

# Contents

<b>1</b>	<b>Abstract:</b>	<b>3</b>
<b>2</b>	<b>Introduction</b>	<b>3</b>
<b>3</b>	<b>Exploratory Data Analysis</b>	<b>3</b>
3.1	Outlier Check & Transformations . . . . .	3
3.2	Variable Elimination . . . . .	3
3.3	Addressing Multicollinearity: Correlation Heat Map for Visual Data Exploration . . . .	4
3.4	Post-Correlation Heat Map Variable Elimination . . . . .	4
3.5	Addressing Multicollinearity: Correlation Matrix for Numerical Analysis . . . . .	5
<b>4</b>	<b>Quadratic Discriminant Analysis</b>	<b>5</b>
4.1	Bartlett's Test . . . . .	6
4.2	Quadratic Discriminant Analysis: Internal Cross-Valdiation and Model Development . .	6
4.3	Quadratic Discriminant Analysis: External Cross-Valdiation . . . . .	6
4.4	Quadratic Discriminant Analysis: Internal vs. External Cross-Validation . . . . .	7
<b>5</b>	<b>Logistic Model Development Using Ordinary Least Squares</b>	<b>8</b>
5.1	Logistic Regression: Model Selection . . . . .	8
5.2	Logistic Regression: Internal Cross Validation and Model Development . . . . .	8
5.3	Logistic Regression: External Cross Validation and Model Development . . . . .	9
5.4	Logistic Regression: Internal vs. External Cross-Validation . . . . .	9
5.5	Logistic Regression: Fitted Model . . . . .	9
<b>6</b>	<b>Shot Made Odds and Probabilities Analysis Using Logistic Regression</b>	<b>10</b>
<b>7</b>	<b>Conclusion</b>	<b>11</b>
<b>8</b>	<b>Appendix A: Figures and Tables</b>	<b>12</b>
8.1	Receiver Operator Characteristic (ROC) Curves . . . . .	12
8.2	Fitted Logistic Regression Model . . . . .	12
8.3	Action Type Frequency Table . . . . .	14
8.4	Variance Inflation Factor . . . . .	15
<b>9</b>	<b>Appendix B: Source Code</b>	<b>16</b>
9.1	Loading Libraries, Importing Data . . . . .	16
9.2	Check for Missing Values . . . . .	16
9.3	Basketball Shot Location Map . . . . .	16
9.4	Shot Location Imputing, toInt Conversion . . . . .	16
9.5	Variable Eliminations, First-Pass . . . . .	16

9.6	Correlation Heat Map . . . . .	17
9.7	Variance Inflation Factor Table . . . . .	18
9.8	Variable Eliminations, Second Pass . . . . .	18
9.9	Correlation Matrix . . . . .	18
9.10	QDA Bartlett Approximation . . . . .	19
9.11	Partitioning Training, Testing Data . . . . .	20
9.12	a Priori Analysis . . . . .	20
9.13	Quadratic Discriminant Analysis: Internal CV and Model Development . . . . .	20
9.14	Quadratic Discriminant Analysis: External CV . . . . .	21
9.15	Quadratic Discriminant Analysis Confusion Matrix Table . . . . .	21
9.16	Predictions from Quadratic Discriminant Analysis . . . . .	22
9.17	Logistic Model Development Using Ordinary Least Squares . . . . .	22
9.18	Logistic Regression: Internal CV . . . . .	23
9.19	Logistic Regression: Internal vs. External CV . . . . .	25
9.20	Logistic Regression: Fitted Model . . . . .	26

# 1 Abstract:

*This project investigates the correlation between multiple potential explanatory variables and Kobe Bryant's ability to make a shot while playing for the NBA team Los Angeles Lakers using data gathered from 1996-2015.*

## 2 Introduction

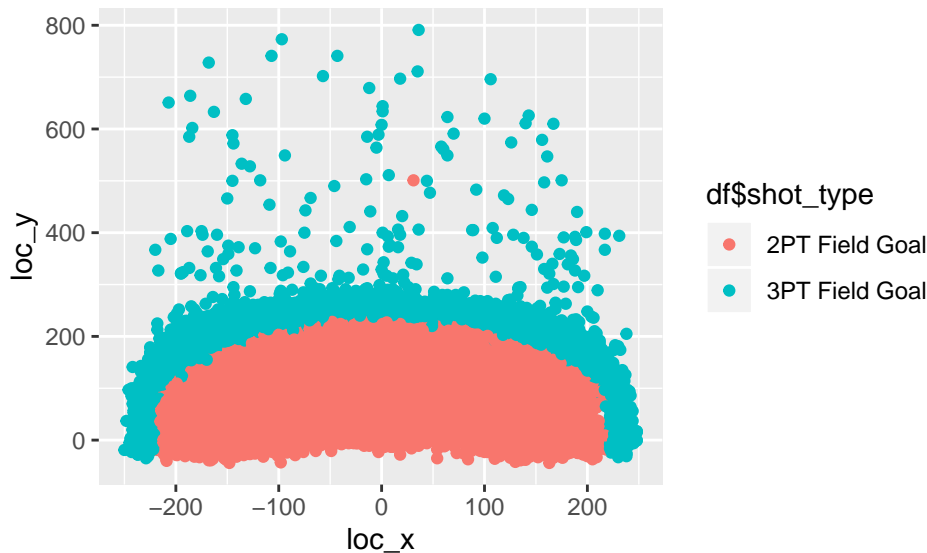
Kobe Bryant, a professional basketball player, spent his entire 20-year career with one team, the LA Lakers as a guard. Bryant started with the Lakers just after high school and went on to win five NBA Championships and 18 All-Star titles. Kobe Bryant's career average points scored per game is 25-points during the regular season and 26-points during playoffs. In addition to these statistics, using the **KobeData** provided, we were able to determine the odds of Kobe making a shot decreases with respect to his distance from the hoop and the probability of him making a shot decreases after taking into consideration his distance from the hoop and analyze the relationship distance and making a shot during the playoffs. During model selection we considered Area Under the Curve (AUC), Mis-Classification Rate, Sensitivity, Specificity and objective / loss function when making model comparisons.

## 3 Exploratory Data Analysis

### 3.1 Outlier Check & Transformations

First, we performed a brief outlier check, which included a graphical analysis of all shots taken, by `loc_x` and `loc_y`. This graphical analysis indicated a 2PT (2-point) Field Goal was recorded from the 3PT (3-point) range. Upon inspection of other attributes - such as `action_type` and `shot_zone_range` - we verified this shot to be a member member of the 3-point level of `shot_type`. Under the assumption shots from beyond the 300 inch mark are more likely to have been incorrectly recorded as 2 points rather than an incorrectly recorded location y, we modified our programming to transform all shots where `loc_y > 300` to be recoded as 3PT Field Goal.

Additionally, `action_type` was recoded where shots types were categorized into `short` and `not_short`. However, modeling performance decreased because the parameter became less descriptive. Further transformations included using indicator variables for `shot_made_flag`, depending on the phase of analysis and modeling.



### 3.2 Variable Elimination

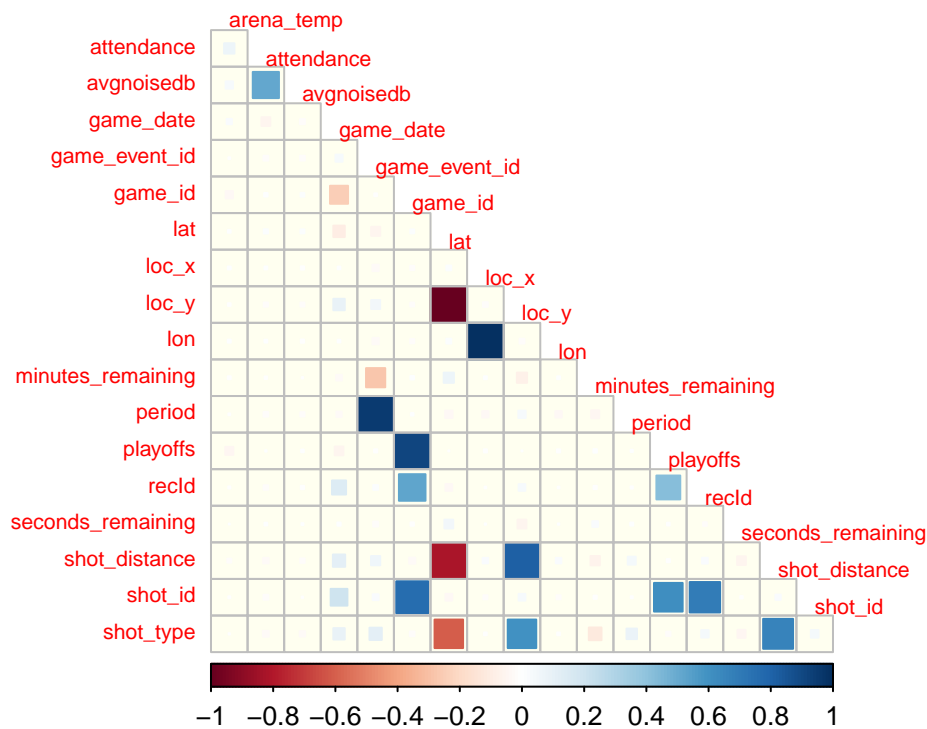
Next, we removed one-level factors. These will never change so are not useful to the model; including can cause issues with model sensitivity since linear trajectories will be down-weighted. Therefore, their significance will be

