

# Homework #3: Crime Prediction

Data 621 Business Analytics and Data Mining

*Aadi Kalloo, Nathan Lim, Asher Meyers, Daniel Smilowitz, Logan Thomson*

*Due July 3, 2016*

## Contents

<b>Data Exploration</b>	<b>2</b>
<b>Data Preparation</b>	<b>8</b>
<b>Model Creation</b>	<b>9</b>
Model 1: Bayesian Information Criterion . . . . .	9
Model 2: Mallow's $C_p$ . . . . .	10
Model 3: Transformed BIC . . . . .	11
Model 4: Significant BIC . . . . .	12
Model 5: Best GLM . . . . .	13
Model 6: Transformed Best GLM . . . . .	14
Model 7: Full Model . . . . .	15
<b>Model Selection and Prediction</b>	<b>17</b>
Model Comparison . . . . .	17
10-fold Cross Validation . . . . .	17
<b>Appendix A</b>	<b>18</b>
<b>Appendix B – Index-wise Results from Predictive Model</b>	<b>19</b>
<b>Appendix C – R Code</b>	<b>20</b>

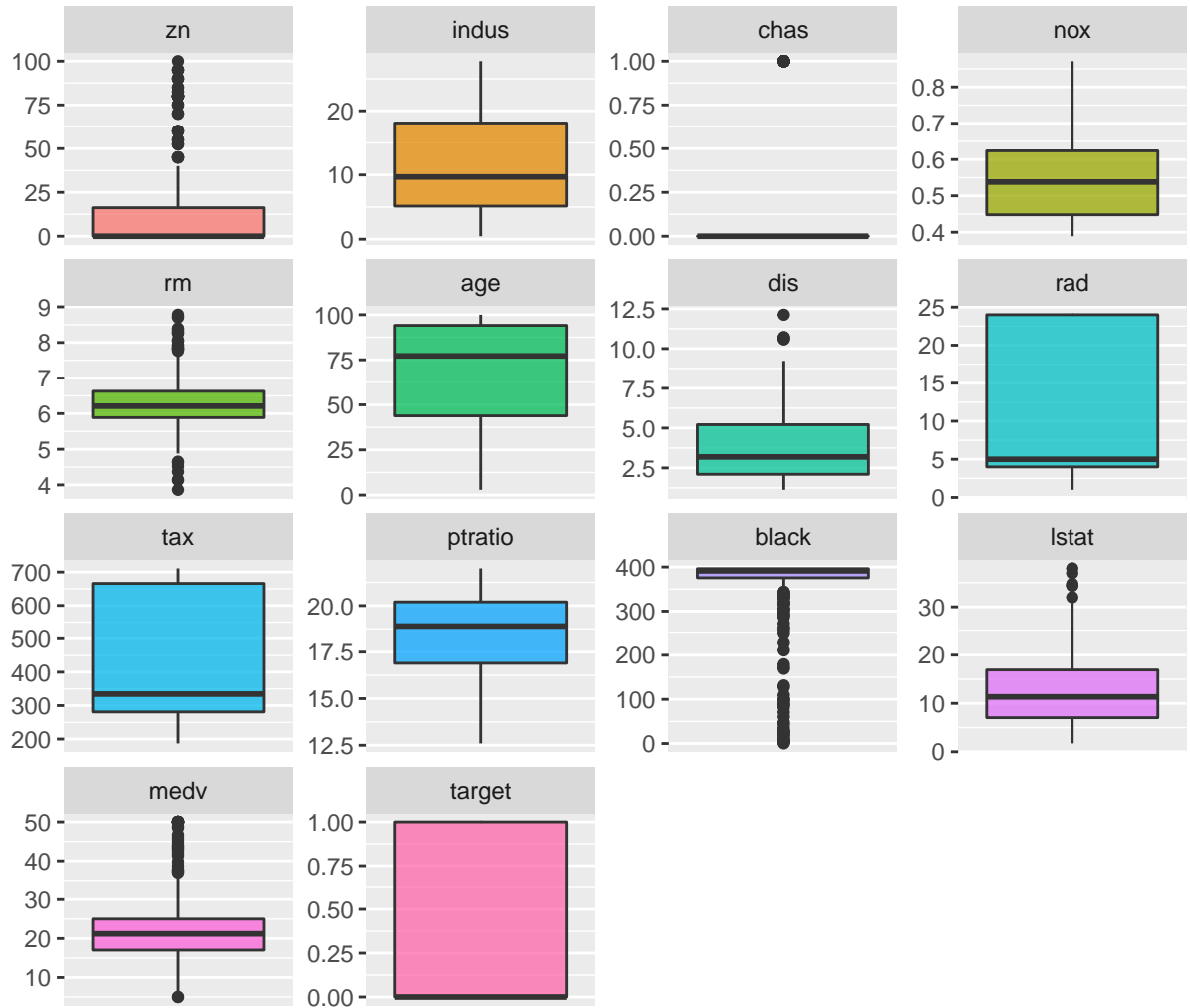
