

Will Kobe Bryant Make His Next Shot: Quadratic Discriminant Analysis and Logistic Regression using R

Paul Adams
Reannan McDaniel
Jeff Nguyen

Master of Science in Data Science, Southern Methodist University, USA

Contents

1	Abstract:	3
2	Introduction	3
3	Exploratory Data Analysis	3
3.1	Outlier Check	3
3.2	Variable Elimination	3
3.3	Addressing Multicollinearity: Correlation Plot for Visual Data Exploration	3
3.4	Post-Correlation Heat Map Variable Elimination	4
3.5	Addressing Multicollinearity: Correlation Matrix for Numerical Analysis	4
4	Quadratic Discriminant Analysis	4
4.1	Bartlett's Test	5
5	Quadratic Discriminant Analysis: Internal Cross-Valdiation and Model Development	5
6	Quadratic Discriminant Analysis: External Cross-Valdiation and Model Development	5
7	Quadratic Discriminant Analysis: Internal vs. External Cross-Validation	7
8	Logistic Model Development using Ordinary Least Squares	8
8.1	Logstic Regression: Model Selection	8
8.2	Logistic Regression: Model Selection Fit-Statistics	9
9	Appendix A: Images and Tables	12

1 Abstract:

1.0.1 *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

[illegible]

3 Exploratory Data Analysis

3.1 Outlier Check

3.1.1 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 Gloal.

3.2 Variable Elimination

3.2.1 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 is may not be significant, it is not conducive to model quality.

3.3 Addressing Multicollinearity: Correlation Plot for Visual Data Exploration

3.3.1 To address multicollinearity among quantitative predictor variables, a correlation heat map was created for visual inspection of correlation. Please see Appendix A to view this correlation heat map.

3.4 Post-Correlation Heat Map Variable Elimination

- 3.4.1 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 make 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.

3.5 Addressing Multicollinearity: Correlation Matrix for Numerical Analysis

- 3.5.1 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.

4 Quadratic Discriminant Analysis

- 4.0.1 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`, `avgnoisedb`, controlling collinearity by eliminating a member of each collinear pair prior to model development.
- 4.0.2 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.

4.1 Bartlett's Test

- 4.1.1 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.
- 4.1.2 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. Therefore, modeling data is broken into a 75% training / 25% testing data split for internal cross-validation. The objective of internal cross-validation is to develop a model using 75% of the data and test it on the remaining 25% in order to assess model fit statistics. Typically, following internal cross-validation, external cross-validation is performed.

5 Quadratic Discriminant Analysis: Internal Cross-Validation and Model Development

- 5.0.1 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.

6 Quadratic Discriminant Analysis: External Cross-Validation and Model Development

- 6.0.1 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).
- 6.0.2 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.
- 6.0.3 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.
- 6.0.4 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.
- 6.0.5 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.

7 Quadratic Discriminant Analysis: Internal vs. External Cross-Validation

- 7.0.1 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.

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

8 Logistic Model Development using Ordinary Least Squares

8.0.1 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.

8.1 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:

8.1.1 Forward Selection Model:

8.1.2 $shot_{made}flag\ action_{type} + attendance + shot_{zone}range + arena_{temp} + game_{event}id + season + seconds_{remaining} + shot_{zone}basic + loc_y + minutes_{remaining} + loc_x + game_{date}$

8.1.3 Backward Elimination Model:

8.1.4 $shot_{made}flag\ recId + action_{type} + game_{event}id + loc_x + loc_y + minutes_{remaining} + playoffs + season + seconds_{remaining} + shot_{type} + shot_{zone}basic + shot_{zone}range + game_{date} + shot_{id} + attendance + arena_{temp} + avgnoisedb$

8.1.5 Stepwise Regression Model:

8.1.6 $shot_{made}flag\ recId + action_{type} + game_{event}id + loc_x + loc_y + minutes_{remaining} + playoffs + season + seconds_{remaining} + shot_{type} + shot_{zone}basic + shot_{zone}range + game_{date} + shot_{id} + attendance + arena_{temp} + avgnoisedb$

8.2 Logistic Regression: Model Selection Fit-Statistics

8.2.1 Based on the fit-statistics generated from each model selection method, the backwards and stepwise models are identical in fit-statistics. They both had the lowest AIC at 25095.19, and the lowest residual deviance at 24905.19. The stepwise model will be used for Cross Validation.

