

EXPERIMENT REPORT



Student Name: Yasaman Mohammadi
Project Name: advmla-2023-spring

Date:
25/08/2023
Deliverables:
notebook name: Mohammadi_Yasaman-24612626-
week2_XGBOOST

1.EXPERIMENT BACKGROUND

1.a. Business Objective

This project aims to create a model that predicts whether a college basketball player will be drafted into the NBA based on their current season's statistics. As a result, NBA teams can make better draft selections and allocate resources more efficiently.

By making accurate predictions, draft decisions can be improved, players can be developed, and teams can perform better. Missed opportunities, wasted resources, and potential damage to the model's credibility can all be attributed to incorrect predictions. Draft choices, player development strategies, and fan engagement are all impacted by the model's success.

1.b. Hypothesis

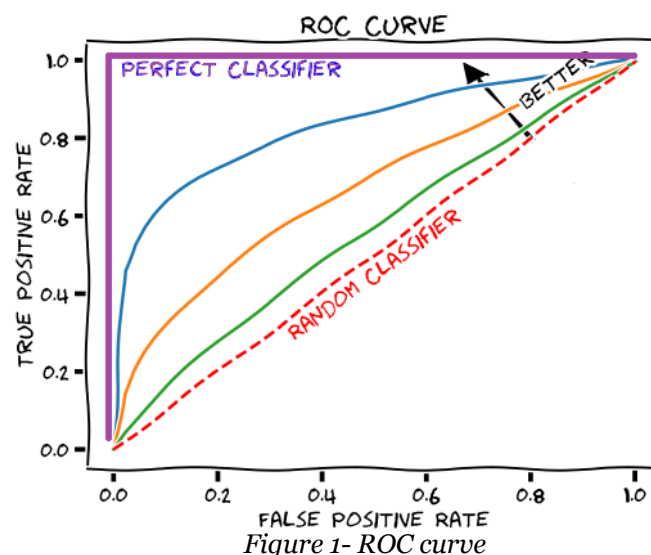
Do college basketball player statistics provide a sufficient basis for predicting their likelihood of being drafted into the NBA?

Accurate predictions can lead to improved draft selections, resulting in a better roster for the team and an overall improvement in performance.

Data-driven predictions can enhance sports analysis, media coverage, and fan discussions. As a result, it can generate excitement and stimulate informed discussion about the potential outcomes of a draft.

1.c. Experiment Objective

The experiment is expected to result in a predictive model with a high AUROC score that accurately predicts whether a college basketball player will be drafted into the NBA based on their current season's statistics. AUROC scores should be significantly higher than random chance, indicating that the model can reliably distinguish between players who are likely to be drafted and those who are not.



Different scenarios might happen based on AUROC curve performance. The area under the receiver operating characteristic (AUROC) is a performance metric for evaluating classification models. The larger the area under this curve, the better the performance.

2.EXPERIMENT DETAILS

2.a. Data Preparation

the steps taken for preparing the data for previous experiment.		Explanation
1	Data collection	Building a machine learning model begins with the collection of data from a variety of sources. In this case, the data was collected from a Kaggle competition.
2	Explore Dataset	A general understanding is derived from the data in this step.
3	Data cleaning	Data should be preprocessed to remove any irrelevant, missing, or corrupted data points and any discrepancies in the data.
4	Data visualization	Visual understanding of data distribution is gained by plotting different charts in both Python and Tableau.
5	Feature engineering	In feature engineering, new features are created, or existing features are transformed to improve the model's performance. This experiment used a label encoder and one hot encoding to transform categorical data into numerical data to consider their influence on the machine learning model.
6	Removing outliers	An outlier is a data point significantly different from the rest. Outliers can adversely affect machine learning models and they need to be removed. (Identifying them by using box plots)
7	Splitting the dataset	The dataset was randomly split into training, validation, and testing. The training set is used to train the model, the validation set is used to tune the model's hyperparameters and prevent overfitting, and the test set is used to evaluate the final performance of the model. The test size parameter specifies the percentage of data allocated to the test set, in this case, 20%. The random state parameter is used to ensure the reproducibility of the split. For many datasets, the 80/20 split is a good rule of thumb since it provides enough data to train the model while still leaving a considerable amount for testing. The ratio may vary depending on the size and complexity of the dataset, as well as the specific problem being addressed. As part of this project, there was a test set without drafted sections and the objective was to predict a probability for the drafted variables at the end.
8	Scaling the dataset	The data is crucial to prevent sensitivity to some features, so all features have the same scale. Since some machine learning algorithms are sensitive to the scale of input features, this is important.

Table 1- the steps taken for preparation of the data for previous experiment.

