will classify the wine. We use classification techniques to check further the correlation of fundamental features like

- fixed acidity
- volatile acidity
- citric acid
- residual sugars
- chlorides
- free sulphur dioxide
- density
- pH
- sulphates
- alcohol
- total sulphur dioxide
- quality

**Implementation:**

1. Downloaded wine quality dataset from Kaggle
2. Combined both white wine and red wine datasets
    - White wine -1
    - Red wine - 0
3. Normalized the features so that equal distribution of variance will be there for all features
4. Trained these features to Logistic Regression, Gaussian SVM, Neural Network Classifiers and obtained confusion matrices, accuracy, ROC Curves
5. Tested the data using above trained models, result is the test accuracy

- Principal Component Analysis, or PCA, is a dimensionality-reduction method that is often used to reduce the dimensionality of large data sets, by transforming a large set of variables into a smaller one that still contains most of the information in the large set.
- Because smaller data sets are easier to explore and visualize and make analysing data much easier and faster
- PCA is keeping enough components to explain 95% variance. After training, 9 components were kept.
- Explained variance per component (in order): 27.5%, 22.3%, 14.4%, 8.8%, 6.6%, 5.6%, 4.8%, 4.6%, 3.1%, 2.1% (variances of least important components hidden)

# Classifiers:

**Logistic Regression:**

Logistic regression is a classification algorithm, used when the value of the target variable is categorized in nature. Logistic regression is most commonly used when the data in the problem has binary output, so when it belongs to one class or another, or is either 0 or 1.

It is important to understand that logistic regression should only be used when the target variables fall into discrete categories and that if there's a range of continuous values the target value might be, logistic regression should not be used.
Some examples are
- Predicting if an email is spam or not
- Tumor is malignant or benign
- Mushroom is poisonous or edible
- Here Wine is White or Red

**Gaussian SVM:**

An SVM classifies data by finding the best hyperplane that separates all data points of one class from those of other class. Gaussian Kernel is a way of measuring the similarity between two training examples in the SVM
The gaussian kernel basically assigns a score proportional to the nearness of the query points; This means that highly varying terrains can be represented.

**Linear Discriminant:**

Linear Discriminant Analysis easily handles the case where the within-class frequencies are unequal and their performances has been examined on randomly generated test data. This method maximizes the ratio of between-class variance to the within-class variance in any particular data set thereby guaranteeing maximal separability.
The use of Linear Discriminant Analysis for data classification is applied to classification problem in speech recognition. We decided to implement an algorithm for LDA in hopes of providing better classification compared to Principal Components Analysis. The prime difference between LDA and PCA is that PCA does more of feature classification and LDA does data classification.
In PCA, the shape and location of the original data sets changes when transformed to a different space whereas LDA doesn't change the location but only tries to provide more class separability and draw a decision region between the given classes. This method also helps to better understand the distribution of the feature data.
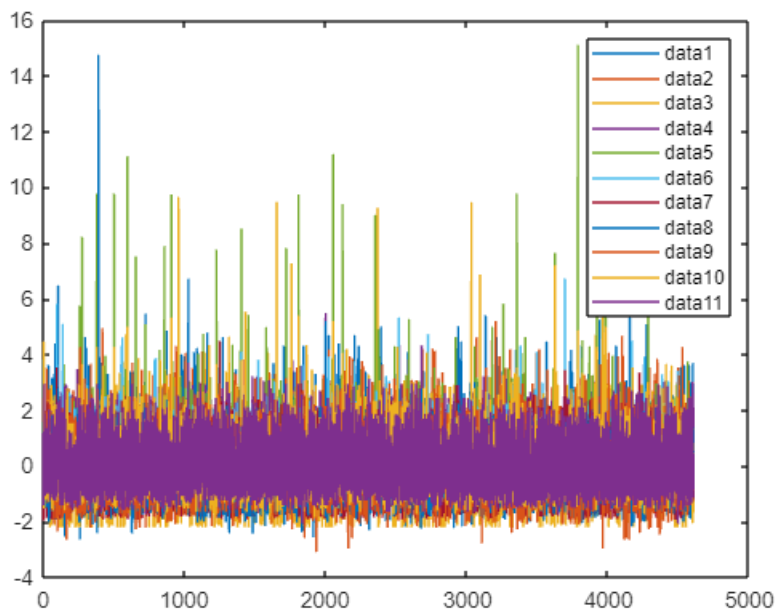
**Neural Network:**

Neural networks are complex models, which try to mimic the way the human brain develops classification rules. A neural net consists of many different layers of neurons, with each layer receiving inputs from previous layers, and passing outputs to further layers.