## Data Exploration

The supplied data set contains various pieces of information about neighborhoods in a major U.S. city. The data set is made up of 466 rows, with 14 variables; 13 predictors and 1 response variable (**target**), which is a binary categorical variable. A summary of each variable is below:

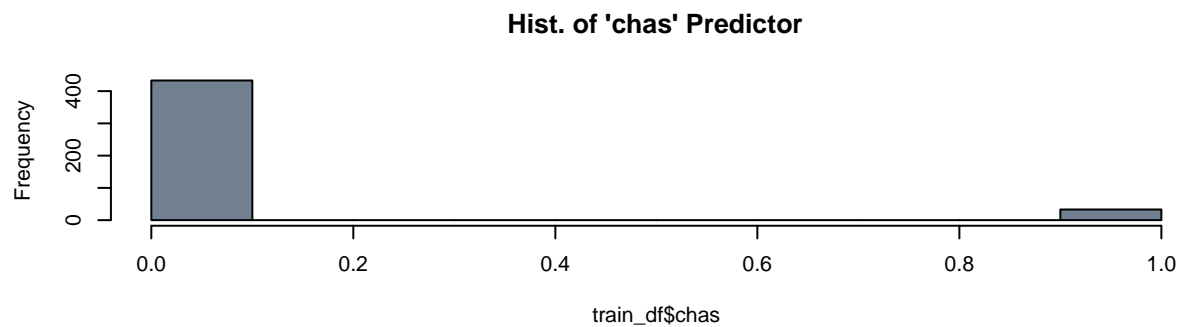
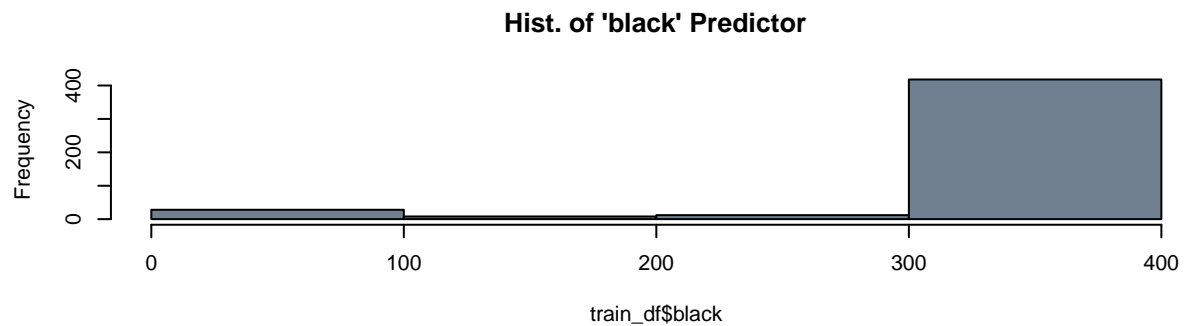
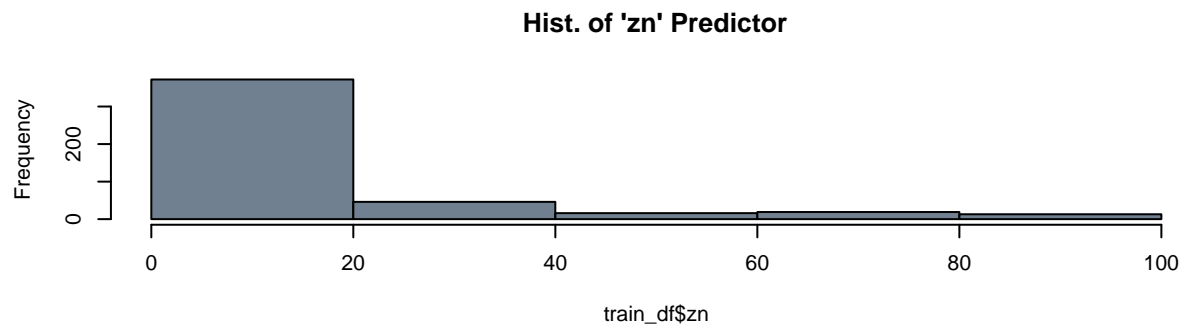
	MEAN	MEDIAN	IQR	SKEWNESS	CORRELATION TO TARGET
<b>zn</b>	11.58	0	16.25	2.18	-0.43
<b>indus</b>	11.11	9.69	12.96	0.29	0.6
<b>chas</b>	0.07	0	0	3.34	0.08
<b>nox</b>	0.55	0.54	0.18	0.75	0.73
<b>rm</b>	6.29	6.21	0.74	0.48	-0.15
<b>age</b>	68.37	77.15	50.22	-0.58	0.63
<b>dis</b>	3.8	3.19	3.11	1	-0.62
<b>rad</b>	9.53	5	20	1.01	0.63
<b>tax</b>	409.5	334.5	385	0.66	0.61
<b>ptratio</b>	18.4	18.9	3.3	-0.75	0.25
<b>black</b>	357.1	391.3	20.63	-2.92	-0.35
<b>lstat</b>	12.63	11.35	9.89	0.91	0.47
<b>medv</b>	22.59	21.2	7.98	1.08	-0.27
<b>target</b>	0.49	0	1	0.03	1

