DATA Science project

Ibrahima Ndiaye

Football Data Challenge on Kaggle
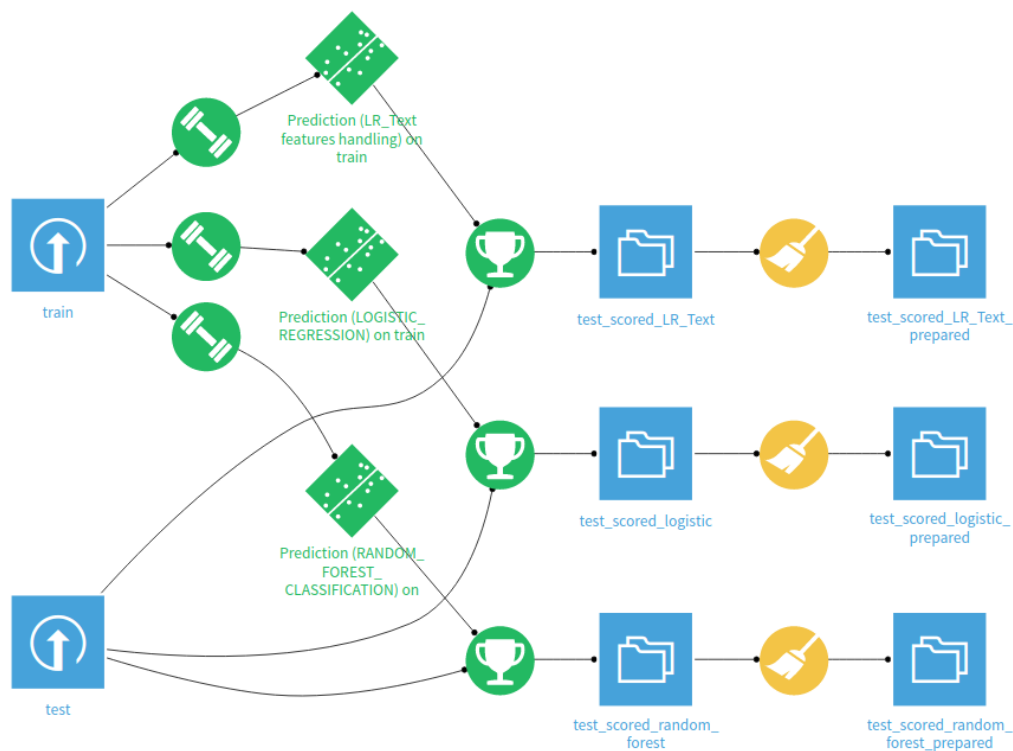
**Project description**

The football data challenge is a kaggle competition found on their website. We are provided with data of the past 8 years of football games and their results in the Italian football league called Serie A .The objective of this project is to be able to predict the result of future football matches with the data provided.
My goal is to be able to predict the result of the games with an accuracy of at least 50%. If the goal is achieved, the project will be considered a success.

I have done 3 submissions on kaggle to be able to know what the score will be in the hopes of getting a score of at least 0.5

This is what my flow looks like at the end of my project:



**The dataset**

The football challenge dataset is the following:
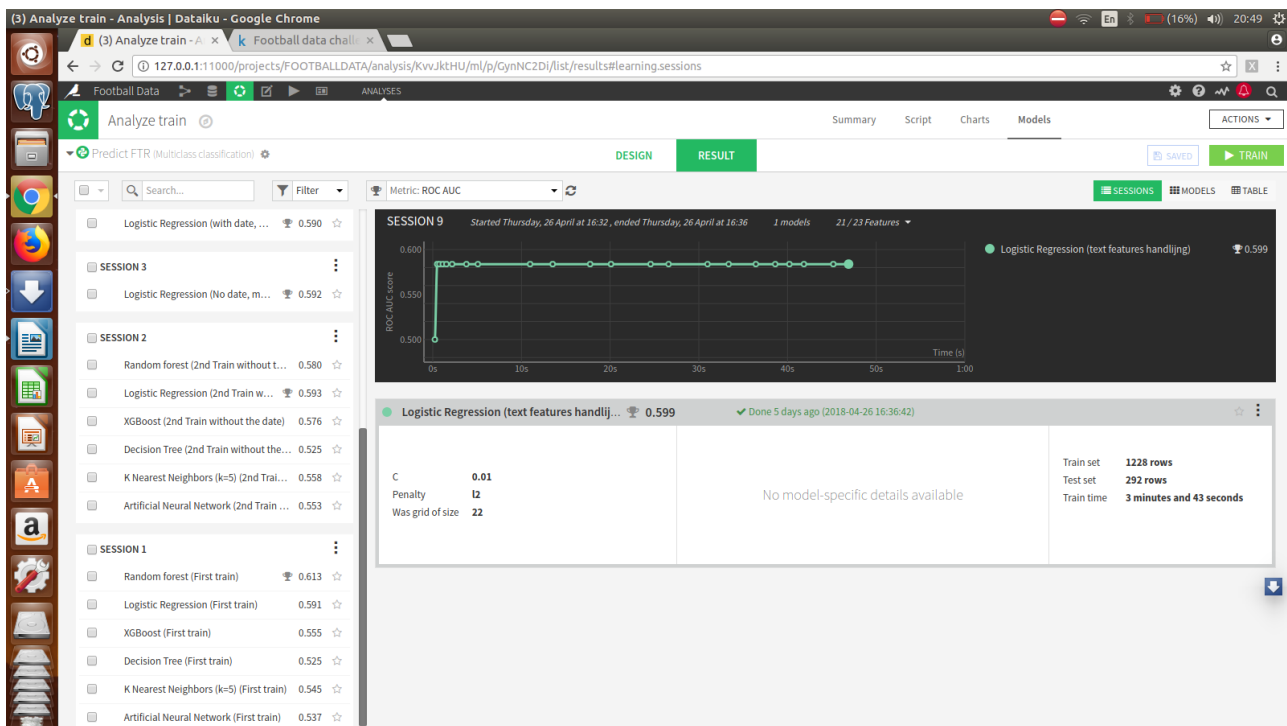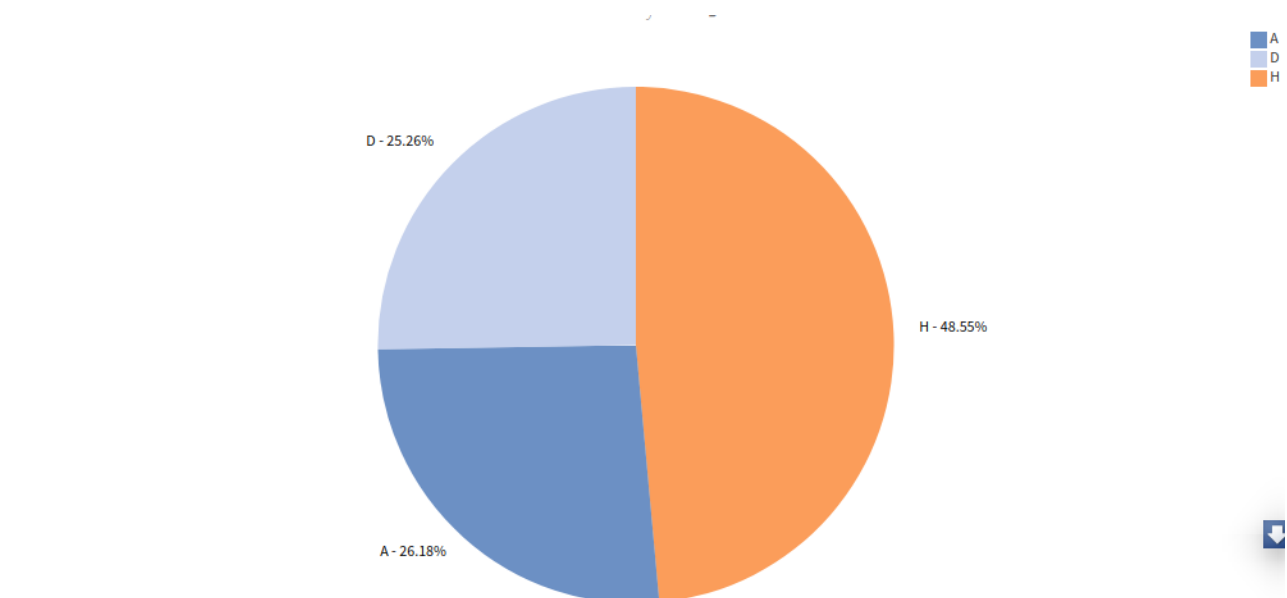
**Train dataset**

- ID = An anonymous ID unique to a given match
- Date = Date of the match
- FTR = Full Time Result (H=Home Win, D=Draw, A=Away Win). This is the target variable

