

Explanation document:

Before I've started to analyze I had to organize all the results in an excel file (Attached).

First step:

Data Explanation:

1. Choosing classes- the target in this project is to calculate the probability the Warriors will win the game. The classes are divided in to two categories: win and loss. I've checked the labels quantity, so we can decide the score for the model– the data is balanced (figure 1), so I'll maximize accuracy.
2. Estimate the features - I noticed that every feature has a float value except the name of the team (feature game), so I dropped this feature.
Also, I decided to ignore the features 'FG Made-Attempted 1', 'FG Made-Attempted 2', '3PT Made-Attempted 1', '3PT Made-Attempted 2', 'FT Made-Attempted 1', 'FT Made-Attempted 2', since the features are pointing the total points for each team - which will tell us who is the winning team.
For each feature I did a statistic description and a box plot - I saw that the feature 'Flagrant Fouls 1' has no impact because this column is a zero column, and I've also succeeded to identify outliers.
The feature 'win' will be the target of the model, so I've separated it from the features.
3. Data Distribution- I performed tests of the various data distributions (for each feature) using graphs (KDE plots) that help us to understand the behavior of the data relative to each category. These distributions can help us understand which features will be most effective in distinguishing between the two classes (win and loss) - the features in which the greatest difference is between the distributions. In this case, the distributions are relatively similar. We also can see that there is
4. Correlation matrix- In this section I examined whether there is a correlation between the various features. A high correlation indicates that we should consider removing features, as they do not explain new information. After the test we found that the parameters are not correlated to each other, the correlations found to be less than 0.4.

second step:

building a decision tree:

I performed a decision tree in order to see what features the tree would choose to be meaningful for separating the win.

Third step:

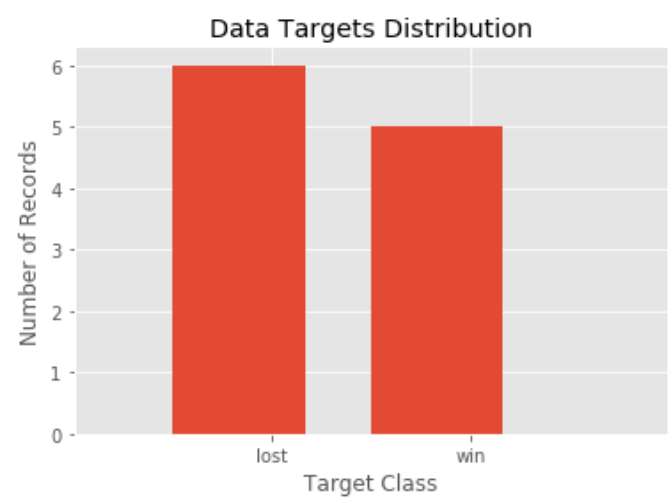
Model:

the model I decided is Random Forest. I think it's most suitable, because this model is using ensembles of trees, where each tree in the ensemble is grown in accordance with a random feature. Final predictions are obtained by aggregating over the ensemble. As the base constituents of the ensemble are tree-structured predictors, and since each of these trees is constructed using an injection of randomness. RF model have succeeded to predict the target in many domains, so with more data it can do an excellent work in the project domain.

For building the model we'll look at the important variables that contribute to the win. After we have a look at the important predictor variables, we need to check the win probability predictions for the final match.

I read that Artificial Neural Network is also suitable for this project, because I didn't learn about this model I didn't want to choose this model, but from what I read this model.

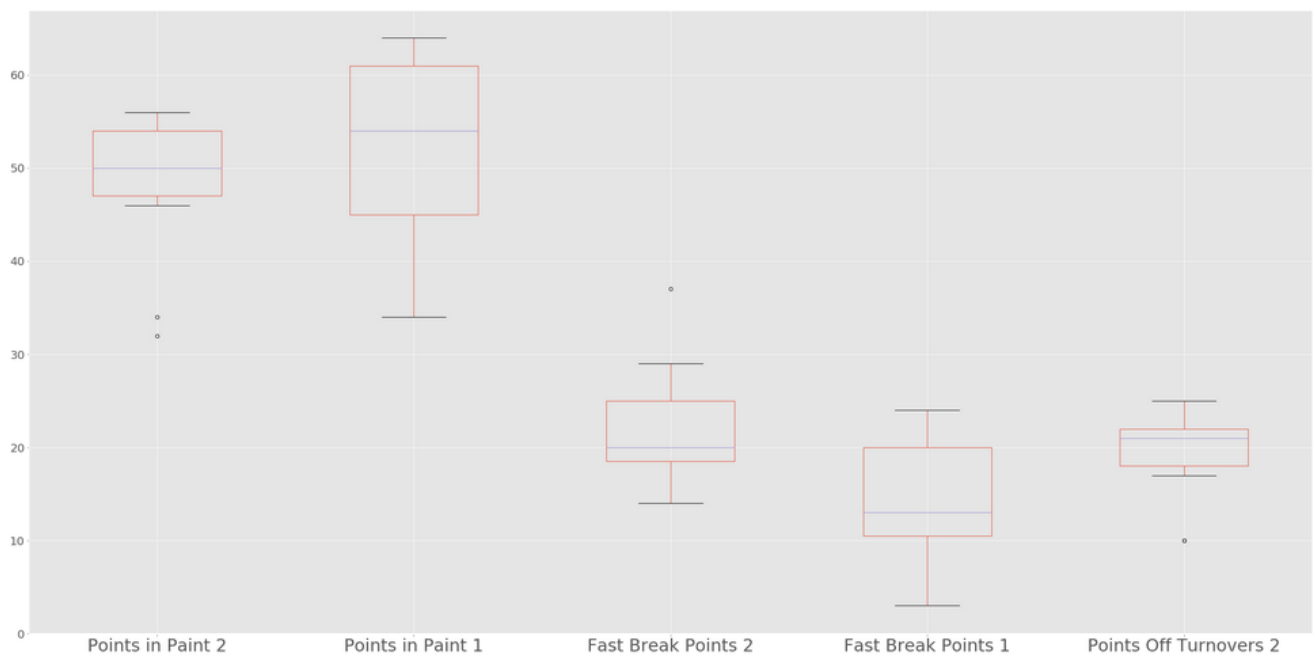
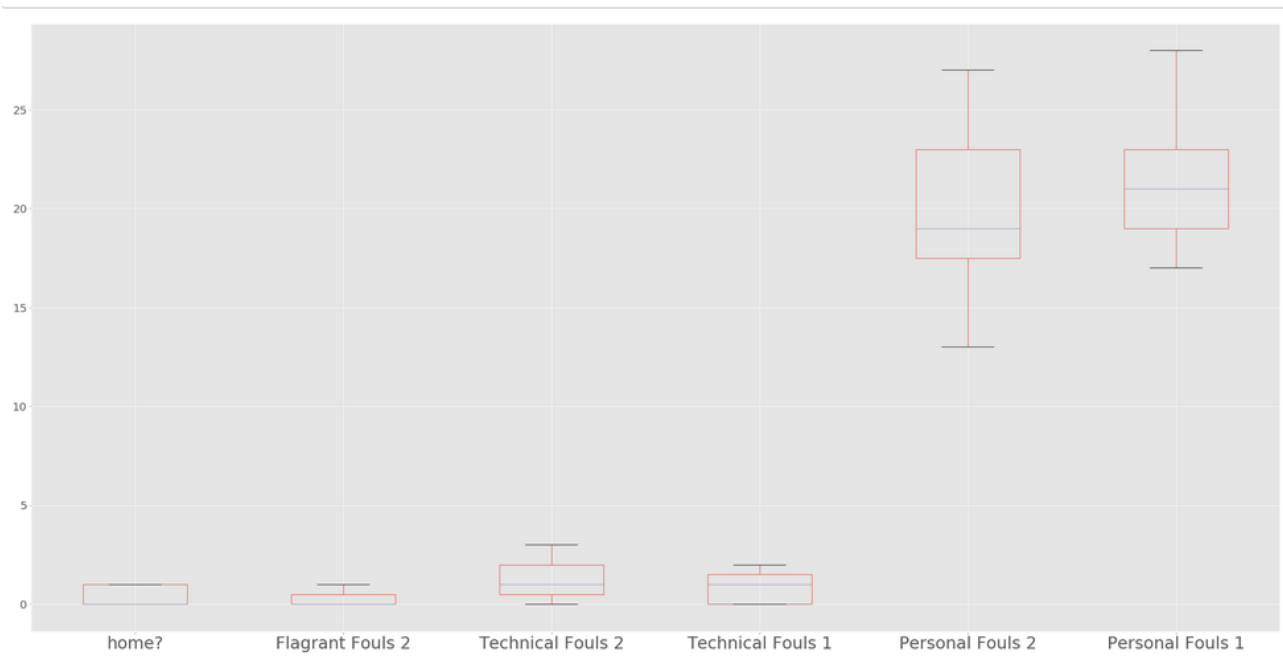
attachments photos:

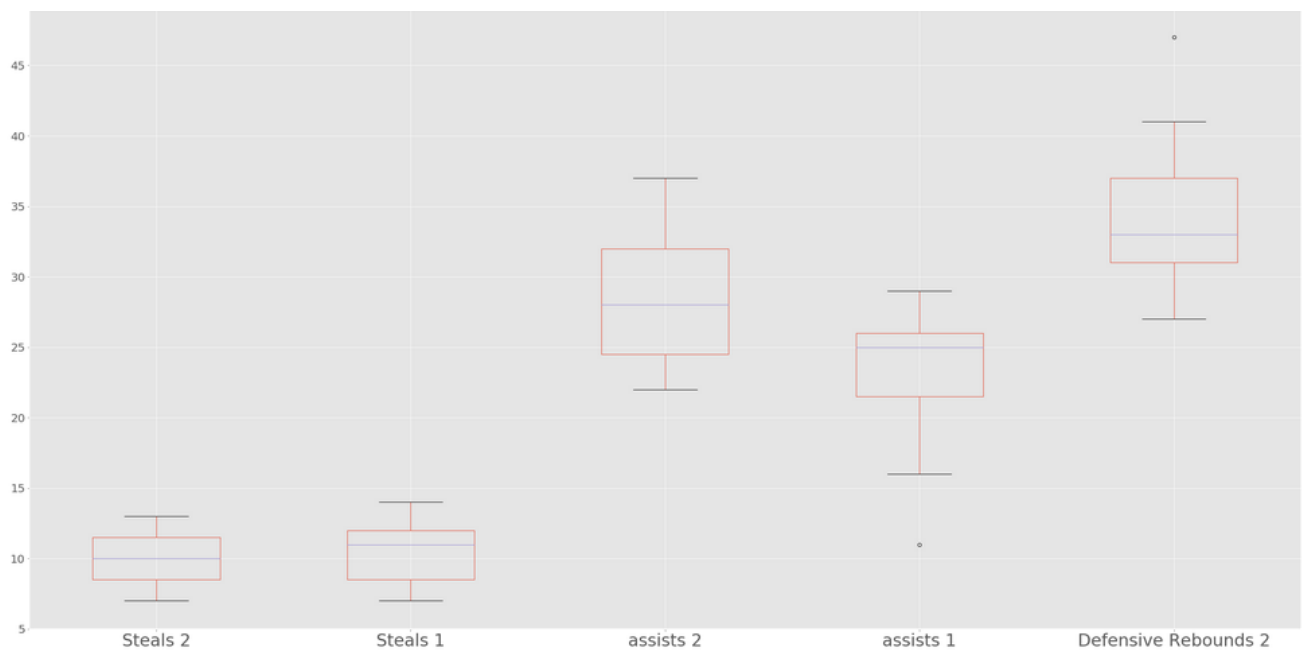
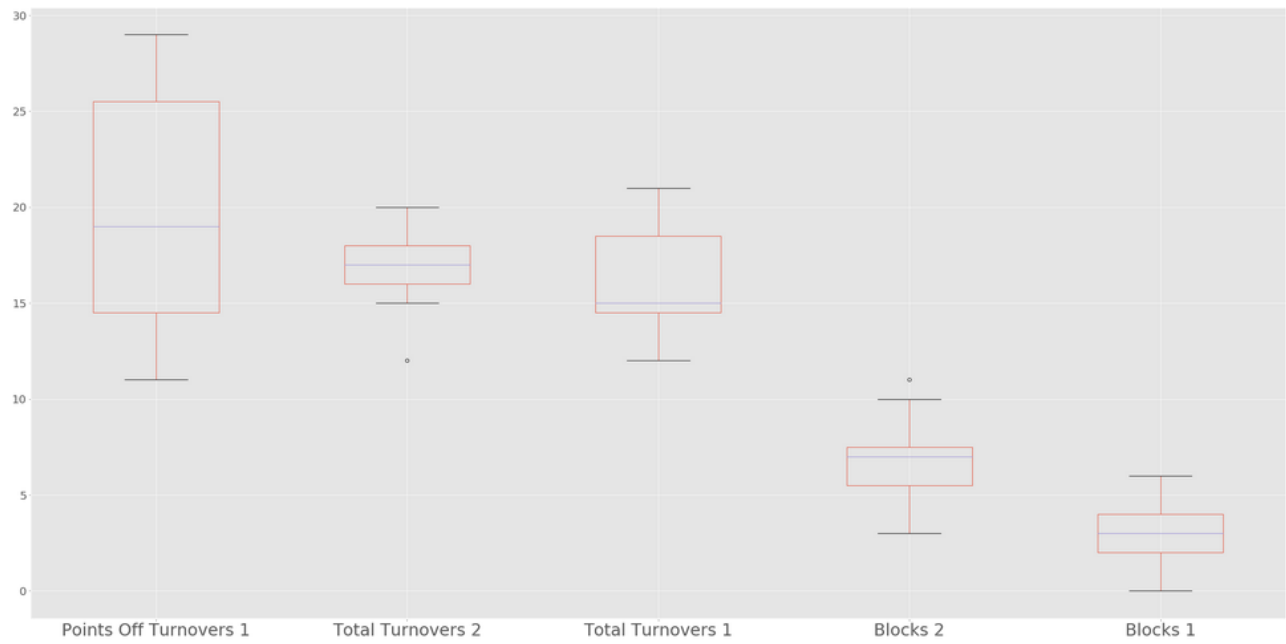


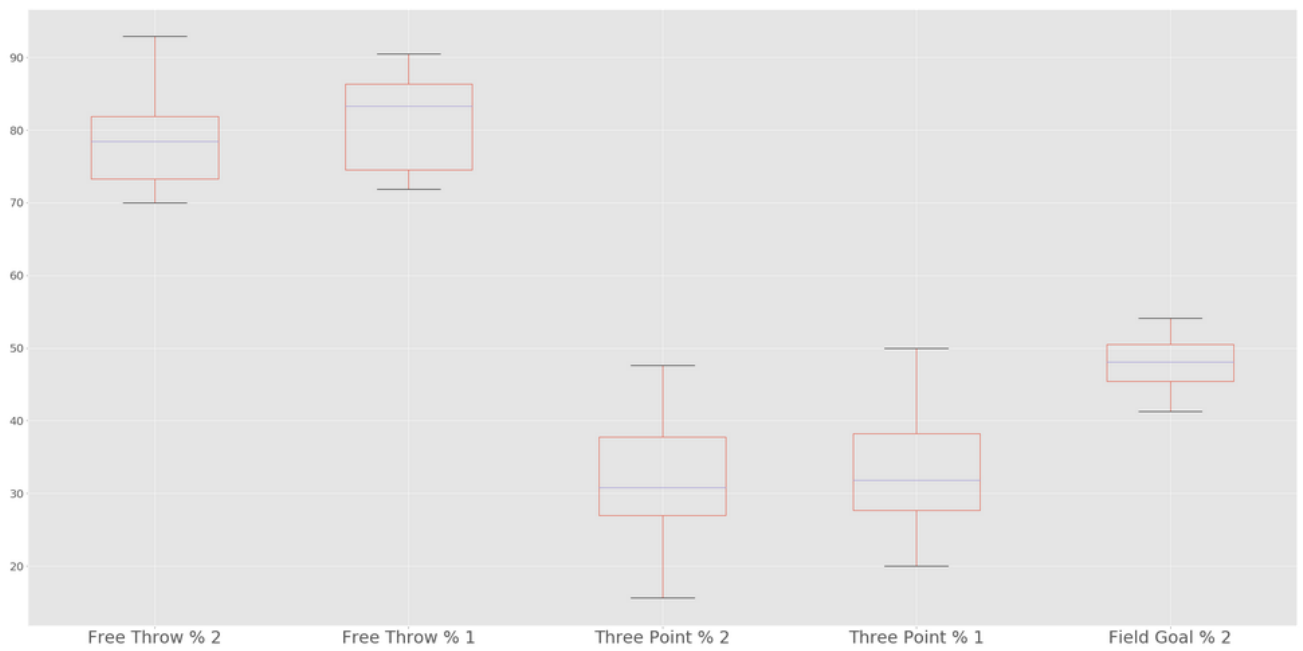
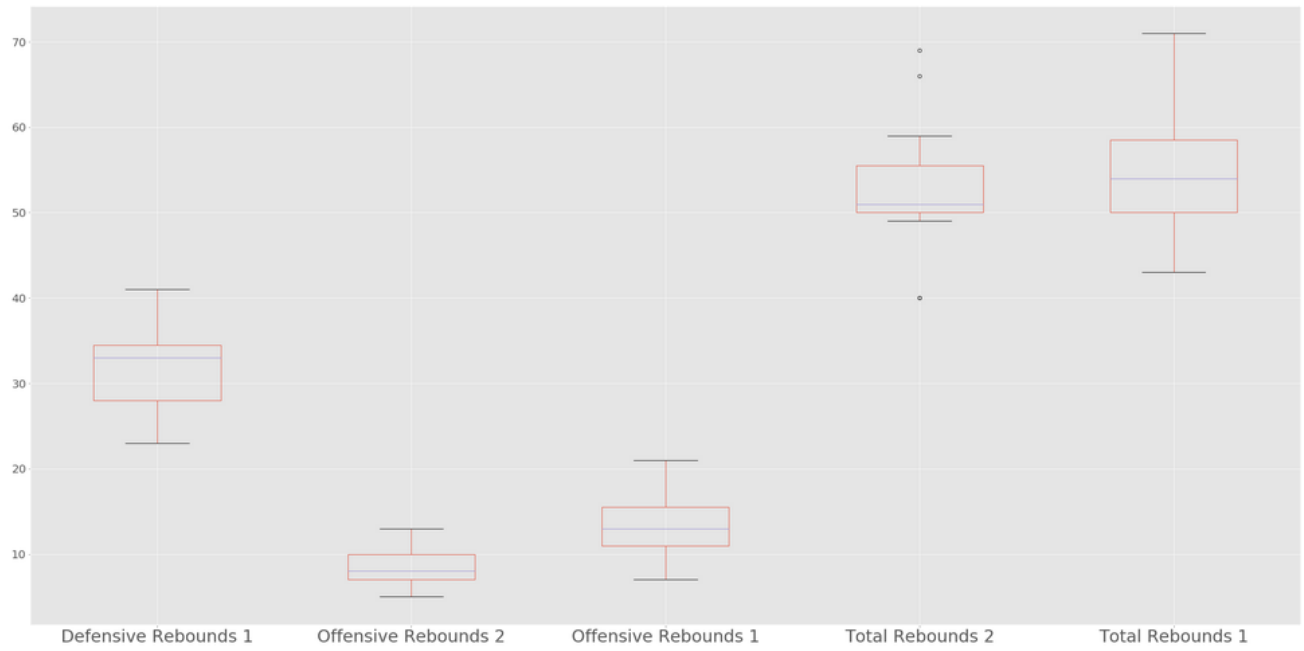
We can declare that the classes are well-balanced.