Selection.Type	AIC	Residual.Deviance
Forwards	25091.70	24909.70
Backwards	25095.19	24905.19
Stepwise	25095.19	24905.19

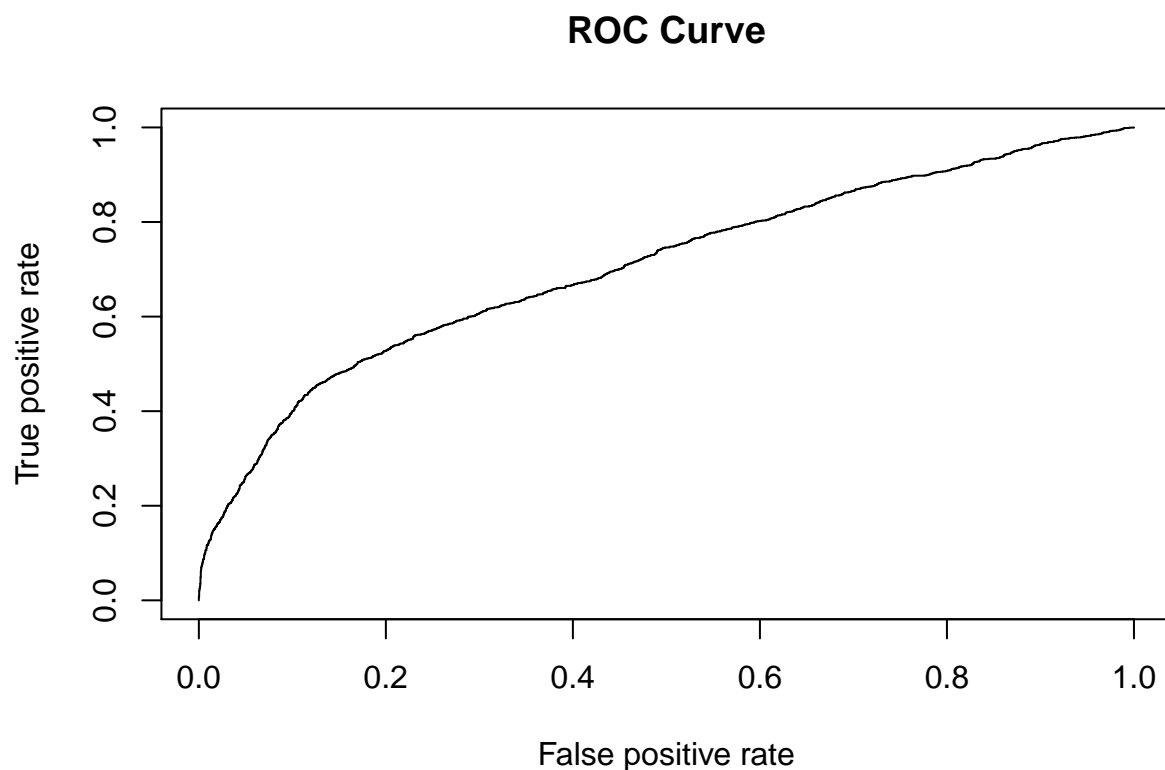
##Logstic Regression: Internal Cross Validation and Model Development ### Based on the AIC and RMSE statistics generated by model selection, the stepwise model has the lowest AIC SOMEVALUE at and RMSE at SOMEVALUE. The stepwise model will then undergo internal cross validation. The data is split into 75% for training and 25% for testing.

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.

##Logstic Regression: External Cross Validation and Model Development

8.2.2 As with the “Quadratic Discriminant Analysis: External Cross-Valdiation 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, Log Loss. When looking at model performance high values for:Accuracy, Precision, Sensitivity, Specificity,AUC are desireable; and low values for Misclassification Rate, Log Loss are desireable. For descriptions of the mentioned terms please refer to the “Quadratic Discriminant Analysis: External Cross-Valdiation and Model Development” section for more information.

```
## [1] 0.3150367
```



