## **EXPERIMENT REPORT**

| Student Name | Archit Murgudkar   |
|--------------|--|
| Project Name | AT1 - Kaggle Competition Week1   |
| Date         | 18/08/2023   |
| Deliverables | AT1_PartA-Binary Classification using<br>Random Forest<br>Model - Random Forest Classifier |

| 1. EXPERIMENT BACKGROUND  |  |  |
|---------------------------|--|--|
|                           |  |  |
| 1.b. Hypothesis           | Features provided will help classify if a certain player gets picked up to play for NBA by training a Random Forest Classifier. There are numerous features present in the dataset which can help predict the dependent variable named, "drafted".   |  |
| 1.c. Experiment Objective | A supervised machine learning classification model predicting whether a college basketball player will make into the dream NBA league depending on the independent variables which are his game skills. The objective is to develop the model and its performance better where overfitting and underfitting to be dealt with if present using hyperparameter tuning and the model performing on the unseen data to help classify the result. |  |

## 2. EXPERIMENT DETAILS

## 2.a. Data Preparation

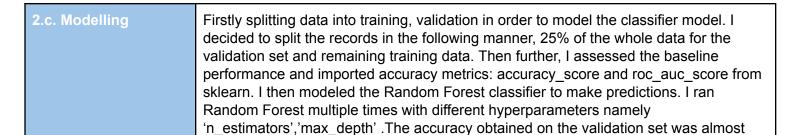
In data preparation, I started with data cleaning first, I made a copy of dataset to perform data cleaning, so that we always have our original dataset if we want to refer and make the changes in the duplicate dataset. Further, I did a check on finding null/missing values in the records. The feature 'rec\_rank' and 'pick' comprised a greater amount of null values, so I decided to totally drop off those columns. Adding any value like mean/median/0 etc. would create extreme bias in the predictions. Similarly 35 rows had 31 fields missing, and imputing so many fields in all these rows can create biases in the dataset. After this, I deleted the following features due to valid reasons as follows:

- 1. num" feature as it is a player specific record
- 2. "year" feature as we want our model to learn generalise patterns
- 3. "ht" feature as it consists false information
- 4. "player\_id" as it consists specific information
- 5. "type" as it consist of the same value in all records and would not help in making any decision.

I imputed the mean in the rest of the columns having missing values. "Yr" had redundant values and nulls which I replaced with the mode of it. Further, I got rid of the duplicates from the dataset

## 2.b. Feature Engineering

Further, I performed exploratory data analysis on the fields in the dataset. I plotted a pie chart to check the distribution of target variable, "drafted" and surprisingly, it was found out that there was a high imbalance in the distribution of target variable. From the heatmap, I understood that the features did not strongly correlate much with the target feature, "porpag" with the correlation of 0.255was the highest among all, and "mid\_ratio" ranked lowest. As our machine learning classifier requires data to be in the numeric format before feeding it, I decided to perform encoding on the features with object data type. Due to high cardinality in the "team" variable, I performed encoding through frequency mapper technique. Another important factor is ensuring that we provide data to the model which is on the same scale, which will help our model make accurate predictions. For this, I performed data scaling by importing the 'StandardScaler' library from sklearn. After accomplishing all this, the data was finally suitable to be fed to the classifier algorithm and carry out modeling.



data set.

equal to that achieved from the training set. So I decided to run this model on the test

| 3. EXPERIMENT RESULTS      |   |  |
|----------------------------|---|--|
|                            |   |  |
| 3.a. Technical Performance | Using the Random Forest model on the unseen data, i.e, the testing set, the score of the performance metric, roc_auc_score came out to be 0.7606.   |  |
| 3.b. Business Impact       | The Random Forest model achieved an accuracy rate of 76.06% on the roc_auc_score performance metric. So the predictions made by the model to classify if a player will get drafted is 76.06 %accurate.  |  |
| 3.c. Encountered Issues    | Major issue I thought of was the high imbalance present in the distribution of the target variable and high cardinality in the features which led to this accuracy level. Also, alot of null and redundant values were present in the variables of the dataset. |  |

| 4. FUTURE EXPERIMENT               |  |  |
|------------------------------------|--|--|
|                                    |  |  |
| 4.a. Key Learning                  | The model would have been more accurate if the features in the dataset were more strongly correlated with the target variable and if the distribution of the target fields would have been more balanced.  |  |
| 4.b. Suggestions / Recommendations | There is a high need to collect more records which would contribute in balancing the dependent variable. Also, more features need to be collected which may have a stronger relation with the target variable which in turn will help the model predict with higher accuracy based on stronger evidence. |  |