

Regression Analysis of Overwatch League Data

Andrew Rodgers

1 Introduction

The goal of this project is to perform a regression on the Overwatch League (OWL) data set. The data is a subset of a stat sheet from recent matches in the 2019 season. Excluded data points are players who did not appear in a match, thus the excluded points are rows of zeros. The quantitative variables measured are Kills(K), Deaths(D),and Ultimates(U). The categorical variable is Role(Tank,Offense,Support). The purpose is to create a model to predict the number of Points earned, given these values for any given player in a future Overwatch League match.

2 Analysis

2.1 Basic MLR Regression

A good place to start the analysis is with the regular, basic regression. From **Figure 1a**, we can see the p-values are around 0, so we can say the model is useful by the F-test. In **Figure 1(b)**, we can see most of the p-value for *Role_Support* is

(a)
ANOVA
table
for
the model
with
indicator
variables

Figure 1

2.2 Variable Selection(Optimal Subsets)

In addition to Role, the data set includes a second categorical variable, Team. We exclude this predictor in our analysis simply because there are 20 responses. If we used a predictor with this many responses, we would not have enough data for each group to make an accurate conclusion. An Extra Sum of Squares Test confirms that Team does not contribute to the model. Before we can choose the best set of predictors for the model, we need to create indicator variables for the Role categorical variable, Tank is used as the reference group. In **Figure 2(a)**, we can clearly see in the plot of K versus Points, different roles have significantly differing slope, so we need interaction terms for K. The other plots do not have clearly different slopes, so we will use the additive model for these. With the new interaction terms, we can add the predictors $K*Role_Support$ and $K*Role_Offense$ to our model. Looking at **Figure 2(b)**, we can see the best model is the full model, to little surprise.

(a)

Best
 subsets
 regressions
 Points
 versus
 Predictors

Figure 2

2.3 Stepwise Regression

This section intentionally left blank.

2.4 Multicollinearity

There will always be some correlation between the predictors because of the way the game works. For example, Kills and Deaths should be inversely related, and Ultimates and Kills should be related. There is multicollinearity present between the predictors, as demonstrated by **Figure 3(a)**. In **Figure 3(b)**, we can see almost all of the VIF's are above 5, which is another indication of multicollinearity.

(a)

Correlation
 Matrix
 for
 predictors

Figure 3

2.5 Outliers

Looking at the 4-in-1 plots, there is only one real outlier in our regression. Calculating Cook's Distance for the data, we see this point (row 72) has a significantly higher Cook's Distance, 0.138, and a much larger DFITS, 1.227, than is expected.

2.6 Transformations

If we exclude the single outlier, the data is very nice, with an excellent symmetry. Looking at histograms of the predictors, we can see all of the distributions are roughly normal or symmetrical. Scatter plots show no clear nonlinear patterns. Thus, we have no reason to suspect our data needs any transformations.

2.7 Assumptions

Figure 4 : 4-in-1 plot

All of the usual assumptions for regressions are met in this case. The 4-in-1 plot shows no correlated errors, normally distributed residuals, and constant variance. We assume the data is linearly related, so we have the correct model function.

3 Conclusions

To little surprise, using the full model with the interaction terms gives us the best model.

$$\begin{aligned} Points = & 4.77 + 3.937 * K - 1.5 * D + 1.850 * U - 9.65 * Role_Offense \\ & + 7.68 * Role_Support + 1.109 * K * Role_Support - 2.03 * K * Role_Offense \end{aligned} \quad (1)$$

We found no evidence any transformations were needed to the data, or that any subset of the predictors would give a better model. In this case there is multicollinearity, which is due to the game rules and scoring system.

The full model works best because of the game rules and scoring system. In fact, the actual Points calculation is a linear combination of Kills, Deaths, and Ultimates, along with 2 other statistics not provided to us. The coefficients in this calculation differ between Roles. The game is designed so that more Ultimates leads to more Kills, which leads to less Deaths; explaining most of the multicollinearity.

Looking at the predicted R-squared of our model, we can say it is fairly good at predicting values based on new data. Applying the model to a small set of test data, our model overestimated Points, but not by a significant margin, likely due to random chance.

While performing the analyses, I was intrigued by the lack of effect that Team had on the model. Maybe it was because there are 20 teams and there would be many fewer observations per team. The more likely answer is that the consistency of skills at the professional gaming level is so high that the most likely variables to affect the model would be hard to quantify; such as teamwork, play style, or something else during the match. Based on this data alone, it does not appear that the best and worst teams would have very different models, but perhaps a larger set of observations would indicate this.