None of our predictors have missing values, so there will be no need to impute any variables. However, some of the predictors have a good amount of skew to the distribution of cases, and may need to be transformed to find a best fit model.

## Distribution of Predictor and Target Variables

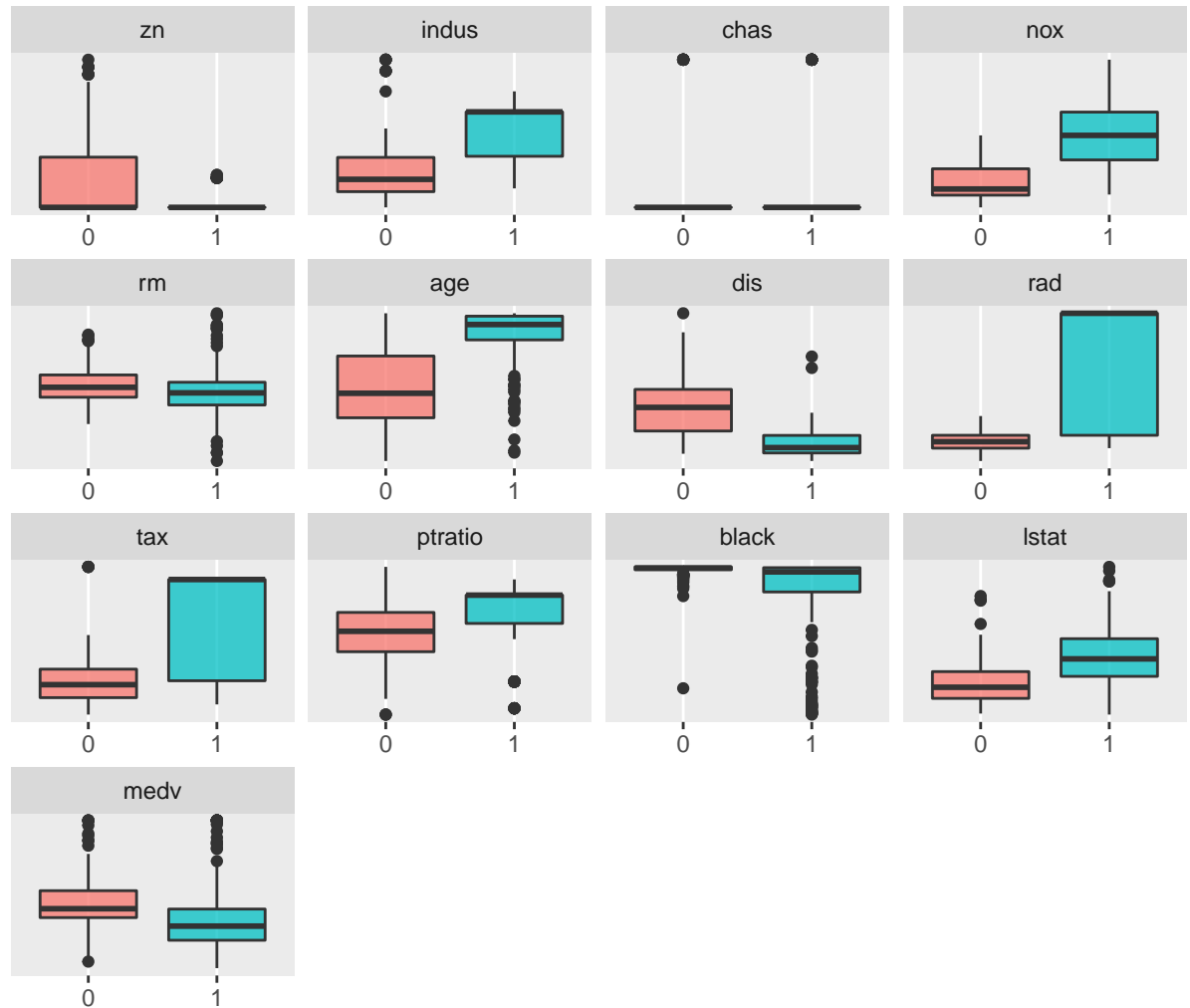


Looking at the box plots above, we can see that the **zn** and **black** predictors are strongly skewed. **chas**, indicating if a suburb borders the Charles River, is also strongly skewed, but like our response