

Classification trees, random forests

Marko Medved

I. PART 1: IMPLEMENTING MULTINOMIAL LOGISTIC REGRESSION AND ORDINAL LOGISTIC REGRESSION

A. Implementation details

The first model implemented is multinomial logistic regression. In this approach, we calculate linear predictors for each class, excluding the reference class. Afterward, we apply the softmax function to compute the probability for the true class. To optimize the model's weights, we use the maximum likelihood estimation (MLE) method, leveraging the *fmin_l_bfgs_b* optimizer with numerical gradients.

An important aspect of this implementation is the inclusion of an intercept column. This is necessary because, in some cases, all feature values may be zero. Without the intercept, this would lead to zero probabilities for all classes. By including the intercept, we ensure that the model can still generate non-zero probabilities for all classes, even when the features are zero.

During the prediction phase, we apply the softmax function to the linear predictors for each class at each data point. This gives us the predicted probabilities for all classes, which are then reported as the model's predictions.

Next, ordinal logistic regression was implemented. The main difference with multinomial logistic regression is the assumption that the target categories have a natural ordering. This assumption allows us to significantly reduce the number of parameters. Instead of calculating separate predictors for each class, we use just one linear predictor along with thresholds (which are fewer than the number of classes) to determine the boundaries between classes.

During the model building phase, we again use the MLE approach, as we did with multinomial logistic regression, and the same optimizer. Instead of explicitly modeling the thresholds, we model the differences between the thresholds. To ensure numerical stability, we constrain these differences (the deltas) to be greater than 1 during the optimization process.

This reduction in parameters, combined with the ordered nature of the categories, makes ordinal logistic regression more efficient for datasets with ordered categorical targets.

B. Testing our implementation on the basketball dataset

We evaluated both classifiers on the provided dataset to predict shot type based on various features. To preprocess the data, we One-Hot encoded the categorical variables and dropped one category from each feature to avoid multicollinearity. Furthermore, we identified strong correlations between certain features. Specifically, the "TwoLegged" and "No Movement" features, as well as the player type indicators "Guard" and "Forward," exhibited high correlation, which we can observe on Figure 1. As a result, we removed the "TwoLegged" and "Forward" features to prevent convergence problems and reduce redundancy.

The model evaluation was conducted using 10-fold cross-validation, and the uncertainty in the results was calculated under the assumption of asymptotic normality. The results are presented in Table I. As expected the multinomial version outperformed the ordinal one, since there is no clear natural order for shot types in basketball.

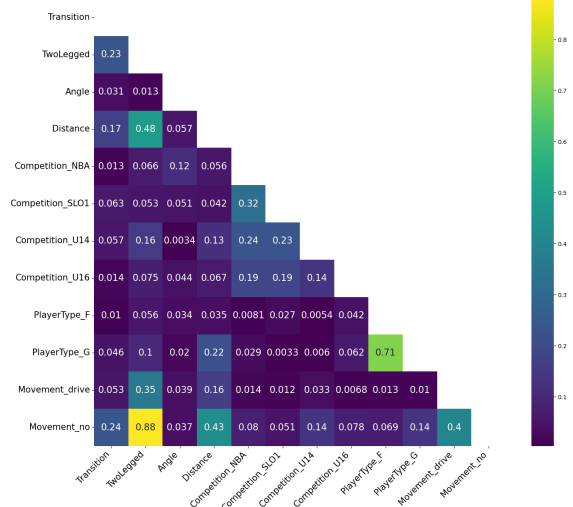


Figure 1. Absolute correlation of features in the basketball dataset.

Model	Accuracy	Log loss
Multinomial	0.7353 \pm 0.0062	0.667 \pm 0.012
Ordinal	0.7138 \pm 0.0064	0.985 \pm 0.018

Table I

COMPARISON OF RESULTS USING MULTINOMIAL AND ORDINAL LOGISTIC REGRESSION

II. PART 2.1: APPLICATION OF THE MULTINOMIAL REGRESSION

In this analysis, we explored the relationships between shot type and various features. To gain insight, we directly compared the coefficients from the linear predictors. These coefficients represent how the log-odds of a given shot type change when a particular feature is adjusted, while holding other features constant.

To estimate the uncertainty of these coefficients, we employed bootstrap sampling. We repeatedly fitted the model to bootstrap samples, calculated the mean of the betas as our best estimate, and used the standard deviation as the measure of uncertainty. We used the "tip-in" shot type as our reference class.

The results of our analysis are displayed in Figure 2. From the plot, we observe that some features exhibit high uncertainty compared to the value of coefficients in certain shot types, suggesting that these features may not contribute significantly to differentiating those shot types from the reference class. On the other hand, several features show consistently high coefficients across multiple shot types.