An alternative way of thinking about a neural net is to think of it as one massive function which takes inputs and arrives at a final output

The neurons within the network interact with the neurons in the next layer, with every output acting as an input for a future function. Every function, including the initial neuron receives a numeric input, and produces a numeric output, based on an internalized function, which includes the addition of a bias term, which is unique for every neuron. That output is then converted to the numeric input for the function in the next layer, by being multiplied with an appropriate weight. This continues until one final output for the network is produced
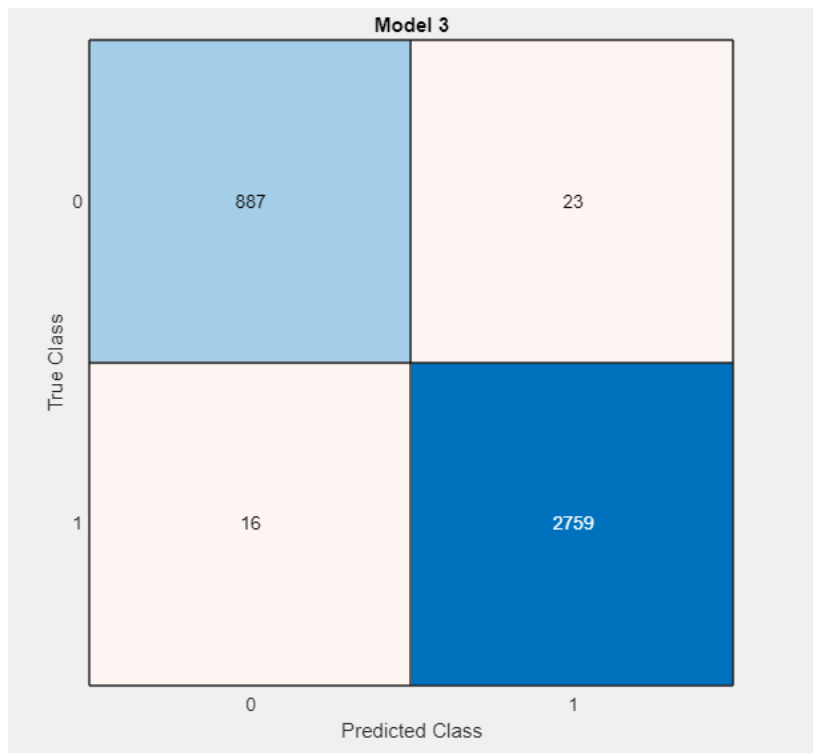
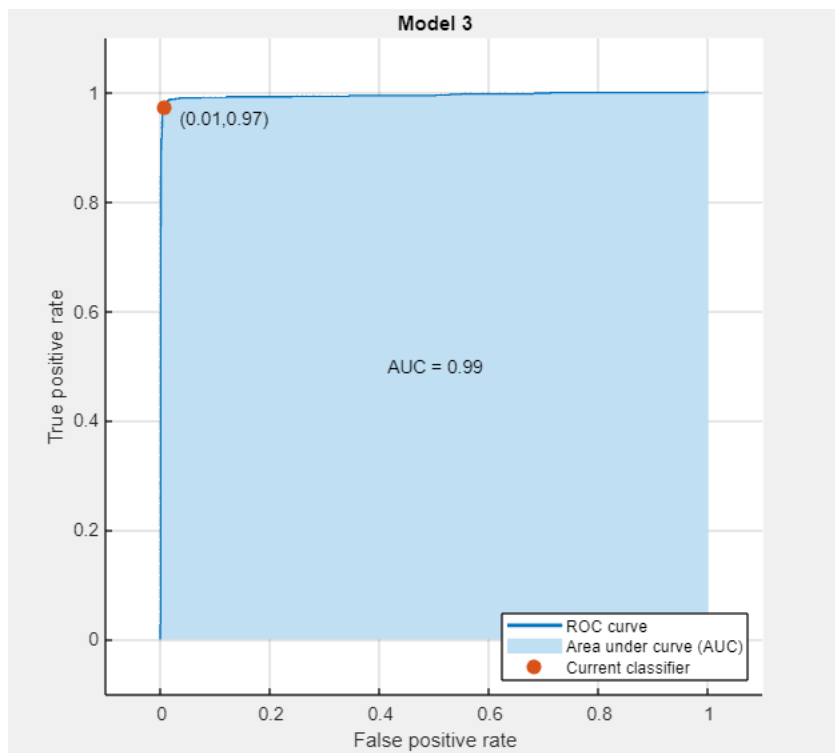## Results:



**Logistic Regression:**

1. Validation Accuracy – 99.3%
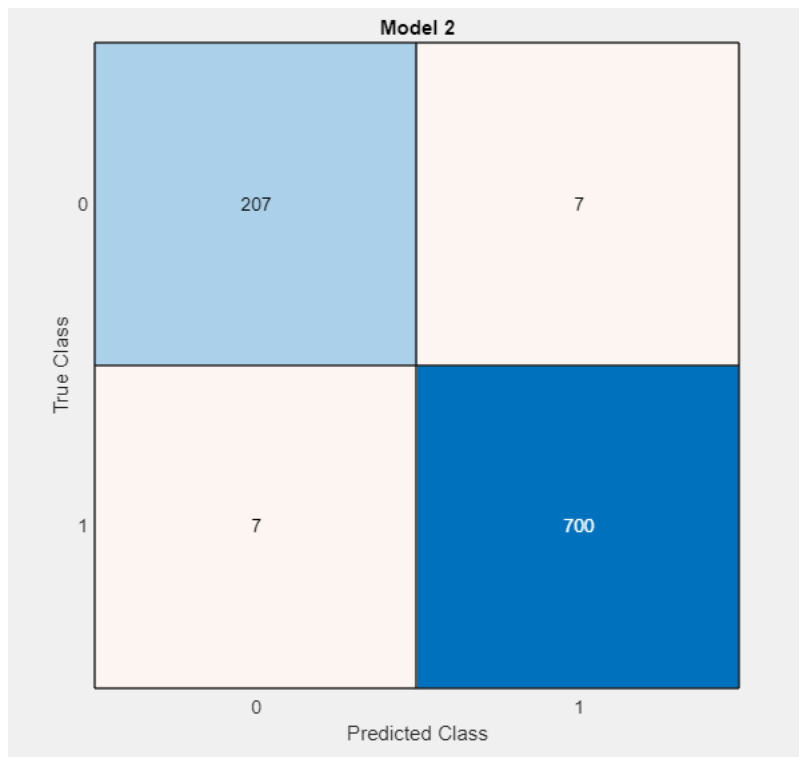2. Test Accuracy – 99%

**Training:**

Confusion Matrix:
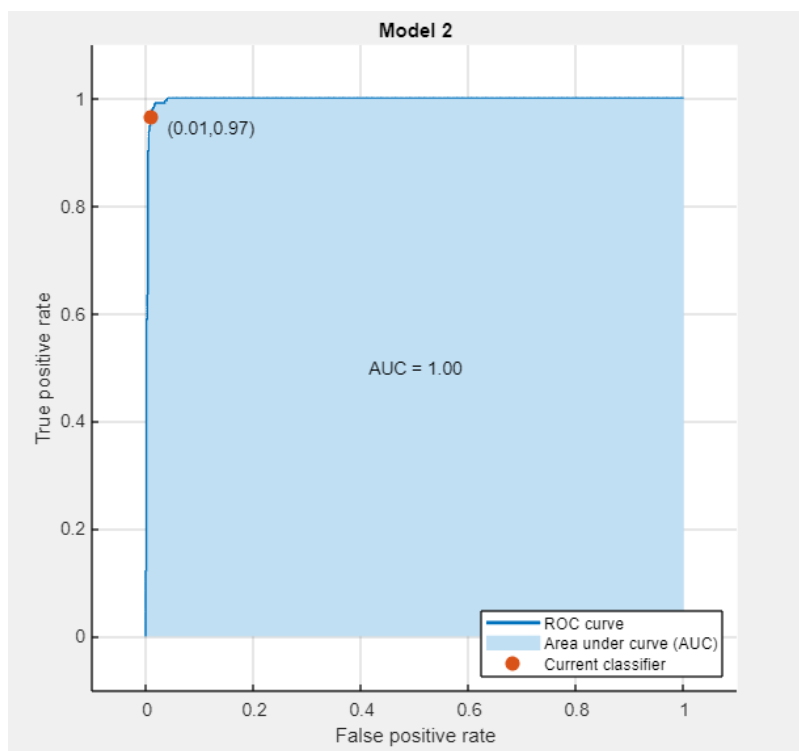


ROC Curve:

**Testing:**
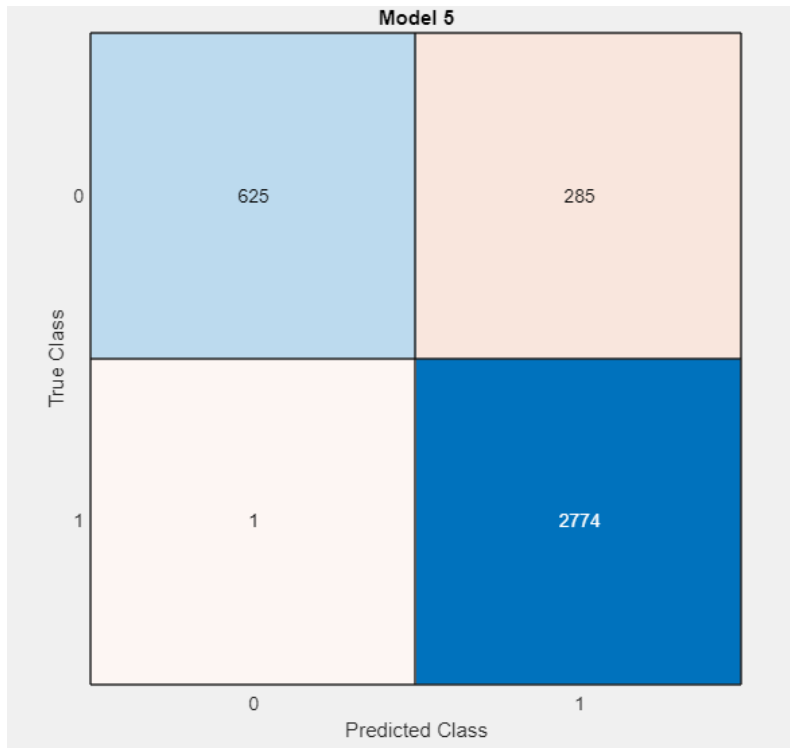
Confusion Matrix:



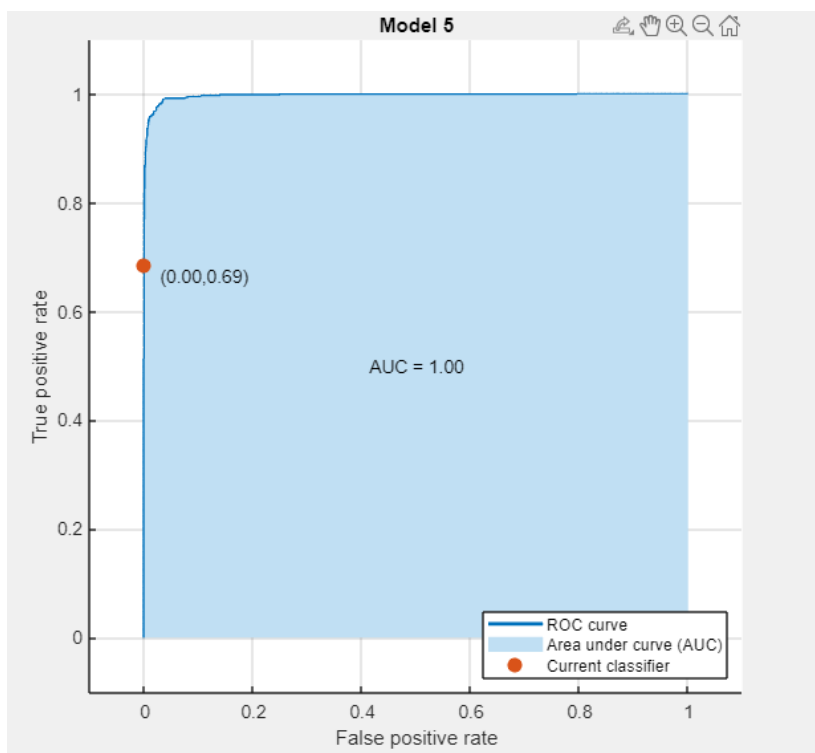ROC Curve:

**Gaussian SVM:**

1. Validation Accuracy – **92.5%**
2. Test Accuracy – **89%**
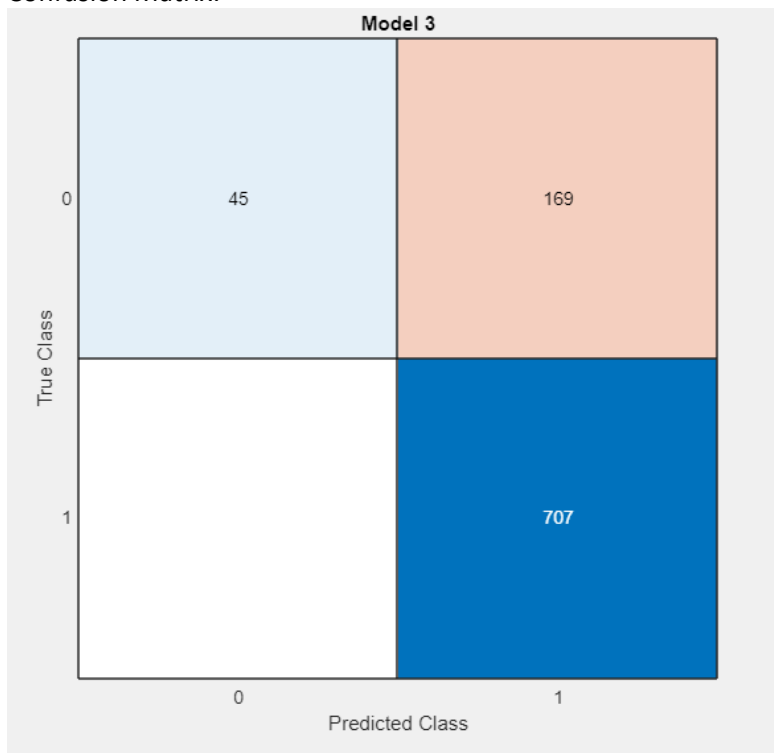
**Training:**

Confusion Matrix:
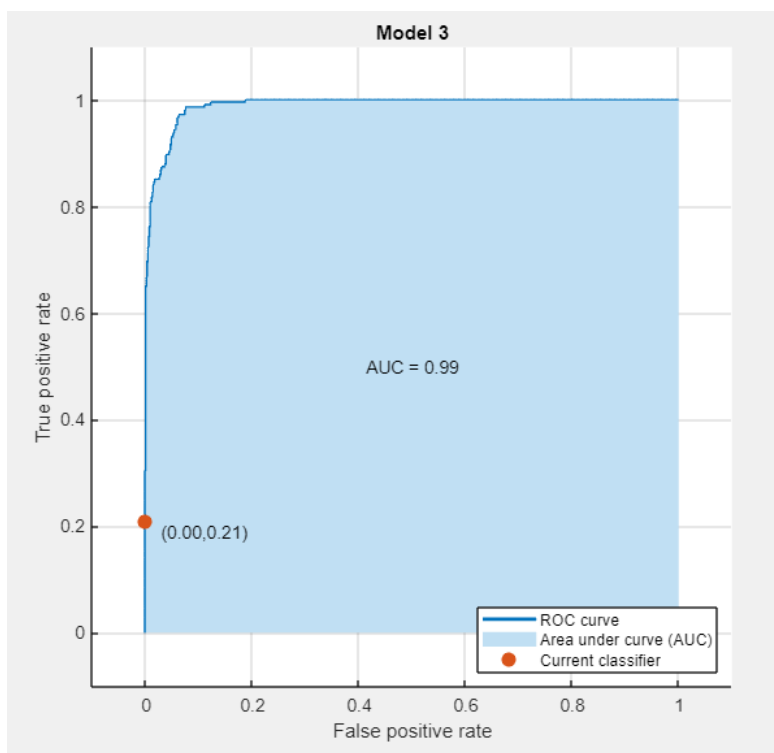


ROC Curve:

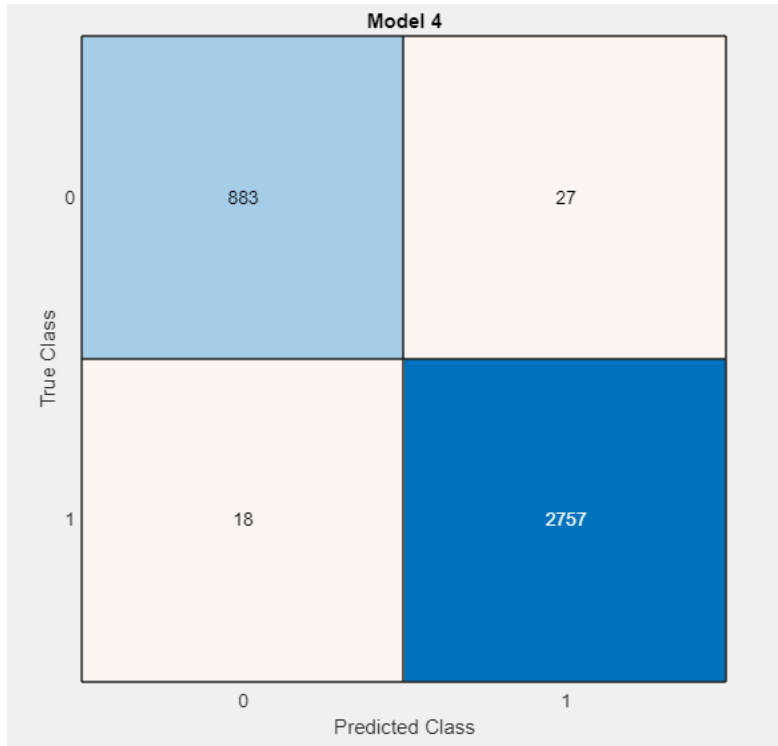**Tested:**

Confusion Matrix:



ROC Curve:

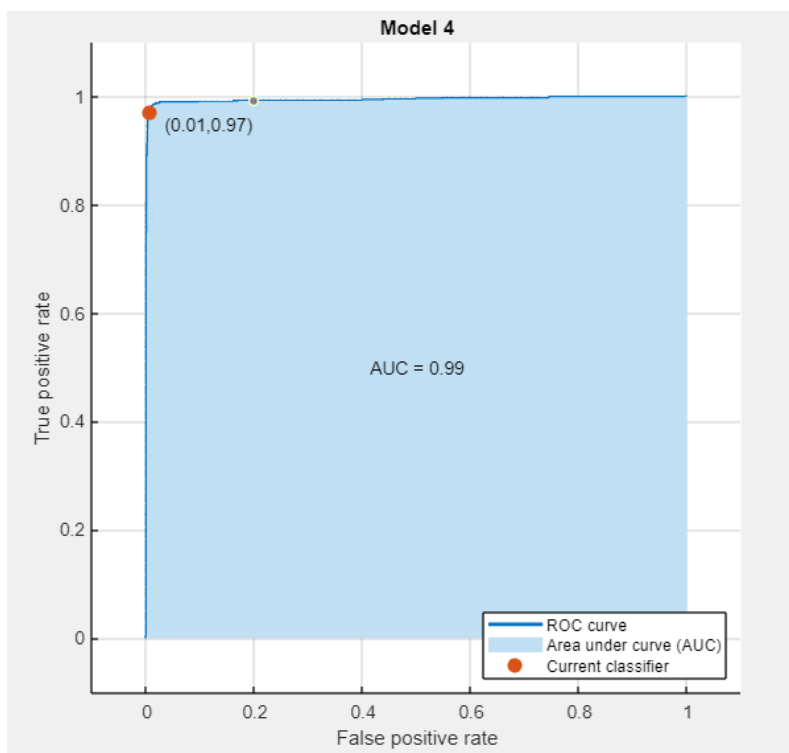**Linear Discriminant:**
1. Validation Accuracy— 97.5%
2. Test Accuracy – 98.3%

**Training:**

Confusion Matrix:
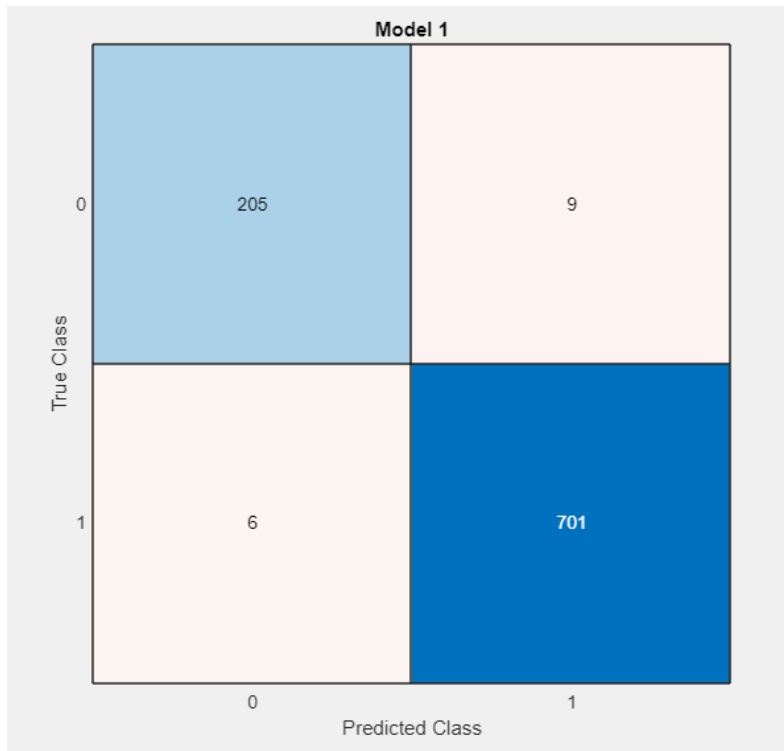


ROC Curve:

**Testing:**

Confusion Matrix:



ROC Curve:

**Neural Networks:**
1. Validation Accuracy – **99.3%**
2. Test Accuracy – **99.6%**

**Training:**
Confusion Matrix:



ROC Curve:

**Tested:**

Confusion Matrix:



ROC Curve:

# Wine Quality:

- o   We will predict the quality of wine on the basis of giving features.
- o   Downloaded wine quality dataset from Kaggle
- o   Dataset has the fundamental features which are responsible for affecting the quality of wine
- o   Debited correlations of fundamental features and quality