Description and Box Plot:

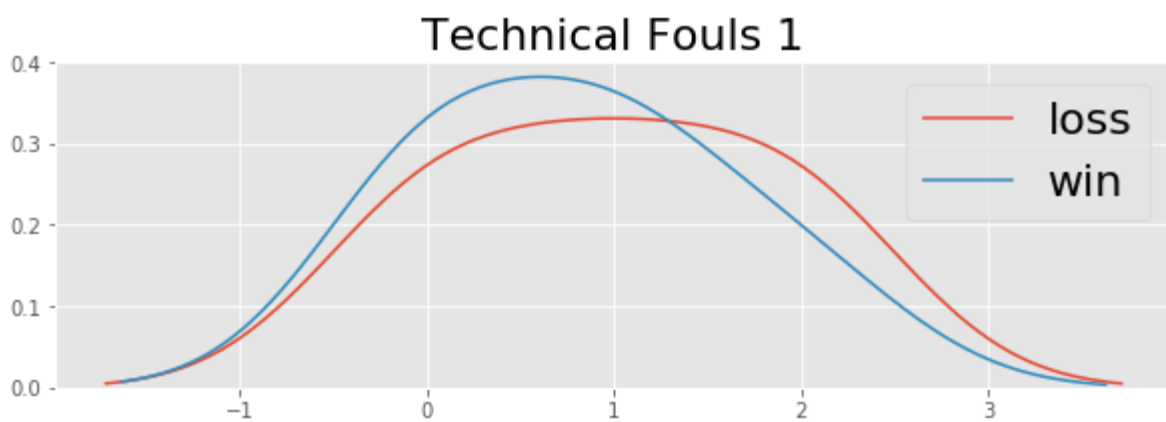
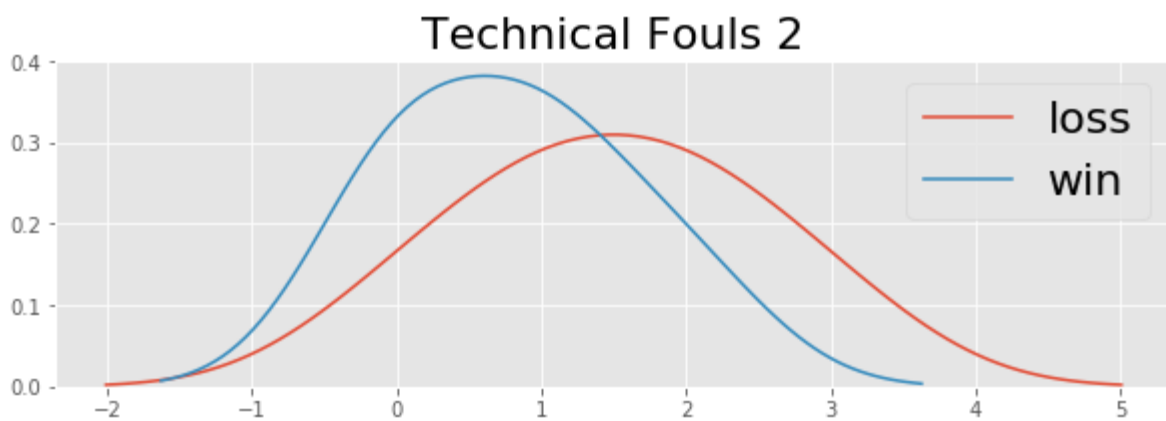
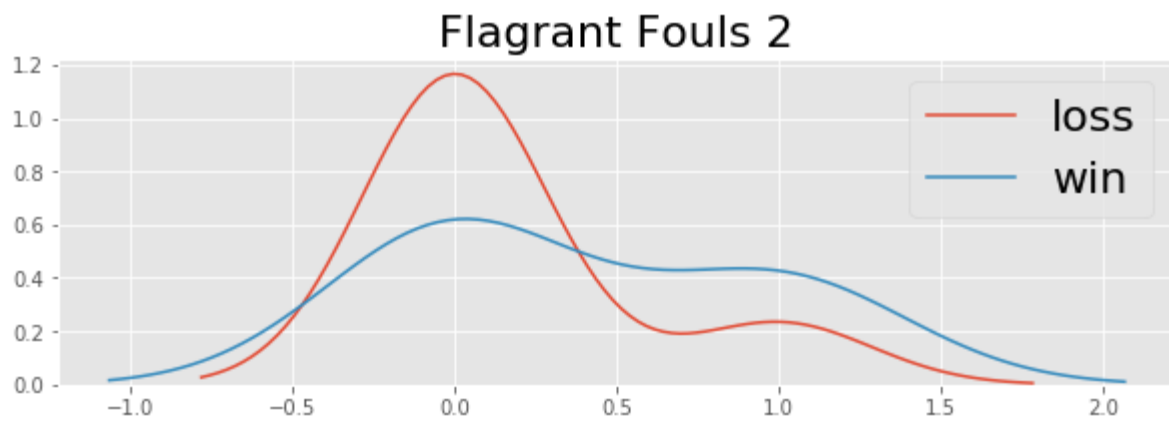
	home?	win?	Flagrant Fouls 2	Flagrant Fouls 1	Technical Fouls 2	Technical Fouls 1	Personal Fouls 2	Personal Fouls 1	Points in Paint 2	Points in Paint 1	Fast Break Points 2	Fast Break Points 1	Points Off Turnovers 2	Points Off Turnovers 1
count	11.000000	11.000000	11.000000	11.0	11.000000	11.000000	11.000000	11.000000	11.000000	11.000000	11.000000	11.000000	11.000000	11.000000
mean	0.454545	0.454545	0.272727	0.0	1.181818	0.909091	20.000000	21.545455	47.818182	52.181818	22.181818	14.363636	19.181818	19.181818
std	0.522233	0.522233	0.467099	0.0	0.981650	0.831209	4.219005	3.777926	7.972681	10.215852	6.554665	7.201010	4.996362	4.996362
min	0.000000	0.000000	0.000000	0.0	0.000000	0.000000	13.000000	17.000000	32.000000	34.000000	14.000000	3.000000	10.000000	10.000000
25%	0.000000	0.000000	0.000000	0.0	0.500000	0.000000	17.500000	19.000000	47.000000	45.000000	18.500000	10.500000	18.000000	18.000000
50%	0.000000	0.000000	0.000000	0.0	1.000000	1.000000	19.000000	21.000000	50.000000	54.000000	20.000000	13.000000	21.000000	21.000000
75%	1.000000	1.000000	0.500000	0.0	2.000000	1.500000	23.000000	23.000000	54.000000	61.000000	25.000000	20.000000	22.000000	22.000000
max	1.000000	1.000000	1.000000	0.0	3.000000	2.000000	27.000000	28.000000	56.000000	64.000000	37.000000	24.000000	25.000000	25.000000



