# Machine learning

Milestone1 - Report

## Data Preprocessing

1. preprocessing techniques
   a. Missing values:

      Dataset contained Some missing values so we dropped any row that contained at least one missing value, but we have done this after choosing the features that we are going to use.

      ```
      data.dropna(how='any',inplace=True)
      ```

   b. Text values:

      Since computers can't deal with string values, so we need to convert any text data to numbers, so we have changed Prime_genre column to numeric values using One-Hot-Encoding technique.

      ```
      one_hot = pd.get_dummies(data['prime_genre'])
      data = data.join(one_hot)
      data = data.drop(['prime_genre'],axis=1)
      ```
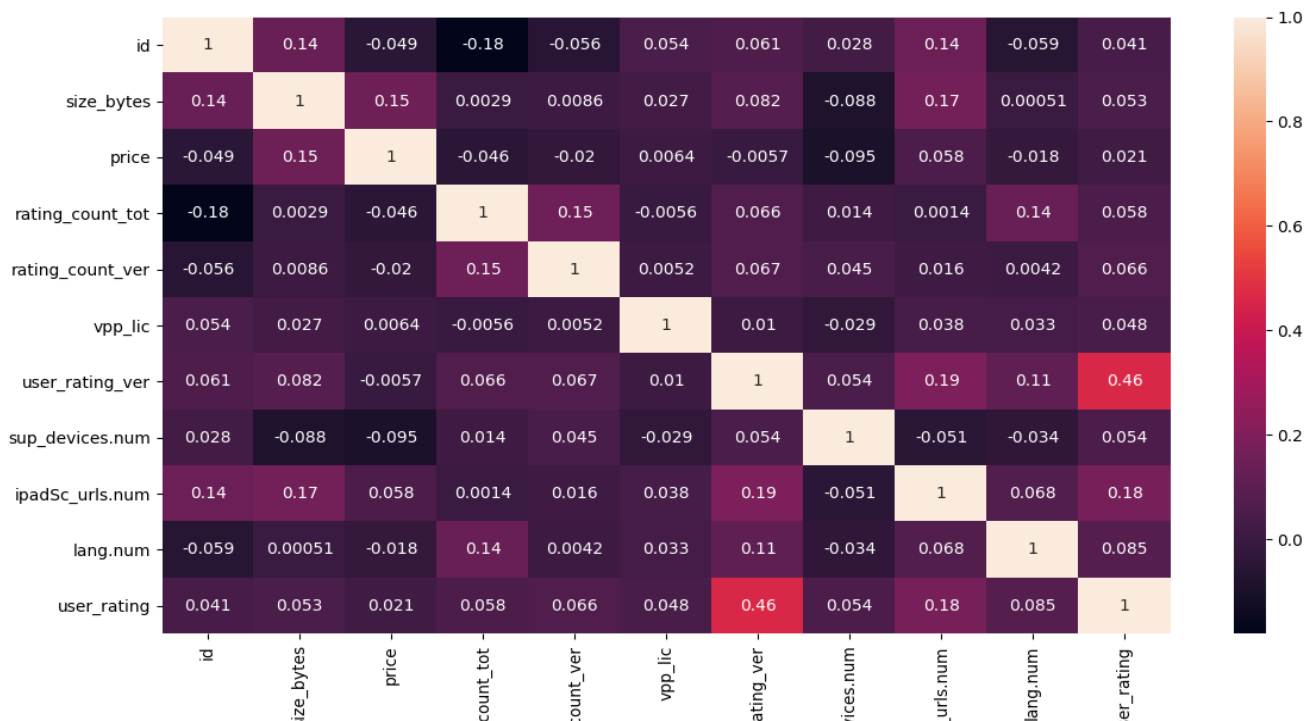
   c. Feature scaling:

      Features values differ very much in ranges between them so som attributes would affect the output more than other attributes only because it has higher ranges, so it  is important to apply feature scaling. We applied min-max scaling on two columns ['rating_count_tot','rating_count_ver']

      ```
      scaler = MinMaxScaler()
      data['rating_count_tot'] =
      scaler.fit_transform(np.array(data['rating_count_tot'])
      .reshape(-1,1))
      ```

```
data['rating_count_ver'] =
scaler.fit_transform(np.array(data['rating_count_ver'])
.reshape(-1,1))
```

