

EXPERIMENT REPORT

Student Name	Trung Kien Hoang
Project Name	Kaggle Competition: The NBA draft
Date	1/09/2023
Deliverables	<Experiment 3> <AdaBoost Classifier with hyperparameters> <Data preparation: Label encoding(yr) + Converter(ht)> Github link: https://github.com/KenUTS/adv_mla_assignment_1.git

1. EXPERIMENT BACKGROUND

Provide information about the problem/project such as the scope, the overall objective, expectations. Lay down the goal of this experiment and what are the insights, answers you want to gain or level of performance you are expecting to reach.

1.a. Business Objective

The main objective of this project for the business is to build a prognostic model capable of evaluating the probability of a collegiate basketball player getting selected for the National Basketball Association (NBA) draft, using their performance information from the current season. The outcomes of the model will be used by diverse stakeholders, including NBA teams, players, sport commentators and fans. Suppose the model accurately forecasts which player will be chosen. In that case, NBA teams can use the model's predictions to help them decide how to pick players in the draft. By finding possible high-performing players, they could improve their team performances. However, if the model's expectations are wrong, teams might not pick the best players with skills in the draft. The accuracy model has the potential to assist players in making informed decisions between NBA draft and their college careers. However, selecting an incorrect model has the potential to lead players to make misguided decisions in their future careers. For sport commentators and fans, the model engages in more comprehensive discussions pertaining to players. However, this method has the potential to generate inaccurate feedback and therefore disappoint supporters in the event that their favorite player is not selected.

1.b. Hypothesis	<p>Hypothesis: The statistical data of collegiate basketball players during their current season have the potential to enhance the likelihood of their selection in the National Basketball Association (NBA) draft.</p> <p>Question: Are there any elements such as height of players or year of study that contribute to an increased likelihood of being selected in the NBA draft?</p> <p>It is worth examining due to its potential to assist players in their preparation for the selection process.</p>
1.c. Experiment Objective	<p>This experiment is implemented by the results of experiment 2, which demonstrate that the AdaBoost model has the best performance.</p> <p>The projected result of the undertaking is a model that effectively forecasts the likelihood of players being selected in the NBA draft. Based on the requirement of the project, the metric used to assess model performance is AUROC (Area Under ROC). It is supposed to be smaller than 1 and greater than 0.8.</p> <p>Scenario 1: Train the selection of algorithms; for all data sets, if the best model's performance is less than 0.8 of the AUROC score, take one or more actions: Change data preparation, feature engineering, hyperparameter optimization, or algorithm selection.</p> <p>Scenario 2: Train the selection of algorithms; for all data sets, if the performance of a model is greater than 0.8, finish the project and deploy the model to the business.</p>

2. EXPERIMENT DETAILS

Elaborate on the approach taken for this experiment. List the different steps/techniques used and explain the rationale for choosing them.

2.a. Data Preparation

1. Adding a new column for shots missed at or near the rim, missed two point shots that were not made at or near the rim and missed dunks for training and testing data.
2. Drop player number as it is not meaningful information, columns with large numbers of missing values such as Rec_Rank and pick columns.
3. Drop ratio columns where it is insufficient due to the players may have drunk made and missing equal to 0 for instance.
4. Fill all missing values of numerical columns with median values.

2.b. Feature Engineering

Because, the shots made and missed could be zero, so the ratio between the shots made and shot missed could be undefined if the number of shots missed is zero. Therefore, for this experiment, a new column for shots missed at or near the rim, missed two point shots that were not made at or near the rim and missed dunks are added for training and testing data. Furthermore, Polynomial Transformation is used to capture the complex among selected numerical features.

Label encoding(yr) + Converter(ht):

1. According to many sources, it has been suggested that the classification of students based on their year of study follows a certain pattern. Freshmen are often referred to as first-year students, sophomores as second-year students, juniors as third-year students, and seniors as fourth-year students, indicating a hierarchical ranking based on their progression through their academic programme. Hence, the use of a label encoding tool is considered the most effective approach for converting categorical data into numeric variables.
2. It seems that there is an issue in the height column, where either Kaggle or Excel has mistakenly converted players' heights into dates. To illustrate, let us use the measurement of 6 feet 11 inches, denoted as 6 - 11 while this measurement was recorded into a date format, namely 11/06/2023. Consequently, a function is established with the purpose of converting such records.

2.c. Modelling

In this experiment, the AdaBoost model combines multiple weak learners to create a strong classifier. It is trained with multiple hyperparameters set in training, validating, and testing dataset.

A function is created to test each hyperparameter sets from given:

'n_estimators' hyper parameters are trained at 100, 200, 300, 400 and 500.

'learning_rate' hyper parameters are trained at 0.1, 0.4 and 1.

3. EXPERIMENT RESULTS

Analyse in detail the results achieved from this experiment from a technical and business perspective. Not only report performance metrics results but also any interpretation on model features, incorrect results, risks identified.

3.a. Technical Performance

The best hyperparameter set is `n_estimators=400` and `learning_rate=0.1`
AUROC score at training dataset is 0.992
AUROC score at validating dataset is 0.987

It can be seen that AdaBoost Classifier is the best model with highest accuracy predictions and very slightly overfitting between training and validating the dataset.

3.b. Business Impact

The elevated AUROC score signifies the model's proficiency in properly discerning collegiate players who possess a high likelihood of being selected in the draft. NBA teams may use this knowledge to enhance their draft plans, players can develop a strategic approach towards their careers, and fans experience heightened enthusiasm for their favorite players. The consequences of inaccurate outcomes, such as an NBA team, may include a decline in team performance and a loss of competitive edge. A player with a lower anticipated probability may risk losing a chance to be picked. A subpar model might also diminish the level of anticipation among fans.

3.c. Encountered Issues

There was an error as "no module of src" when i tried to import a function from `src/data`. To solve this problem, a path to home folder is set up by using `sys.path.insert(0, '..')`.

4. FUTURE EXPERIMENT

Reflect on the experiment and highlight the key information/insights you gained from it that are valuable for the overall project objectives from a technical and business perspective.

4.a. Key Learning

The present algorithm performs flawlessly in predicting future NBA players. The hyperparameter settings are examined and an optimal set is found to reduce overfitting.

4.b. Suggestions / Recommendations

The chosen algorithm and hyperparameters set for this experiment is deemed sufficiently effective in completing the project.

