**PUNE INSTITUTE OF COMPUTER TECHNOLOGY
DHANKAWADI, PUNE**


DATA MINING AND WAREHOUSING MINI-PROJECT REPORT
ON


**"PREDICTING SOCCER GAME RESULTS USING VARIOUS
MODELS"**

**SUBMITTED BY**

Nachiket Erlekar    41434
Ashay Koradia      41429

**Under the guidance of**
Prof. Deepali Kadam

DEPARTMENT OF COMPUTER ENGINEERING
**Academic Year 2020-21**

# Contents

# 1   Problem Statement

Consider a labeled dataset belonging to an application domain. Apply suitable data preprocessing steps such as handling of null values, data reduction, discretization. For prediction of class labels of given data instances, build classifier models using different techniques (minimum 3), analyze the confusion matrix and compare these models. Also apply cross validation while preparing the training and testing datasets.

# 2  Abstract

Classification is a form of data analysis that extracts models describing importantdata classes. Such models, called classifiers, predict categorical (discrete, unordered)class labels. For example, we can build a classification model to categorize bank loanapplications as either safe or risky. Such analysis can help provide us with a better understanding of the data at large. In this project we use multiple classification models toanalyse the outcome of Soccer game played between various teams. Use apply suitabledata preprocessing steps.We then compare performance of classification models to findwhich one is the best

# 3   Hardware and Software Requirements

## 3.1   Hardware Requirements

1. 500 GB HDD

2. 4GB RAM

3. Monitor

4. Keyboard

## 3.2   Software Requirements

1. 64 bit Open Source Operating System like Ubuntu 18.04

2. Python 3

3. Google Colab

4. Different Libraries

5. Libararies like sklearn, pandas, matplotlib

# 4   INTRODUCTION

We have been provided with the data regarding various aspects of the home team, the opposition team and their supporters for a number of Soccer games.
The Data fields are

1. Id − Unique id given to each game.

2. game_seq − Sequence of the game in the history of FIH (International Soccer Federation).

3. season_end − Year in which the corresponding season ended.

4. date − Date on which the game was played.

5. season_game_seq − Sequence of the game in the corresponding season.

6. playoff − Whether the game is a playoffs game.

7. team_id − Unique id for the home team.

8. Elo − Elo rating for the home team before the game.

9. opp_team_id − Unique id for the opposition team.

10. opp_Elo − Elo rating for the opposition team before the game.

11. win_equivalent − Equivalent number of wins for the home team in a season.

12. bet_ratio − Fraction of bets placed on the home team.

13. home_crowd − Number of supporters for the home team.

14. opp_crowd − Number of supporters for the opposition team.

15. total_crowd − Total number of attendees for the game.

16. game_result − Win or loss for the home team (Win - 1, Loss - 0).

The train set contains 45000 records while the test set contains 13107 records. We drop the date column from our analysis. The null entries are as follows

| Attribute | Null Count |
|---|---|
| Elo | 9197 |
| opp_Elo | 7006 |
| win_equivalent | 12263 |

Table 1: Null Counts

We fill the null Elo and opp_elo entries with the mean value of Elo and opp_Elo attribute respectively i.e. 1501.184  1501.837

# 5 OBJECTIVE

- To understand data preprocessing

- To perform classification on dataset and predict labels for test dataset.

