# ECE356 - Database Systems

# Lab 4 - Data Mining

## Group 3 report

Our general approach to this lab is summarized in the following steps:

1. Write the SQL to extract the dataset from the Lahman2016 database. The columns in the datasets following this order: playerID; feature1; feature2; …; featureN; classification.
2. Use the command "mysql -u <uid> -p -h marmoset04.shoshin.uwaterloo.ca db356_<uid> -B < lab4.sql | tr '\t' ',' > res.csv" and the corresponding database password to export the datasets(features) from the database to desired directory in .csv format. (in this case, it is res.csv)
3. Write the Python Decision Tree to train the dataset in the exported .csv file. The Python Decision Tree will output the report for the first trial for both gini and entropy and then the accuracy and predictions for the rest 5 iteration of trials.

All of our choice of datasets(features) during this lab will be tested using the above approach to determine the quality of the model.

## Part 1. Analysis and Results:

The first set of features we chose is the playerId, all the columns that represents the players' stats in game and the votedBy, votes and classification (HallOfFame.inducted) which is listed below:

```
playerid, b_G, b_AB, b_R, b_H, b_2B, b_3B, b_HR, b_RBI, b_SB, b_CS, b_BB,
b_SO, b_IBB, b_HBP, b_SH, b_SF, b_GIDP,p_W, p_L, p_G, p_GS, p_CG, p_SHO,
p_SV,p_IPouts, p_H, p_ER, p_HR, p_BB, p_SO, p_BAOpp, p_ERA, p_IBB, p_WP,
p_HBP, p_BK, p_BFP, p_GF, p_R, p_SH, p_SF, p_GIDP, votedBy, votes and
classification
```

Where b and p means the column is from Batting or Pitching table respectively.

The reason we started with this set of features is because the data in the Batting and Pitching table record all the life time behaviors of players during their matchups. This are the most important stats that can be used to determine whether a player is good and will be likely to nominate or elected into the Hall of Fame. The features votedBy and votes is the fact that how a player is elected for Hall of Fame

With this set of features, the models trained from the Python Decision Tree has the following results:

| Model | Accuracy | F1-score |
|---------|--------------------|--------------------|
| Gini | 0.9786776859504132 | 0.7994736842105263 |
| Entropy | 0.9669421487603306 | 0.7706024096385543 |

We use both Accuracy and F1-score to analyze the quality of our model because sometimes, accuracy is not a good parameter to determine whether the model is good or not. F1 Score is the weighted average of Precision and Recall. F1 is usually more

useful than accuracy, especially if the features are unevenly distributed. With the consideration of F1 score, we can have a better analysis. We can see here, both the Accuracy and the F1 score from this model is very high. This shows that this model is Overfitting. Overfitting usually means that a model has trained the data too well. During the learning, there are dominating/bias features acting as noise which will negatively impact the performance of the model on new data. In this case the dominating/bias features are votedBy and votes because even without other features, the model can still achieve a high Accuracy and the F1 score. This will make the model not able to predict new data and negatively impact the models quality. The votedBy and votes are simply the fact of why a player is elected. These two features are highly biased and will not contribute to the model's learning ability, so they should be removed.

Now we have:

```
playerid, b_G, b_AB, b_R, b_H, b_2B, b_3B, b_HR, b_RBI, b_SB, b_CS, b_BB,
b_SO, b_IBB, b_HBP, b_SH, b_SF, b_GIDP,p_W, p_L, p_G, p_GS, p_CG, p_SHO,
p_SV,p_IPouts, p_H, p_ER, p_HR, p_BB, p_SO, p_BAOpp, p_ERA, p_IBB, p_WP,
p_HBP, p_BK, p_BFP, p_GF, p_R, p_SH, p_SF, p_GIDP, and classification
```

And the models trained from the Python Decision Tree has the following results:

| Model | Accuracy | F1-score |
|---|---|---|
| Gini | 0.7829752066115702 | 0.3356626506024096 |
| Entropy | 0.7657446808510639 | 0.3510743801652892 |

We can see here, both the Accuracy and the F1 score from this model is very low. The model is not doing good for predictions. One of the reasons which might cost this problem is that some of features from the set have no weights during a prediction of elections meaning that those features are not crucial facts which determine whether a player is good or not. These features will also have a negative impact on the model's learning ability. This means that we need all the individual weight for each features.

