

ML Final Report

Mr. O'clock insanely practices

December 22, 2024

1 Introduction

This final project is aim to use machine learning to predict HTMLB outcome. We have explore through many algorithms and data preprocess, and result in 0.58989 accuracy in predicting same season outcome(stage1) and 0.55473 in predicting a whole new season(stage2). We are going to break it down what have we done with every model and preprocess and what are the results.

2 Data Preprocess

Dealing with NAN

There are many NAN in all kinds of data. To deal with them, we can either fill in 0, fill in mean of other non-NAN rows, or fill in mean of the data of same team.

Team Attribute and Pitcher Attribute

We try to add 60 boolean columns: `is_home_team_{team name}` and `is_away_team_{team name}` which indicate that if home team or away team equal to team name. Same method can be used on pitcher name.

However we try it on random forest an observed that these attribute may have low importance. We guess that perhaps this kind of data is too scatter to make it important

Alternatively, we can labelize string columns, this may work in algorithm based on decision tree.

Drop std and skew

Since we think std and skew has little relationship with winning or lose. We try to drop all of those column.

The result is that there are no difference in accuracy. However, training speed is faster using this data.

3 Models

KNN

PLA

Linear Regression

Logistic Regression

Logistic regression is a machine learning technique that predicts the probability of an outcome based on the weighted contribution of predictors. The sigmoid function ensures the probabilities are valid, and the coefficients are adjusted during training to maximize how well the model explains the observed data.

This is one of the very first models we have tried on, and the best one among all trained logitic regression models used the parameters as follows:

```
model = LogisticRegression(solver='liblinear', max_iter=1000, random_state=42)
```

where the logistic regression model is included in the package *sklearn.linear_model*.

We split the provided training data into 80% for training and 20% for validation. The above model could reach a training accuracy of 63.42% and validation accuracy 54.52%. The public score we got on kaggle is 56.9% for stage 1, and 57.64% for stage 2. For private, we got an accuracy of 57.1% for stage 1 and 53.85% for stage 2. In stage 2, a larger gap emerges between the public and private scores, which we attribute to using the same training technique for both stages. We did not account for the influence of time, an especially critical factor in stage 2, which should have been incorporated into our model. Overall, it's an efficient model but assumes a linear relationship which might not always be the case especially in complex problems like this.

SVM

Random Forest

CATboost

Catboost is an algorithm based on Gradient Boosting Tree, which is similar to XGboost and Light-GBM. However, Catboost can deal with string data type directly, which is a very good characteristic for this problem since we think string attribute like team name or pitcher name are important for training and cannot be dropped.

In all of the following experiment, i use mean of column to replace numerical column's NAN, and use 'NAN' to replace all NANs in string column.

In the initial tries, we use parameters as following

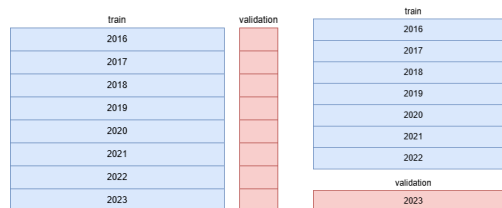
```
iterations = 500, depth = 5, learning_rate = 0.05, loss_function = 'Logloss'
```

This is our first algorithm that reaches 0.58 in public test. However, it is still not good enough to pass benchmark.

To find a better set of parameters, we use optuna to help us find it. And to mimic to stage1 / stage2 situation we will met, we use following different strategy to cut training and validation test for optuna.

stage1 Since we are predicting the last half season given each year's first half season, we use the last 0.2 of each season's data as validation data, and other as training data, and feed these to optuna to find the best parameters.

stage2 We are predicting a whole new year. Hence we take 2023's data as validation data, and other as training data, and feed these to optuna to find best parameters.



When using optuna, we search in following range

```
params = {
    "iterations": trial.suggest_int("iterations", 100, 1000),
    "depth": trial.suggest_int("depth", 4, 10),
    "learning_rate": trial.suggest_float("learning_rate", 1e-3, 0.3, log=True),
    "l2_leaf_reg": trial.suggest_float("l2_leaf_reg", 1, 10),
    "bagging_temperature": trial.suggest_float("bagging_temperature", 0, 1),
    "random_strength": trial.suggest_float("random_strength", 0.1, 10),
    "loss_function": "Logloss",
}
```

The best parameters lead to a 0.57985 accuracy in stage1. And 0.54493 in stage2. Which is not quite stable between two stages.

Package Reference

```
optuna, catboost
```

Blending and Bagging

4 Conclusion