2. Analysis
   a. Applied correlation between features to see which features have the most effect on the output (User_Rating) and decide which features we are going to use in our regression models



   b. After Applying correlation it was observed that the user_rating_ver has the most effect on the target value

3. Features Used

| Feature | Description |
|---|---|
| rating_count_tot | User Rating counts (for all version) |
| rating_count_ver | User Rating counts (for current version) |
| user_rating_ver | Average User Rating value (for all version) |
| prime_genre | Primary Genre |
| sup_devices.num | Number of supporting devices |
| lang.num | Number of supported languages |
| ipadSc_urls.num | Number of screenshots showed for display |

# Regression Techniques

As we know that regression models are used to describe relationships between variables by fitting a line to the observed data, Regression allows you to estimate how a dependent variable changes as the independent variable changes.

So ,we used two techniques, Multiple linear regression and Polynomial Regression.

We also tried **enhancing** the results using **L2** regularization.

- **Multiple Linear Regression:**

  "Y=bo+b1x1+b2x2+...........+BnXn"

  Multiple Linear regression are based on the assumption that there is a linear relationship between both the dependent and independent variables or predictor variable and target variable, and it also assumes that there's no major correlation between the independent variables.

  Multi linear regression can be linear and nonlinear, and it has one y and one or more x.

  After Implementation of the multiple linear regression

  We got a mean square error between  0.37652360484071495 and 0.35761225751504806

  The root of mean square error is between

  0.6136151928046721 and 0.5963103640175121

  The R2 score is between

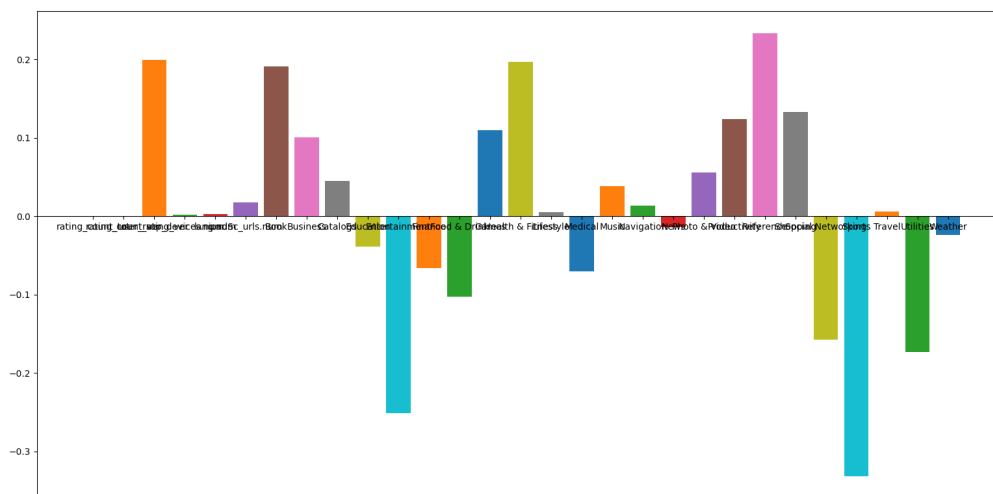  0.22929096545166372 and 0.21844404907891346

