



Prediction tool for mixed martial arts bout outcome

Danil Burov

Fontys University of Applied Sciences

January 14, 2025

Contents

| | | |
|----------|-------------------------------|-----------|
| 1 | Introduction | 2 |
| 2 | Business understanding | 2 |
| 3 | Data understanding | 2 |
| 4 | Data preparation | 3 |
| 4.1 | Deleting columns | 3 |
| 4.2 | One-Hot-Encoding | 5 |
| 4.3 | Scalling | 5 |
| 4.4 | Factorzing | 6 |
| 4.5 | Filling NaN values | 7 |
| 5 | Modelling | 8 |
| 5.1 | Optimization | 10 |
| 6 | Conclusion | 15 |

1 Introduction

The purpose of this document is to report in a comprehensive way all the work that has been done throughout my fifth semester regarding the creation of an AI tool for predicting the outcome of MMA bouts. To do so, I will explain each step that I took in order to deliver the final outcome of the project.

2 Business understanding

The primary objective of this project is to develop a predictive tool that provides insights into the outcome of a bout between two fighters. Rather than offering a binary prediction (0 or 1, representing true or false outcomes), the tool generates probabilistic predictions on a scale between 0 and 1. This allows users to understand the likelihood of various outcomes, such as the probability of a fighter submitting their opponent, achieving a knockout, or winning by decision. By delivering granular insights, the tool empowers users to make data-driven assessments of match outcomes.

One notable risk associated with the project is its potential misuse for betting purposes. The tool could be exploited to gain unfair advantages in gambling scenarios. To mitigate this concern, the tool will be kept private and its access limited to legitimate and approved stakeholders.

Beyond the potential concerns, the tool holds immense value as an analytical resource for fighters and their coaching teams. Fighters can leverage the tool to analyze their strengths and weaknesses relative to their opponents. For instance, by understanding the likelihood of success in specific fighting techniques, fighters can tailor their training regimens to exploit their opponent's vulnerabilities and address their own shortcomings.

3 Data understanding

The dataset selected for this project is sourced from Kaggle and contains data on every UFC fight from mid-2010 to the present. Each row in the dataset represents an individual bout, making it well-suited for predictive modeling and analysis. This dataset was chosen for several reasons: it is regularly updated on a weekly basis, ensuring that the data remains current, and it already includes pre-processed statistics such as the difference in metrics between the two fighters in each bout. Additionally, it provides key variables such as betting odds, fighter rankings, and fight statistics, as well as the target variable for this project, the bout "Winner."

Despite its strengths, including over 6,000 entries and 118 features, there are still several limitations to consider regarding the accuracy and completeness of the model. Many critical factors that influence fight outcomes are not included

in the dataset due to their subjective or difficult-to-quantify nature. For example:

Mental State: A fighter’s psychological readiness, confidence, and focus leading up to the fight can significantly affect their performance.

Health Status: Factors such as injuries, illnesses, or the general health of the fighter during training camp or on fight day are not accounted for.

Weight Cutting and Rehydration: Many fighters undergo extreme weight cuts before weigh-ins, only to rehydrate and regain substantial weight before the fight. The extent of their weight recovery and its impact on their physical condition during the bout is a crucial factor that is not captured in the dataset.

Moreover, many day-to-day factors influencing a fighter’s performance are inherently subjective and nearly impossible to gather systematically. Fighters themselves may not be fully aware of how these aspects—such as their recovery quality, nutrition, or emotional state—affect their performance in a bout. This lack of granular data presents a challenge in achieving high accuracy for predictive models.

While these limitations highlight gaps in the dataset, they also underscore the opportunity for future research and data collection in combat sports. If more detailed physiological and psychological data could be integrated into such datasets, it would pave the way for more accurate and robust predictive models in the domain of mixed martial arts.