- HomeTeam = Home Team

- AwayTeam = Away Team

- B365H = Bet365 home win odds

- B365D = Bet365 draw odds

- B365A = Bet365 away win odds

- BWH = Bet&Win home win odds

- BWD = Bet&Win draw odds

- BWA = Bet&Win away win odds

- IWH = Interwetten home win odds

- IWD = Interwetten draw odds

- IWA = Interwetten away win odds

- LBH = Ladbrokes home win odds

- LBD = Ladbrokes draw odds

- LBA = Ladbrokes away win odds

- VCH = VC Bet home win odds

- VCD = VC Bet draw odds

- VCA = VC Bet away win odds

- WHH = William Hill home win odds

- WHD = William Hill draw odds

- WHA = William Hill away win odds

With all of the information provided above, our goal is to be able to predict what the full time result of future games will be on the test data set which will have all of the information, except the full time result which we will of course predict.

In order to get varying results, I have made 3 different training models which were individually submitted to Kaggle to see what the best score would be.

Before diving into the procedure of how we have completed our project, there was one important analysis to be made. If we look at the image below, we get a percentage distribution of the results of the games at full time. What is impressive is the amount of games which were won by the home team as opposed to the draws and away team wins. The home team wins about 48.55 percent of the games that they play and this could be due to any number of factors. Football experts usually put it down to the fact that they are playing on their own soil which they are used to and also the support of the home team's fans and the atmosphere they provide. This gives the home team an advantage when playing football games. Also, the same could be said about the away team which are not used to the field of the home team, and may perhaps be intimated by the home team's fans and the overwhelming atmosphere and ambiance created by them which thus affects their performance.

We have done multiple sessions to determine which algorithms and settings would be the best to use in order to get the best prediction results. 9 sessions were made, but we will focus on the main ones which yielded the best results.

Before starting this project, we knew that logistic regression was the most suitable algorithm for this type of situation. But for curiosity purposes and to be able to see the results by myself, I have used 6 different algorithms for the first 2 sessions to better determine which algorithm would be best suited for our goals.

**Algorithms**

- **Random Forest**

- **Logistic Regression**
- **XGBoost**
- **Decision Tree**
- **K Nearest Neighbors**
- **Artificial Neural Network**

These different algorithms provided different results, and the 2 best algorithms were the Random Forest and the Logistic Regression algorithms.
In our first session, the Random Forest was the winner with an ROC AUC score of 0.613.
After playing around with the different features of the Logistic Regression algorithm, I wasn't able to get more than 0.599 and that is only by changing the settings in features handling to Text.
Other than that, the normal score for Logistic Regression is 0.593

## Submissions to Kaggle

In this section, I will explain how I obtained the different scores in the kaggle competition.



## Submission 1 - Random Forest Algorithm

In this first submission, I have used the Random Forest Algorithm as it gave me a pretty good score when I was working with different algorithms in the model prediction section. Also I was curious to find out how it would fare once submitted to Kaggle.

I have left the default settings for the Random Forest algorithm and it gave me a ROC AUC score of  0.613 which was the highest achieved in DSS.



After scoring this last model, we got 4 extra columns :

- Probability that the home team wins – proba_H
- Probability that the away team wins – proba_A
- Probability that the score is a draw - proba_D
- The predicted result of the match – prediction

In order to submit our results to Kaggle, we had to follow the format of having only 2 columns ; the match ID and the FTR which is the predicted result of the match. It was thus necessary to prepare the dataset to conform to this format.

All we had to do was to keep the "ID" column and the "prediction" column, and then proceed by renaming the "prediction" column to "FTR" :

This final format was then submitted to Kaggle and the scores obtained were a private score of 0.49508 and a public score 0.43606.

## Submission 2 – Logistic Regression Algorithm



The second submission was with a Logistic Regression algorithm that yielded at ROC AUC score of 0.593, which was decent but not as good as the previous one achieved with the Random forest.

In the features handling we could see that the defaut category variable type was kept, as we wanted to first see what kind of result we would be getting by just keeping all of the default settings.



After scoring this second model, we got 4 extra columns :

- Probability that the home team wins – proba_H
- Probability that the away team wins – proba_A
- Probability that the score is a draw - proba_D

- The predicted result of the match – prediction

In order to submit our results to Kaggle, we had to follow the format of having only 2 columns ; the match ID and the FTR which is the predicted result of the match. It was thus necessary to prepare the dataset to conform to this format.