# 6  Scope

We select dataset of soccer games of various seasons. We try to apply many models and compare which one is the best model amongst them.
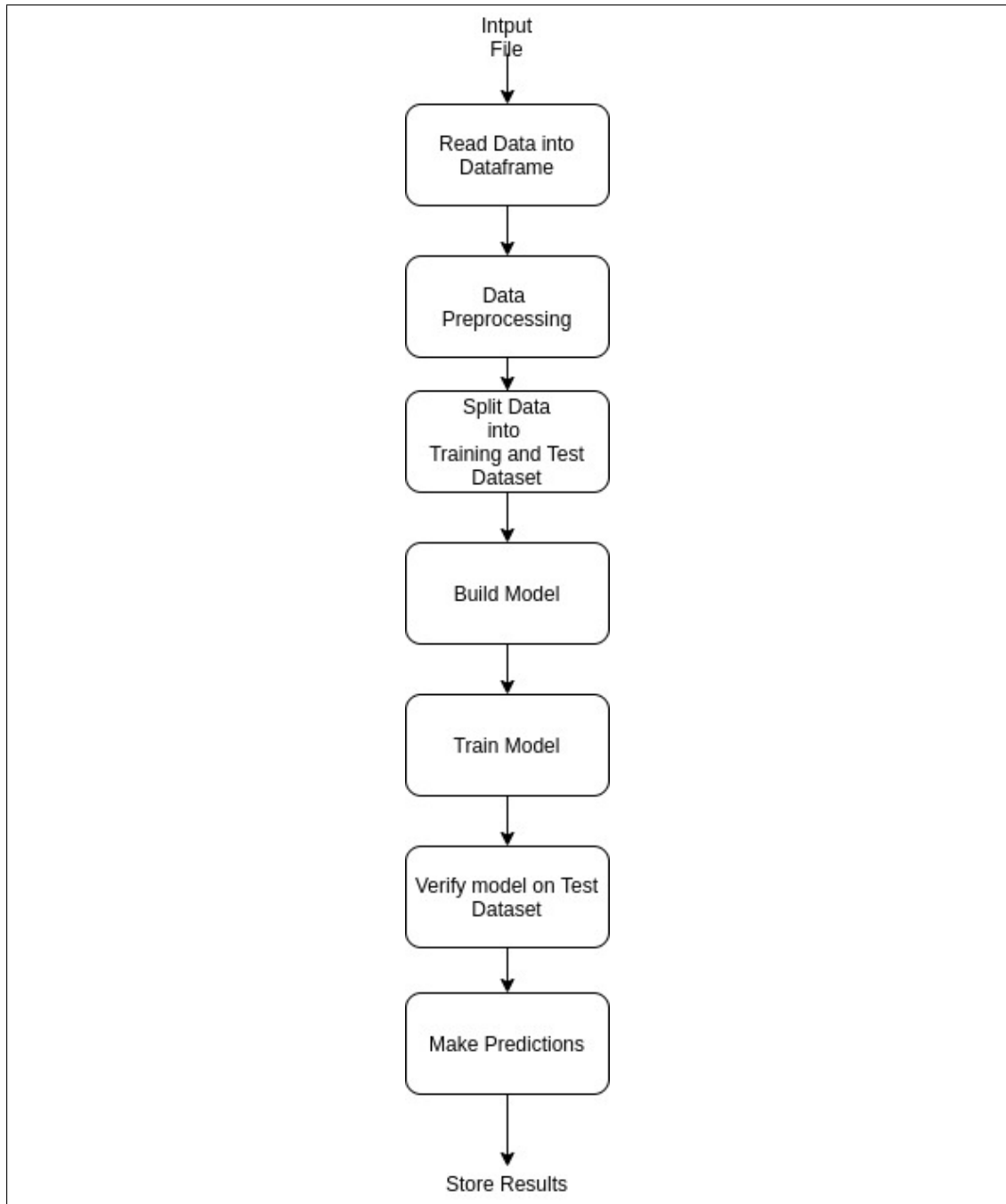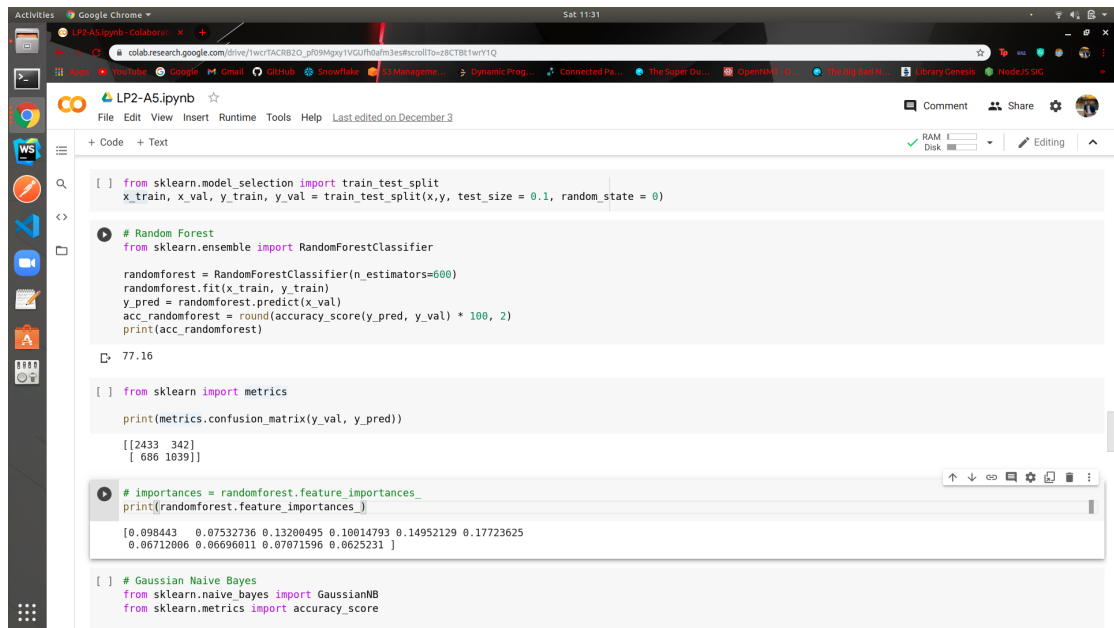
# 7 System Architecture