4 Data preparation

In order to explain the data preparation done regarding the ufc dataset I will reference most of my explanations to the appendix. This was done to keep a neatly organized text. I will be using snippets from my Jupyter notebook. All these features are displayed here:

4.1 Deleting columns

When the dataset was chosen it was known that many of the features will not be used for the purpose of the tool. All these features were dropped for a reason and in the upcoming subsection it will be explained why. Initially the number of features was 118. After dropping the unnecessary features, the number was

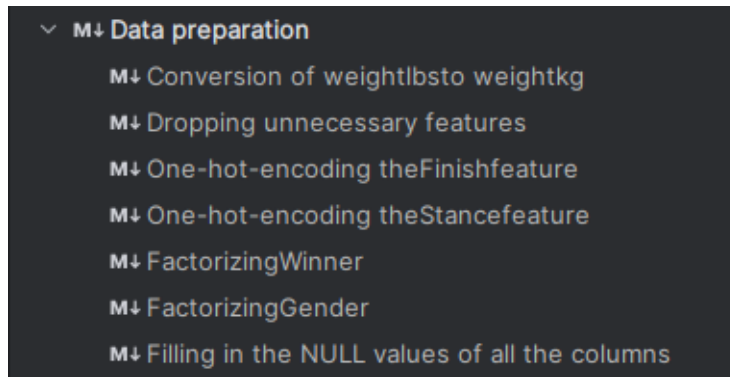


Figure 1: Overview of data preparation processing steps

reduced to 49.

These features were rankings, streak orientated, or keeping track of the fighters' decisions and win/lose streak, and other miscelanous features that are irrelevant.

The code can be found below:

```

1 features_to_drop = [
2
3     # Blue features
4     "BlueFighter", "BlueExpectedValue", "BMatchWCRank", "BWFlyweightRank", "BWFeatherweightRank",
5     "BWStrawweightRank", "BWBantamweightRank", "BHeavyweightRank",
6     "BLightHeavyweightRank", "BMiddleweightRank", "BWelterweightRank",
7     "BLightweightRank", "BFeatherweightRank", "BBantamweightRank",
8     "BFLyweightRank", "BPFPRank", "BlueCurrentLoseStreak", "BlueCurrentWinStreak",
9     "BlueDraws", "BlueLongestWinStreak", "BlueLosses", "BlueTotalRoundsFought", "BlueTotalTitleBouts",
10    "BlueWinsByDecisionMajority", "BlueWinsByDecisionSplit", "BlueWinsByDecisionUnanimous",
11    "BlueWinsByKO", "BlueWinsBySubmission", "BlueWinsByTKODoctorStoppage",
12    "BlueWins",
13
14    # Red features
15    "RedFighter", "RedExpectedValue", "RMatchWCRank", "RWFlyweightRank", "RWFeatherweightRank",
16    "RWStrawweightRank", "RWBantamweightRank", "RHeavyweightRank",
17    "RLightHeavyweightRank", "RMiddleweightRank", "RWelterweightRank",
18    "RLightweightRank", "RFeatherweightRank", "RBantamweightRank",
19    "RFLyweightRank", "RPFPRank", "RedCurrentLoseStreak", "RedCurrentWinStreak", "RedDraws", "RedLongestWinStreak", "RedLosses",
20    "RedTotalRoundsFought", "RedTotalTitleBouts", "RedWinsByDecisionMajority",
21    "RedWinsByDecisionSplit", "RedWinsByDecisionUnanimous", "RedWinsByKO",
22    "RedWinsBySubmission", "RedWinsByTKODoctorStoppage", "RedWins",
23
24    # Other
25    "Date", "Location", "Country", "TitleBout", "WeightClass", "EmptyArena", "BetterRank", "FinishDetails", "FinishRoundTime",
26 ]
27
28 ufc_data_master = ufc_data_master.drop(columns = features_to_drop)
29
30 ufc_data_master.columns
31 [7]

```

Figure 2: Unnecessary features dropped

4.2 One-Hot-Encoding

One-hot-encoding is used to transform categorical values into numerical values. In this case the 'Finish' feature had 5 different types of finishes. When I one-hot-encode the feature I get 5 different columns with all the possible finishes and every bout (each row in the dataset) has a true or false to the according finish. It is not a good idea to factorize with that many different options because the model could develop bias towards the bigger numbers.