lessened by the constant state of the additional parameters. While this may not be significant, it is not conducive to model quality.

### 3.3 Addressing Multicollinearity: Correlation Heat Map for Visual Data Exploration

To address multicollinearity among quantitative predictor variables, a correlation heat map (below) was created for visual inspection of correlation. Red corresponds to negative correlation (where an increase in a predictor causes a decrease in the value of its collinear counterpart) whereas blue corresponds to positive correlation.

**Correlation among Predictor Variables**



### 3.4 Post-Correlation Heat Map Variable Elimination

Following our correlation heat map, we decided to eliminate some collinear terms. However, some of the collinearity is useful to capture the instances where the terms are unique. For example, **combined\_shot\_type** (factor variable) is collinear with **shot\_distance** (quantitative variable), but it also accounts for the method Kobe may use to take a shot. For example, distance may be relatively the same between 10 and 11 feet, but the factor levels used to derive their **short** or **far** indications may differ. This difference could be whether Kobe makes a potentially more accurate heel-planted shot or if he is forced to lean forward and take a riskier shot at basket; the difference in distance may only be one foot, but the difference in technique could measure significant relative to the odds of success. Because this is a healthy level of confusion for the model, some multicollinearity was accepted.

### 3.5 Addressing Multicollinearity: Correlation Matrix for Numerical Analysis

After deselecting the most obvious collinear terms through visually inspection of the correlation plot, a correlation matrix for analyzing the remaining results. Collinear quantitative data was preliminarily removed following correlation plot analysis to desaturate the model to an extent that allows more distinction among significance measures for terms in the correlation matrix. Following the removal of predictor variables after visually inspecting the correlation heat map, we analyzed a correlation matrix. However, the matrix itself did not identify any remaining collinearity at a threshold of correlation necessitating removal of like-terms. Consequently, no further predictor variables are removed. Please refer to the table below for a list of the remaining top 10 collinear terms following first-pass variable removal. While this risks over-fitting, we believed the variables were different enough to capture useful information for the model.

The Variance Inflation Factor (VIF) scores across all remaining variables is another test applied in consideration for multicollinearity among explanatory variables. A threshold of roughly 1 or higher is considered a strong level of collinearity. Both VIF and correlation matrix statistics were considered during model development.

#### 3.5.1 Top 10 Multicollinear Terms, Correlation Matrix:

Correlation Predictor Variable	Correlation Response Variable	Correlation	p-Value
loc_y	shot_distance	0.81812	p < 0.0001
recId	shot_id	0.69172	p < 0.0001
shot_distance	shot_type	0.66861	p < 0.0001
playoffs	shot_id	0.61299	p < 0.0001
loc_y	shot_type	0.60662	p < 0.0001
attendance	avgnosedb	0.51092	p < 0.0001
recId	playoffs	0.42527	p < 0.0001
game_date	shot_id	0.20932	p < 0.0001
recId	game_date	0.14773	p < 0.0001
game_event_id	shot_type	0.11238	p < 0.0001

#### 3.5.2 Top 10 Multicollinear Terms, Variance Inflation Factor:

Variable	Variance Inflation Factor
shot_id	39939311.92
recId	39937863.50
shot_distance	3.59
loc_y	3.15
shot_type	1.90
playoffs	1.76
game_event_id	1.61

## 4 Quadratic Discriminant Analysis

As requested within the requirements of this study, a Linear Discriminant Analysis must be assessed and provided. Discriminant analysis is an operation that compares a categorical response variable against measures of quantitative predictor variables. As a result, analysis for this section is performed on the numerical predictors, which include recId, game\_event\_id, game\_id, loc\_x, loc\_y, minutes\_remaining, seconds\_remaining, shot\_distance, shot\_made\_flag, shot\_type, game\_date, shot\_id, attendance, arena\_temp, avgnosedb, controlling collinearity by eliminating a member of each collinear pair prior to model development.

Linear Discriminant Analysis requires a linear boundary between the predictor variables, respective of the response. If the boundary between predictors and response is not linear, Quadratic Discriminant Analysis

(QDA) must be used. Wilks' Lambda distribution is used to assess the nature of boundary linearity, which is a required understanding to develop a well-fit discriminant classification model. However, because of the large dimensions of the data set analyzed in this study, an approximation of Wilks' Lambda must be used, rather than Wilks' Lambda itself. Bartlett's Test is an approximation of Wilks' Lambda that can be used for models with large dimensions by applying a measure against the Chi-Square distribution. This method is applied herein to assess linearity.

Also to note of importance, a priori knowledge was set to the default proportions when performing Quadratic Discriminant Analysis.

## 4.1 Bartlett's Test

The result of the Bartlett's test returned statistically significant results, indicating the null hypothesis of linearity must be rejected in favor of the alternate, which is that the discriminant boundary is non-linear. Consequently, we proceed with a model based on Quadratic Discriminant Analysis to provide predictive responses from a discriminant model. However, we proceed with caution, as the quadratic version of the discriminant analysis is at greater risk for over-fitting to the data than Linear Discriminant Analysis as the boundary is required to conform more closely to the data rather than to the mean of the data. This was also taken into consideration when assessing the results of the Logistic Regression model development that occurs afterward. Bartlett's Test of this data set yielded a significant p-value, where  $p < 0.0001$ , indicating that the proportion of distribution beyond the derived test statistic is beyond that which could be explained by chance. Therefore, we must reject the null hypothesis that the boundary for analysis is linear; the boundary is non-linear. Thus, an analysis using Quadratic Discriminant Analysis is applied. Please refer to the table Bartlett Test's Wilks' Lambda Approximation below.

### 4.1.1 Bartlett Test's Wilks' Lambda Approximation

	Statistics
Chi Square Statistic	1037.24251
Degrees of Freedom	14
Wilks' Lambda	0.9511
p-Value	$p < 0.0001$

## 4.2 Quadratic Discriminant Analysis: Internal Cross-Valdiation and Model Development

Following removal of significant levels of multicollinearity from the dataset and partitioning into a 75% training / 25% testing split, internal cross-validation is performed. The specifics of this test involves 25 folds of the data - meaning the 75% training data is divided into 25 partitions. The model is then trained on 1/25th of the original 75%, then tested against the remaining 24/25ths, 1/25ths at-a-time. This test is repeated 5 times, with each repeat involving a different random partitioning of the 25 specified folds of the data. Finally, the model developed using the 75% training split is then applied to the 25% testing split and predictions are measured against the actuals of that split to develop model statistics such as Accuracy, Misclassification, Precision, Sensitivity and Specificity.

This is internal cross-validation and its objective is to assist in identifying a model that can perform well across different data sets of the same source by using partitions of data from the same set to simulate an environment where those partitions are actually data from different sets. This is typically performed prior to external cross-validation.

## 4.3 Quadratic Discriminant Analysis: External Cross-Valdiation

After building a model using internal cross-validation, which applied 5 repeated internal cross-validations across the 25 folds of training data, a confusion matrix was constructed and analyzed. Next, we applied the model developed using the 75% training split to make predictions against the entire portion of data that includes values

for `shot_made_flag` in order to assess how closely the model can predict against the entire data set compared to the actuals.

