

DSA4211 Group 13 Project Assignment

Predicting the winners of UFC fights

Ragnhild Laursen & Jeppe Sjørup & Tay Zhi Yao Joel & SONG Xiaohan

October 2019

Abstract

In this report, we built models mainly with two methods, one was logistic regression and the other one was Neural Network, to predict the winner in 1078 fighting competitions from Ultimate Fighting Championship(UFC) with 159 feature data for each fighting prediction.

Before the analysis, we split the data into training and testing set, and the train-test split was 7:3. For feature selection, we only kept the non-constant features and removed those features which were highly correlated, so the number of features that joint into the analysis was reduced from 159 to 135.

For Logistic Regression, we firstly did ordinary logistic regression without regularization, the model accuracy on the test set was 68.1%. Then we added L1 regularization into the model and applied cross-validation to choose the best λ . The model accuracy was therefore increased slightly to 68.6%, and the selected λ was 0.017. Furthermore, since LASSO regression was used, many predictors' coefficients were shrunk to 0, and the dimension of the features was therefore decreased to 26.

For Neural Network, we attempt to implement a feed-forward neural network on the 135 features and output the probabilities for each observation being 1 of the K classes (Here $K = 2$). We configured our model with 5 layers in all, one input layer, one output layer and three hidden layers, where we used Exponential linear units (ELU), Sigmoid and Tanh function as activation functions respectively. For model training, we used a validation split of 20%, which means we use 80% of the training data for model training and the else for validation. Furthermore, we use Adam optimizer and categorical loss entropy as our loss function and applied an early stopping callback function with the patience of 5 epochs to prevent overfitting. Lastly, we set our epoch number as 40, and batch size as 128. After implementing the Neural Network model on the testing set, we finally got a test accuracy of 68.9%.

In the end, we discussed the respective pros and cons of these two models. From the perspective of prediction accuracy, the Neural Network looks better, however, sometimes, it might be justifiable to forgo the benefits of accuracy that the neural network has to offer in return for the interpretability of the logistic regression model depending on the use case.

Key words: *Logistic Regression LASSO Neural Network Early Stopping*

1 Introduction

In this report we look at predicting winners of fights from Ultimate Fighting Championship (UFC) based on data available prior to the start of the fight. In our investigation we use data hosted on [kaggle.com](https://www.kaggle.com) and compiled by Rajeev Warriar [1]. The data contains information about 3592 UFC fights conducted from 1993 to 2019.

UFC is an professional American mixed martial arts organization that was founded in 1993. They host events throughout the world, where fighters from different combat sports compete in one versus one fights.[2] The fights are split up in different weight classes from straw weight (105-115 lbs) to super heavyweight (265+ lbs), with each weight class having a reigning champion. Each fight has a pre-determined maximum number of rounds, which is normally 3 or 5 with a duration of 5 minutes per round. The fights can end prematurely either by a knockout (KO), submission, or a technical knockout (TKO), in which case the fighter that gets knocked out would be declared the loser and the other fighter would be declared the winner. If no stoppage occurs the winner is selected by judges using a points system after the fight has ended.[3] The two fighters are distinguished by the color of their shorts. Each fight has red fighter and a blue fighter. Traditionally the fighter given the red shorts is considered the favorite before the fight, therefore historically the win percentage of the red fighter is 66.26% as derived from the dataset.

2 Dataset

2.1 Data Structure

The data consist of 159 features containing various statistics known before the fight starts. The data has been processed so that each feature only contains information known before the fight takes place. Out of the 159 features, there is 142 features describing the two fighters, i.e. 71 features for each fighter. This contains their fighting stance, height, weight, age, matches won, average strikes and take downs among many other descriptive statistics obtained from previous fights. The remaining 17 features describes the format of the fight, including the weight class, gender and the maximum number of rounds. All of the feature are quantitative as they are either dummy, discrete or continuous variables. There is only one feature, which is a factor variable of two factors indicating whether the fight was a title or not, which is changed to a dummy variable.

The response is also a factor variable indicating whether the *Blue* or *Red* fighter won. The response is changed into a dummy variable, where 1 indicates red fighter and 0 indicates the blue fighter.