One-hot-encoding was used not only to encode the 'Finish' feature but also the 'Stance' feature. 'Stance' had three possible values 'Southpaw', 'Orthodox', 'Switch'.

The code for one-hot-encoding the 'Finish' feature can be found:

```
One-hot-encoding the 'Finish' feature
One-hot-encoding is used to transform categorical values into numerical values. In my case the 'Finish' feature had 5 different types of finishes. When I one-hot-encode the 'Finish' feature I get 5 different columns with all the possible finishes and a match has a true or false to the according finish. It is not a good idea to use factorizing with that many different options because the model could develop bias towards a bigger number.
Issue -> The issue it solves is the models I am using need to have a numerical value instead of categorical.
Solution -> One-hot-encoding is used instead of factorizing to reduce the model bias.

ufc_data_master = pd.get_dummies(ufc_data_master, columns = ['Finish'], prefix='Finish')
ufc_data_master.head()
[10]
```

Figure 3: One-Hot-Encoding 'Finish' feature

The code for one-hot-encoding the 'Stance' feature can be found below:

```
One-hot-encoding the 'Stance' feature
I have three different stances that a fighter could have. I am one-hot-encoding it to make the categorical data to numerical, since the models I am using require numerical data.

1 ufc_data_master = pd.get_dummies(ufc_data_master, columns = ['RedStance'], prefix = 'RedStance')
2 ufc_data_master = pd.get_dummies(ufc_data_master, columns = ['BlueStance'], prefix = 'BlueStance')
3
4 ufc_data_master.head()
[11]
```

| | RedOdds | BlueOdds | Winner | Gender | NumberOfRounds | BlueAvgSigStrLanded | BlueAvgSigStrPct | BlueAvgSubAtt | BlueAvgTDLanded | BlueAvgTDPct |
|---|---------|----------|--------|--------|----------------|---------------------|------------------|---------------|-----------------|--------------|
| 0 | 558.0 | -888.0 | Blue | MALE | 5 | 4.39 | 0.57 | 0.0 | 0.44 | |
| 1 | 250.0 | -310.0 | Blue | MALE | 3 | 3.24 | 0.57 | 1.3 | 6.39 | |
| 2 | 205.0 | -250.0 | Red | MALE | 3 | 2.10 | 0.37 | 1.0 | 1.97 | |
| 3 | 360.0 | -470.0 | Blue | FEMALE | 3 | 2.89 | 0.48 | 1.0 | 2.68 | |
| 4 | -485.0 | 370.0 | Red | MALE | 3 | 5.35 | 0.57 | 0.0 | 0.92 | |

Figure 4: One-Hot-Encoding 'Stance' feature

4.3 Scalling

In general all the features in the dataset need to be scaled to standard deviation. That is the reason why all the features that had percentage were scaled to a decimal to the first zero.

Additionally I converted the 'Weight' feature from LBS to KG, because it is more convenient for me.

The code can be found below:

```
Conversion of weight (lbs) to weight (kg)  
This conversion is a personal preference, I feel more confident working with kilos instead of pounds  
  
1 ufc_data_master["BlueWeightKg"] = ufc_data_master["BlueWeightLbs"] * 0.453592  
2 ufc_data_master["RedWeightKg"] = ufc_data_master["RedWeightLbs"] * 0.453592  
3  
4 ufc_data_master = ufc_data_master.drop(columns=["BlueWeightLbs", "RedWeightLbs"])  
[9]
```

Figure 5: Converting the 'Weight' feature from lbs to kg

4.4 Factorizing

For some features one-hot-encoding is not exactly needed. In this case it was decided to factorize values that can have only two possible values. However in some cases it could be better to even use one-hot-encoding with these features as well because some models can deal with more features more efficiently. The reason these features were factorized is because in the beginning the project was supposed to train a classification model where the output is either 0 or 1, however the idea changed but factorizing is still good enough for the model that is currently being used.

Since the dataset contains both genders the feature that was factorized is the 'Gender' where the 0 is man and the 1 is a woman.