For instance, the "Movement_no" coefficient is notably influential. The likelihood of no movement would intuitively be high for the "above head" shot. However, this comparison is made against the "tip-in" shot, where the only movement in the dataset is "Movement_no"—which is also true for "dunk"

and "hook shot" categories. In contrast, "Movement_drive" has a clear, intuitive impact. For example, a drive is less likely to be associated with an "above head" shot but is much more common for "layup" and "other" shots. For the "tip-in", "hook shot", and "dunk", the likelihood of a drive is relatively similar to the reference class.

The "Distance" feature also shows an interesting pattern: the "tip-in"(reference) shot is associated with the shortest distances, followed by the "dunk", "layup", and "hook shot". Both "other" and "above head" shots, on the other hand, have higher coefficients, which aligns with their nature (mostly for the above head), as they often occur from farther distances.

When examining the "Competition_U14" feature, we see that "tip-in" is less likely compared to all shot types except for the "dunk". This makes sense, as younger players (under 14) are unlikely to perform tip-ins and even less likely to dunk.

Finally, the intercept plays a significant role across all shot types. It represents the baseline log-odds of selecting a particular shot type when all other features are zero.

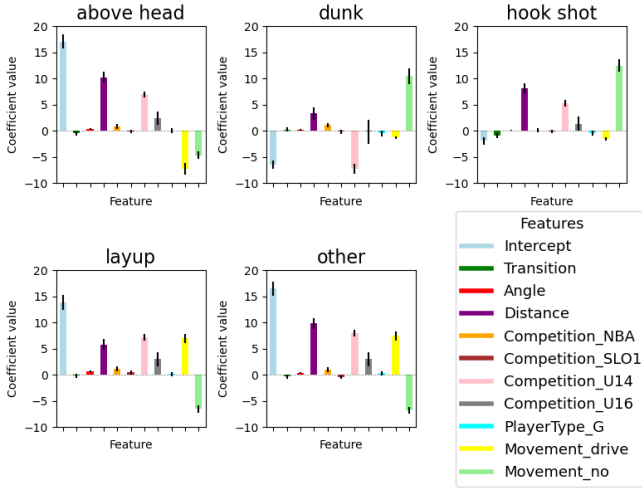


Figure 2. Comparison of relationships of features and different shot types

III. PART 2.2: APPLICATION OF THE ORDINAL REGRESSION

In this section, we designed a data-generating process (DGP) where ordinal logistic regression outperforms multinomial logistic regression. The process was designed as follows: We begin by assigning a class to each data point. For each class, we generate the corresponding feature values by sampling from a uniform distribution spanning from $i - 0.5$ to $i + 0.5$, where i is the current class. This ensures that the data is ordinal in nature. To introduce some variability and reduce the correlation between features, we then add noise sampled from a standard normal distribution.

We tested the DGP using a dataset with 100 training data points and 100 testing data points. The results of this experiment can be seen in Table II.

IV. PART 3: GLM DIAGNOSTICS

In this section, we implemented linear regression under the assumption that the residuals follow a normal distribution, utilizing the least squares method. We applied this model to

Model	Accuracy	Log Loss
Multinomial Logistic Regression	0.580 ± 0.050	0.942 ± 0.100
Ordinal Logistic Regression	0.630 ± 0.049	0.823 ± 0.070

Table II

COMPARISON OF ACCURACY AND LOG LOSS BETWEEN MULTINOMIAL AND ORDINAL LOGISTIC REGRESSION MODELS.

predict basketball shot distance based on the shot angle. The results of the fitted linear regression can be seen in Figure 3.

Upon examining the residuals, we observe that their distribution is bi-modal, rather than normal as we initially assumed. This is not entirely surprising, given that the distribution of shot distances, as shown in Figure 4, is also bi-modal. Additionally, the correlation between shot angle and shot distance is very low (-0.056), which further contributes to the similar distribution of residuals across all data points.