Applying the model to the entire dataset as **external cross-validation** provides the model an opportunity to test against different data and more closely simulate a real-life scenario than internal cross-validation. Internal and external cross-validation is performed for later Logistic Regression models as well. Following external cross-validation of both models, the metrics are compared to determine the best model (Quadratic Discriminant Analysis versus Logistic Regression).

A confusion matrix is a table of results from cross-validation. Some key metrics provided by a confusion matrix include **Accuracy**, **Precision**, **Sensitivity** and **Specificity**. **Accuracy** is the number of all correct predictions divided by the number of all predictions. **Precision** is the ratio of the number of correctly classified positive predictions divided by the number of all positive predictions. **Sensitivity** (also called **Recall**) is the number of correctly classified positive predictions divided by all positive actuals - this is similar to precision, except that sensitivity measures against actual values. **Specificity** is the number of correctly classified negative predictions divided by all negative actuals. Simplistically, sensitivity is the true positive rate whereas specificity is the true negative rate. Higher Accuracy, Precision, Sensitivity, and Specificity is desirable.

Another important component for cross-validation is the **Misclassification Rate**. The Misclassification Rate is a descriptor of how often a model is wrong. This value is equal to the total number of False Positives plus the False Negatives divided by all predictions. A lower misclassification rate is desirable.

In addition to the misclassification rate, Accuracy, Precision, Sensitivity, Specificity and Misclassification Rate, the Logarithmic Loss function is applied to measure . A lower logarithmic loss value is desirable as logarithmic loss increases as predicted probability diverges from the actual response values and conversely decreases as predicted probability moves converges toward the actual response values.

Two final metrics used in this analysis are the **Area Under the Curve (AUC)** and **Receiver Operating Characteristic (ROC) curve**. The ROC bounds an area the area which the AUC describes. As a discrimination threshold changes, the ROC visually represents the correct diagnostic ability of a binary classification model and is a plot of the true positive against the false positive rate at those varied thresholds. As the AUC describes the area under this curve, a higher AUC is more desirable than a lower AUC. As mentioned, these metrics will be analyzed when comparing internal to external cross-validation to ensure consistency as well as between the QDA and Logistic Regression models.

## 4.4 Quadratic Discriminant Analysis: Internal vs. External Cross-Validation

Using the two confusion matrix output tables immediately below, the performance across internal and external cross-validations of the QDA model can be compared. As indicated in those figures, the model performed highly similarly across both cross-validation techniques, indicating the model is consistent and reasonably fit, after controlling for the variables selected for modeling.

### 4.4.1 Confusion Matrix Results for Quadratic Discriminant Analysis

	Internal CV Statistics	External CV Statistics
Sensitivity	0.51431	0.51187
Specificity	0.66761	0.66993
Precision	0.56104	0.55695
Accuracy	0.59826	0.59917
Misclassification Rate	0.40174	0.40083
Logarithmic Loss	0.70127	0.70127
Area Under the Curve	0.40904	0.59090

## 5 Logistic Model Development Using Ordinary Least Squares

Logistic Regression is a classification technique that is best suited for dichotomous response variables - in the case of the Kobe data the response is '0' for shot missed, or '1' for shot made. Compared to discriminant analysis techniques, multiple explanatory variables, interactions, and categorical variables can be used, allowing for a potentially more descriptive model. For this type of regression, coefficients are in log-odds where each coefficient needs to be exponentiated to yield odds ratios - this is done for ease of interpretation. Logistic Regression can also be used to generate predictions that yield the probability of an observation having the desired traits of the response variable as occurring or not.

### 5.1 Logistic Regression: Model Selection

A preliminary, manual variable elimination process was performed during the analysis of multicollinear terms in preparation for model development. Below we perform logistic regression using Ordinary Least Squares (OLS). In preparation for the model development, a starting model and a finishing model must be developed to provide the scope of variable selection. These initial models are used by forward, backward, and stepwise model selection methods are used to help select a combination of variables that result in the lowest Residual Deviance and/or AIC. The selection method that generates the model with the lowest AIC/Residual Deviance is then used for internal cross validation to further tune the model which allows for better prediction. Below are models that each selection method generated:

**Forward Selection Model:**  $\text{shot\_made\_flag} = \text{action\_type} + \text{attendance} + \text{arena\_temp} + \text{game\_event\_id} + \text{season} + \text{seconds\_remaining} + \text{minutes\_remaining} + \text{loc\_y} + \text{game\_date} + \text{loc\_x}$

**Backward Elimination Model:**  $\text{shot\_made\_flag} = \text{recId} + \text{action\_type} + \text{game\_event\_id} + \text{loc\_x} + \text{minutes\_remaining} + \text{season} + \text{seconds\_remaining} + \text{shot\_distance} + \text{game\_date} + \text{shot\_id} + \text{attendance} + \text{arena\_temp}$

**Stepwise Regression Model:**  $\text{shot\_made\_flag} = \text{recId} + \text{action\_type} + \text{game\_event\_id} + \text{loc\_x} + \text{minutes\_remaining} + \text{season} + \text{seconds\_remaining} + \text{shot\_distance} + \text{game\_date} + \text{shot\_id} + \text{attendance} + \text{arena\_temp}$

Based on the fit-statistics generated from each model selection method, the backwards and stepwise models are identical in fit-statistics with an AIC at 25167.77, and the residual deviance at 25001.77. The forward model selection out-performs both backwards and stepwise models with an AIC of 25166.48, but has a higher a residual deviance at 25001.48. Compared to the forward selected model, the backwards and stepwise models have lower residual deviances, although their residual deviances are very close in value to the forward model, their AIC values are larger compared to the forward selected model. Due to the model selection process, the forward selection model does not contain `shot_distance`; in order to meet the project requirements the stepwise model will be used. The stepwise model has a slightly higher AIC compared to the backwards selected model, but has a lower residual deviance. Residual deviance is a statistic used for goodness-of-fit, where a lower value is desirable. We chose to explore this route because the stepwise model has a AIC value that is similar to the backwards model, but has lower residual deviance.

Selection Type	AIC	Residual Deviance
Forwards	25168.23	25004.23
Backwards	25167.77	25001.77
Stepwise	25169.75	25001.75

### 5.2 Logistic Regression: Internal Cross Validation and Model Development

After identifying that the forward selected model is the best candidate based on it's AIC and residual deviance, its features are tuned using an internal cross validation. A training set is generated by randomly selecting 75% of the observations, with the remaining 25% serving as the validation (test) set. The model is tuned using 25 folds and is repeated 5 times, this process is the same as described in section 6, "Quadratic Discriminant Analysis: External Cross-Valdiation and Model Development." Prior to the internal cross validation process, action types that rarely



occur in the data are recoded to similar, but differ action types, i.e. “Running Tip Shot” (1 observation) to “Tip Shot” (many observations). If levels within the evaluation data exist but are not present in the training data, a trained model will have difficulty making predictions. Infrequently occurring action types are recoded to similar action types to avoid this situation.

### 5.3 Logistic Regression: External Cross Validation and Model Development

External cross-validation is used to evaluate a model after it has been tuned in the internal cross-validation process. External cross validation confusion matrix statistics, ROC/AUC, and misclassification rates can be compared to the internal cross-validation statistics to help assess performance. As with the “Quadratic Discriminant Analysis: External Cross-Validation and Model Development” section, confusion matrices are used to assess model performance where we will be focusing on **Accuracy**, **Precision**, **Sensitivity**, **Specificity**, **Misclassification Rate**, **AUC**, and **Log Loss**. When looking at model performance high values for **Accuracy**, **Precision**, **Sensitivity**, **Specificity**, **AUC** are desirable; and low values for **Misclassification Rate**, and **Log Loss** are desirable. For descriptions of the mentioned terms please refer to the “Quadratic Discriminant Analysis: External Cross-Validation and Model Development” section for more information.

### 5.4 Logistic Regression: Internal vs. External Cross-Validation

The fit statistics between internal and external cross-validation can be compared to assess model performance. Both internal and external cross-validation methods yielded fit statistics very similar in value. **Sensitivity** performed well, while the other fit statistics are at a moderate performance. This suggests that the model is well-fit.