Another feature that has only two values is the 'Winner' feature, which has either red or blue, depending on who won the fight (Red and Blue are universal colours used to indicate the corner that fighter fought from.). When factorized, 0 corresponds to Red and 1 corresponds to Blue.

Code for factorizing the 'Gender' feature can be found below:

```
Factorizing 'Gender'

The 'Gender' feature was factorized for the same reason as the 'Winner' feature. I had only two possible values 'Man' or 'Woman', now the 0 is man ** and 1 is woman.

Using one-hot-encoding is not best practise here because the feature is binary, always either 0 or 1. When there are more than 2 values then one-hot-encoding should be considered. As well as, one-hot-encoding would have created more columns which just adds up to the training time in this case.

ufc_data_master['Gender'], gender_mapping = pd.factorize(ufc_data_master['Gender'])

ufc_data_master.head()

[12]
```

Figure 6: Factorizing the 'Gender' feature

Code for factorizing the 'Winner' feature can be found below:

```
Factorizing 'Winner'

The 'Winner' feature had either 'Red' or 'Blue' as a winner, in order to normalize it I factorized the feature providing me with a binary value 0 or 1, where 0 is Red and 1 is Blue

This was needed because my models could not use the string type of the feature to model the data

ufc_data_master['Winner'], winner_mapping = pd.factorize(ufc_data_master['Winner'])

ufc_data_master.head()

[13]
```

Figure 7: Factorizing the 'Winner' feature

4.5 Filling NaN values

The dataset has **415296** total entries from all features from which **19601** are NULL values. There were three options how to deal with all the different values. First was to fill in the missing values with 0 or specific string text wherever is needed. Second was to delete the rows that had these NULL values. However, the condition for this is the rows that are missing values to be less than 10 percent. And the last possible action was to develop a ML model that feels in these values. This option was not developed.

Most of the features did not have a significant impact on the accuracy of the model that is why they were filled with 0. Most of the rows that had NULL features were both from the red and blue corner (e.g. The 'RedOdds' feature and 'BlueOdds' feature in the row were both NULL so filling them with 0 just fixes the empty value)

Code for filling the NaN values can be found below:


```

The rows that had NaN values were dropped because they are <10% of the whole dataset.

1 ufc_data_master = ufc_data_master.dropna(subset=['FinishRound'])
2 ufc_data_master = ufc_data_master.dropna(subset=['TotalFightTimeSecs'])
[14]

1 null_values_finish_rounds = ufc_data_master['FinishRound'].isnull().sum()
2 null_values_finish_round_time = ufc_data_master['FinishRoundTime'].isnull().sum()
3 null_values_total_fight_time = ufc_data_master['TotalFightTimeSecs'].isnull().sum()
4
5 print(f"Number of null values in 'FinishRound': {null_values_finish_rounds}")
6 print(f"Number of null values in 'FinishRoundTimeSecs': {null_values_finish_round_time}")
7 print(f"Number of null values in 'TotalFightTimeSecs': {null_values_total_fight_time}")
[15]

```

Figure 8: filling the NaN values

5 Modelling

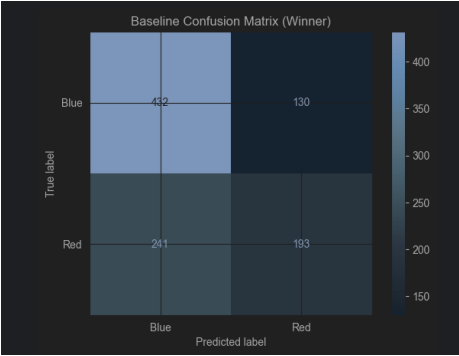
Choosing a model for the data was difficult given the reason that the project idea was to not predict if a fighter is going to win (0 or 1), but rather see the likelihood of a fighter winning in three different outcomes against the opponent. Knowing this, logistic regression seemed like a good fit for the project.

Logistic regression under the hood is a classification model that uses the Sigmoid function. Although, given its classification nature, it provides the probability output between 0 and 1.

