# EXPERIMENT REPORT

| | |
|---|---|
| **Student Name** | Mark Brackenrig (12964298) |
| **Project Name** | Assignment 1B |
| **Date** | 14/02/2021 |
| **Deliverables** | notebooks/brackenrig_mark-12964298-week2-SMOTE.ipynb<br>models/brackenrig_mark_12964298_week2_SMOTE.sav<br>github: https://github.com/roger-yu-ds/assignment_1 |

| 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 goal of this project is to predict whether an NBA 'rookie' will have a career that spans at least 5 years. A rookie is an NBA player who is in the first year of their career.<br><br>This model has numerous applications – for example, a collectibles investor could use this prediction to inform purchases of 'rookie' basketball cards while they are still cheap. |
| **1.b. Hypothesis** | The previous experiment highlighted that using an ensemble model of two different classifiers (random forest and logistic regression) saw a slight improvement in AUC. The main weakness of this model was that it struggled to predict the minority class.<br><br>Based on the results of the previous experiment, I think that using an up sampling method would improve results by allowing the model to better learn patterns in the minority class. |

| 1.c. Experiment Objective | The objective is to test whether training the model on a dataset with an up-sampled minority class will improve the AUC on the test dataset. |
|---|---|

## 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 | A shared function `download_data` was used to retrieve data from the Kaggle API, unzips the files and saves them into the raw data directory. This data contains two CSV files – 'train' and 'test'.

The training dataset is cleaned by removing the ID columns 'Id' and 'Id_old'. The test dataset is split into 'X_test' (without the id column) and 'test_id' – a pandas series of the Ids.

Next, the training set is split into the independent and dependent variables using the shared function `separate_target`.

Lastly, the training data is then transformed using SMOTE (Synthetic Minority Oversampling Technique). SMOTE randomly selects one of the K nearest neighbours to synthetically generate a data point between the two data points in the feature space. SMOTE will create new data points until the classes are represented evenly in the sample.

There was no notable data cleaning performed on the dataset as it was a generally clean dataset with no missing values, or obvious erroneous entries.

Unlike experiment 1, the models were trained on the entire training dataset rather than 80% of the training data. This decision was made as cross validation was used, so generating a new validation set was unnecessary. |
|---|---|

| | |
|---|---|
| **2.b. Feature Engineering** | Feature engineering was the same as experiment 1. The baseline model was used to conduct the test – which was a VotingClassifier model consisting of a linear regression model and a random forest model.<br><br>The logistic regression model applied the StandardScaler method for normal standardization, and Principal Components Analysis was used to create orthogonal independent variables. This was used due to the high levels of collinearity identified in the dataset (see appendix for correlation matrix). No feature engineering was used for the random forest model. |
| **2.c. Modelling** | The model uses a voting classifier with a random forest classifier and a logistic regression classifier as inputs. The voting classifier uses the 'soft' voting method which classifies the result based on the highest average probability of the underlying classifiers.<br><br>The Hyperparameter tuning was kept the same as in experiment one. In this experiment, testing the original model was tested against the original model trained on the oversampled training set.<br><br>CalibrationClassifierCV was used to apply probability calibration on training of the new sample. Since the model included synthetic data in the training, it was important to calibrate the probabilities generated.<br><br>CalibrationClassifierCV was also used to train the model without calibrating against the original training data, but using cross-validation using the SMOTE dataset. |

---

| 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** | BASE MODEL<br>AUC (Training): 0.713<br>AUC (Test – Kaggle): 0.709<br><br>(On raw training data)<br>Precision (Negative Class / Positive Class): 0.54/0.84 |

| | |
|---|---|
| | Recall (Negative Class / Positive Class): 0.07/0.99<br>F1 (Negative Class / Positive Class): 0.13/0.91<br><br>SMOTE<br>AUC (Training): 0.720<br>AUC (Test – Kaggle): 0.711<br><br>(On raw training data)<br>Precision (Negative Class / Positive Class): 0.3/0.9<br>Recall (Negative Class / Positive Class): 0.59/0.72<br>F1 (Negative Class / Positive Class): 0.39/0.80<br><br>SMOTE – Calibrated on Original Data<br>AUC (Training): 0.720<br>AUC (Test – Kaggle): 0.711<br><br>(On raw training data)<br>Precision (Negative Class / Positive Class): 0.65/0.83<br>Recall (Negative Class / Positive Class): 0.01/1<br>F1 (Negative Class / Positive Class): 0.02/0.91<br><br>SMOTE – Calibrated using CV on Synthetic Data<br>AUC (Training): 0.723<br>AUC (Test – Kaggle): 0.712<br><br>(On raw training data)<br>Precision (Negative Class / Positive Class): 0.3/0.9<br>Recall (Negative Class / Positive Class): 0.59/0.72<br>F1 (Negative Class / Positive Class): 0.4/0.80 |
| **3.b. Business Impact** | The best performing model was the SMOTE model using cross validation calibration on the synthetic data. The model has improved its ability to predict the minority class marginally. Unfortunately, the minor improvement is not large enough to realistically improve the business applications of the model.<br><br>In saying this, small systematic improvements to the model will eventually create business value by improving predictions. By using synthetic data, we were able to improve on the AUC of the model marginally. To a certain extent, this has overcome the issue encountered in the first experiment where there wasn't enough data in the minority class. |