Data Preparation Steps for the Experiment:

1. **Data Cleaning:** To address missing values, a tailored approach was taken. Categorical columns were imputed with the mode rather than discarding all instances with missing values, while numerical columns were filled with zeros. For example, the "Rec_Rank" column, which reflects the player's high school recruiting rank and is closely tied to the target, underwent imputation to retain its significance.
2. **Feature Engineering:** One hot encoding and label encoding began after gaining insights from a Tableau dashboard. The roles of the "team" and "conf" columns in influencing the target were evident, prompting us to encode them numerically using label encoding. However, further investigation revealed that these columns needed a distinct ordinal relationship to use this approach. Consequently, these columns were excluded from the analysis. Solely, the "year" column underwent label encoding due to its relevance.
3. **Enhanced Feature Integration:** Building upon the previous experiment, we reconsidered the "ht" column. In this instance, columns containing month information were transformed into numerical representations of month numbers. Additionally, a conversion to height in feet format was performed, supplementing the dataset with more comprehensive information.

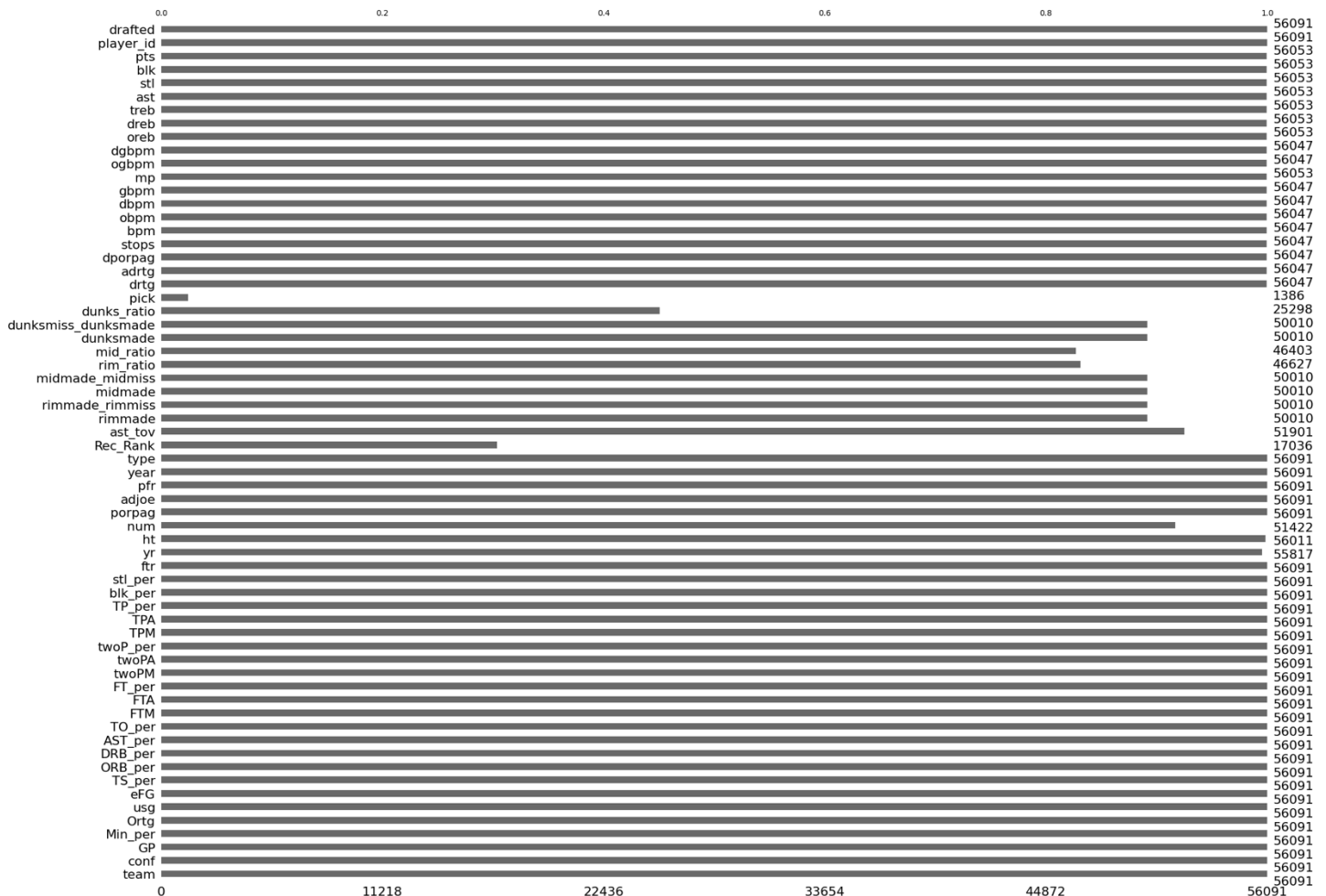


Figure2-The mnso bar chart for fiding missing values

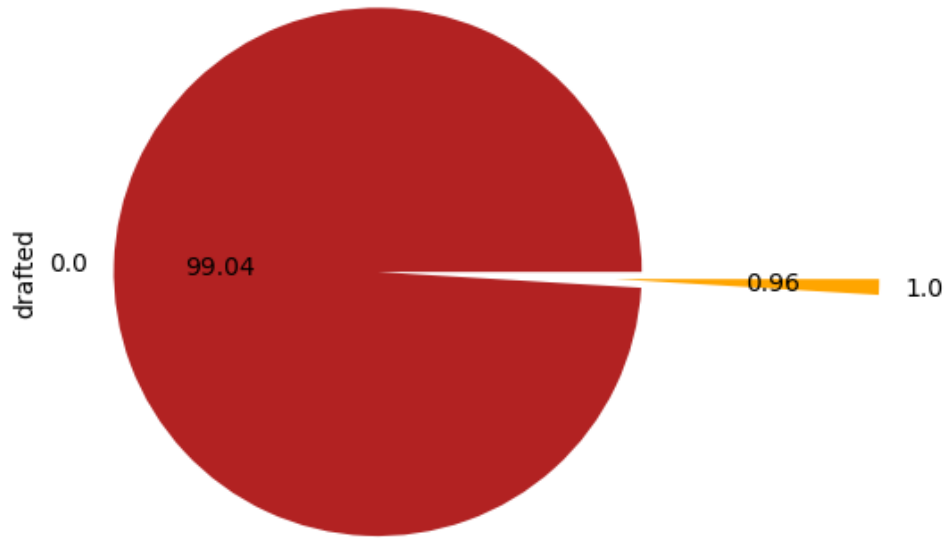


Figure 3-Pie chart of drafted distribution

2.b. Feature Engineering

One hot encoding and label encoding began after gaining insights from a Tableau dashboard. The roles of the "team" and "conf" columns in influencing the target were evident, prompting us to encode them numerically using label encoding. However, further investigation revealed that these columns needed a distinct ordinal relationship to use this approach. Consequently, these columns were excluded from the analysis. Solely, the "year" column underwent label encoding due to its relevance.

2.c. Modelling