All we had to do was to keep the "ID" column and the "prediction" column, and then proceed by renaming the "prediction" column to "FTR" :





We now have the correct format for submission to kaggle. This second submission achieved a private score of 0.44590 and a public score of 0.45245. In my opinion this was a pretty balanced score and most probably reflected the most correct score out of all the submissions.

## Submission 3 - Logistic Regression Algorithm with Text Variable Type



The ROC AUC score that we have obtained using the Logisctic Regression algorithm while using the Text Variable Type in the features handling section was 0.599 as we can see above.



This last submission was made using the Text Variable Type instead of Categorical Variable type which was used by default.

After scoring this last model, we got 4 extra columns
- Probability that the home team wins – proba_H
- Probability that the away team wins – proba_A
- Probability that the score is a draw - proba_D
- The predicted result of the match – prediction

In order to submit our results to Kaggle, we had to follow the format of having only 2 columns ; the match ID and the FTR which is the predicted result of the match. It was thus necessary to prepare the dataset to conform to this format.

All we had to do was to keep the "ID" column and the "prediction" column, and then proceed by renaming the "prediction" column to "FTR" :
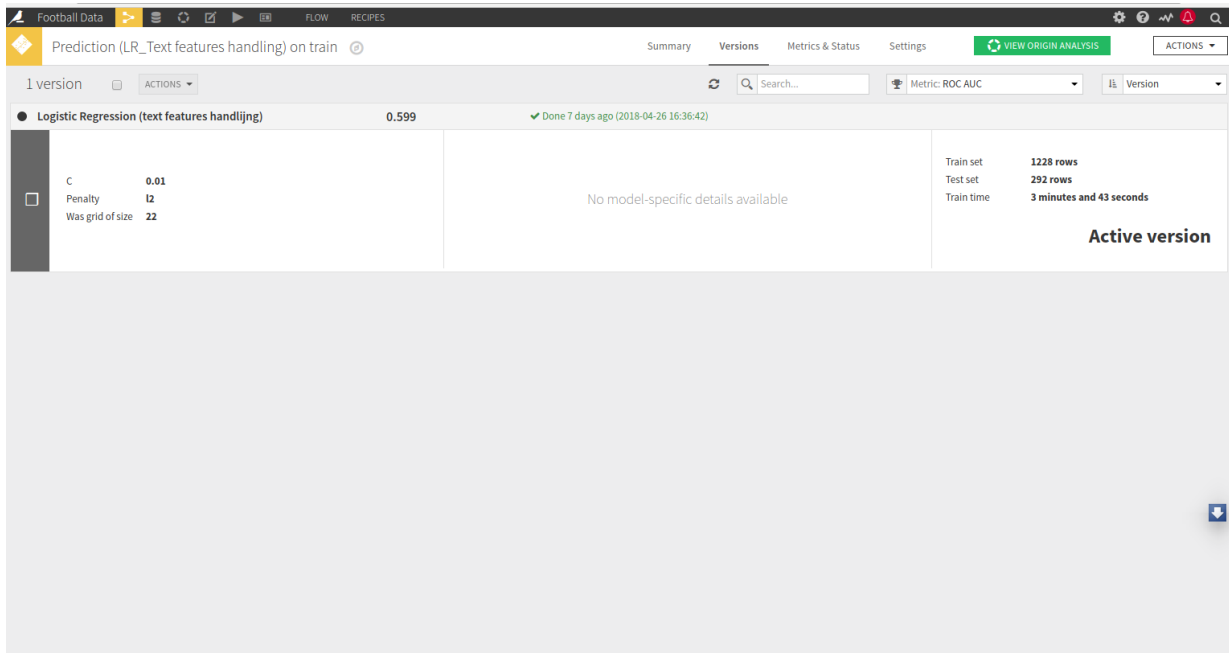
My private score was 0.40327 and my public score was 0.43606 which was the worst score I have achieved out of all 3 submissions.

**CONCLUSION**

This Football data challenge has helped me put into practice the different tools and knowledge that we have acquired about Data Science and Dataiku DSS during our course. It helped me analyze the results and to be able to tell which algorithm was the best to use and for what reasons. Unfortunately, I was not able to add addition features to further enhance and enrich my data as I had planned to do in the beginning of the project. This was due to my inability to find the appropriate data to work on. I wish to work in future projects where I will be able to add features on my own, in order to achieve better results, and more importantly to achieve the objectives which have been set. Even though the final score is not too far from the targeted score (~ 0.45 and 0.5 respectively), I would have been more satisfied to attain my initial goals.