- **L2 regularization**

  Ridge Regression added a term in ordinary least square error function that regularizes the value of coefficients of variables. This term is the sum of squares of coefficient multiplied by the parameter The motive of adding this term is to penalize the variable corresponding to that coefficient not very much correlated to the target variable. This term is called L2 regularization.
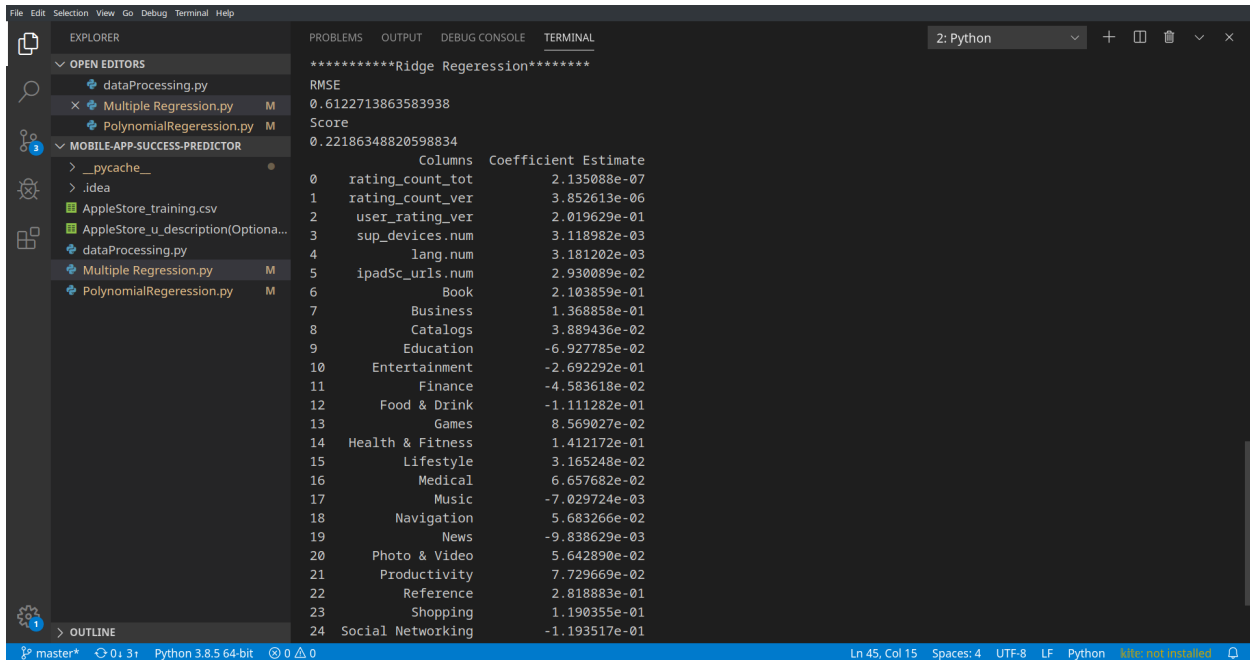
## alpha = 1



## alpha = 10

As we can observe from the above plots that alpha helps in regularizing the coefficient and make them converge faster.

The score value is increased by ( 0.002 to 0.05 )

And the error value is decreased by ( 0.01 to 0.03 )

```
File  Edit  Selection  View  Go  Debug  Terminal  Help

EXPLORER                                  PROBLEMS   OUTPUT   DEBUG CONSOLE   TERMINAL                          2: Python                    +  ⬚  🗑  ⌄  ×

∨ OPEN EDITORS                            ***********Ridge Regeression********
    ⬥ dataProcessing.py                  RMSE
  ×  ⬥ Multiple Regression.py    M        0.6122713863583938
    ⬥ PolynomialRegeression.py  M        Score
∨ MOBILE-APP-SUCCESS-PREDICTOR            0.22186348820598834
  > __pycache__              ●                     Columns  Coefficient Estimate
  > .idea                               0     rating_count_tot          2.135088e-07
  ▦ AppleStore_training.csv             1     rating_count_ver          3.852613e-06
  ▦ AppleStore_u_description(Optiona... 2       user_rating_ver         2.019629e-01
  ⬥ dataProcessing.py                   3       sup_devices.num         3.118982e-03
  ⬥ Multiple Regression.py    M         4             lang.num          3.181202e-03
  ⬥ PolynomialRegeression.py  M         5        ipadSc_urls.num        2.930089e-02
                                        6                 Book          2.103859e-01
                                        7             Business          1.368858e-01
                                        8              Catalogs         3.889436e-02
                                        9             Education        -6.927785e-02
                                        10       Entertainment        -2.692292e-01
                                        11             Finance        -4.583618e-02
                                        12        Food & Drink        -1.111282e-01
                                        13               Games         8.569027e-02
                                        14     Health & Fitness        1.412172e-01
                                        15           Lifestyle         3.165248e-02
                                        16             Medical         6.657682e-02
                                        17               Music        -7.029724e-03
                                        18          Navigation         5.683266e-02
                                        19                News        -9.838629e-03
                                        20        Photo & Video        5.642890e-02
                                        21         Productivity        7.729669e-02
                                        22           Reference         2.818883e-01
                                        23            Shopping         1.190355e-01
  > OUTLINE                              24     Social Networking       -1.193517e-01
master*  ↻ 0↓ 3↑  Python 3.8.5 64-bit  ⊗ 0 ⚠ 0                          Ln 45, Col 15   Spaces: 4   UTF-8   LF   Python   kite: not installed   🔔
```