Intput
File

Read Data into
Dataframe

Data
Preprocessing

Split Data
into
Training and Test
Dataset

Build Model

Train Model

Verify model on Test
Dataset

Make Predictions

Store Results

Figure 1: System Architecture

# 8 Test Cases



Figure 2: Output for Random Forest Classifier



Figure 3: Output for GaussianNB and Logistic regression classifier

Figure 4: Output for SVC and Linear SVC



Figure 5: Output for Decision tree and Gradient boosting classifier

# 9    Result

The Accuracy for Various models are:

| Model | Accuracy |
|---|---|
| DecisionTree | 73.24 |
| RandomForest | 77.16 |
| GaussianNB | 66.56 |
| LogisticRegression | 65.47 |
| SVC | 61.67 |
| LinearSVC | 61.4 |
| GradientBoostingClassifier | 83.62 |

Table 2: Accuracy of vaious Models

We see that Gradient Boost Classifier gives the best score. We then use this model to perform training and testing of the model. After training, the model gives an accuracy of 83.62 %.
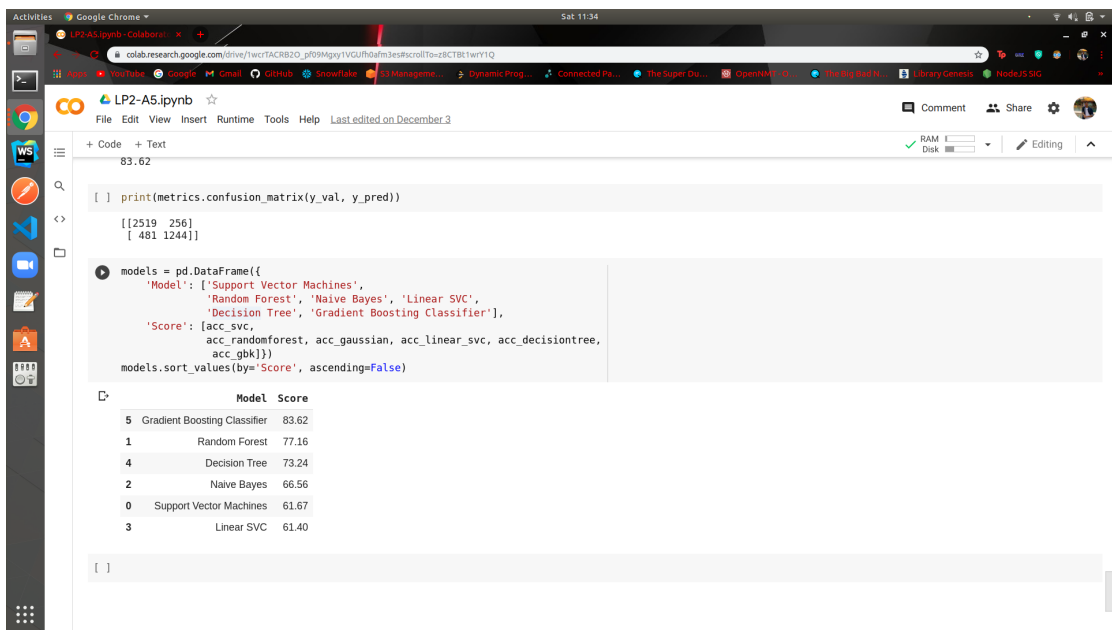


Figure 6: Comparison of various models

# 10   Conclusion

We have analysed the Soccer game dataset and performed data pre-processing steps.We have experimented multiple classification models and found out the best performer amongt them. We presented classification of soccer game results to predict the win/loss using Gradient Boost Classifier. We report a classification accuracy of 83.62

# References

[1] https://scikit-learn.org/stable/modules/generated/sklearn.ensembl

[2] https://scikit-learn.org/stable/modules/generated/sklearn.ensembl

[3]  https://www.kaggle.com/c/datawiz19round1/data