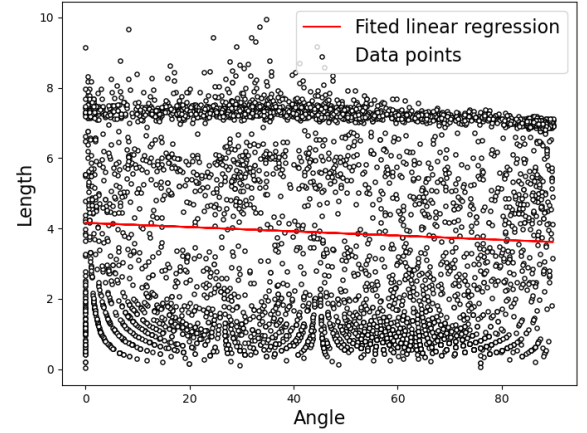


Figure 3. Linear regression, predicting shot length from shot angle

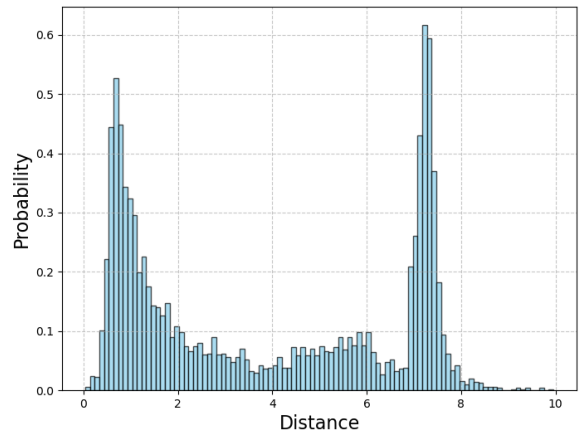


Figure 4. Distribution of the distance of basketball shots

As part of our diagnostic checks, we plotted the Q-Q plot (Figure 5) to assess whether the assumption of normally distributed residuals holds. In the Q-Q plot, we compare the

quantile residuals (we choose the quantile residuals because they are approximately normally distributed as long as we choose the correct distribution from the exponential family) to the theoretical quantiles of a standard normal distribution. If the residuals were normally distributed, the points would align closely with the reference line. However, as shown in the plot, the residuals deviate significantly from this line, indicating that the normality assumption is not justified. Again, we can see the consequences of the distribution of our residuals being bi-modal.

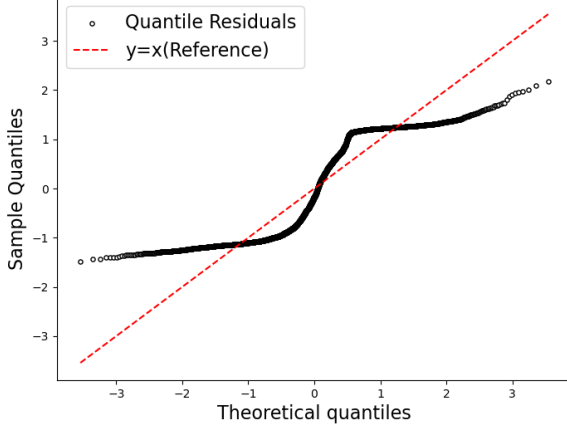


Figure 5. Quantile-quantile plot

The second plot is the quantile residuals vs fitted values plot. We chose to use quantile residuals because they tend to have approximately constant variance (if we choose the correct distribution), which makes them suitable for assessing homoscedasticity. This plot helps us verify whether the variance of the residuals changes across the range of fitted values.

In Figure 6, we observe that there is no apparent trend or pattern in the variance of the residuals as the fitted values increase. This suggests that the residuals are homoscedastic. Here this would not necessarily mean that we chose the correct distribution. The homoscedasticity comes from the fact that the two variables are not really correlated so we have approximately the same distribution of residuals at every point and consequently the same variance.

The last plot, shown in Figure 7, represents Cook's distance for each data point. Cook's distance measures the influence of individual data points on the regression model, taking into account both the features and the target values. A high Cook's distance indicates that a data point has a significant impact on the model, potentially altering the estimated regression coefficients.

In the plot, we observe that none of the points have a particularly high Cook's distance. As a rule of thumb, data points with a Cook's distance close to 1 are considered highly influential. Since all the points lie way below this threshold, we can conclude that there are no particularly influential data points that would unduly affect the model. Therefore, it is safe to assume that the regression results are not being significantly impacted by any single data point.

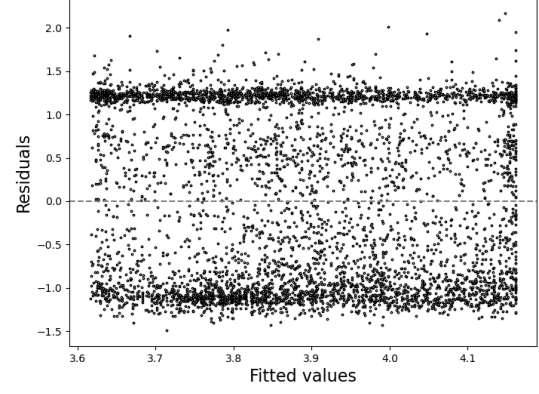


Figure 6. Residuals vs fitted values

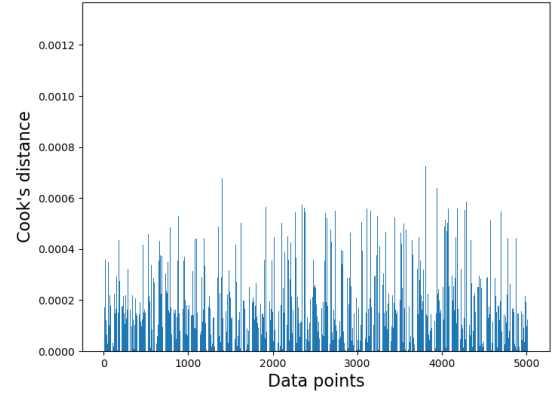


Figure 7. Cook's distance across data points