- **Polynomial  Regression**

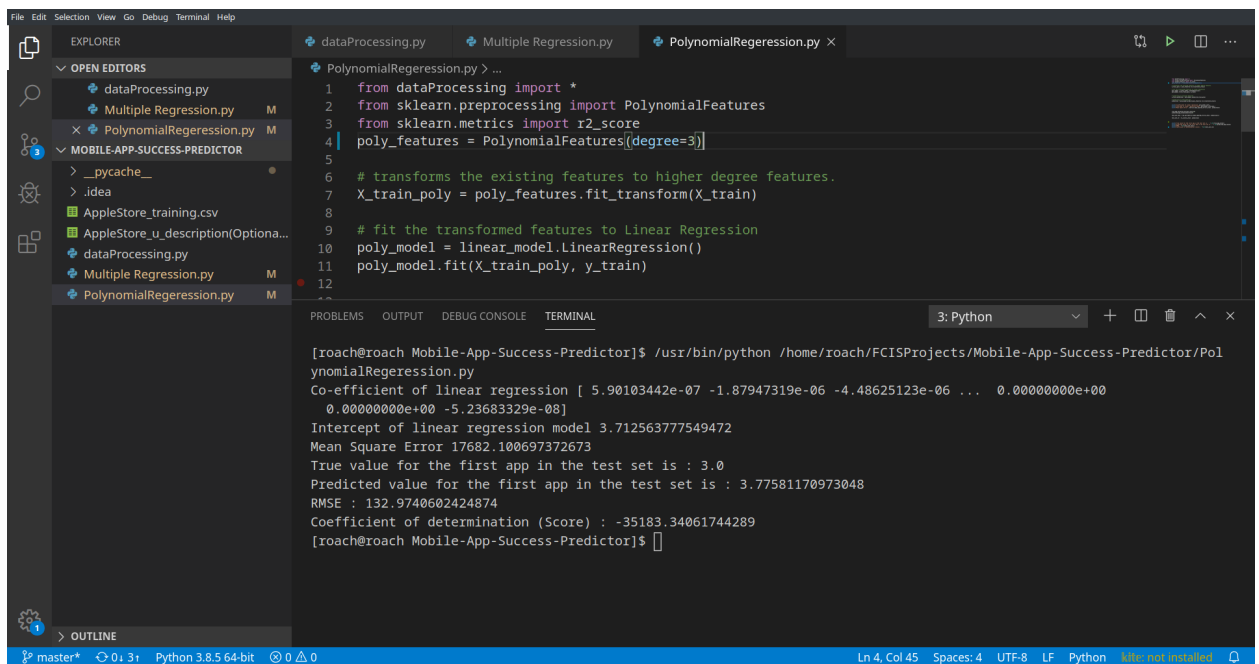  "$Y = b_o + b_1X + b_2X^2 + \ldots\ldots + b_nX^n$"

  Polynomial regression is one of the types of linear regression in which the relationship between the independent variable x and dependent variable y is modeled as nth degree polynomial regression fits a nonlinear relationship between the value of x and the corresponding conditional mean of y, donated $E(y|x)$.

  Polynomial regression provides the best approximation of the relationship between the dependent and independent variable.