Above is the he individual weight for each features from python feature_importances. We can see there, the first a few features have the weight 0.0. This means those features are meaningless for determining whether a player is good or not and should be eliminated from the set. Once all of those zero-weight features get eliminated, re-run the program and there will be new zero-weight features appears. Do the elimination again and again and finally we have reached this set (final set we used):

```
playerid, p_HBP, p_SHO, b_AB, p_L, p_SV, p_GF, p_CG, p_R, b_HBP, b_GIDP,
b_BB, b_H, b_SH, b_IBB, b_2B, b_3B, b_SB, b_RBI, b_SO, b_G, p_W, b_R,
classification
```

And the models trained from the Python Decision Tree has the following results:

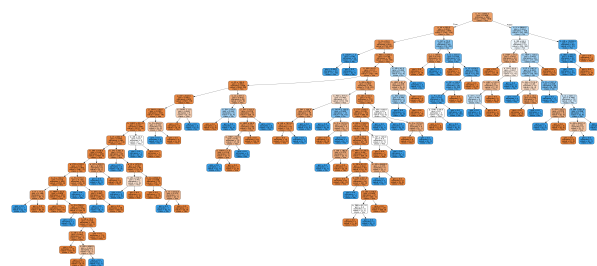| Model | Accuracy | F1-score |
|---|---|---|
| Gini | 0.8512396694214877 | 0.6086956521739131 |
| Entropy | 0.859504132231405 | 0.6136363636363636 |

We can see there, the Accuracy and the F1 score from this model improved a lot. The results are acceptable. Some sample output for all 5 iterations:
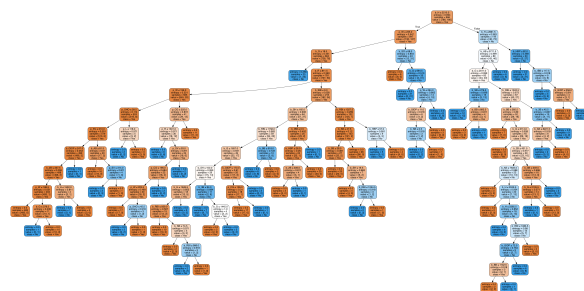


The Accuracy and the F1 score after iterations becomes high. This is because the 20% of the test data is randomly selected for the whole set. After the very first trial, the Accuracy and the F1 score is expected to increase because some of the test data are likely coming from the training data. This also show that the quality of our model are acceptable and can be used for prediction.