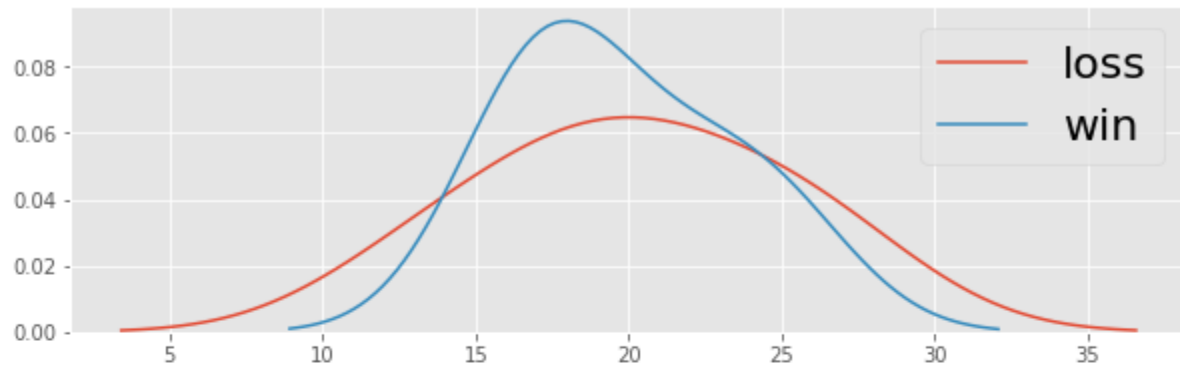




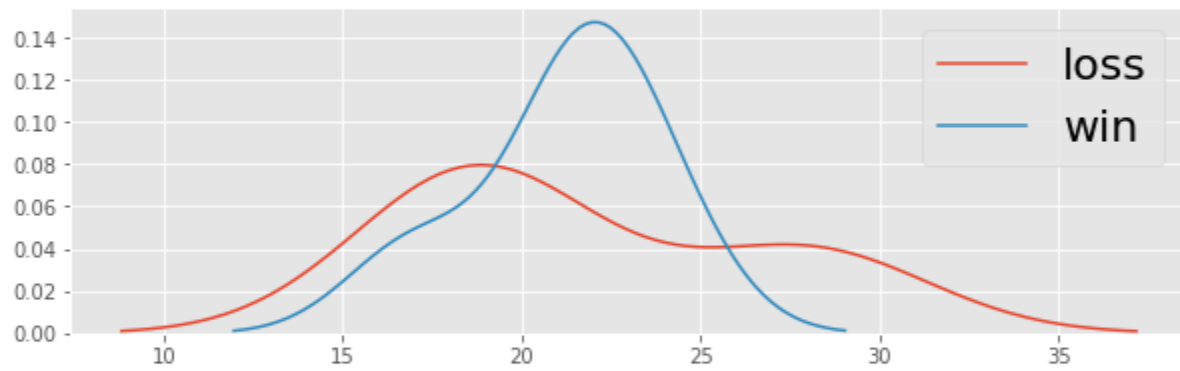
Distribution of each feature within the classes