→ <u>First : Degree of 3</u>

Produced very high root mean square error

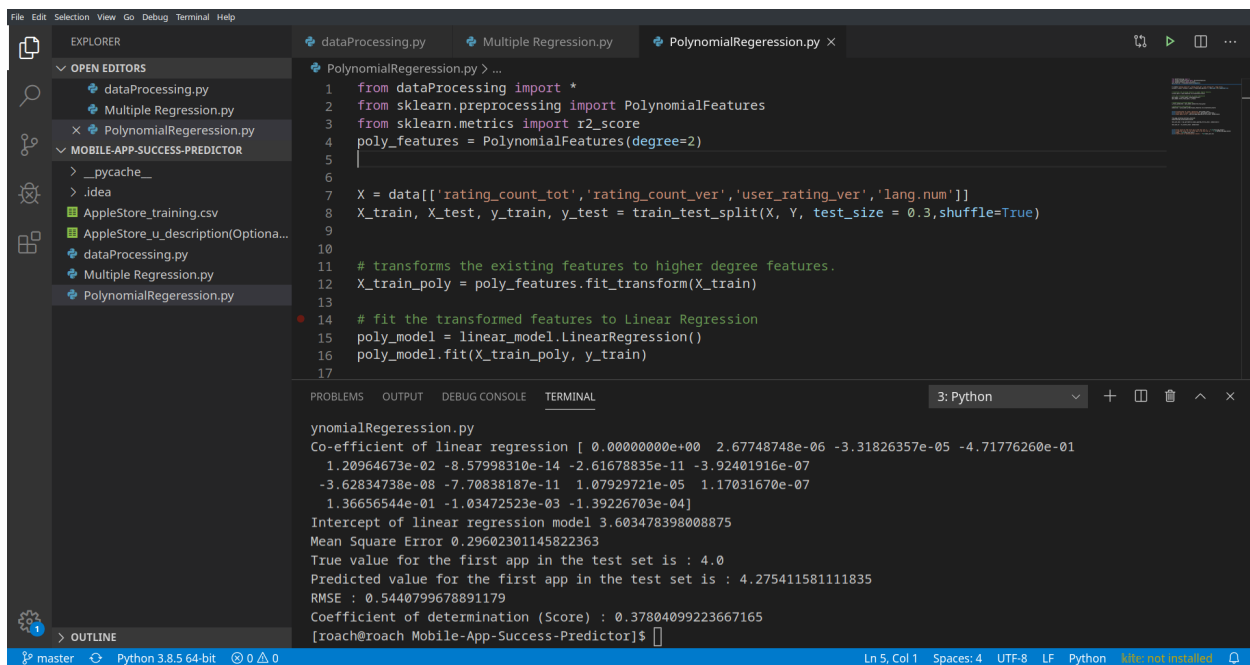Resulting in overfitting.

→ <u>Second : Degree of 2</u>

Produced lower root mean square error and higher R2 score than multiple linear regression.

Providing the better solution with this data.

Features used :

[rating_count_tot, rating_count_ver, user_rating_ver, lang.num]



## Time of training

The time of training was between 7-10 seconds

Except the polynomial of degree 3 which took around 2-3 minutes.

## Conclusion

- ➢ Time of training increases with the increase of the complexity.
- ➢ In higher polynomial degrees overfitting can occur.
- ➢ L2 regularization increases the accuracy of training.
- ➢ Data preprocessing and good feature selection leads to better results with better accuracy.

# Machine learning

## Milestone2 - Report

## **Data Preprocessing**

1. preprocessing techniques

    a. Missing values:

        Dataset contained Some missing values so we dropped any row that contained at least one missing value, but we have done

```
#drop null rows
data.dropna(how='any',inplace=True)
```

        this after choosing the features that we are going to use.

    b. Text values:

        Since ML Algorithms can't read text value data, we need to convert any text data to numbers, so we have changed Prime_genre column to numeric values using One-Hot-Encoding technique.

```
#Onehot encoding X values
dummy=pd.get_dummies(Classify_X['prime_genre'],prefix="Genre",drop_first=False)
Classify_X=pd.concat([Classify_X,dummy],axis=1)
Classify_X=Classify_X.drop(['prime_genre'],axis=1)
```

c. Feature scaling:

Features values differ very much in ranges between them so som attributes would affect the output more than other attributes only because it has higher ranges, so it  is important to apply feature scaling.