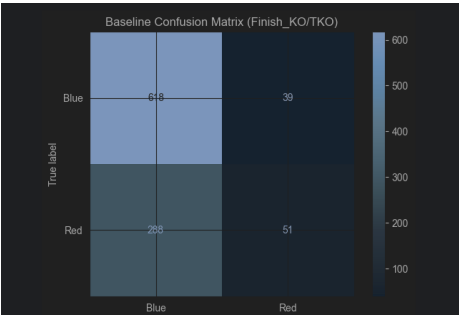
In this project, four different models were trained with the data. This is due to the fact that the target variable needed to be existent in the dataset. In the beginning, an imaginary target variable was created, which was a combination of all the targets: 'Winner', 'Finish_KO/TKO', 'Finish_SUB', 'Finish_U-DEC'. For this project, other finishes were not taken into account for simplicity reasons. Although further development could be the prediction of the rest finish outcomes of a match, e.g., 'Doctor stoppage'.

Four different baseline models were trained for the purpose of the project. Each one having the according target variable. The winner prediction model has a baseline accuracy of **0.62%**, the KO/TKO has **0.67%**, submission has **0.80%**, and the unanimous decision model has **0.66%** accuracy. However, looking into the confusion matrixes of all the models it was easily established that the model has a big bias towards the blue corner. Below will be displayed the confusion matrixes of all the baseline models.

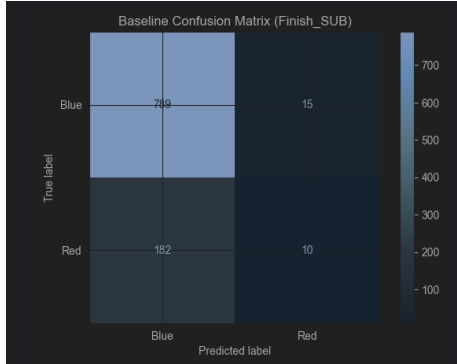
Winner baseline model



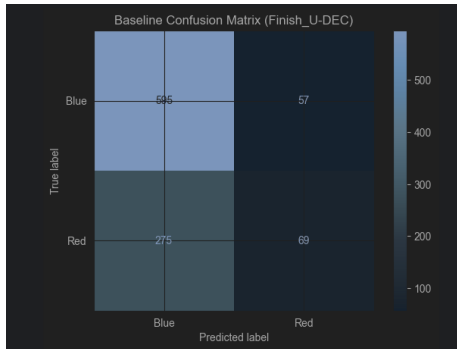
Finish_KO/TKO baseline model



Finish_SUB baseline model



Finish_U-DEC baseline model

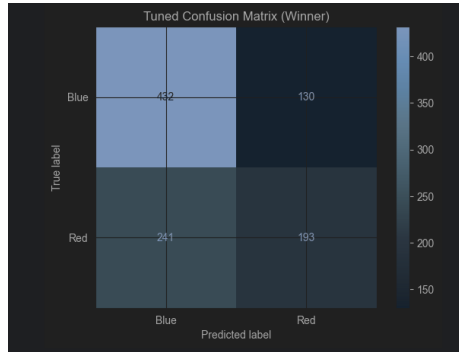


It is very clear that the models are not really predicting well. All of them seem to have some sort of bias towards the blue fighter and they are mainly predicting if the blue fighter will win and they cannot predict that well if a red fighter will win.