| 3.c. Encountered Issues | While the SMOTE technique improved predictions, using synthetic data in model training needs to be used with caution. By using synthetic data, you are essentially 'guessing' what the classes should look like, and allow your model to train itself on this guess.<br><br>A way around this could be using a bootstrap resampling method to under sample the majority class. While this would mean the model is trained on less data, it also avoids the issue of generating synthetic data.<br><br>While this experiment focused on resampling the data, feature engineering, and base classifiers have not been explored in great detail. It is entirely possible that we are not using the most appropriate features or base classifiers for this task. More exploration is needed in this area. |
| --- | --- |

| 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 | Our hypothesis was correct in assuming that the original model struggled to predict the minority class, and therefore an oversampling method improved the model.<br><br>Our learnings thus far have shown that ensembling different classifiers and reweighting the training data has resulted in improving the prediction on the test set.<br><br>We still need to complete more exploration into base classifiers and feature engineering as these have not been explored in great detail. |
| --- | --- |
| 4.b. Suggestions / Recommendations | Future experiments should revisit comparing base classifiers. There are potentially many more complex classifiers not yet considered in our modelling process which could improve on the final model we use.<br><br>The only feature engineering used to date has been PCA. Various other techniques could be employed (such as Non-negative matrix factorization or t-SNE) to account for non-linearity and improved information preservation.<br><br>However, regardless of the modelling or feature engineering methods employed, there is an inherit limitation based on the |

dataset. Since we are trying to predict whether an NBA rookie will have a 5 year long career, it is entirely possible that there a level of randomness in the data. For example, career ending injuries to a star prospect would be a significant outlier in our dataset. While this is not possible in this context, in a professional context finding more data (such as reason for career ending) could be used to improve the model. In this context, we could potentially look at excluding or downvoting outliers to improve the model.
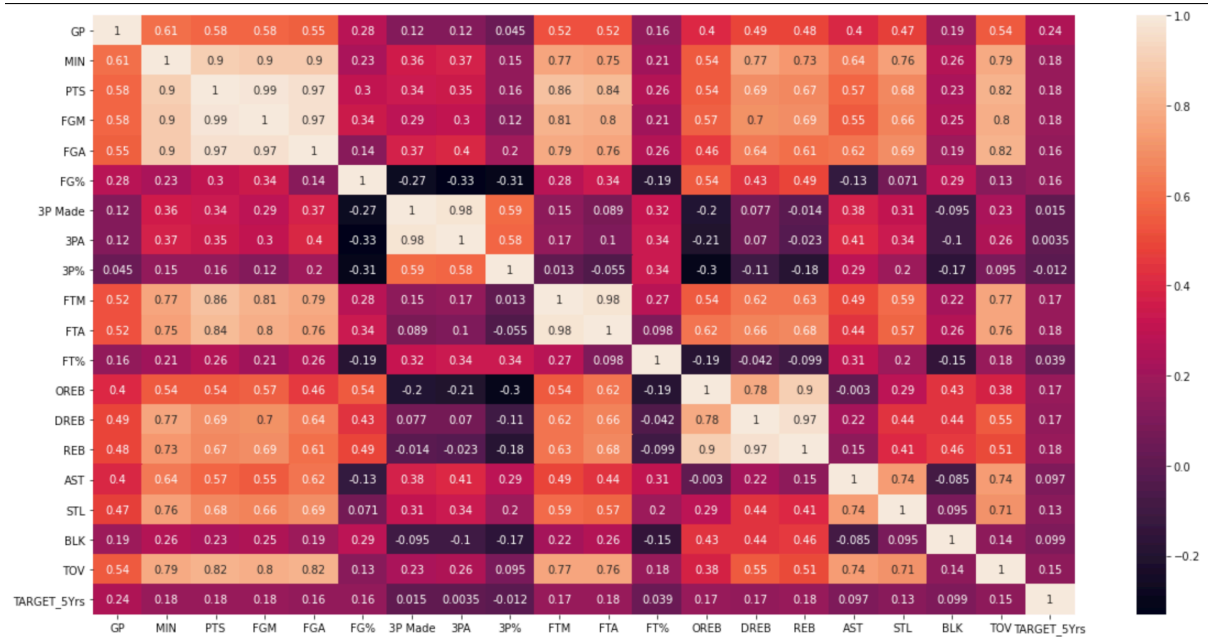
# Appendix

## Correlation Matrix



*Figure 1: Correlation Matrix Heatmap*