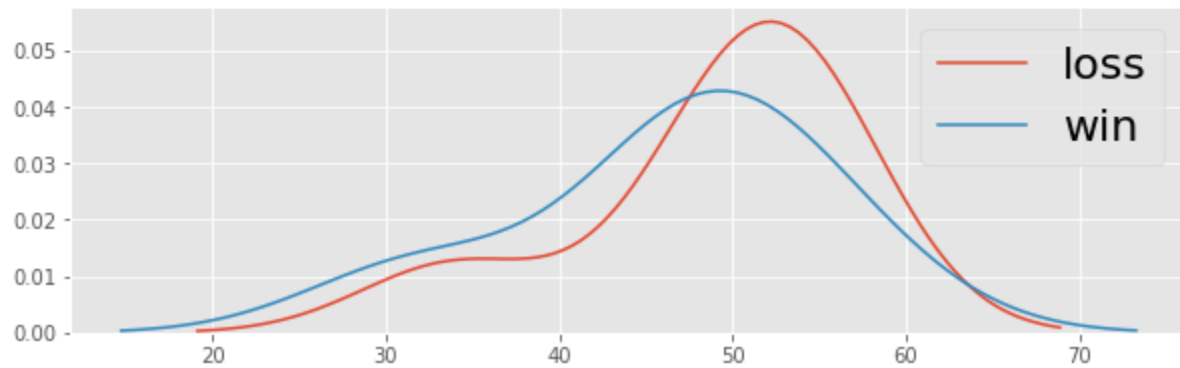
Personal Fouls 2



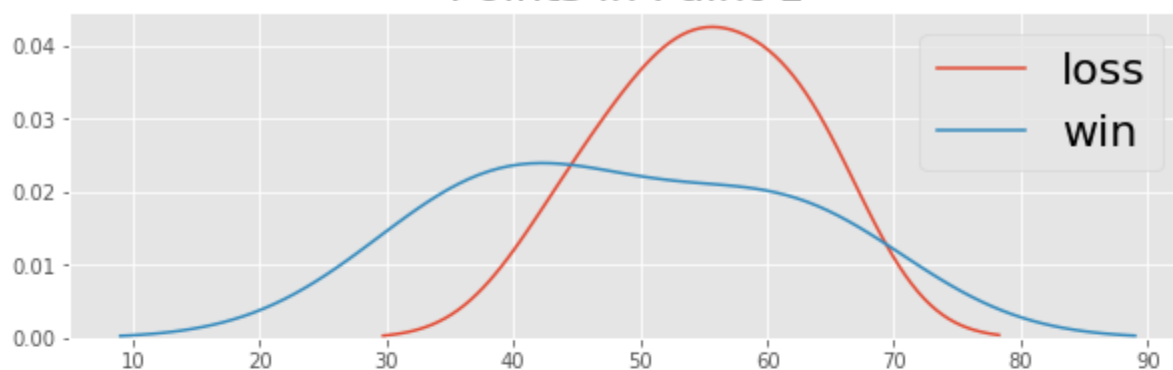
Personal Fouls 1



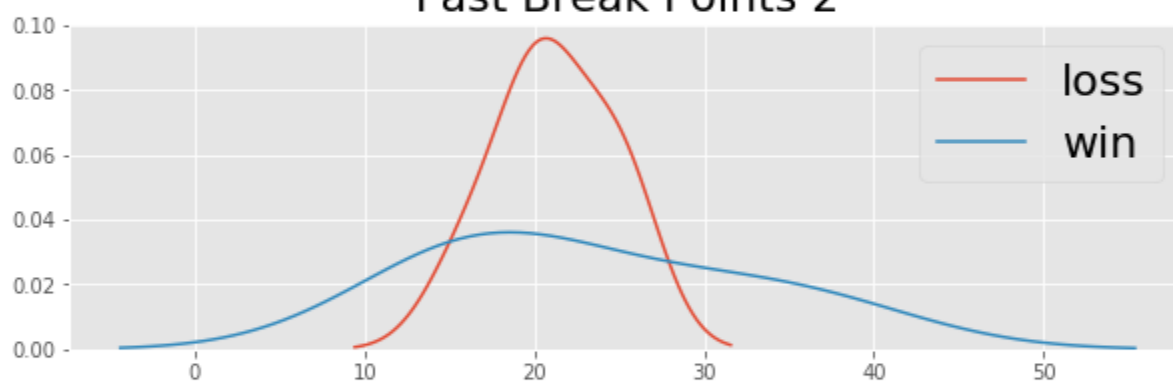
Points in Paint 2



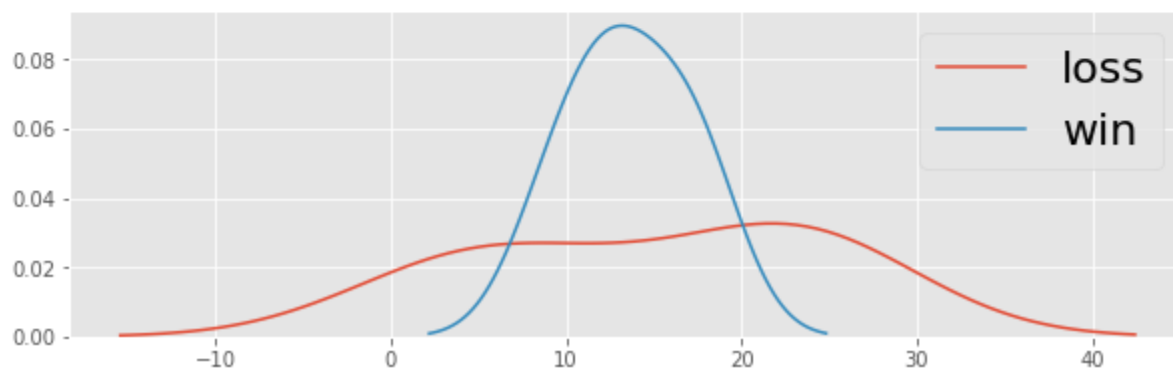
Points in Paint 1



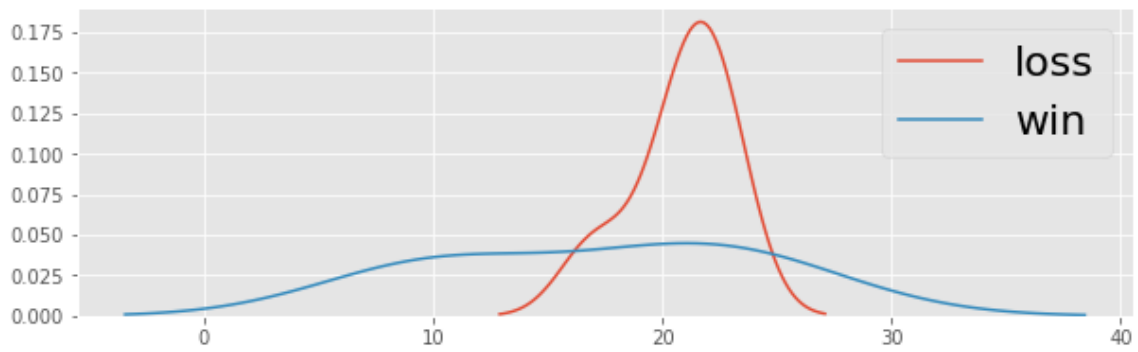
Fast Break Points 2



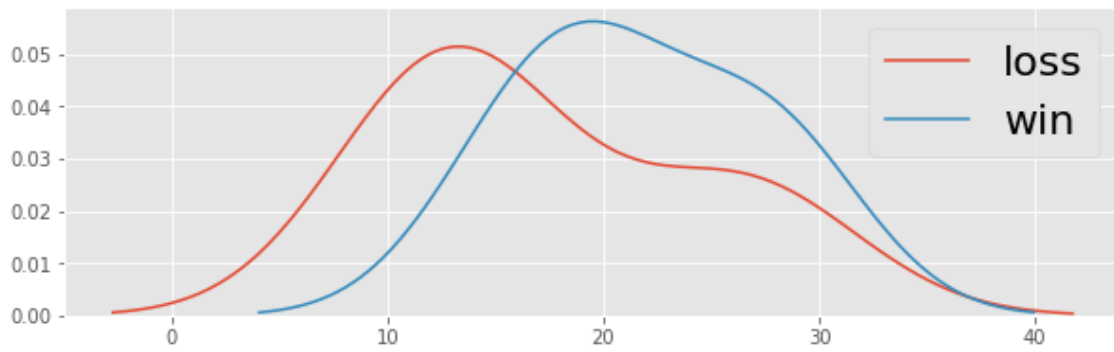
Fast Break Points 1



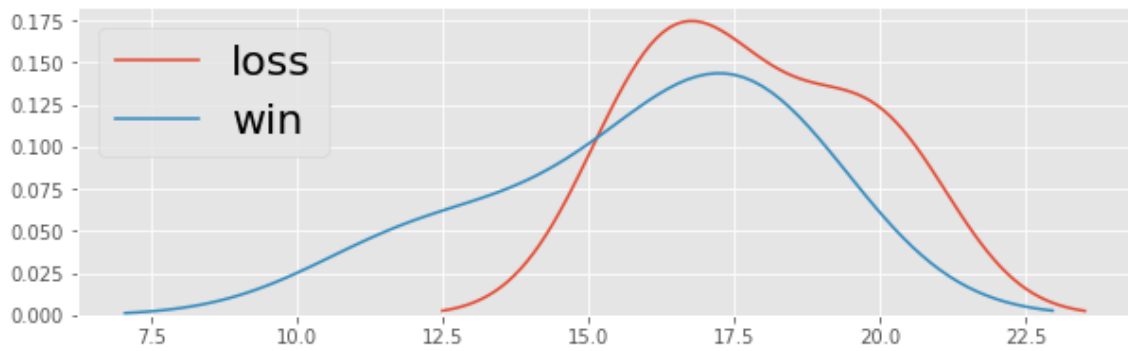
Points Off Turnovers 2