variable, it is a two-level categorical predictor and the skew shows because 92% of the values are 0.



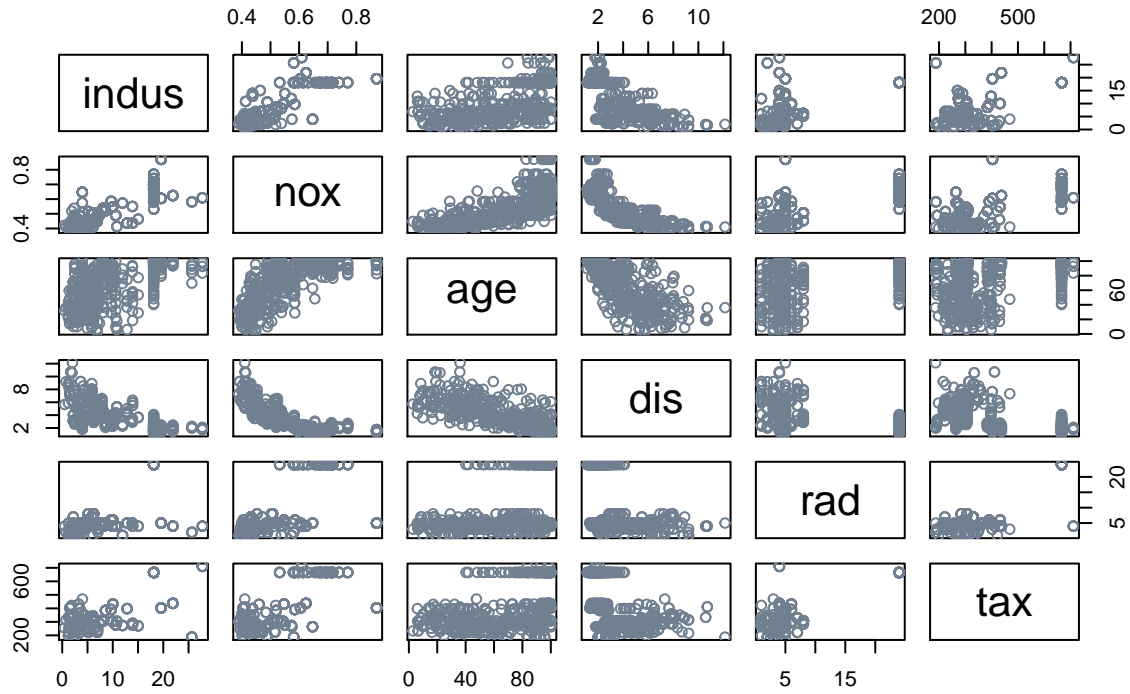
Since our response variable `target` is a two-level factor, we can take a look at a plot of each predictor, subset by `target` and see the relationship between the predictor and our response variable. Right away, the difference in means for some of the variables is quite apparent. These variables also have a higher correlation to `target`, and may be some of the more significant predictors in our models.