	Internal CV Statistics	External CV Statistics
Sensitivity	0.86998	0.86449
Specificity	0.46260	0.46223
Precision	0.66693	0.66478
Accuracy	0.68786	0.68440
Misclassification Rate	0.31214	0.31214
Logarithmic Loss	0.74483	0.74483
Area Under the Curve	0.70215	0.70452

### 5.5 Logistic Regression: Fitted Model

The stepwise selection model was tuned using k-fold internal cross-validation. The **playoffs** coefficient was added back into the model as the project requirements ask for interpretation of that particular coefficient. Additionally the “Hosmer and Lemeshow goodness of fit (GOF) test” was run on the fitted model to assess if the null hypothesis of: the model does not provide information to predict the outcome variable. Based on the evidence the model that was fit for this project can reject the null, where the model can provide information to predict the outcome variable ( $x = 17.361$ ,  $df = 8$ ,  $p\text{-value} = 0.02656$ ). The stepwise selection model generated is listed below where Logit coefficients for each feature can also be found in the “Fitted Logistic Regression Model” portion of the appendix:

$$\text{shot\_made\_flag} = \text{recId} + \text{action\_type} + \text{game\_event\_id} + \text{loc\_x} + \text{minutes\_remaining} + \text{season} + \text{seconds\_remaining} + \text{shot\_distance} + \text{game\_date} + \text{shot\_id} + \text{attendance} + \text{arena\_temp} + \text{playoffs}$$

Several notable features that this model generated are several of the **action\_types**. The odds ratios of Kobe making a shot are revealed once the coefficients are exponentiated. In the table below we can see that the odds ratios are all very large. This suggests that when Kobe performs one of the listed levels of **action\_type**, his chances of making the shot are very high. The frequency of the **action\_type** can also be ascertained by looking at the standard errors from the model output. Small standard errors for a coefficient suggest that there are a large number of observations with that attribute with a similar outcomes, and a large standard error would suggest that there are fewer observations with a similar outcome. When moving from the base **action\_type** of “Alley Oop Dunk Shot” to “Driving Bank Shot” for every one unit increase in “Driving Bank Shot” the logit value for **shot\_made\_flag** increases by approximately 10.5 units; in terms of odds ratios, compared to using the “Alley Oop Dunk Shot,”

Kobe's chances of making the shot when using the "Driving Bank Shot" is at a factor of approximately 37080. The "Driving Bank Shot" also has a standard error of approximately 535, this value is higher than the other action types suggesting that Kobe uses the "Driving Bank Shot" infrequently but rarely misses this shot type when he chooses to use it. When a frequency table of the `action_type` for the whole dataset is created, the total number of times the "Driving Bank Shot" occurs is 1 - compared to 5 observations for "Jump Shot." The pattern of high odds ratios and high standard errors exist for the top 5 shots; this suggests that Kobe rarely uses these `action_types` but almost always makes the shot if that particular `action_type` is used.

	Coefficient (logit)	Odds Ratio
(Intercept)	13.422543	675050.90
'action_typeDriving Bank shot'	10.520842	37080.34
'action_typeReverse Slam Dunk Shot'	10.298657	29692.71
'action_typeHook Bank Shot'	10.213455	27267.61
'action_typeDriving Floating Bank Jump Shot'	10.196184	26800.73
'action_typeTurnaround Finger Roll Shot'	9.954396	21044.54

## 6 Shot Made Odds and Probabilities Analysis Using Logistic Regression

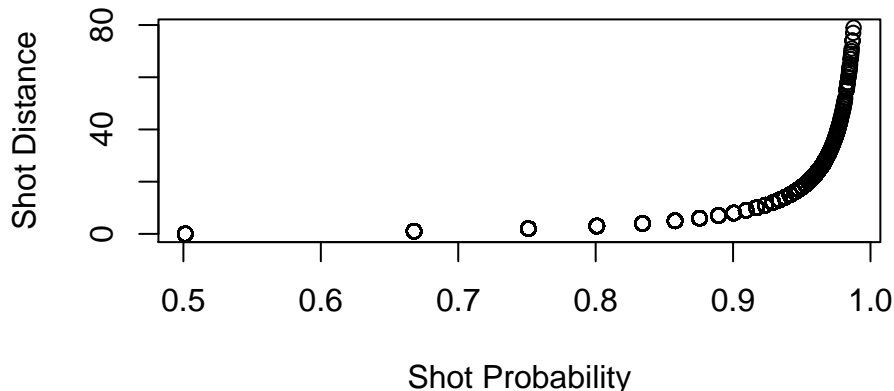
When analyzing odds for basketball many are interested in the spread, or by how much will the favored team win or lose. There are many other odds played in professional sports, but we are focused on the odds Bryant making a shot or not relative to his distance from the hoop during regular season games and the playoffs.

There is sufficient evidence to suggest that `shot_distance` from the hoop is not significant ( $p\text{-value} < 0.053303$ ) with respect to all other covariates. Although the  $p\text{-value}$  for `shot_distance` is barely outside  $\alpha = 0.05$ , it still provides practical information on whether or not Kobe makes his shots.

For every additional unit in `shot_distance` from the hoop, the estimated odds ratio of him making the shot increase by a factor of  $e^{0.00514126} = 1.0053956042$ . This is for every additional unit increase in `shot_distance`, the estimated odds of Bryant making the shot increase by a factor of  $0.0053956042 * (1 - 1.0053956042)$ . A 95% confidence interval for the multiplicative change is  $(-0.0006427041, 0.01143028475)$  holding all other variables constant.

We also were curious if the odds of Bryant making a shot further away from the hoop had a difference during the playoffs. The shot distance is still significant but there is virtually no difference in the odds if the shot was made or not in a regular season game.

### Shot Probability over Distance



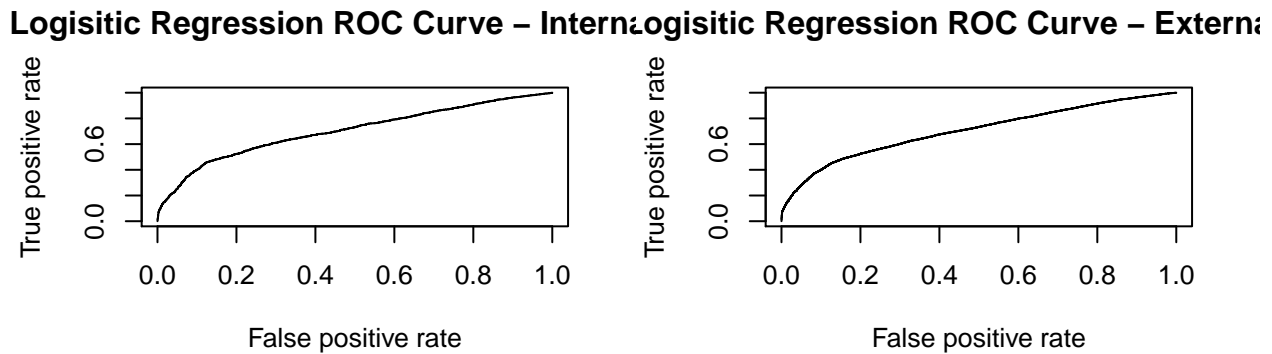
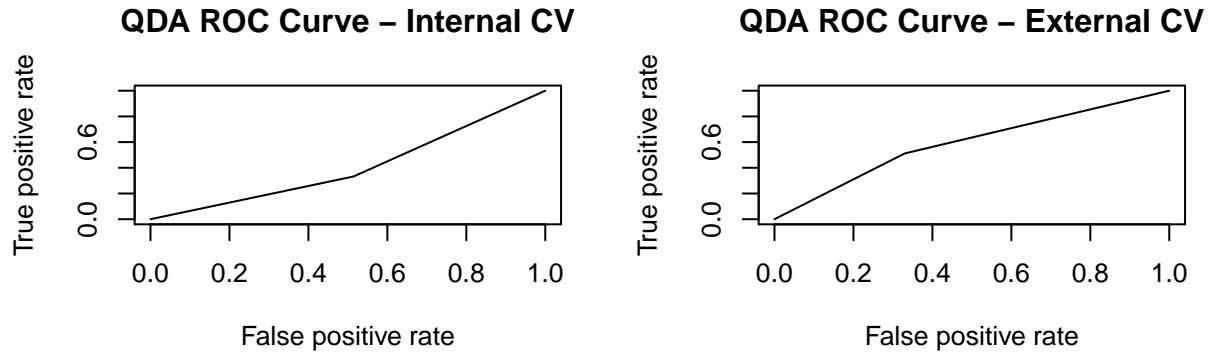
## 7 Conclusion

Analysis of the Kobe Bryant Shot Made data shows Kobe is consistent with his shot made by type, location on the court, distance from the hoop and whether he played a regular season game or a playoff game. Based on the analysis performed, a Logistic Regression Model produced better results than a Quadratic Discriminant Model after taking all variable and model analysis into consideration. These statistics are represented in Kobe Bryant's 20-year successful career.