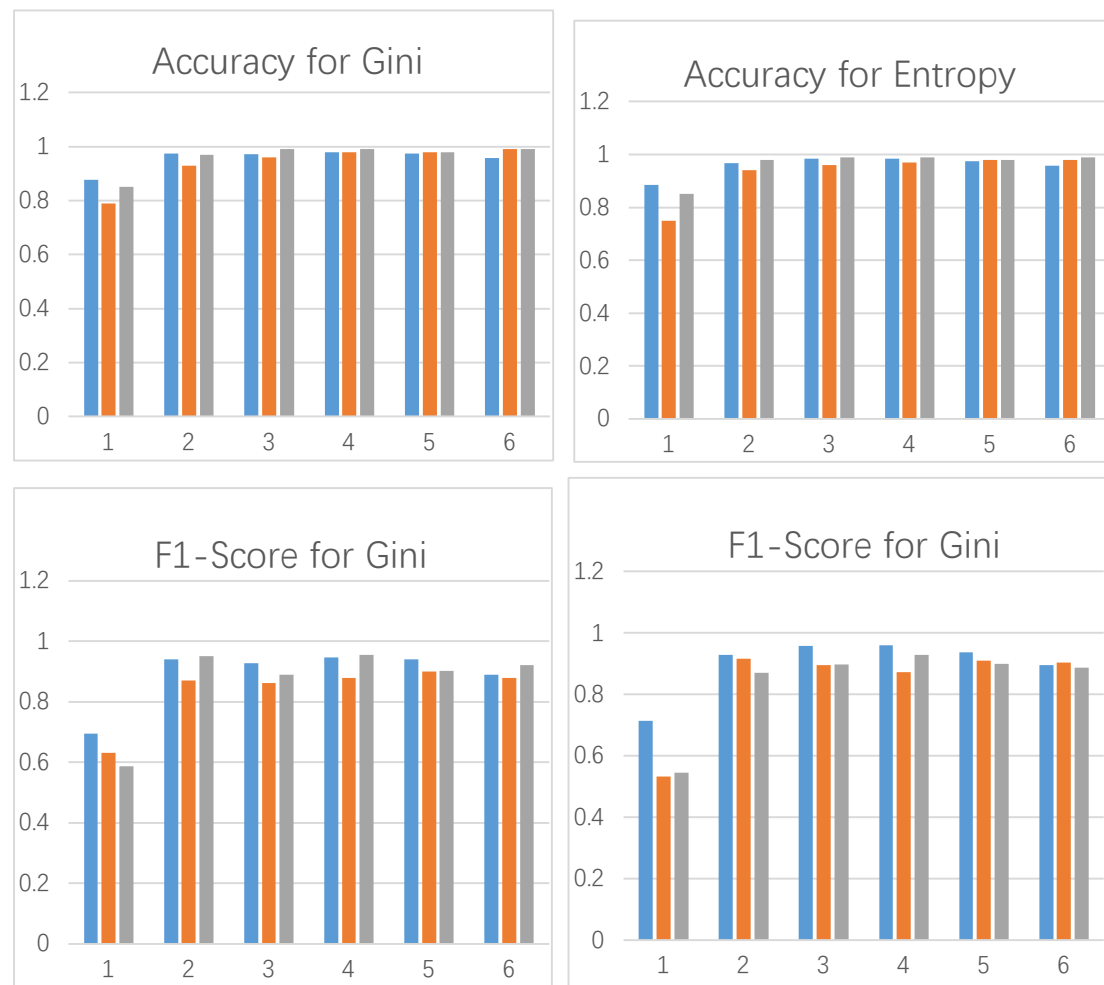The sample of decision trees are shown below:
For Gini:



For Entropy:



Blue indicates inducted = 0 and Red indicates inducted = 1

**Part 2. Comparison:**

For comparisons, the easiest way is to visualize the similarities and differences between all 3 choices of feature sets and their resulting models on their Accuracy and the F1 score in bar charts.



For all the plots above, blue bars are the values for the first choice (include votedby and votes), red bars are the values for the second choice (exclude votedby and votes) and grey bars are the values for the final choice (based on feature importance).

As we can see here, for gini accuracy, on the very first iteration, the second choice and the final choice starts low and first choice starts high. For later iteration, all three converged to a high value around 99%. However the first choice is lower compare to the rest because of the model's quality.

For gini F1 score, it follows a similar trend as the gini accuracy.

For entropy accuracy, on the very first iteration, the second choice and the final choice starts low and first choice starts high. For later iteration, all three converged to a high value around 95%. However the first choice are the final choice usually higher than the second choice. This might due to the fact that the second choice has meaningless features.

For entropy F1 score, it follows a similar trend as the entropy accuracy. However, on the $6^{th}$ iteration, all three choices have a similar f1 score around 94%. This might be costed by the uncertain percentage of changes during each trial.