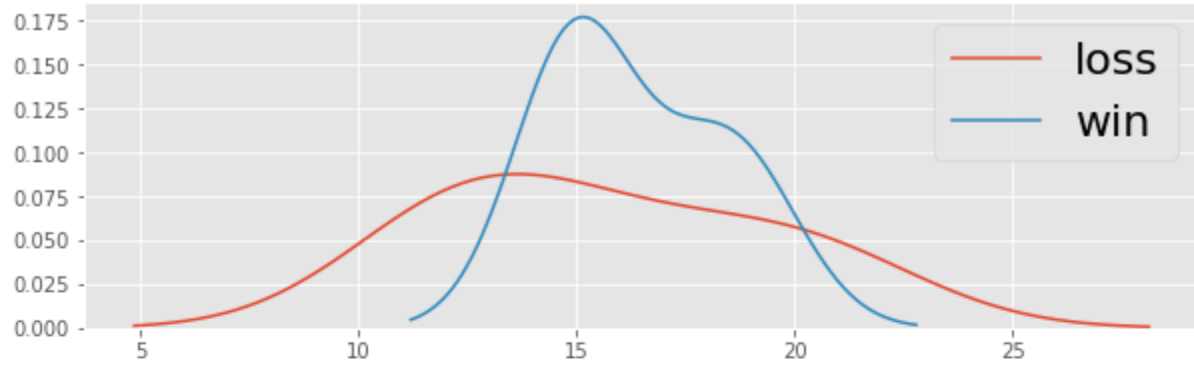
Points Off Turnovers 1



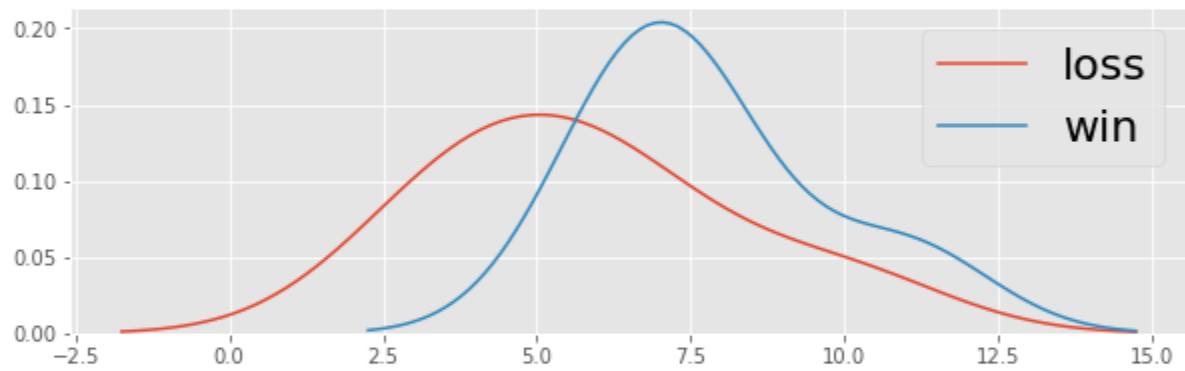
Total Turnovers 2



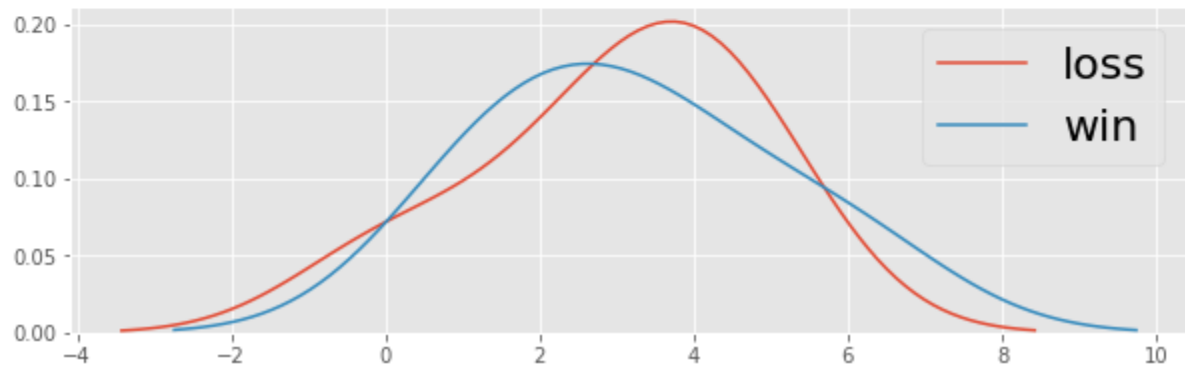
Total Turnovers 1



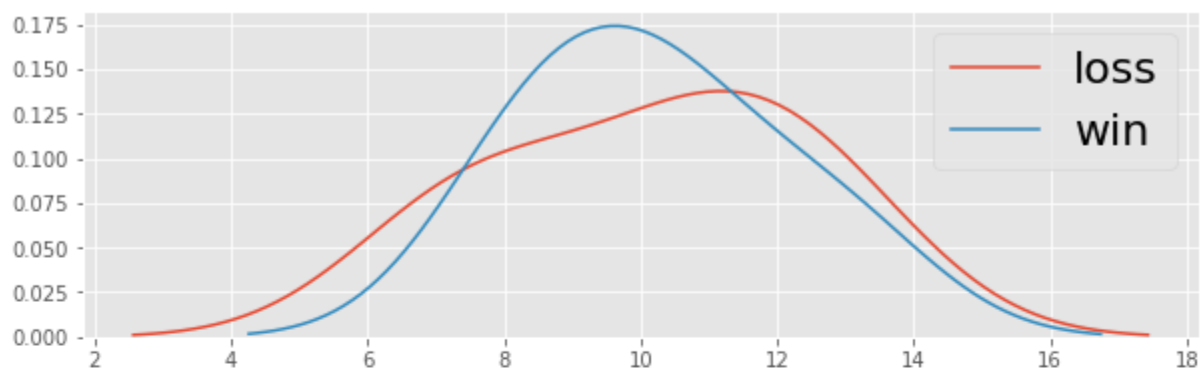
Blocks 2



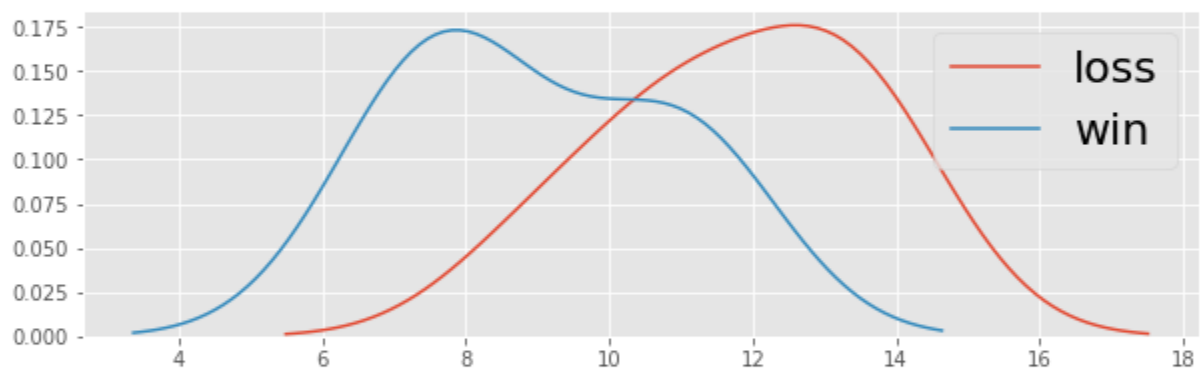
Blocks 1



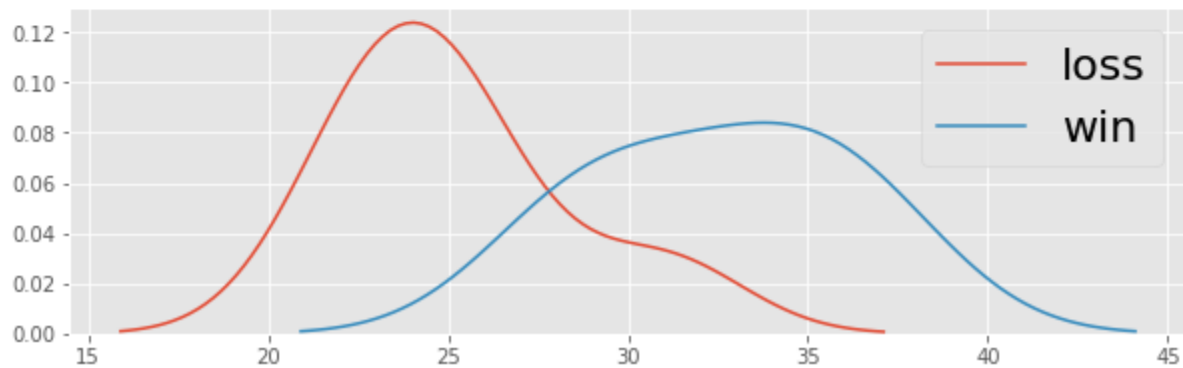
Steals 2



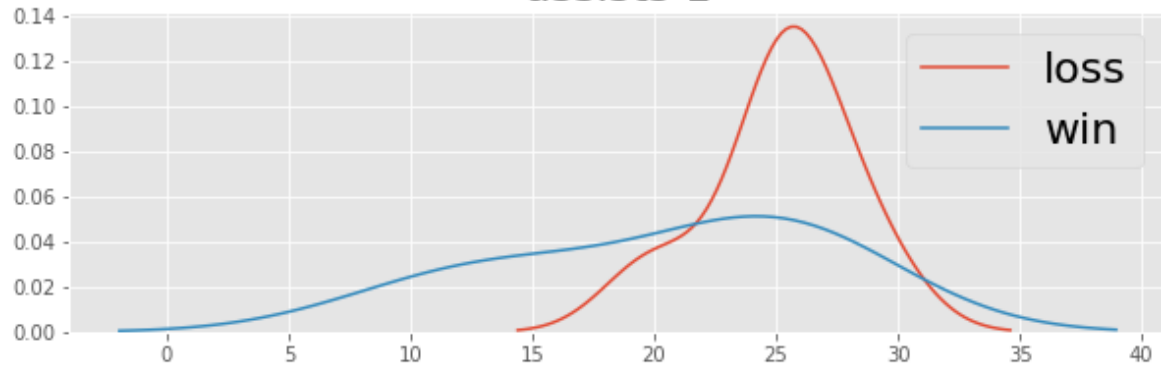
Steals 1



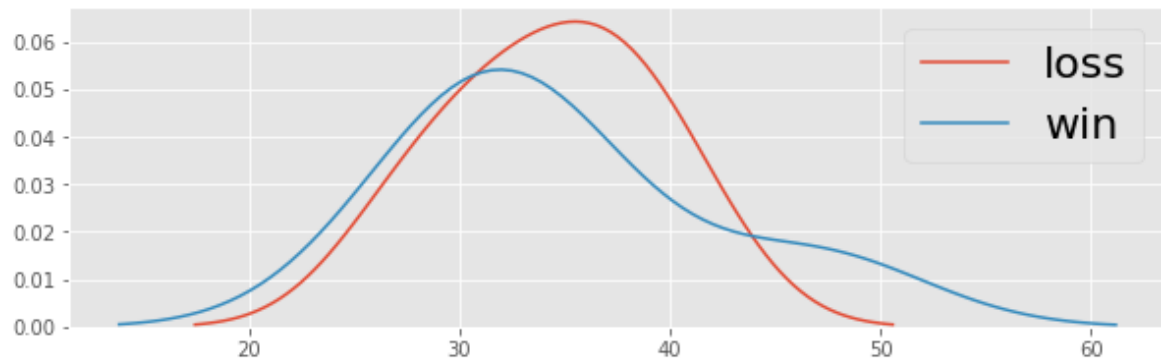