## Data Visualization:

**Histograph:**

# Quality vs Fundamental Features



quality vs alcohol

quality vs fixed acidity

quality vs density

quality vs chlorides

quality vs total sulphur dioxide

quality vs citric acid

quality vs pH

quality vs sulphides

**Correlation:**

Here we used statistical method which is used to evaluate the strength of bonding of the relationship between two quantitative variables.



```
 1  |      fixed acidity  volatile acidity  citric acid  residual sugar  chlorides  \
 2  0               7.0             0.270         0.36            20.7      0.045
 3  1               6.3             0.300         0.34             1.6      0.049
 4  2               8.1             0.280         0.40             6.9      0.050
 5  3               7.2             0.230         0.32             8.5      0.058
 6  4               7.2             0.230         0.32             8.5      0.058
 7  ...             ...               ...          ...             ...        ...
 8  6492            6.2             0.600         0.08             2.0      0.090
 9  6493            5.9             0.550         0.10             2.2      0.062
10  6494            6.3             0.510         0.13             2.3      0.076
11  6495            5.9             0.645         0.12             2.0      0.075
12 ∨ 6496           6.0             0.310         0.47             3.6      0.067
13  |
14  |      free sulfur dioxide  density    pH  sulphates  alcohol  quality  \
15  0                     45.0  1.00100  3.00       0.45      8.8        6
16  1                     14.0  0.99400  3.30       0.49      9.5        8
17  2                     30.0  0.99510  3.26       0.44     10.1        6
18  3                     47.0  0.99560  3.19       0.40      9.9        6
19  4                     47.0  0.99560  3.19       0.40      9.9        5
20  ...                    ...      ...   ...        ...      ...      ...
21  6492                  32.0  0.99490  3.45       0.58     10.5        5
22  6493                  39.0  0.99512  3.52       0.76     11.2        6
23  6494                  29.0  0.99574  3.42       0.75     11.0        9
24  6495                  32.0  0.99547  3.57       0.71     10.2        5