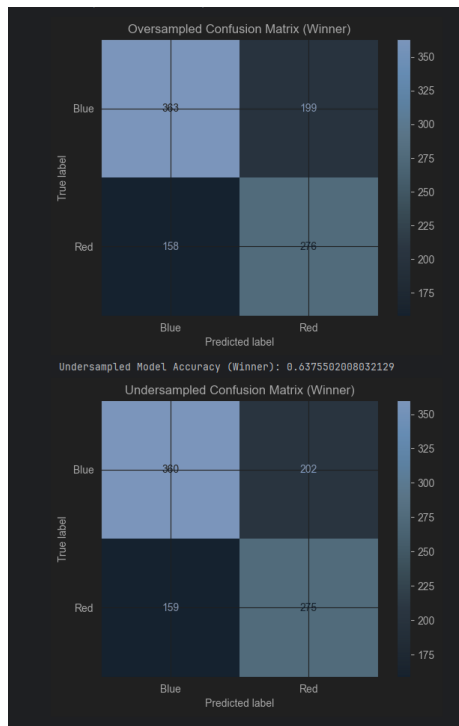
5.1 Optimization

Logistic regression was used to train four models with four different target variable. However on all of them this bias towards the blue corner was persistent. In order to mitigate this bias, balancing the set was needed. Hyperparameter tuning, oversampling and undersampling were performed on the datasets and then the data was trained again using Logistic regression. For all models it seemed that oversampling worked best for balancing the set, however regarding accuracy it wasn't always the best pick. Below I will showcase all the balancing results for all four models.

Hyperparameter Tuning Score - Winner

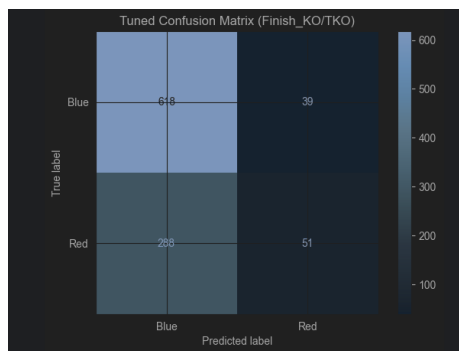


Oversampling/Undersampling Score - Winner

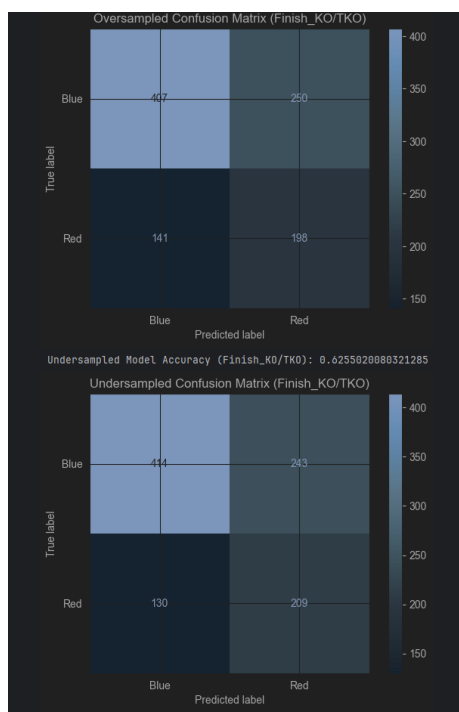


In this model, hyperparameter tuning and balancing the set results in a better prediction after all. It is seen that in the oversampling technique the confusion matrix shows a better prediction when it comes to red winners, however with a trade-off for more false predictions on the blue side.