assists 2



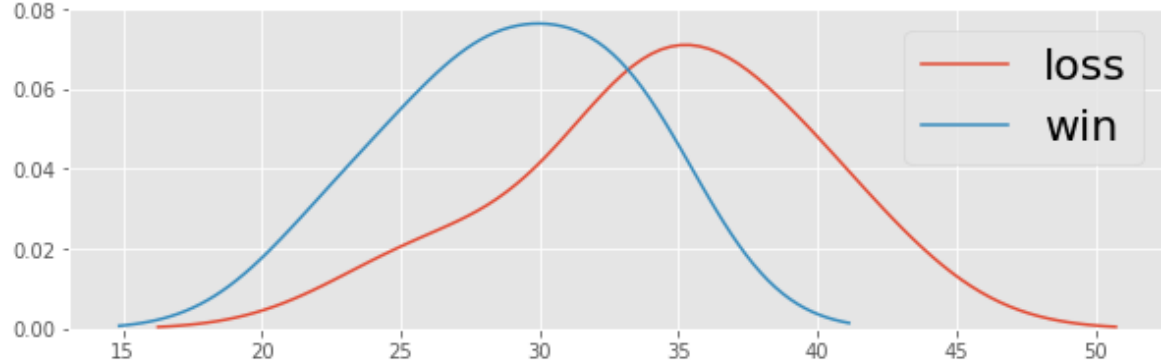
assists 1



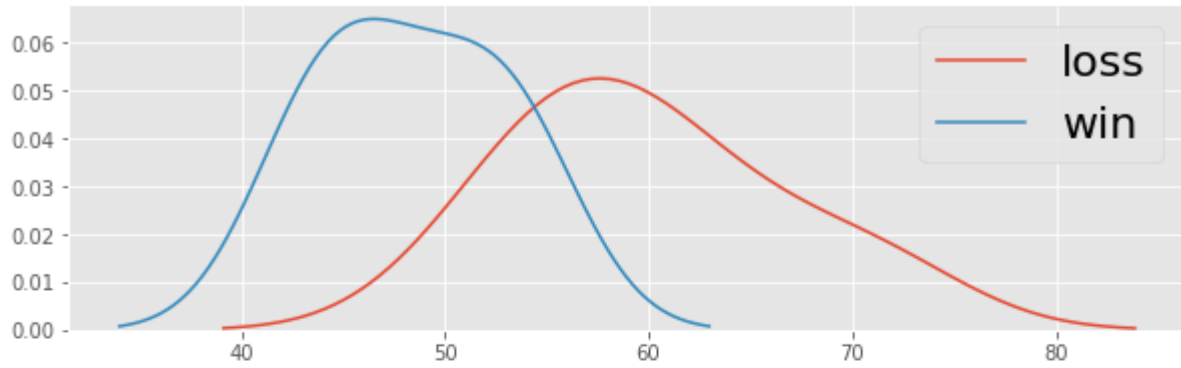
Defensive Rebounds 2



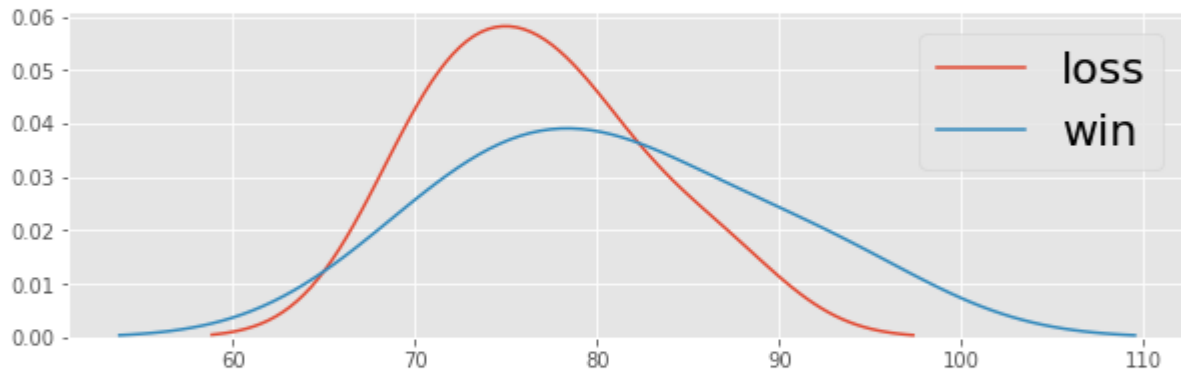
Defensive Rebounds 1



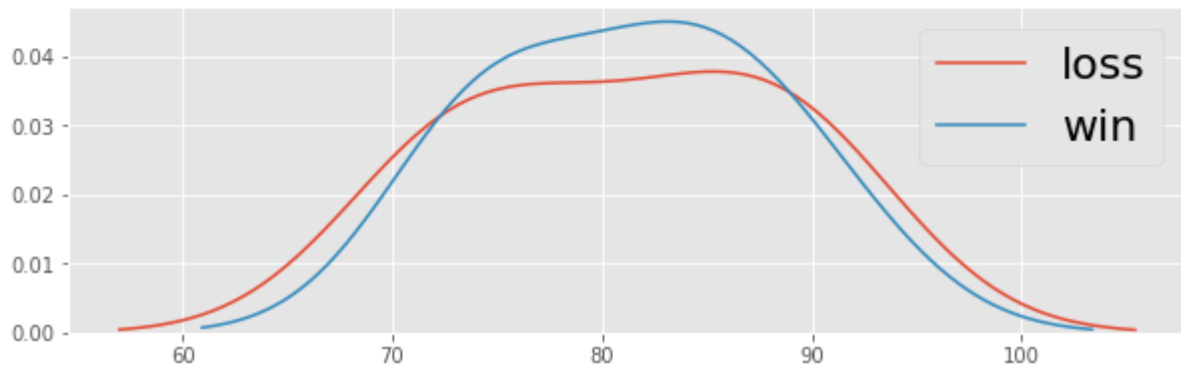
Total Rebounds 1



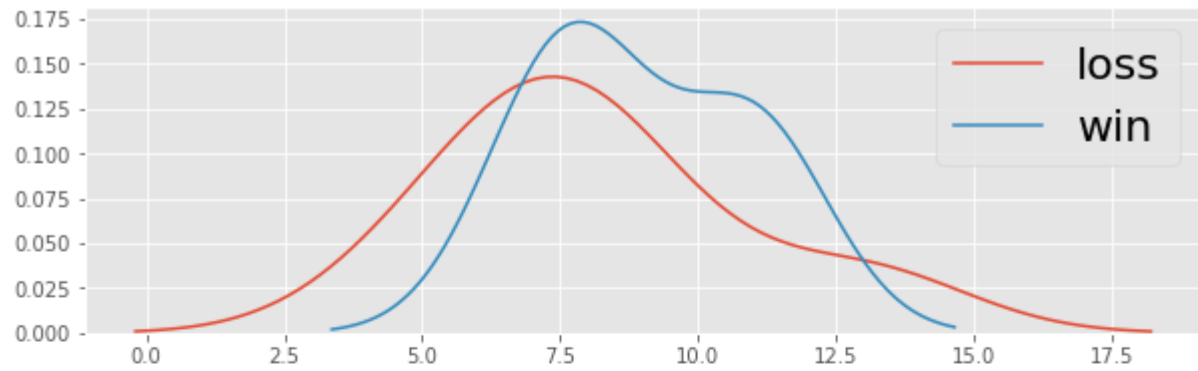
Free Throw % 2



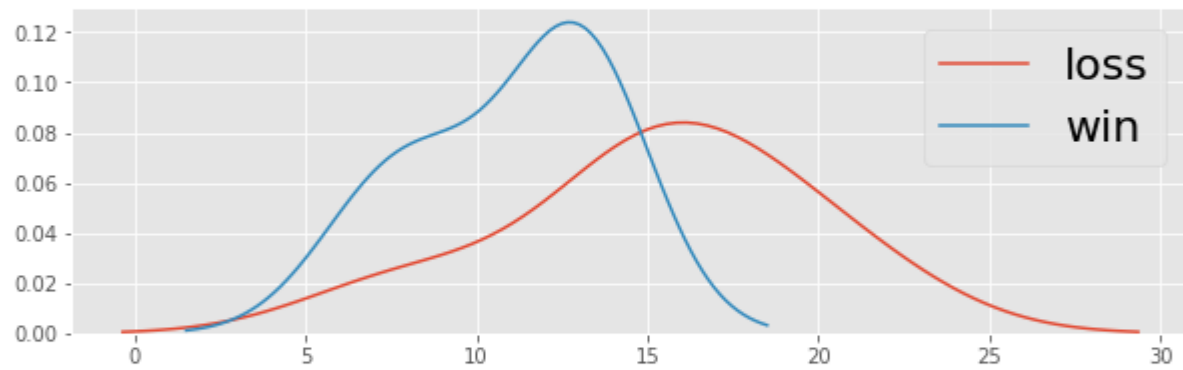
Free Throw % 1



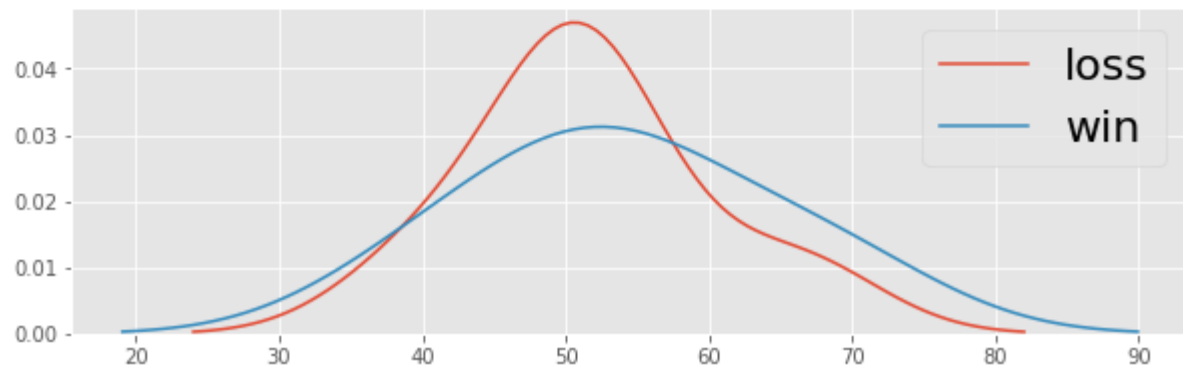