We applied min-max scaling on two columns ['rating_count_tot','rating_count_ver']

```
#feature scaling
scaler = MinMaxScaler()
data['rating_count_tot'] = scaler.fit_transform(np.array(data['rating_count_tot']).reshape(-1,1))
data['rating_count_ver'] = scaler.fit_transform(np.array(data['rating_count_ver']).reshape(-1,1))
#print("************* Scaling **************")
```

## Features Used

| Feature | Description |
|---|---|
| rating_count_tot | User Rating counts (for all version) |
| rating_count_ver | User Rating counts (for current version) |
| prime_genre | Primary Genre |
| sup_devices.num | Number of supporting devices |
| lang.num | Number of supported languages |
| ipadSc_urls.num | Number of screenshots shown for display |

# Classification Techniques

 As we know that classification and regression are two major prediction problems which are usually dealt with Data Science and Machine Learning.

Classification is the process of finding or discovering a model or function which helps in separating the data into multiple categorical classes discrete values.

Regression is the process of finding a model or function for distinguishing the data into continuous real values instead of using classes or discrete values. It can also identify the distribution movement depending on the historical data

In classification , data is categorized under different labels according to some parameters given in input and then the lapels are predicted for the data.

## • Logistic Regression

It's a classification algorithm used to assign observations on a set of discrete classes , some examples of classification like email spam or not spam.

• Simple logistic regression =logistic regression with one

predictor variable.

• Multiple logistic regression =logistic regression with multiple

predictor variables.

• And it also transforms its output using the logistic sigmoid

function to return a probability value.

# Logistic Regression Results:

Bar Graph For Logistic regression:



**Logisitc Regeression**

## Support Vector Machine

As we know SVM is supervised machine learning which can be used for both classification or regression problems. SVM is most used with linearly separable data.

It mostly used in classification problems, in this algorithm we plot each data item as a point in n-dimensional space (n-> number of features)

If the data isn't linear we perform a transformation using Kernel functions.

In this model we used SVM with C=0.4 and Polynomial Kernel Function of degree 3.

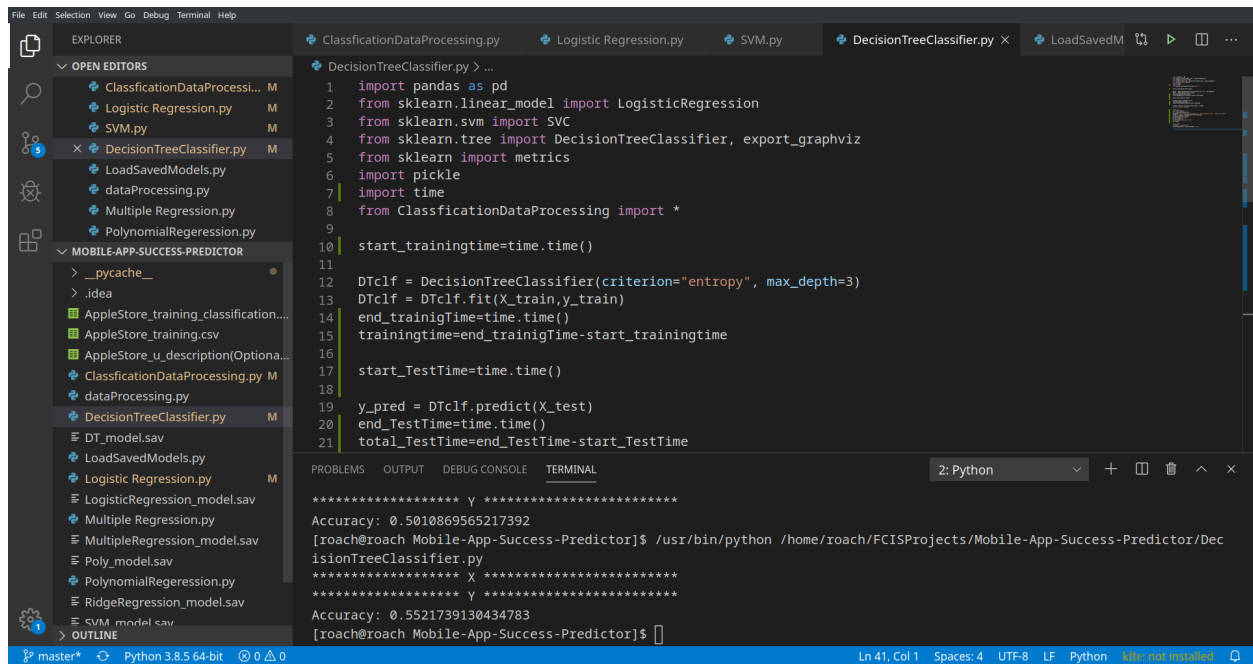## SVM Results:

<u>SVM Plotted Bar Graph:</u>



## **Decision Tree**

Decision Tree Analysis is a general, predictive modelling tool that has applications spanning a number of different areas. In general, decision trees are constructed via an algorithmic approach that identifies ways to split a data set based on different conditions. It is one of the most widely used and practical methods for supervised learning.