In this experiment, XGBoost was utilized to construct a model capturing the interplay between independent variables and the probability of a binary outcome. Notably, this model boasts higher complexity than the preceding Polynomial Logistic Regression approach.

Hyperparameter		Explanation
1	N estimator	<p>indicate how many trees are present in the forest. Generally, the larger the number of trees, the better the ability to learn the data. To find the optimal value, a grid search was used to find the optimal number of trees. However, adding too many trees can significantly slow down the training process.</p> <p>With grid search, a n estimator of 200 to 350 was selected for this experiment, and 250 was determined to be the best value.</p>
2	Max depth	<p>In XGBOOST, the maximum depth is calculated as the longest path between the root node and the leaf node. With grid search, a maximum depth of 15 to 20 was selected for this experiment, and 17 was determined to be the best value. To prevent overfitting, higher values are not given</p>
3	Gamma (Minimum loss reduction)	<p>Gamma specifies the minimum loss reduction required for a split. As a result, the algorithm becomes more conservative. Depending on the loss function, the values may vary.</p> <p>With grid search, eta values between 0.01 to 0.8 were selected for this experiment, and 0.05 was determined to be the best value. A larger gamma will result in a more conservative algorithm, so higher values were not used.</p>
4	eta (learning rate)	<p>To prevent overfitting, step size shrinkage is used in the update process. At each step of the boosting process, the weights of new features are achieved, and eta shrinks the weights of new features to make the boosting process more conservative. With grid search, eta values between 0.1 and 0.3 were selected for this experiment, and 0.2 was determined to be the best value.</p>

Table 2-Hyperparameters tuning explanation

It was decided to use stratified k-fold cross-validation due to the imbalanced distribution of the target feature in the primary data. Therefore, training and test data in each fold will reflect the imbalanced distribution of the target feature in the primary data.