Offensive Rebounds 2



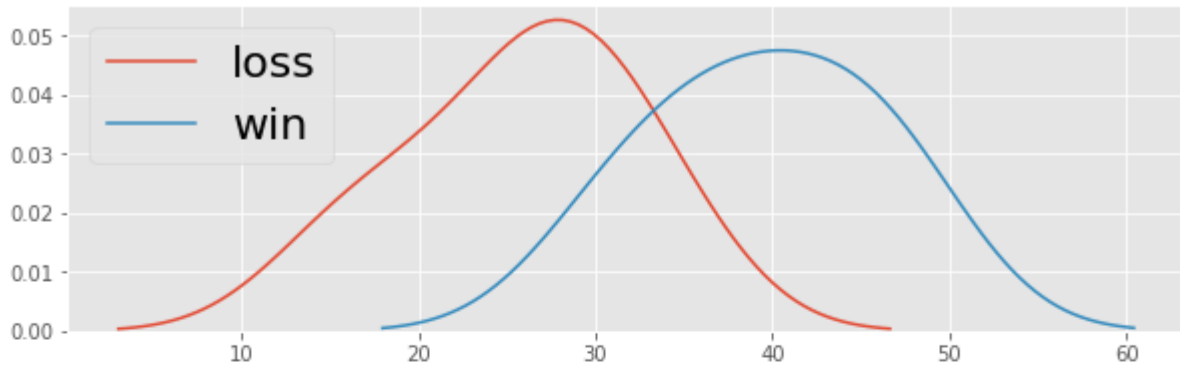
Offensive Rebounds 1



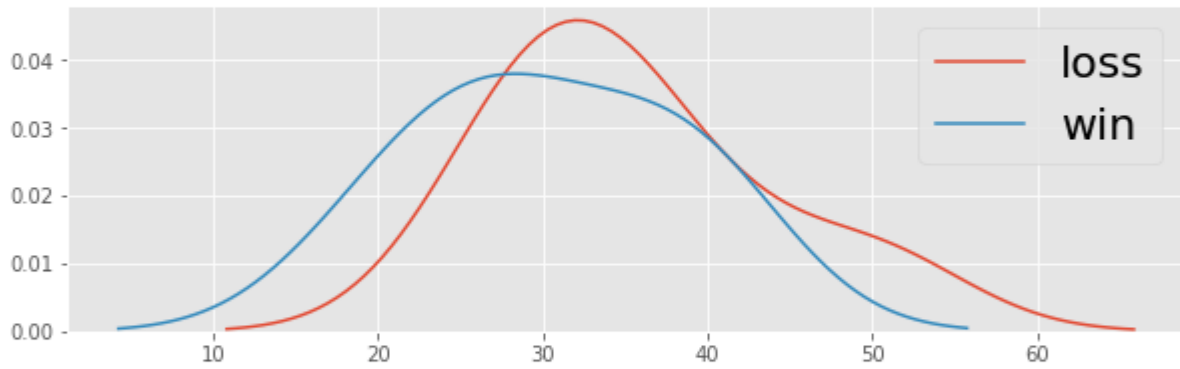
Total Rebounds 2



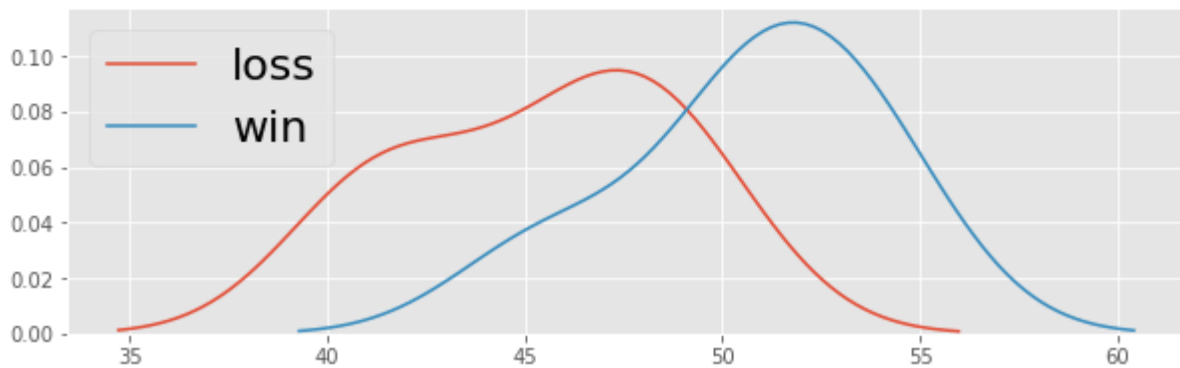
Three Point % 2



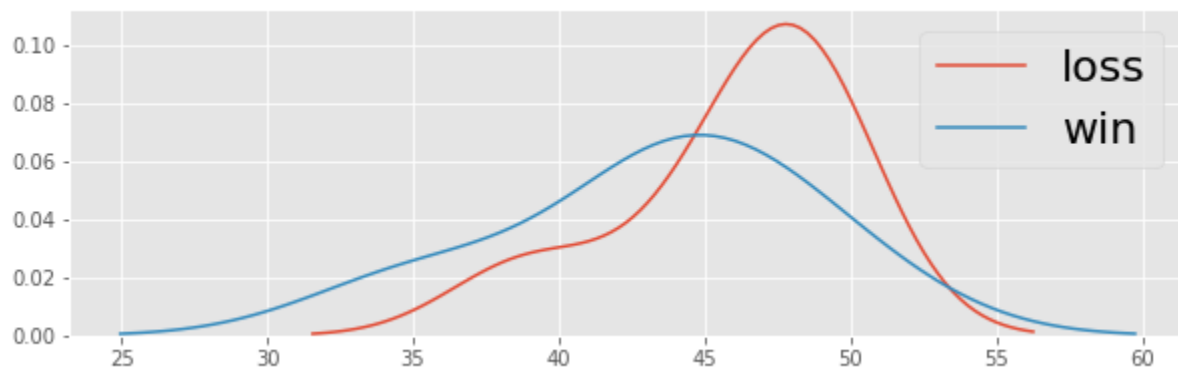
Three Point % 1

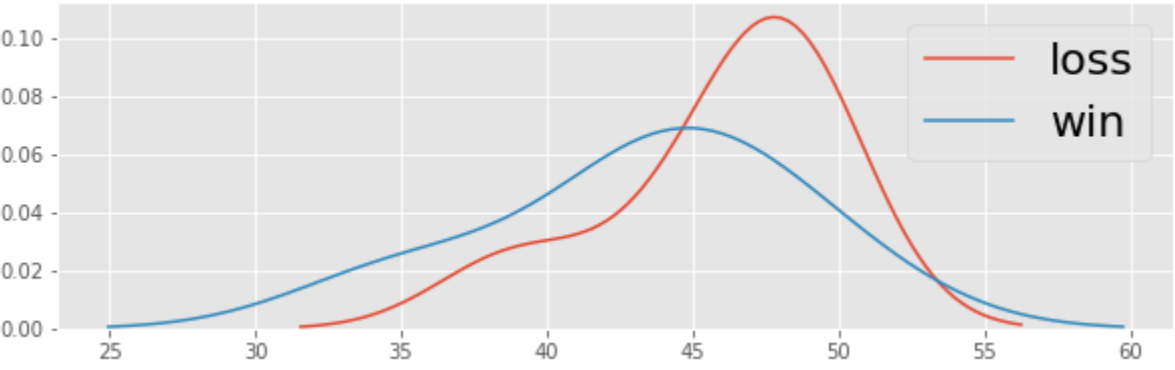


Field Goal % 2

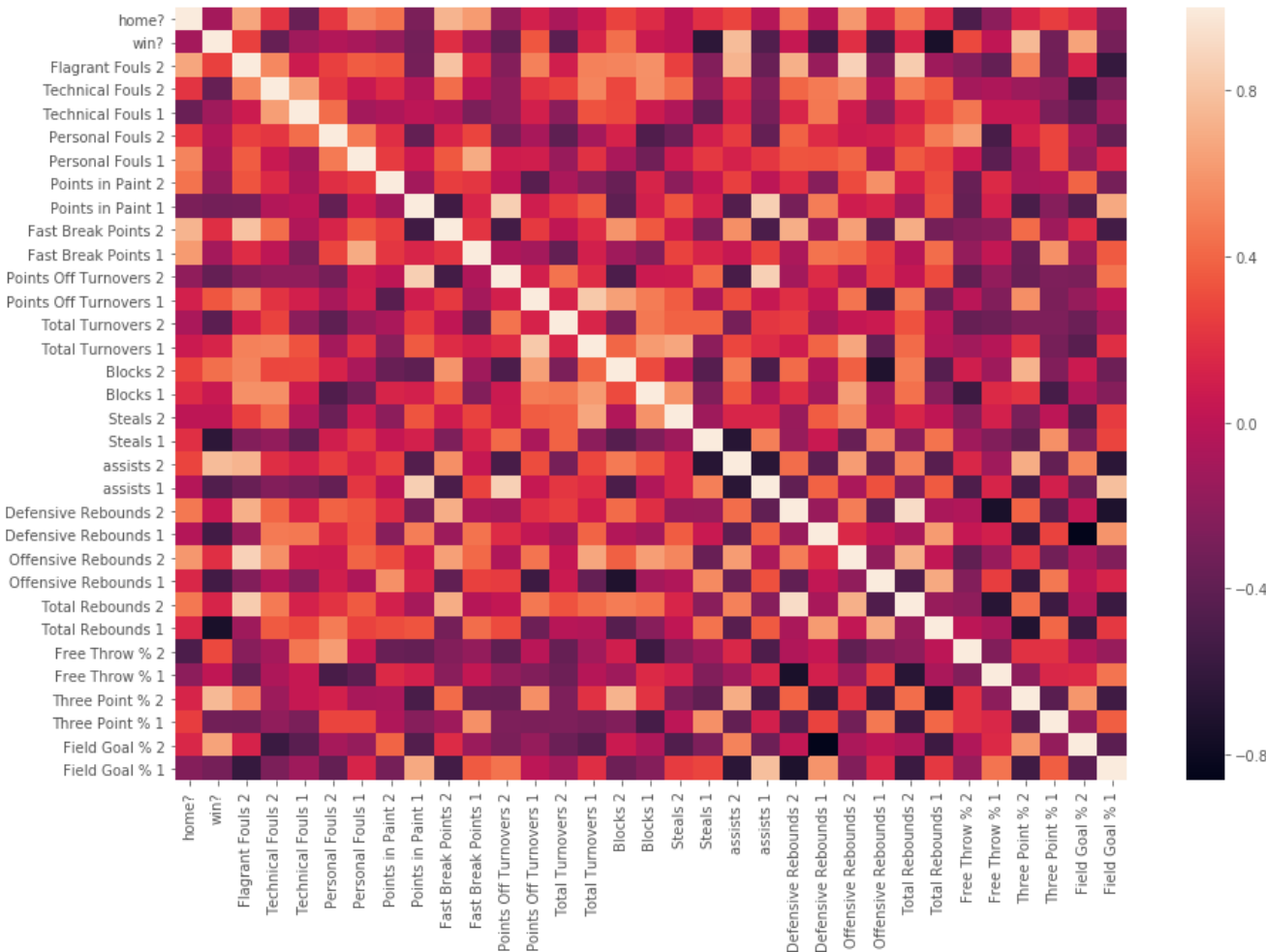


Field Goal % 1





Correlation Matrix:



Feature importance - Deseccion Tree