## Distribution of Predictors by Target

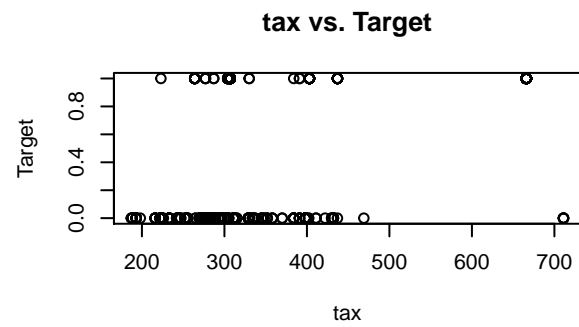
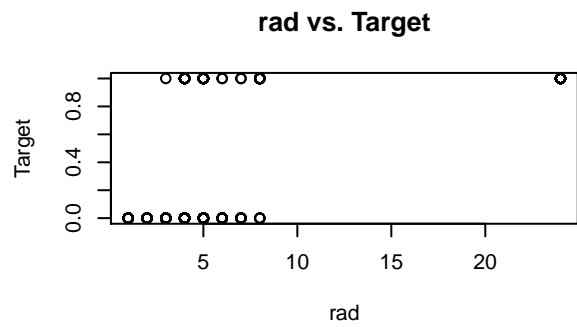
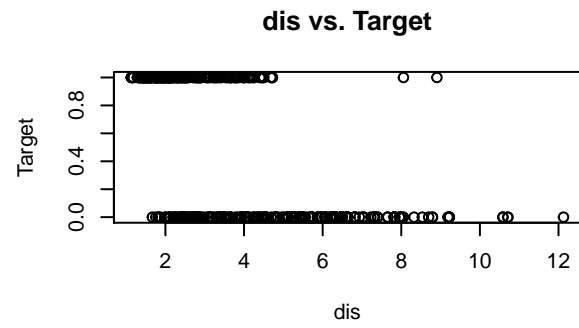
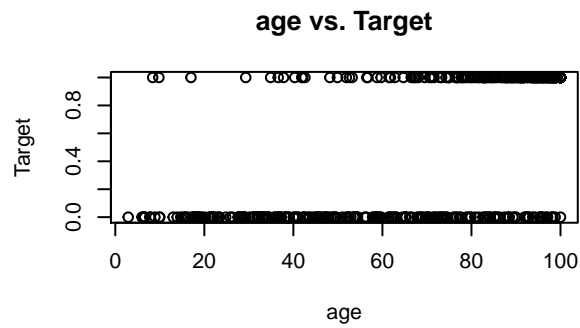
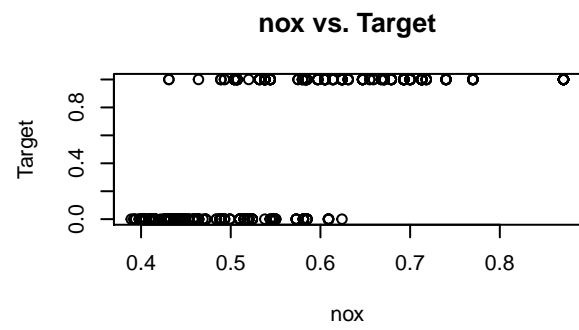
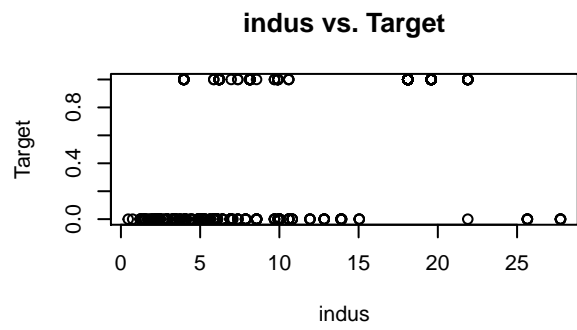