## 8 Appendix A: Figures and Tables

### 8.0.1 Appendix A contains content related to the paper

### 8.1 Receiver Operator Characteristic (ROC) Curves



### 8.2 Fitted Logistic Regression Model

	Coefficient (logit)	Odds Ratio
(Intercept)	13.4225434	675050.9018004
recId	0.0136368	1.0137302
‘action_typeAlley Oop Layup shot‘	-2.3205263	0.0982219
‘action_typeCutting Layup Shot‘	-1.6513466	0.1917915
‘action_typeDriving Bank shot‘	10.5208422	37080.3407838
‘action_typeDriving Dunk Shot‘	0.5245757	1.6897417
‘action_typeDriving Finger Roll Layup Shot‘	-1.1149420	0.3279343
‘action_typeDriving Finger Roll Shot‘	-1.8462292	0.1578312
‘action_typeDriving Floating Bank Jump Shot‘	10.1961843	26800.7262598
‘action_typeDriving Floating Jump Shot‘	-3.4111350	0.0330037
‘action_typeDriving Hook Shot‘	-2.7840394	0.0617884
‘action_typeDriving Jump shot‘	-4.0797951	0.0169109
‘action_typeDriving Layup Shot‘	-2.2080771	0.1099118
‘action_typeDriving Reverse Layup Shot‘	-1.9126942	0.1476820
‘action_typeDriving Slam Dunk Shot‘	0.1321462	1.1412751

(continued)

	Coefficient (logit)	Odds Ratio
‘action_typeDunk Shot‘	-2.2672636	0.1035953
‘action_typeFadeaway Bank shot‘	-0.7136965	0.4898302
‘action_typeFadeaway Jump Shot‘	-3.0468033	0.0475106
‘action_typeFinger Roll Layup Shot‘	-2.3265883	0.0976283
‘action_typeFinger Roll Shot‘	-3.1228802	0.0440302
‘action_typeFloating Jump shot‘	-2.2637304	0.1039619
‘action_typeFollow Up Dunk Shot‘	-1.0661569	0.3443293
‘action_typeHook Bank Shot‘	10.2134549	27267.6128042
‘action_typeHook Shot‘	-3.8692331	0.0208744
‘action_typeJump Bank Shot‘	-2.0944755	0.1231348
‘action_typeJump Hook Shot‘	-2.0526109	0.1283992
‘action_typeJump Shot‘	-4.1200063	0.0162444
‘action_typeLayup Shot‘	-3.8516668	0.0212443
‘action_typePullup Bank shot‘	-2.9546883	0.0520949
‘action_typePullup Jump shot‘	-2.2243604	0.1081366
‘action_typePutback Dunk Shot‘	-2.7899501	0.0614243
‘action_typePutback Layup Shot‘	-2.6921832	0.0677329
‘action_typeReverse Dunk Shot‘	0.0771918	1.0802493
‘action_typeReverse Layup Shot‘	-2.8497340	0.0578597
‘action_typeReverse Slam Dunk Shot‘	10.2986570	29692.7136467
‘action_typeRunning Bank shot‘	-1.5553045	0.2111251
‘action_typeRunning Dunk Shot‘	-1.5628760	0.2095326
‘action_typeRunning Finger Roll Layup Shot‘	-2.8121239	0.0600773
‘action_typeRunning Finger Roll Shot‘	-16.9093575	0.0000000
‘action_typeRunning Hook Shot‘	-1.8963490	0.1501157
‘action_typeRunning Jump Shot‘	-2.2497152	0.1054292
‘action_typeRunning Layup Shot‘	-2.8535071	0.0576418
‘action_typeRunning Pull-Up Jump Shot‘	-16.5617800	0.0000001
‘action_typeRunning Reverse Layup Shot‘	-4.5986325	0.0100656
‘action_typeSlam Dunk Shot‘	0.8467543	2.3320653
‘action_typeStep Back Jump shot‘	-2.5585379	0.0774179
‘action_typeTip Shot‘	-4.1454923	0.0158356
‘action_typeTurnaround Bank shot‘	-1.9936634	0.1361956
‘action_typeTurnaround Fadeaway shot‘	-3.0459246	0.0475523
‘action_typeTurnaround Finger Roll Shot‘	9.9543965	21044.5414148
‘action_typeTurnaround Hook Shot‘	-3.4711028	0.0310827
‘action_typeTurnaround Jump Shot‘	-2.9466004	0.0525179
game_event_id	-0.0002398	0.9997602
loc_x	0.0001497	1.0001498
minutes_remaining	0.0122044	1.0122792
‘season1997-98‘	-0.0578076	0.9438315
‘season1998-99‘	0.3796978	1.4618428
‘season1999-00‘	0.5389535	1.7142120
‘season2000-01‘	0.7859562	2.1945043
‘season2001-02‘	0.8929332	2.4422828
‘season2002-03‘	0.9829625	2.6723614
‘season2003-04‘	0.9639388	2.6220037
‘season2004-05‘	1.4038822	4.0709735
‘season2005-06‘	1.6091042	4.9983318

(continued)

	Coefficient (logit)	Odds Ratio
‘season2006-07‘	1.7103964	5.5311536
‘season2007-08‘	1.7773374	5.9140887
‘season2008-09‘	2.0093797	7.4586889
‘season2009-10‘	1.9749527	7.2062790
‘season2010-11‘	2.1296754	8.4121361
‘season2011-12‘	2.2033561	9.0553535
‘season2012-13‘	2.4169619	11.2117448
‘season2013-14‘	1.6501936	5.2079882
‘season2014-15‘	2.4556923	11.6544994
‘season2015-16‘	2.5550064	12.8713824
seconds_remaining	0.0025664	1.0025697
shot_distance	0.0053921	1.0054066
game_date	-0.0004337	0.9995664
shot_id	-0.0136309	0.9864616
attendance	0.0001672	1.0001673
arena_temp	0.0374076	1.0381160
playoffs	-0.0373788	0.9633111

### 8.3 Action Type Frequency Table

Var1	Freq
Alley Oop Dunk Shot	76
Alley Oop Layup shot	59
Cutting Layup Shot	6
Driving Bank shot	1
Driving Dunk Shot	196
Driving Finger Roll Layup Shot	47
Driving Finger Roll Shot	52
Driving Floating Bank Jump Shot	1
Driving Floating Jump Shot	3
Driving Hook Shot	13
Driving Jump shot	19
Driving Layup Shot	1335
Driving Reverse Layup Shot	67
Driving Slam Dunk Shot	38
Dunk Shot	176
Fadeaway Bank shot	22
Fadeaway Jump Shot	693
Finger Roll Layup Shot	21
Finger Roll Shot	23
Floating Jump shot	75
Follow Up Dunk Shot	10
Hook Bank Shot	5
Hook Shot	61
Jump Bank Shot	223
Jump Hook Shot	16
Jump Shot	12712
Layup Shot	1734

(continued)

Var1	Freq
Pullup Bank shot	10
Pullup Jump shot	318
Putback Dunk Shot	3
Putback Layup Shot	9
Reverse Dunk Shot	52
Reverse Layup Shot	276
Reverse Slam Dunk Shot	15
Running Bank shot	35
Running Dunk Shot	14
Running Finger Roll Layup Shot	5
Running Finger Roll Shot	4
Running Hook Shot	28
Running Jump Shot	620
Running Layup Shot	42
Running Pull-Up Jump Shot	1
Running Reverse Layup Shot	6
Slam Dunk Shot	264
Step Back Jump shot	93
Tip Shot	121
Turnaround Bank shot	50
Turnaround Fadeaway shot	299
Turnaround Finger Roll Shot	1
Turnaround Hook Shot	8
Turnaround Jump Shot	739