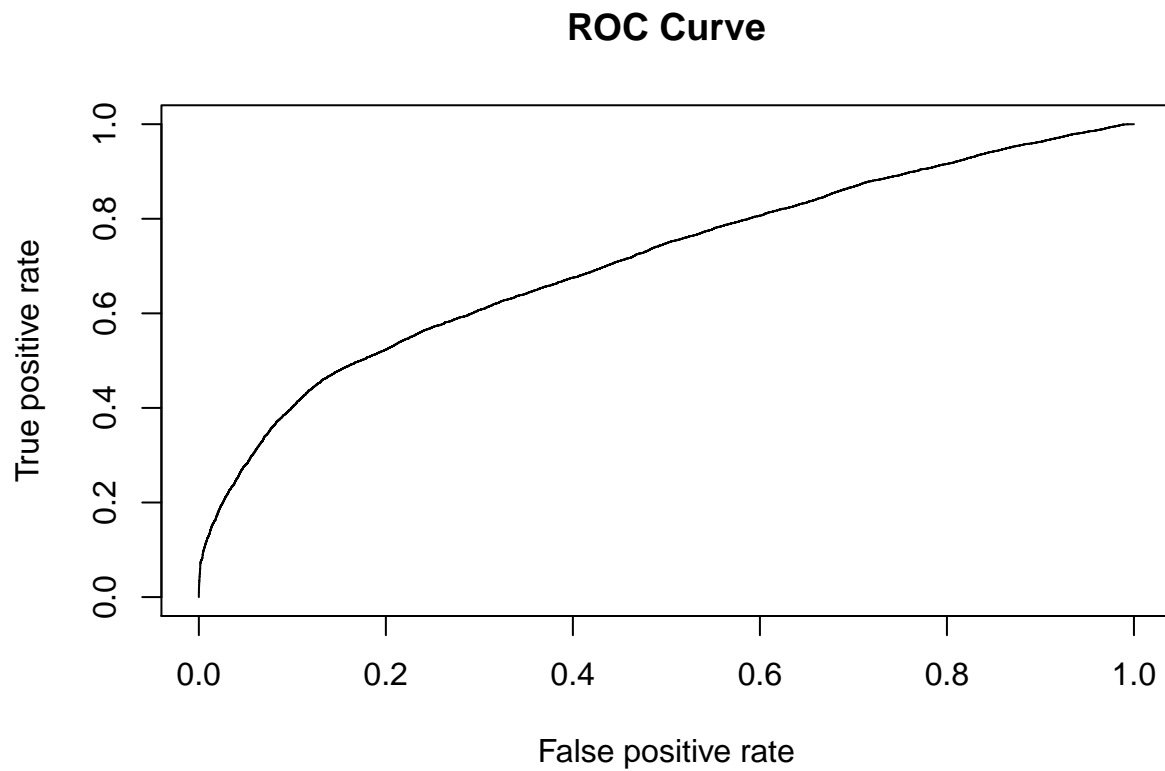
```
## [1] 0.7055852
```

```
## [1] 0.754176
```

```
## [1] 0.3150367
```

Table 1: (#tab:Confusion Matrix Tables Logistic)Internal Cross-Validation Confusion Matrix

Internal CV Statistics	
Sensitivity	0.51431
Specificity	0.66761
Precision	0.56104
Accuracy	0.59826
Misclassification Rate	0.40174
Logarithmic Loss	0.70127
Area Under the Curve	0.40904



```
## [1] 0.7093543
```

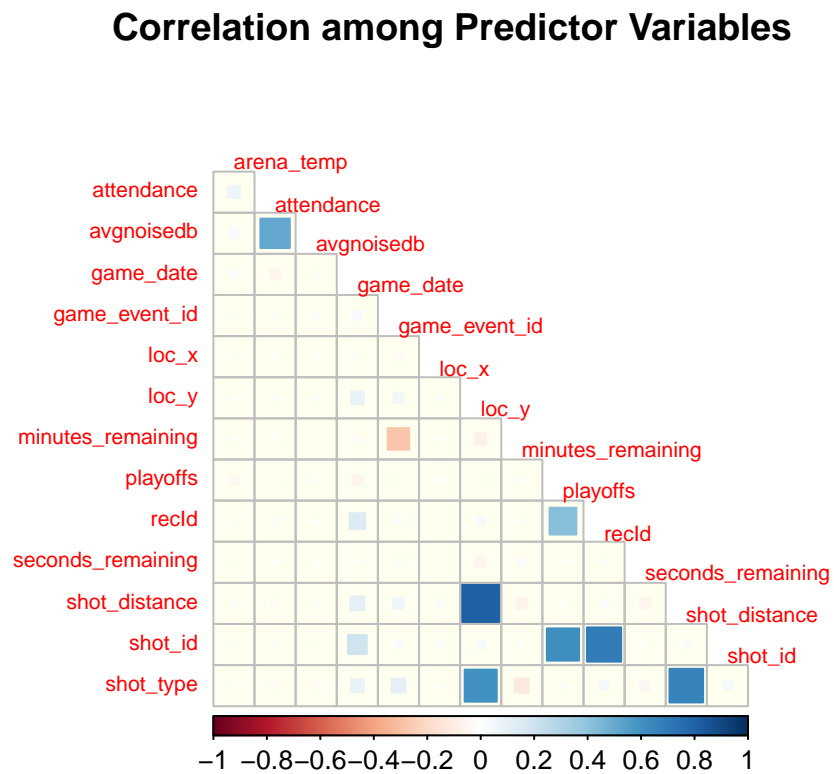
```
## [1] 0.7580061
```

Table 2: (#tab:Confusion Matrix Tables Logistic)External Cross-Validation Confusion Matrix