The predictors with the highest correlation to **target** are **indus**, **nox**, **age**, **dis**, **rad**, and **tax**, and all are related to each other in some form when we consider industrial zoning and the effects it has on a surrounding area. In fact, looking at these variables against each other, we see that they are highly correlated as well, evidence that multicollinearity exists.

## Highly Correlated Predictors



Scatterplots are shown below for those variables with higher correlations and larger differences in means when subset by the **target** response variable. Plots for all predictors against the **target** variable are included in the appendix.



## Data Preparation

As stated in Part 1, there are no missing values in the training data set, so we do not need to impute any of the variables, or exclude any cases in the data set. The predictor **black**, which measures the proportion of blacks by town appears to already have been transformed (adding a constant value, squaring, then multiplying by 1000). Undoing this transformation simply results in data that is even more skewed against the response variable, so the transformation was done to reduce this.

Combining other predictors or using ratios did not seem to create any new meaningful predictors that would improve our models, however due to the left- and right-skewness of some of our predictors, transformations like log, square root, or  $1/X$  reduce the amount of skew and increase the correlation of our predictors to the response.

In the table below, we can see the effect of each of the transformations on the correlation to the **target** variable. Many of the transformations do not help in improving the correlation, but the log transformation of **nox** and **rad** may fit our models better, as well as the square root of **indus**.

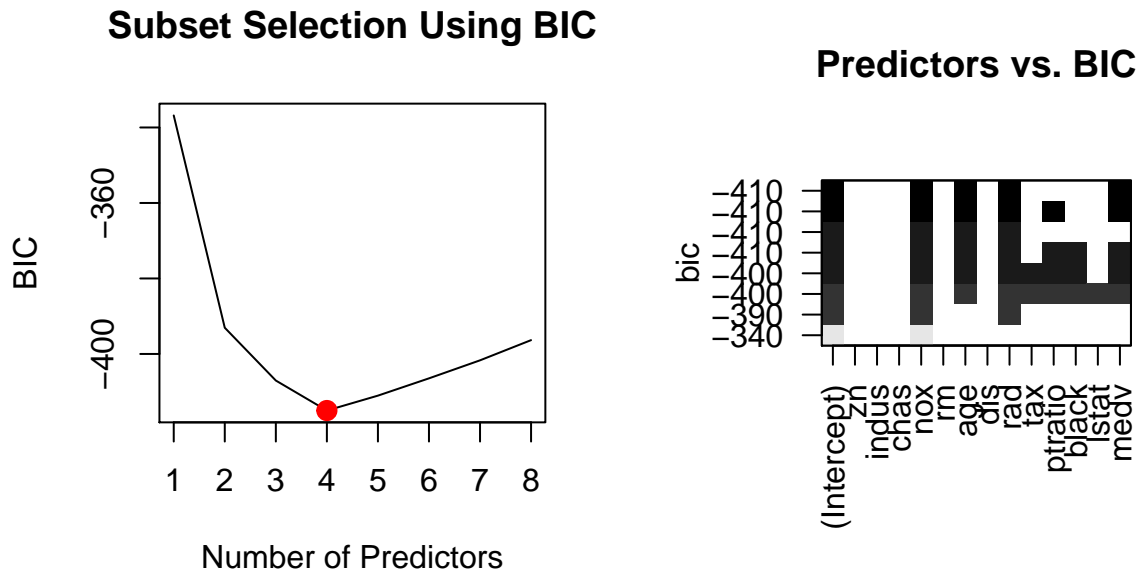
	No Transform	LOG	SQ.ROOT	1/X
<b>zn</b>	-0.43	NA	-0.47	NA
<b>indus</b>	0.6	0.6	0.62	-0.44
<b>chas</b>	0.08	NA	0.08	NA
<b>nox</b>	0.73	0.75	0.74	-0.75
<b>rm</b>	-0.15	-0.17	-0.16	0.19
<b>age</b>	0.63	0.54	0.6	-0.33
<b>dis</b>	-0.62	-0.66	-0.64	0.64
<b>rad</b>	0.63	0.64	0.64	-0.49
<b>tax</b>	0.61	0.61	0.61	-0.6
<b>prratio</b>	0.25	0.22	0.24	-0.19
<b>black</b>	-0.35	-0.27	-0.32	0.08
<b>lstat</b>	0.47	0.44	0.47	-0.32
<b>medv</b>	-0.27	-0.36	-0.32	0.38
<b>target</b>	1	NA	1	NA