3. EXPERIMENT RESULTS

3.a. Technical Performance

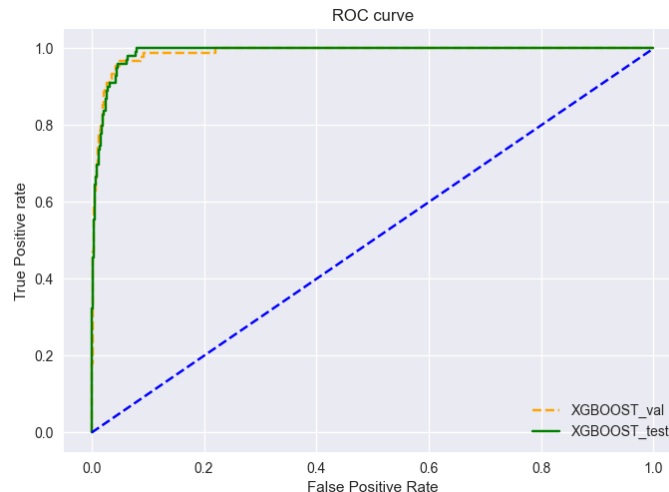


Figure 5-AUROC curve of XGBOOST model

The higher the AUC, the better the model will be at distinguishing between positive and negative classes. This ROC curve indicates good performance, but there is room for improvement.

The polynomial regression model achieved a final score of 0.95891 on the Test set, as confirmed by submitting the results on Kaggle.

After implementing the earlier adjustments on EDA, the AUC score of the polynomial regression model on the Test set significantly improved from 0.9888 to 0.99.

Subsequently, upon submitting the XGBOOST results on Kaggle, there was a slight enhancement in the final score from 0.95891 to 0.9612.

It may be possible to achieve better results using other models or feature engineering.

3.b. Business Impact

The success of the model affects several aspects of the NBA teams' operations, including draft choices, player development, resource allocation, fan engagement, and team performance. Errors in results can have a wide range of negative impacts, including wasted resources, missed opportunities, credibility damage, and long-term competitive disadvantages. For the model to provide meaningful value to NBA teams and stakeholders, accurate predictions are critical.

3.c. Encountered Issues

There might be better solution to imputing missing data to achieve higher scores.

4.FUTURE EXPERIMENT

4.a. Key Learning

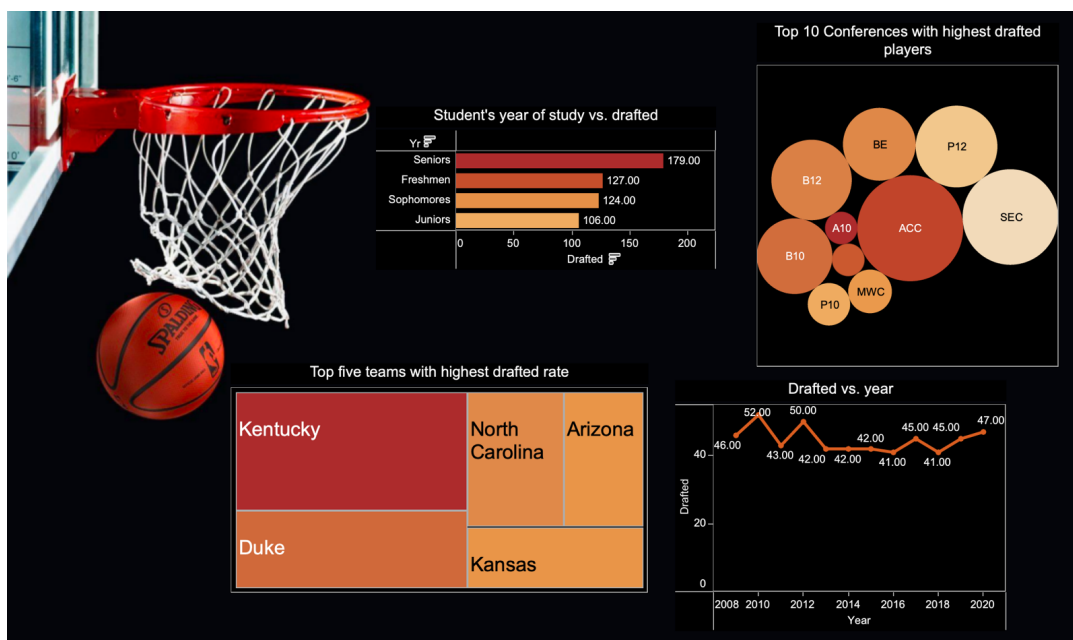


Figure 6- Tableau Dashboard

An interactive dashboard in Tableau has been created to visualize the correlation between drafted players' rates and teams, seasons, and conferences.

Moreover, we regarded XGBoost as a more intricate model than the previous one, slightly improving predictions. However, it's worth noting that even better machine learning models, like ADABOOST, could lead to achieving even higher scores.

4.b. Suggestions / Recommendations

In future experiments, different machine learning models with hyperparameter tuning will be used. In order to deploy the final solution into production, the user interface must be refined so that predictions are easy to access, ethical biases must be addressed proactively, and ongoing model adjustments should be monitored in real-time.

Private repo link:

https://github.com/JYasimo/Kaggle_competition_NBA_league/tree/main