A summary over the response and the four other features for the whole data are depicted in table 1. The feature starts with R if it describes the red fighter and B if it describes the blue fighter. In detail, `R_age` is the age of the red fighter, `R_losses` is the total number of losses through the career of the red fighter, `R_avg_opp_SIG_STR_landed` is the average significant strikes landed from previous opponents of the red fighter and `B_Height_cms` is the height in cm of the blue fighter.

	Winner	R_age	R_losses	R_avg_opp_SIG_STR_landed	B_Height_cms
min.	0	19	0	0	152.4
median	1.0000	30	2	25	180.3
mean	0.6626	29.93	2.352	27.09	179.4
max.	1	46	14	151	210.8

Table 1: A summary over the response and four features from the data matrix.

2.2 Data Processing

Before analysing the data, first, we randomly assign 70% of the observations to a training set and treat the rest as our test set, which we will use to validate our models. To keep our test set independent from our training set, we only process our data using estimates based on the training set.

Afterwards, we remove the constant columns and highly correlated columns, i.e. a correlation larger than 0.95, only based on the columns from the training set. The same corresponding columns are removed from the test set. The processing above removed 23 features from our data, which leaves us with 135 features in our data.

Last, we scale the data based on the scales found from the training data. This includes finding the mean and standard error of each feature in the training data and then scaling each observation in the data using these parameters.

When referring to the data in the analysis below, we are always working with processed data i.e. correlated and constant features removed and the data is scaled. This also includes the split between training and test set, which remain the same through the report.

3 Methods and Results

The goal of our analysis is twofold. Primarily we want predict the winner of fights. This implies that we should perform binary classification of our data. The baseline of this classification is predicting the favorite (red) as the winner. Secondly we want to identify important features related to winning fights.

In our investigation of the data we have chosen logistic regression and neural network. For logistic regression we first fit a model based on all the features. Afterwards we try to fit a new logistic regression model, where we regularize our parameters using lasso. As we use lasso, the remaining coefficients also represent the most important features. The logistic regression has a linear decision boundary, so afterwards we try to predict neural network which has a highly non-linear decision boundary to see if this gives a better prediction.

3.1 Logistic Regression

Logistic regression is a soft classifier that models the probabilities of a given data point belonging to each of K classes. In the case of the UFC data set we want to model whether the winner is *Blue* or *Red* i.e. 0 or 1, which means $K = 2$. Letting $Y \in \{0, 1\}$ denote the

class and $x \in \mathbb{R}^d$ a given data point of features, then we model the probabilities by

$$P(Y = 1(Red)|X = x) = \frac{\exp(\beta_0 + \beta_1^T x)}{1 + \exp(\beta_0 + \beta_1^T x)}$$

$$P(Y = 0(Blue)|X = x) = \frac{1}{1 + \exp(\beta_0 + \beta_1^T x)}$$

where $\beta_0 \in \mathbb{R}$, and $\beta_1 \in \mathbb{R}^d$ are the parameters that determine the model. The model is fit via maximum likelihood and parameters can be found via the Newton Raphson-algorithm as described in [4, p. 120]. The log-likelihood, which we want to maximize is given by

$$\ell(\beta_0, \beta_1) = \sum_{i=1}^n \{y_i(\beta_0 + \beta_1^T x_i) - \log(1 + \exp(\beta_0 + \beta_1^T x_i))\}$$

3.1.1 Non-regularized logistic

As a baseline model we fit the full data set using logistic regression. The confusion table for the prediction can be seen in table 2. The train and test error can then be calculated as

$$\text{train accuracy} = \frac{344 + 1455}{344 + 213 + 502 + 1455} = 0.7156$$

$$\text{test accuracy} = \frac{109 + 626}{109 + 91 + 252 + 626} = 0.6818$$

(a) Training results			(b) Test results		
	Blue	Red		Blue	Red
Blue	344	213	Blue	109	91
Red	502	1455	Red	252	626

Table 2: Confusion tables for logistic regression using all 135 features.

We see that the training accuracy is a bit higher than the test accuracy, when doing ordinary logistic regression. This is an indication that we have overfitted our train data. One way to mitigate this is using regularized logistic regression.