Hyperparameter Tuning Score - Finish_KO/TKO

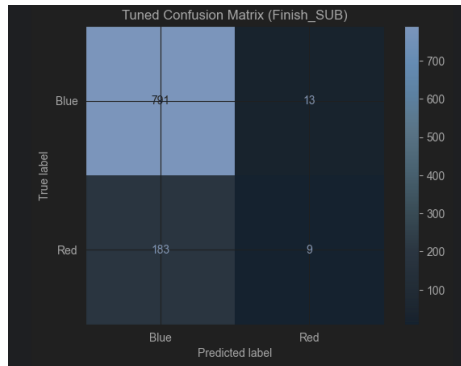


Oversampling/Undersampling Score - Finish_KO/TKO

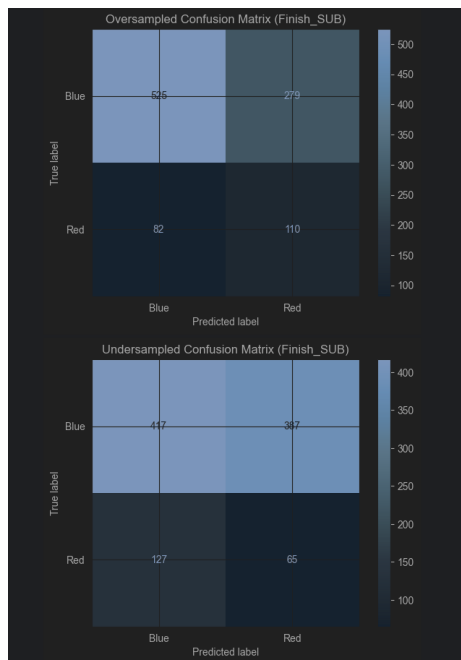


Here the techniques applied do not result in a favourable result. Although, using oversampling did increase the true predictions for red, they are still quite low and the prediction is not reliable.

Hyperparameter Tuning Score - Finish_SUB

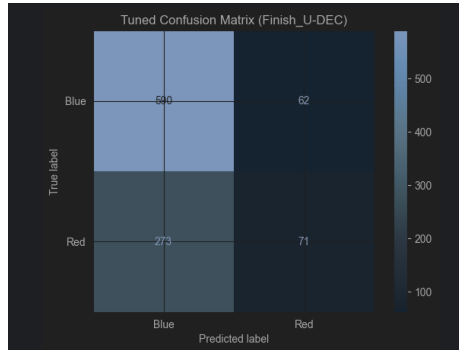


Oversampling/Undersampling Score - Finish_SUB

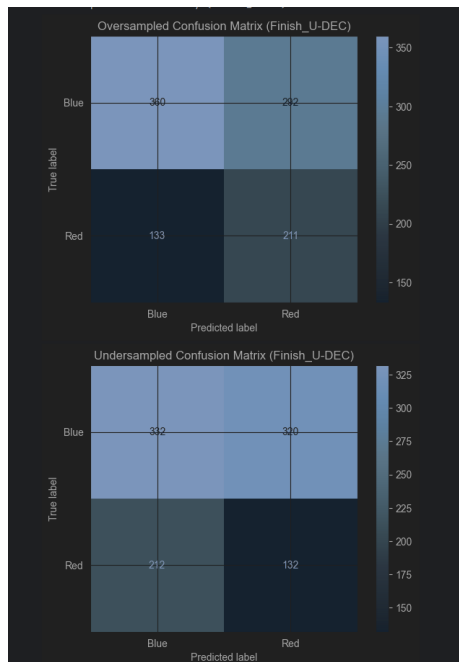


Unfortunately, this is the worst predicting model. There is an exponentially high bias towards the blue fighter and even with balancing of the set the bias is still too high.

Hyperparameter Tuning Score - Finish_U-DEC



Oversampling/Undersampling Score - Finish_U-DEC



This model has the same outcome as the submission model. It just does not perform well at all. Based on the confusion matrixes it is seen that there is an even distribution of true - positivities and true-negatives which is absolutely incorrect.

It is very clear that hyperparameter tuning doesn't really boost the performance of the model. The model is still predicting blue fighters better than red fighters, however the accuracy in some of the cases does increase.

Opposed to oversampling which boosts the fairness of the model on all of the models. This is due to the fact that the oversampling technique populates the minority class (blue) with synthetic data. However this should be taken with a grain of salt because the synthetic data could be duplicates which might influence the models' accuracy and prediction.

6 Conclusion

Overall, the baseline models are very insufficient and if the model is to be used in any way it needs to have a more balanced set due to the difference between the red and blue winners in the dataset. It makes it so that the dataset develops a bias towards a winner. However, there is at least one good model from each outcome. The highest score of accuracy was achieved in the following optimized models: oversampled winner models with accuracy **0.64%**, hyperparameter tuned finish_ko/tko model with accuracy **0.67%**, hyperparameter tuned submission model with **0.80%**, baseline model for the unanimous decision model **0.66%**. It is very difficult to predict an outcome of a match or which fighter will win. There are many attributes that need to be taken into account. There is a possible further continuation of the project, looking into betting odds. I personally think that this could be a good direction of the project.