25 ∨ 6496                 18.0  0.99549  3.39       0.66     11.0        8
26  |
27  |      type_white  best quality
28  0               1             0
29  1               1             1
30  2               1             0
31  3               1             0
32  4               1             0
33  ...           ...           ...
34  6492            0             0
35  6493            0             0
36  6494            0             1
37  6495            0             0
38  6496            0             1
39  |
40  [6497 rows x 13 columns]
```
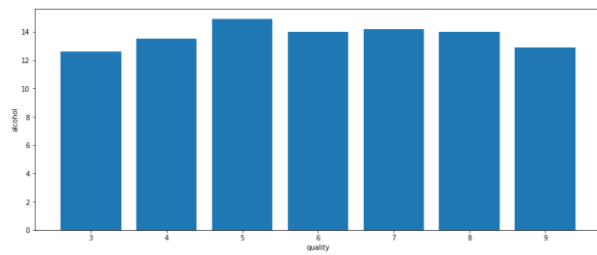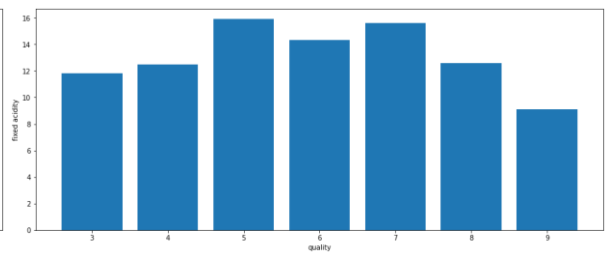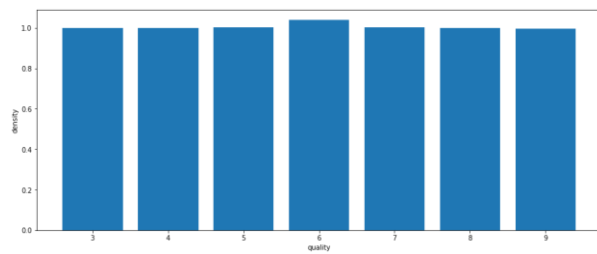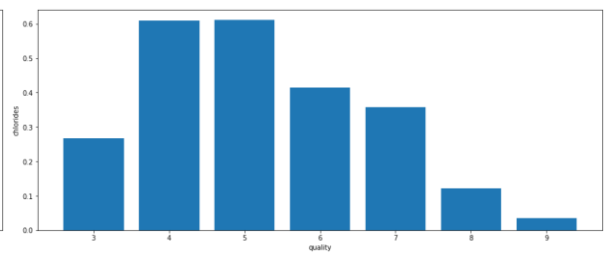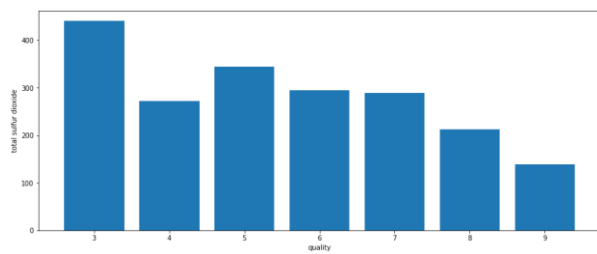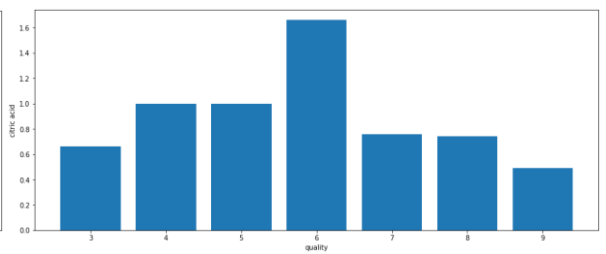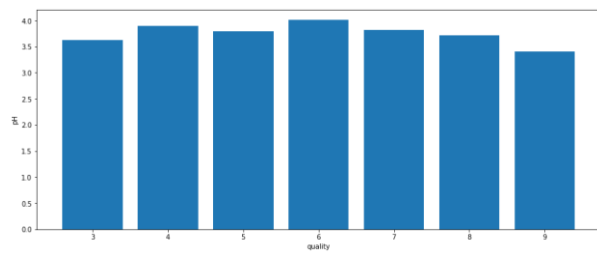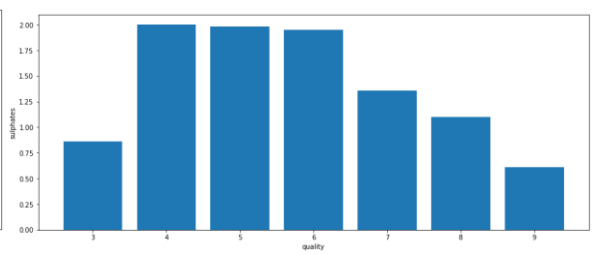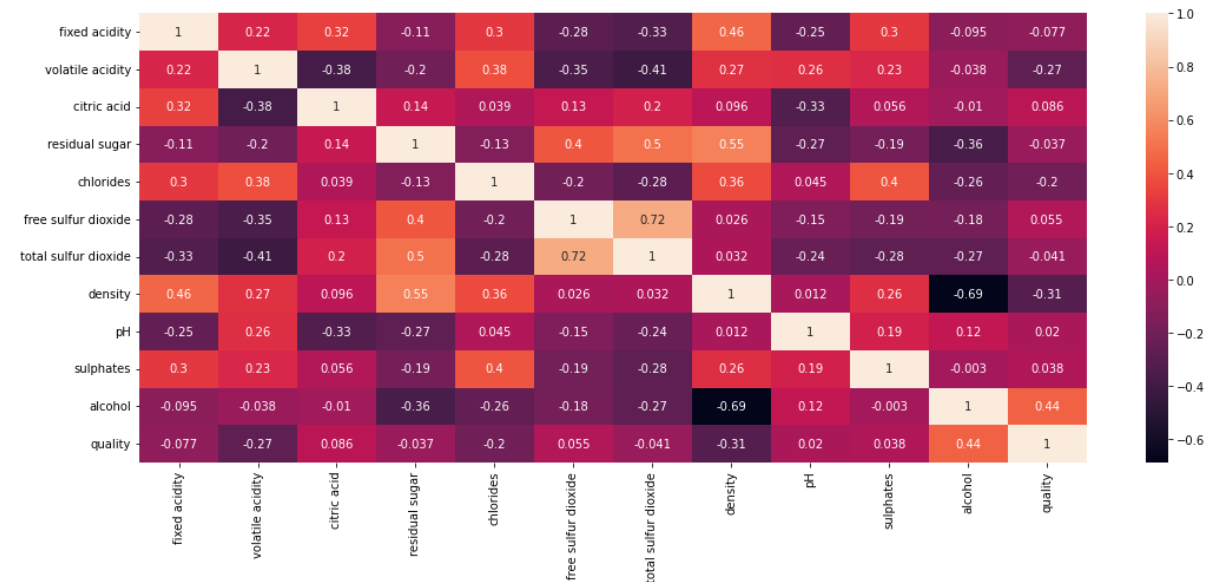
**Conclusions:**

1. Neural Network Classifier classified the wine data more accurately.

2. If total sulphur dioxide is less in wine more will be the quality of wine.

3. More density less will be the quality of the wine

4. Gaussian SVM classifier is less accurate to classify the wine type to white and red


**Contribution:**

1. Kakarla V S S Pavan Teja(S20190020216) – Python code to plot histogram of all quantitative variables, to plot quality vs different variables, to correlate the quantitative variables, Wine Quality.

2. K Sreenivasulu Reddy(S20190020217) – Matlab code to classify the wine data using Gaussian SVM and Neural Network Classifier, Classification App, Generation of Confusion Matrix, ROC Curves, Accuracy, Report.

3. K Litheesh Kumar(S20190020218) – Loading wine data, Matlab code to classify the wine data using Logistic Regression and Linear Discriminant Classifiers, Classification App, Generation of Confusion Matrix, ROC Curves, Accuracy, PPT.


**Reference:**

1. https://archive.ics.uci.edu/ml/datasets/Wine+Quality

2. https://www.geeksforgeeks.org/wine-quality-prediction-machine-learning/

3. https://in.mathworks.com/help/deeplearning/ug/wine-classification.html

4. https://www.javatpoint.com/classification-algorithm-in-machine-learning