3.1.2 Regularized logistic regression

In regularized logistic regression, the weights are penalized in the loss function, restricting them from becoming too large. In LASSO regularization we use the L_1 -norm as a penalization. The optimal estimates of $\beta_0 \in \mathbb{R}$ and $\beta_1 \in \mathbb{R}^d$ are found by maximizing

$$\ell(\beta_0, \beta_1) - \lambda \sum_{j=1}^d |\beta_{1j}|$$

where $\ell(\beta_0, \beta_1)$ is the log-likelihood function of our data. The absolute value of the weights are subtracted in as we want to minimize that value.

In order to fit the model, the parameter λ has to be decided on. In order to make this selection we use crossvalidation. We used 10 folds in our crossvalidation. That is we split the training set in 10 equally sized parts, and fit the data using 9 of those parts, and then predicting on the last set. Then we repeat this until every part has been used as validation set. For choosing the sequence of λ 's we rely on `glmnet`'s default implementation. That is we set λ_{max} equal to the value where no features are included in the model and the next values equal to the point where the new weights are turned non-zero until λ_{min} as the λ in which all weights are non-zero, such that all features are included in the model. The results from cross-validation can be seen in fig. 1. By using the one-standard error rule we obtain an optimal lambda of $\lambda_{1se} = 0.017265$. Using this λ our model reduces to only 26 different features, as the weights for all the other features have been set to 0. Next we want to evaluate this model on the training and test set. The results from this prediction can be seen in table 3 and the training and test errors can be computed from the confusion tables as:

$$\begin{aligned}\text{train accuracy} &= \frac{335 + 1452}{335 + 211 + 516 + 1452} = 0.7108 \\ \text{test accuracy} &= \frac{46 + 693}{47 + 24 + 315 + 693} = 0.6855\end{aligned}$$

(a) Training set			(b) Test set		
	Blue	Red		Blue	Red
Blue	335	211	Blue	46	24
Red	516	1452	Red	315	693

Table 3: Confusion tables for regularized logistic regression, effectively using only 26 different features.

3.1.3 Important features

LASSO regularization also performs feature selection as it sets coefficients to zero. For the cross-validation plotted in fig. 1 there is new weights included for each value of λ . By setting $\lambda = 0.0714$ i.e. higher than the optimal value of lambda, the features with corresponding coefficients different from zero are seen in table 4. The most important features used to predict the winner are about the red fighters age, losses and significant strikes landed from the opponent.

The predicted winner for a data point x is the red fighter, when the $\beta_0 + \beta^T x > 0$. The model has a high intercept (β_0), which implies the baseline prediction is red as the winner. The positive intercept agrees with the fact that red is the most common winner. The features included in table 4 are all features that are negatively correlated with red winning. That is if the red fighter is older, he has a lower chance to win. If he has more losses beforehand he

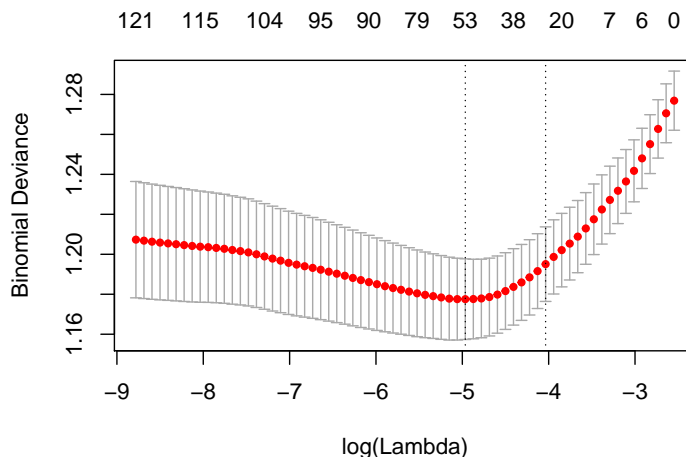


Figure 1: Crossvalidation results for logistic regression. The leftmost vertical line indicates the λ for which the crossvalidation error was lowest and the rightmost vertical line indicates the λ one standard error from the minimum value.

is more likely to lose and if his opponent has landed more significant strikes on average he is more likely to lose. Therefore the winner of the fight seems to be mostly influenced by details of the red fighter. If the red fighter is young, with few losses and average significant strikes landed from the opponent, then the red fighter is always predicted to win in this model, no matter how good the blue fighter is. The summary of the three features without scaling are depicted in table 1.