## 8.4 Variance Inflation Factor

Variable	Variance Inflation Factor
shot_id	39939311.92
recId	39937863.50
shot_distance	3.59
loc_y	3.15
shot_type	1.90
playoffs	1.76
game_event_id	1.61
attendance	1.39
avgnoisedb	1.37
game_date	1.13
minutes_remaining	1.09
shot_made_flag	1.05
arena_temp	1.02
loc_x	1.01
seconds_remaining	1.01

## 9 Appendix B: Source Code

### 9.0.1 *Appendix A contains source code created for this project*

### 9.1 Loading Libraries, Importing Data

```
library(pacman)
p_load(rrcov, MASS, dplyr, purrr, ggplot2, Hmisc, pcaPP, knitr, kableExtra, caret,
cluster, robustbase, ROCR, Metrics, bookdown, ResourceSelection, usdm)

# Reading Data
df <- read.csv("./modelingKobeData.csv", header=T, sep="," , strip.white=T,
stringsAsFactors = F, na.strings=c(""))
df.preds <- read.csv("./predictionKobeData.csv", header=T, sep="," ,
strip.white=T, stringsAsFactors = F, na.strings=c(""))
```

### 9.2 Check for Missing Values

```
apply(df, function(cnt) sum(length(which(is.na(cnt)))))
```

### 9.3 Basketball Shot Location Map

```
shotsTaken <- data.frame(df$loc_x, df$loc_y, df$shot_distance)
colnames(shotsTaken) <- c("loc_x", "loc_y", "shot_distance")

ggplot(shotsTaken, aes(x=loc_x, y=loc_y)) +
  geom_point(aes(colour = df$shot_type))
```

### 9.4 Shot Location Imputing, toInt Conversion

```
df[which(df$loc_y > 300), "shot_type"] <- "3PT Field Goal"
df.preds[which(df.preds$loc_y > 300), "shot_type"] <- "3PT Field Goal"

# Convert the points to integer values
df$shot_type <- ifelse(df$shot_type=="2PT Field Goal", 2, 3)
df.preds$shot_type <- ifelse(df.preds$shot_type=="2PT Field Goal", 2, 3)
```

### 9.5 Variable Eliminations, First-Pass

```
df <- df %>% mutate_if(is.integer, as.numeric) %>%
mutate_if(is.character, as.factor) %>% data.frame()
df <- df %>%
  subset(select=-c(
```



```

team_id, # dropping since this is a uniform distribution of data
team_name, # dropping since this is a uniform distribution of data. Also
collinear with team_id
combined_shot_type, # dropping this in favor of combined_shot_type
shot_zone_area, # this is ambiguous and less descriptive than geospatial data
shot_zone_range,
shot_zone_basic,
matchup # removing in favor of opponent; Kobe only played for LAL so that
will never change
    )
)

df.preds <- df.preds %>%
mutate_if(is.integer, as.numeric) %>% mutate_if(is.character, as.factor) %>%
data.frame()
df.preds <- df.preds %>%
  subset(select=-c(
    team_id, # dropping since this is a uniform distribution of data
    team_name, # dropping since this is a uniform distribution of data. Also
    collinear with team_id
    combined_shot_type, # dropping this in favor of combined_shot_type
    shot_zone_area, # this is ambiguous and less descriptive than geospatial data
    shot_zone_range,
    shot_zone_basic,
    matchup # removing in favor of opponent; Kobe only played for LAL so that
    will never change
  )
)

# create numeric dataframe for correlation plot
df.numeric <- df %>% keep(is.numeric)
df.numeric.preds <- df.preds %>% keep(is.numeric)

```

## 9.6 Correlation Heat Map

```

corr.plot <- corrplot::corrplot(cor(df.numeric %>% subset(select=-c(
shot_made_flag)))
    , title = "Correlation among Predictor Variables"
    , type = "lower"
    , tl.pos = "ld"
    , method = "square"
    , tl.cex = 0.65
    , tl.col = 'red'
    , order = "alphabet"
    , diag = F
    , mar=c(0,0,5,0)
    , bg="ivory1"
    ,tl.srt=.05
)

```

## 9.7 Variance Inflation Factor Table

```
vifFrame <- vif(df.numeric)
vifFrame <- vifFrame[order(-vifFrame$VIF),]
vifFrame.Top7 <- data.frame(head(vifFrame, 7))
```

## 9.8 Variable Eliminations, Second Pass

```
df <- df %>% subset(select=-c(
  lat, # dropping lat because it is collinear with loc_y and shot_distance
  lon, # dropping lon because it is collinear with loc_x and shot_distance
  period, # dropping period in favor of game event id - game event id is
          # more descriptive and continuous
  game_id # dropping playoffs for game_id;
          # game ID can capture playoffs seasonally
))

df.preds <- df.preds %>% subset(select=-c(
  lat,
  lon,
  period,
  game_id
))

df.numeric <- df %>% keep(is.numeric) %>% mutate_if(is.integer, as.numeric)
df.numeric.preds <- df.preds %>% keep(is.numeric) %>% mutate_if(
  is.integer, as.numeric)
```

## 9.9 Correlation Matrix

```
flattenCorrMatrix <- function(cormatrix, pmatrix) {
  ut <- upper.tri(cormatrix)
  data.frame(
    row = rownames(cormatrix)[row(cormatrix)[ut]],
    column = rownames(cormatrix)[col(cormatrix)[ut]],
    cor = (cormatrix)[ut],
    p = pmatrix[ut]
  )
}

options(scipen=999)
options(max.print=100000)

correlation.matrix <- Hmisc::rcorr(as.matrix(df.numeric), type="pearson")
corDF <- data.frame(flattenCorrMatrix(correlation.matrix$r,
correlation.matrix$p))

corDF.ordered <- data.frame(corDF[order(-corDF$cor),])
```

```

cordDF.ordered.Top10 <- data.frame(head(cordDF.ordered, 10), row.names = F)

colnames(cordDF.ordered.Top10) = c("Correlation Predictor Variable",
"Correlation Response Variable", "Correlation", "p-Value")

cordDF.ordered.Top10$Correlation <- round(as.numeric(as.character(
cordDF.ordered.Top10$Correlation)), digits=5)

cordDF.ordered.Top10$p-Value <- ifelse(as.numeric(as.character(
cordDF.ordered.Top10$p-Value)) < 0.0001, "p < 0.0001", as.numeric(as.character(
cordDF.ordered.Top10$p-Value)))

```

## 9.10 QDA Bartlett Approximation

```

dfTrain.numeric <- df.numeric[which(!is.na(df.numeric$shot_made_flag)),]
prediction.Data.numeric <- df.numeric[which(is.na(df.numeric$shot_made_flag)),]

dfTrain.numeric$shot_made_flag <- as.factor(dfTrain.numeric$shot_made_flag)
dfTrain.numeric$shot_made_flag <- ifelse(dfTrain.numeric$shot_made_flag=="1",
"made", "not_made")
dfTrain.numeric <- dfTrain.numeric %>% mutate_if(is.integer, as.numeric) %>%
mutate_if(is.character, as.factor) %>% data.frame()

Bartlett_ChiSq <- rrcov::Wilks.test(shot_made_flag ~ ., data=dfTrain.numeric,
method = "c", approximation = "Bartlett")

# Wilk's Lambda produces significant p-value in Bartlett's test so we need to
use a Quadratic Discriminant Analysis instead of Linear
format(round(Bartlett_ChiSq$p.value, 2), nsmall=4)

# Wilks' Lambda plus degrees of freedom used in Bartlett's chi-squared test
WilksDegreesofFreedom <- rbind(as.numeric(paste0(Bartlett_ChiSq$parameter, sep = " ")))

# p-value from Bartlett's test
Bartlett_ChiSq$p.value
Bartletts_p <- format(round(as.numeric(Bartlett_ChiSq$p.value), 2), nsmall=4)

# Because Bartlett's p-value is less than 0.0001 (indicated above), updating to
shorter form:
Bartletts_p = ifelse(Bartlett_ChiSq$p.value < 0.0001, "p < 0.0001",
Bartlett_ChiSq$p.value)
#Bartletts_p <- "p < 0.0001"

dfBartlett <- data.frame(WilksDegreesofFreedom, Bartlett_ChiSq$wilks, Bartletts_p)
colnames(dfBartlett) <- c("Chi-Square Statistic", "Degrees of Freedom",
"Wilks' Lambda", "p-Value")

dfBartlett$`Chi-Square Statistic` <- round(as.numeric(as.character(
dfBartlett$`Chi-Square Statistic`)), digits=5)
dfBartlett$`Wilks' Lambda` <- round(as.numeric(as.character(
dfBartlett$`Wilks' Lambda`)), digits=5)

bartlettsTest <- data.frame(rbind(dfBartlett$`Chi-Square Statistic`,