A decision tree is a tree-like graph with nodes representing the place where we pick an attribute and ask a question; edges represent the answers the to the question; and the leaves represent the actual output or class label. They are used in non-linear decision making.

<u>In this model </u>we used a decision tree with max depth of 3 which resulted in approximate  accuracy to other models with much less training and test time resulting in much better performance.
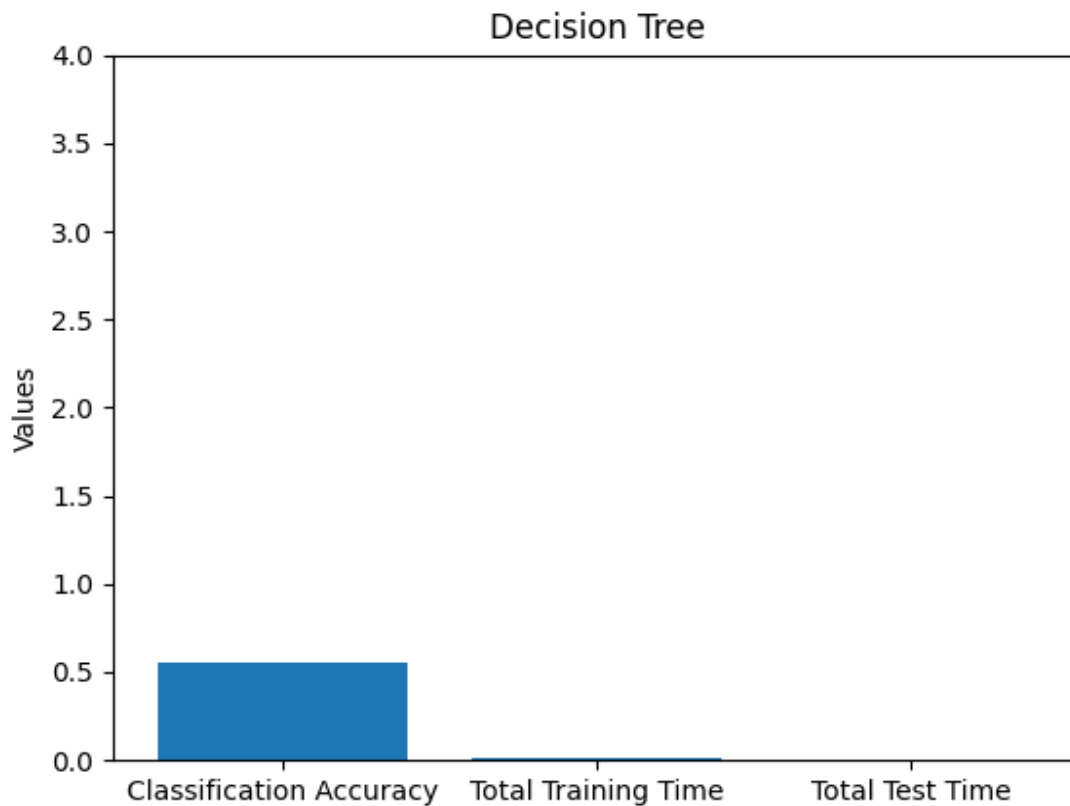
## Decision Tree Results:

Decision Tree Plotted  Bar Graph:



## **Conclusion:**

• SVM Try to maximize the margin between the closest support vectors.

• Logistic Regression maximizes likelihood.

• Data preprocessing and good feature selection leads to better results with

better accuracy.

• Accuracy of the Logistic regression is slightly better than SVM in this

Model but training time is much larger while logistic regression has less

test time.

• Decision Tree accuracy is close slightly less but training and test time are way less than SVM and LR in this model.