Feature	Weight
(Intercept)	0.67919167
R_age	-0.02186342
R_avg_opp_SIG_STR_landed	-0.02715641
R_losses	-0.01652865

Table 4: The intercept and weights different from zero, when setting $\lambda = 0.0714$ in regularized logistic regression.

3.2 Neural Networks

Neural networks using softmax as output function is also a soft classifier like logistic regression, which means it returns probabilities for each observation being in one of the K classes. The difference from logistic regression is that the decision boundary will be highly non-linear in neural networks. The neural network is therefore expected to perform better than logistic regression if the decision boundary between the two winners are actually non-linear.

3.2.1 Feed-forward Neural Network

In this report, we attempt to implement a feed-forward neural network on the 135 features and predict the winner of fights. Information in feed-forward neural networks move only in one direction, from the input nodes into the hidden nodes and finally into the output nodes. As the response is categorical with 2 levels; "Red" and "Blue", we utilized one hot encoding on the response.

3.2.2 Network Shape

The network shapes and dimensions that we used can be seen in table 5.

Layer	Units	Parameters	Activation function
Input	135	0	-
Hidden layer 1	50	6800	Exponential Linear Unit (ELU)
Hidden layer 2	30	1530	Sigmoid
Hidden layer 3	30	930	Hyperbolic tangent (Tanh)
Output	2	62	Softmax

Table 5: The table shows the structure of the neural network using classification. This includes the number of units and parameters for each layer together with the choice of activation function for each of the hidden layers.

We use 135 units in our input layer because the dimensions of our input data is 135. The output layer has 2 units because we utilized one hot encoding during conversion of the output variable. Since the variable Winner is in binary, this will lead to a dimension of 2 after one hot encoding. We use Softmax function to convert our output into probabilities that sum to 1. For each layer the number of parameters can be calculated as size of output (bias terms) plus size of input times size of output (linear combination parameters).

3.2.3 Model Attributes

We use Adam optimizer and categorical loss entropy as our loss function. The metric that we use is accuracy as we are doing categorical analysis. We also utilize an early stopping callback function. It monitors validation loss with a patience of 5. This means that if there is no improvement in the validation loss after 5 epochs, training will be stopped. The motivation to use an early stopping callback function is to prevent overfitting. Lastly, an epoch of 40 and batch size 128 is used with a validation split of 20% during the training of the neural network.

3.2.4 Model Training

As mentioned above, we use a validation split of 20% in training the neural network. This means that we do a 20% split of the train data and we train our neural network using 80% of the train data and validate it on 20% of the train data. This way, we will not be

doing forward looking into our test data during training and we only use the test data for predictions.

3.2.5 Results

The plot of the training/test accuracy and loss can be seen below in fig. 2. As we had used an early stopping callback function, we sought to prevent the neural network from over-fitting and potentially giving us worse predictions. The network converged within the designated epochs, which implies early stopping was engaged and further training would therefore be meaningless.

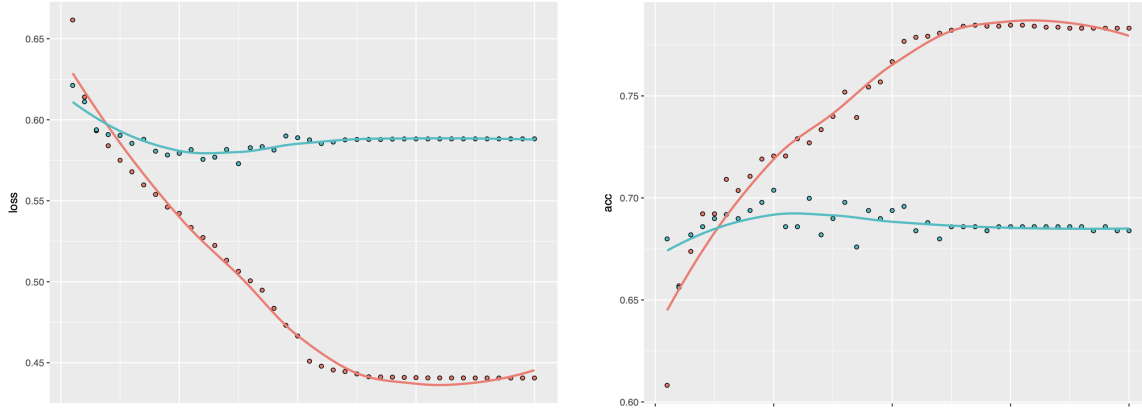


Figure 2: History of training and validation Loss and accuracy throughout training of the model. The red line and blue line indicates training and validation respectively. The error and accuracy stops changing between epochs towards the end as early stopping prevents the parameters from changing.