```

```
dfBartlett$`Degrees of Freedom`,dfBartlett$`Wilks' Lambda`,Bartletts_p))
rownames(bartlettsTest) <- c("Chi Square Statistic","Degrees of Freedom",
"Wilks' Lambda","p-Value")
colnames(bartlettsTest) <- "Statistics"
```

## 9.11 Partitioning Training, Testing Data

```
dfTrain <- df[which(!is.na(df$shot_made_flag)),]
prediction.Data <- df[which(is.na(df$shot_made_flag)),]

### Full Data train/test split for Logistic
test_sample_size <- floor(0.75 * nrow(dfTrain))
set.seed(123)
train_ind <- sample(seq_len(nrow(dfTrain)), size = test_sample_size)
subDF.Train <- dfTrain[train_ind, ] #75% training
subDF.Test <- dfTrain[-train_ind, ] # 25% testing

#### Numeric Data train/test split for QDA
test_sample_size <- floor(0.75 * nrow(dfTrain.numeric))
set.seed(123)
train_ind <- sample(seq_len(nrow(dfTrain.numeric)), size = test_sample_size)
subDF.Train.numeric <- dfTrain.numeric[train_ind, ] #75% training
subDF.Test.numeric <- dfTrain.numeric[-train_ind, ] # 25% testing
```

## 9.12 a Priori Analysis

```
# MASS package used for qda()
#df.numeric <- df.numeric[order(df.numeric$shot_made_flag),]
kobe.qda <- qda(shot_made_flag ~ ., CV=T, data=dfTrain.numeric)

data.frame(mean(kobe.qda$posterior[,1]), mean(kobe.qda$posterior[,2]))

shot_made_flag_Posterior <- rbind("0", "1")
# a Priori noted, used for testing, but default for proportions used
proportion_Posterior <- rbind(mean(kobe.qda$posterior[,1]), mean(
kobe.qda$posterior[,2]))
priori <- data.frame(shot_made_flag_Posterior, proportion_Posterior)
```

## 9.13 Quadratic Discriminant Analysis: Internal CV and Model Development

```
subDF.Train.numeric$shot_made_flag <- as.factor(
subDF.Train.numeric$shot_made_flag)
#subDF.Train.numeric$shot_made_flag <- ifelse(
subDF.Train.numeric$shot_made_flag=="1", "made", "not_made")
subDF.Train.numeric <- subDF.Train.numeric %>% mutate_if(
is.integer, as.numeric) %>% mutate_if(is.character, as.factor) %>%
data.frame()

# k=25 folds, repeat random folding for internal cross-validation 5 times:
```

```

train.Control <- caret::trainControl(method = "repeatedcv",
                                     number = 25,
                                     repeats = 5,
                                     #summaryFunction = twoClassSummary,
                                     summaryFunction = mnLogLoss,
                                     classProbs = T) # classProbs=T to get mnLogLoss (
                                     also for twoClassSummary)

# build the model using the 75% partitioned from the internal dataset (
the set with all shot_made_flag response results):
qda.filtered <- train(shot_made_flag ~ .
                     , data = subDF.Train.numeric
                     , method = "qda"
                     , trControl=train.Control
                     , preProcess = c("center", "scale", "spatialSign")
                     #, preProcess = "spatialSign"
                     , metric = "logLoss"
                     )

# test the model on the 25% partitions from the internal dataset:
internal_cv.predicted.qda <- suppressWarnings(predict(qda.filtered, newdata =
subDF.Test.numeric))

# build a confusion matrix for internal cross-validation to see performance:
confusion_matrix_results.internal<- confusionMatrix(table(
internal_cv.predicted.qda, subDF.Test.numeric$shot_made_flag))

```

## 9.14 Quadratic Discriminant Analysis: External CV

```

external_cv.predicted.qda <- suppressWarnings(predict(qda.filtered, newdata =
dfTrain.numeric))
confusion_matrix_results.external <- confusionMatrix(table(
external_cv.predicted.qda, dfTrain.numeric$shot_made_flag))

```

## 9.15 Quadratic Discriminant Analysis Confusion Matrix Table

```

### Internal Cross-Validation Confusion Matrix
confusionFrame.internal <- data.frame(rbind(round(
Sensitivity.confusion.internal, digits=5),round(
Specificity.confusion.internal, digits=5),round(
Precision.confusion.internal, digits=5),round(
Accuracy.confusion.internal, digits=5),round(
misclassification.QDA.internal, digits=5),round(
logLoss.quadratic, digits=5),round(AUC.internal, digits=5)))

rownames(confusionFrame.internal) <- c("Sensitivity",
"Specificity","Precision","Accuracy",
"Misclassification Rate","Logarithmic Loss","Area Under the Curve")

### External Cross-Validation Confusion Matrix
confusionFrame.external <- data.frame(rbind(round(

```

```

Sensitivity.confusion.external, digits=5),round(
Specificity.confusion.external, digits=5),round(
Precision.confusion.external, digits=5),round(
Accuracy.confusion.external, digits=5),round(
misclassification.QDA.external, digits=5),round(
logLoss.quadratic, digits=5),round(AUC.external, digits=5)))

rownames(confusionFrame.external) <- c("Sensitivity",
"Specificity","Precision","Accuracy",
"Misclassification Rate","Logarithmic Loss","Area Under the Curve")

confusionFrame <- data.frame(confusionFrame.internal, confusionFrame.external)
colnames(confusionFrame) <- c("Internal CV Statistics",
"External CV Statistics")

```

## 9.16 Predictions from Quadratic Discriminant Analysis

```

# Apply the developed model to the external data that needs predictions:
pred.qda.filtered <- suppressWarnings(predict(qda.filtered,
newdata = df.numeric.preds))

df.preds$shot_made_flag <- pred.qda.filtered

#write.csv(df.preds, "Predicted_Results.csv", row.names = F)

```

## 9.17 Logistic Model Development Using Ordinary Least Squares

```

### Initial Models
model.forward.Start <- glm(shot_made_flag~1, family=binomial(link='logit'),
data = dfTrain)
model.Allvar <- glm(shot_made_flag ~ ., family=binomial(link='logit'),
data = dfTrain)

#### Forward Selection
model.Forward <- stepAIC(model.forward.Start, direction = "forward",
trace = F, scope = formula(model.Allvar))

summary(model.Forward)
model.Forward$anova
##### Forward Selection Model Suggestion
forward.glm <- glm(shot_made_flag ~ action_type + attendance + arena_temp +
  game_event_id + season + seconds_remaining + minutes_remaining +
  loc_y + game_date + loc_x + playoffs
  , family=binomial(link='logit')
  , data=dfTrain)

summary(forward.glm)
#####

# Backward Elimination
model.Backward <- stepAIC(model.Allvar, direction = "backward",

```

```

trace = F, scope = formula(model.forward.Start))
summary(model.Backward)
model.Backward$anova
##### Backward Elimination Model Suggestion
back.glm <- glm(shot_made_flag ~ recId + action_type + game_event_id +
  loc_x + minutes_remaining + season + seconds_remaining +
  shot_distance + game_date + shot_id + attendance + arena_temp
  , family=binomial(link='logit')
  , data=dfTrain)

summary(back.glm)
back.glm$aic
#####

# Stepwise Regression
model.Stepwise <- stepAIC(model.Allvar, direction = "both", trace = F)
summary(model.Stepwise)
model.Stepwise$anova
##### Stepwise Regression Model Suggestion
step.glm <- glm(shot_made_flag ~ recId + action_type + game_event_id +
  loc_x + minutes_remaining + season + seconds_remaining +
  shot_distance + game_date + shot_id + attendance + arena_temp + playoffs
  , family=binomial(link='logit')
  , data=dfTrain)

summary(step.glm)
step.glm$aic

# Model Selection Statistics AIC and RMSE
model_stats <- data.frame(cbind(rbind("Forwards","Backwards",
"Stepwise"),rbind(round(forward.glm$aic, digits=2),round(
back.glm$aic, digits=2),round(step.glm$aic, digits=2)),rbind(
round(forward.glm$deviance, digits=2),round(back.glm$deviance,
digits=2),round(step.glm$deviance, digits=2))))
colnames(model_stats) <- c("Selection Type", "AIC", "Residual Deviance")

kable(model_stats,format="latex", booktabs = T) %>%
  kable_styling(latex_options="striped", position = "center")