	External CV Statistics
Sensitivity	0.86318
Specificity	0.46406
Precision	0.66521
Accuracy	0.68450
Misclassification Rate	0.31504
Logarithmic Loss	0.70127
Area Under the Curve	0.59090

9 Appendix A: Images and Tables

9.0.1 Correlation Heat Map



Correlation Predictor Variable	Correlation Response Variable	Correlation	p-Value
arena_temp	arena_temp	0.51092	$p < 0.0001$
attendance	arena_temp	0.51092	$p < 0.0001$
game_date	arena_temp	0.51092	$p < 0.0001$
game_event_id	arena_temp	0.51092	$p < 0.0001$
loc_x	arena_temp	0.51092	$p < 0.0001$
loc_y	arena_temp	0.51092	$p < 0.0001$
minutes_remaining	arena_temp	0.51092	$p < 0.0001$
playoffs	arena_temp	0.51092	$p < 0.0001$
recId	arena_temp	0.51092	$p < 0.0001$
seconds_remaining	arena_temp	0.51092	$p < 0.0001$

Table 3: (#tab:Bartletts Results for Appendix)Bartlett Test's Wilks' Lambda Approximation

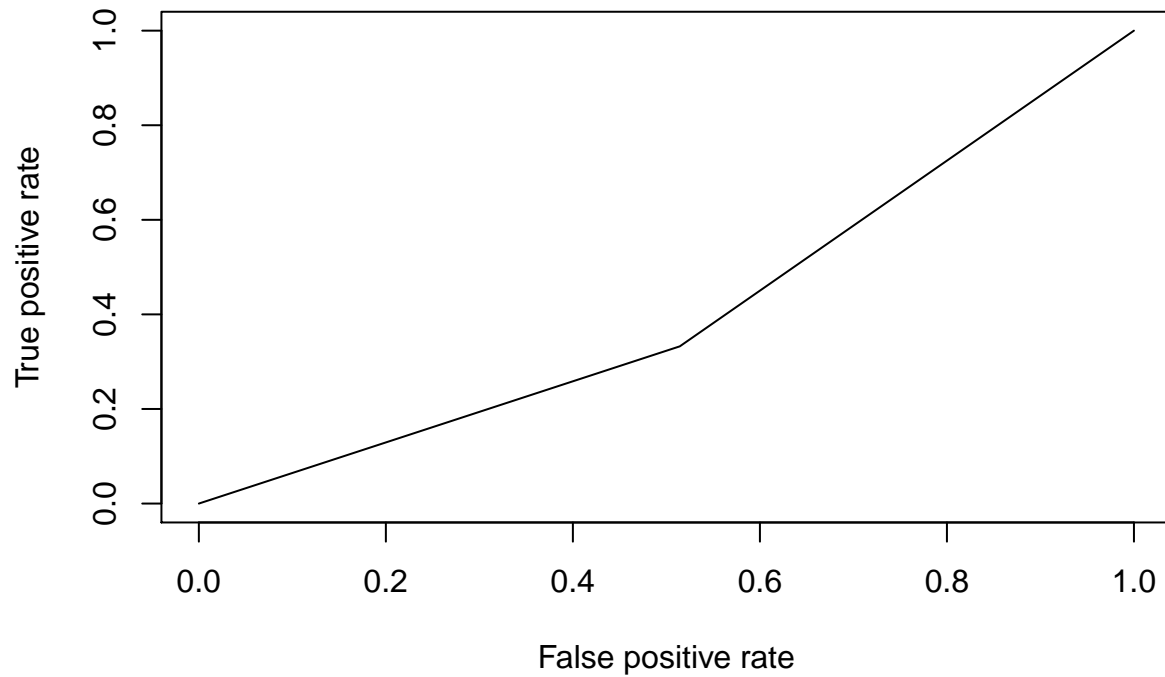
	Metric Output
Chi Square Statistic	1037.24251
Degrees of Freedom	14
Wilks' Lambda	0.9511
p-Value	$p < 0.0001$

9.0.2 Correlation Matrix - Top 10 Collinear Terms

9.0.3 Bartlett Test's Wilks' Lambda Approximation

9.0.4 ROC Curves

ROC Curve – Internal CV



ROC Curve – External CV