The green line in the plot above indicates the validation loss/accuracy and the red line indicates the training loss/accuracy. The x-axis is the number of epochs which is 40 in this case.

As our predictions are numerical probabilities, we use a threshold of 0.5 in deciding which categories they belong to. Having a predicted probability greater than 0.5 would result in the category being assigned to Red, else it would be assigned to Blue.

We construct a confusion matrix below to visualize our predicted directions as compared to the actual directions. The test accuracy is also shown below.

(a) Predicted VS Test

	Blue	Red
Blue	133	107
Red	228	610

$$\text{test accuracy} = \frac{133 + 610}{133 + 107 + 228 + 610} = 0.6892393$$

The results of the neural network leads to an accuracy of 68.9%.

4 Discussion

In this session, we compare the results of the two methods we used and discuss their pros and cons.

According to the result, we got 68.55% for the accuracy of LASSO logistic regression and 68.92% for the accuracy of the Neural Network. So we could see, although the Neural Network outperforms slightly, the prediction impacts of these two methods were generally similar. They both outperform the baseline learner of always predicting "red" as the winner, which yields a test accuracy of 66.26%.

The two methods yields different decision boundaries. The neural network fits a non-linear decision boundary to our data, while logistic regression produces a linear decision boundary. Therefore neural networks can fit a wider range of models, and could be a reason why the network slightly outperforms logistic regression in this case.

Although the Neural Network outperformed, it is a relatively more complex model. Since we have 9322 different parameters feature interpretation becomes difficult. It is hard to analyse what impact changing one of the features will have in the outcome generally. This is because the original features have been merged and transformed many times after passing through all the hidden layers. Hence it would be hard for us to tell which factors are essential.

Also, we need to note that Neural Networks require much more data compared to other traditional machine learning methods to have significant results. In the current data set used, we have about 3592 rows of data which might be insufficient for the neural network, as relatively few iterations of training is needed for the network to converge and as seen in fig. 2 the training loss is converging to 0 when the validation error is stabilizing. This indicated that the model is prone to overfitting, which is indicates that we do not have enough datapoints.

However, In logistic regression, especially with LASSO, the original features were not transformed and the dimension was also reduced significantly from 135 to 26. Therefore, we can easily tell which features are more important, which may be more efficient as we can use the model to guide the fighters in their training.

5 Conclusion

Generally speaking, logistic regression and neural network have their pros and cons in this case. Perhaps it might be justifiable to forgo the benefits of accuracy that the neural network has to offer in return for the interpretability of the logistic regression model depending on the use case.

From the perspective of a trainer training his fighter for the competition, certainly the logistic regression model would be more helpful as it gives a clearer interpretation of the features and provides the fighters with a direction to focus on to increase their odds of winning.

However, from the perspective of perhaps someone that is betting on the game, the focus would be aiming to make the prediction more accurate and accuracy would be imperative here. Hence, the Neural Network would be a more optimal choice.

For further work in this field it could be interesting to look at the soft clustering we produce. From the probabilities the odds can be calculated, and these can be compared to the various odds offered by sports betting companies. It could be interesting to see if there are cases where our models performs very different from real sport odds.

References

- [1] Rajeev Warriar, UFC-Fight historical data from 1993 to 2019,
<https://www.kaggle.com/rajeevw/ufcdata>
- [2] History of the UFC <https://www.ufc.com/history-ufc>
- [3] Unified rules of mixed martial arts
<https://www.ufc.com/unified-rules-mixed-martial-arts>
- [4] Hastie, Tibshirani and Friedman (2009). The Elements of Statistical Learning (2nd edition), Springer-Verlag. 763 pages.