```

## 9.18 Logistic Regression: Internal CV

```

#Cross Validation
#Recategorize rare shot types into more common, but similar shot types
dfTrain$action_type = as.character(dfTrain$action_type)
dfTrain = dfTrain %>%
  mutate(action_type = if_else(action_type == "Running Tip Shot",
  'Tip Shot', action_type))%>%
  mutate(action_type = if_else(action_type == "Tip Layup Shot",
  'Tip Shot', action_type)) %>%
  mutate(action_type = if_else(action_type == "Putback Slam Dunk Shot",
  "Slam Dunk Shot",action_type)) %>%
  mutate(action_type = if_else(action_type == "Running Slam Dunk Shot",

```

```

  "Slam Dunk Shot",action_type))
dfTrain$action_type = as.factor(dfTrain$action_type)

df.preds$action_type = as.character(df.preds$action_type)
df.preds = df.preds %>%
  mutate(action_type = if_else(action_type == "Turnaround Finger Roll Shot",
    'Finger Roll Shot', action_type)) %>%
  mutate(action_type = if_else(action_type == "Putback Slam Dunk Shot",
    'Slam Dunk Shot',action_type)) %>%
  mutate(action_type = if_else(action_type == "Running Slam Dunk Shot",
    'Slam Dunk Shot',action_type))
df.preds$action_type = as.factor(df.preds$action_type)

#K-fold CV
set.seed(100)
Train <- createDataPartition(dfTrain$shot_made_flag, p=0.75, list=FALSE)
training <- dfTrain[ Train, ]
testing <- dfTrain[ -Train, ]

#train for specificity??? option
ctrl <- trainControl(method = "repeatedcv",
  number = 25,
  repeats = 5,
  classProbs = T)

#combined shot type used instead of action type - test set has action types
that are not in the training set

mod_fit <- train(shot_made_flag ~ recId + action_type + game_event_id +
  loc_x + minutes_remaining + season + seconds_remaining +
  shot_distance + game_date + shot_id + attendance + arena_temp + playoffs,
  data=training, method="glm",
  family="binomial",
  trControl = ctrl,
  tuneLength = 5,
  metric = "logLoss")

#####
#Internal Model Performance Metrics
#confusion matrix https://rpubs.com/dvorakt/255527
pred = predict(mod_fit, newdata=testing)
cf = confusionMatrix(table(data=as.numeric(pred>0.5), testing$shot_made_flag))
misclassificationRateInternal = (cf$table[2,1]+cf$table[1,2]) / sum(cf$table)

#ROC/AUC
# Compute AUC for predicting Class with the model
pred <- prediction(pred, testing$shot_made_flag)
perf <- performance(pred, measure = "tpr", x.measure = "fpr")

aucInternal <- performance(pred, measure = "auc")
aucInternal <- aucInternal@y.values[[1]]

#LOG LOSS AND PREDICTION FOR TRAINING DATA
testing$prob = predict(mod_fit, newdata=testing)
loglossTraining = testing %>%
  mutate(logloss = testing$shot_made_flag * log(1-testing$prob) +

```



```

(1-testing$shot_made_flag)*log(1-testing$prob))

#Will generate log loss value
loglossValue = -1/5174 * sum(loglossTraining$logloss)

#####
#External Model performance Metrics
ex.pred = predict(mod_fit, newdata=dfTrain)
ex.cf = confusionMatrix(table(data=as.numeric(ex.pred>0.5),
dfTrain$shot_made_flag))
misclassificationRateEx = (cf$table[2,1]+cf$table[1,2]) / sum(cf$table)

#ROC/AUC
# Compute AUC for predicting Class with the model
expredROC <- prediction(ex.pred, dfTrain$shot_made_flag)
experf <- performance(expredROC, measure = "tpr", x.measure = "fpr")

aucEx <- performance(expredROC, measure = "auc")
aucEx <- aucEx@y.values[[1]]

#LOG LOSS AND PREDICTION FOR External CV
dfTrain$prob = predict(mod_fit, newdata=dfTrain)
loglossTraining = dfTrain %>%
  mutate(logloss = dfTrain$shot_made_flag * log(1-dfTrain$prob) +
    (1-dfTrain$shot_made_flag)*log(1-dfTrain$prob))

#Will generate log loss value
loglossValue = -1/20697 * sum(loglossTraining$logloss)

```

## 9.19 Logistic Regression: Internal vs. External CV

```

##### Logistic Internal Cross-Validation Metrics
lr.SpecSense.confusion.internal <- data.frame(cf$byClass)
lr.AccuracyP.confusion.internal <- data.frame(cf$overall)

lr.Accuracy.confusion.internal <- AccuracyP.confusion.internal[1,]
lr.Sensitivity.confusion.internal <- SpecSense.confusion.internal[1,]
lr.Specificity.confusion.internal <- SpecSense.confusion.internal[2,]
lr.Precision.confusion.internal <- SpecSense.confusion.internal[5,]

##### External Cross-Validation Metrics
lr.SpecSense.confusion.external <- data.frame(ex.cf$byClass)
lr.AccuracyP.confusion.external <- data.frame(ex.cf$overall)

lr.Accuracy.confusion.external <- lr.AccuracyP.confusion.external[1,]
lr.Sensitivity.confusion.external <- lr.SpecSense.confusion.external[1,]
lr.Specificity.confusion.external <- lr.SpecSense.confusion.external[2,]
lr.Precision.confusion.external <- lr.SpecSense.confusion.external[5,]

### Internal Cross-Validation Confusion Matrix
lr.confusionFrame.internal <- data.frame(rbind(round(
lr.Sensitivity.confusion.internal, digits=5),round(
lr.Specificity.confusion.internal, digits=5),round(

```

```

lr.Precision.confusion.internal, digits=5),round(
lr.Accuracy.confusion.internal, digits=5),round(
misclassificationRateInternal, digits=5),round(
loglossValue, digits=5),round(aucInternal, digits=5)))

rownames(lr.confusionFrame.internal) <- c("Sensitivity",
"Specificity","Precision","Accuracy",
"Misclassification Rate","Logarithmic Loss","Area Under the Curve")

### External Cross-Validation Confusion Matrix
lr.confusionFrame.external <- data.frame(rbind(round(
lr.Sensitivity.confusion.external, digits=5),round(
lr.Specificty.confusion.external, digits=5),round(
lr.Precision.confusion.external, digits=5),round(
lr.Accuracy.confusion.external, digits=5),round(
misclassificationRateEx, digits=5),round(loglossValue,
digits=5),round(AUC.external, digits=5)))

rownames(lr.confusionFrame.external) <- c("Sensitivity",
"Specificity","Precision","Accuracy",
"Misclassification Rate","Logarithmic Loss","Area Under the Curve")

lr.confusionFrame <- data.frame(lr.confusionFrame.internal,
lr.confusionFrame.external)
colnames(lr.confusionFrame) <- c("Internal CV Statistics",
"External CV Statistics")

```

## 9.20 Logistic Regression: Fitted Model

```

###HL Test
hltest = hoslem.test(training$shot_made_flag, fitted(mod_fit), g=10)

#Odds Ratios
topshots = as.data.frame(coef(mod_fit$finalMode))
names(topshots)[1] = "Coefficient (logit)"
topshots$`Odds Ratio` = exp(topshots$`Coefficient (logit)`)
topshots = topshots[order(topshots$`Odds Ratio`,decreasing = TRUE),]

#Predictions
df.preds$prob = predict(mod_fit, newdata=df.preds)
logistic_prediction = df.preds %>%select(recId, prob)
#write.csv(logistic_prediction,"logistic_prediction.csv")

```