# Model Creation

## Model 1: Bayesian Information Criterion

The first model created utilizes the Bayesian Information Criterion (BIC) to determine the number of predictors to use and which predictors should be used.



The left plot above shows that the BIC is minimized using 4 predictors. The plot on the right shows that the 4 predictors with the lowest BIC are **nox**, **age**, **rad**, and **medv**. As such, a model is created using these predictors; this model is presented below.

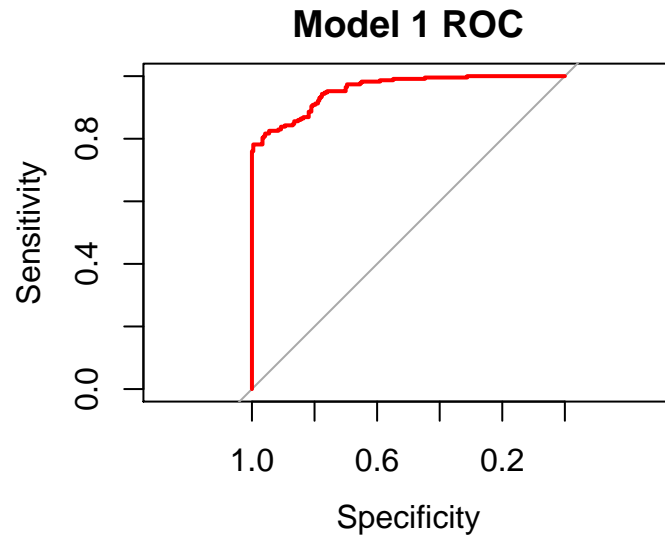
	Estimate	Std. Error	z value	Pr(> z )
<b>nox</b>	23.62	3.936	6.003	1.942e-09
<b>age</b>	0.01824	0.009172	1.989	0.04673
<b>rad</b>	0.4528	0.1093	4.144	3.413e-05
<b>medv</b>	0.04481	0.02319	1.932	0.05338
<b>(Intercept)</b>	-17.63	2.168	-8.131	4.246e-16

(Dispersion parameter for binomial family taken to be 1 )

Null deviance:	645.9 on 465 degrees of freedom
Residual deviance:	232.8 on 461 degrees of freedom

The coefficients in this model indicate that nitrogen oxides concentration has the strongest, as well as the most statistically significant, effect on on the target variable. Age, highway access, and home value all have far weaker effects on the target. All of the estimated coefficients are statistically significant at the  $\alpha = 0.5$  level except **medv** – the p-value for this coefficient is 0.0534.

The receiver operating characteristic (ROC) curve and confusion matrix returned by this model are shown below:

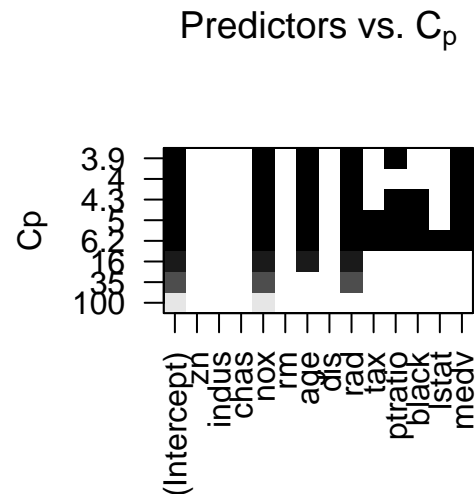
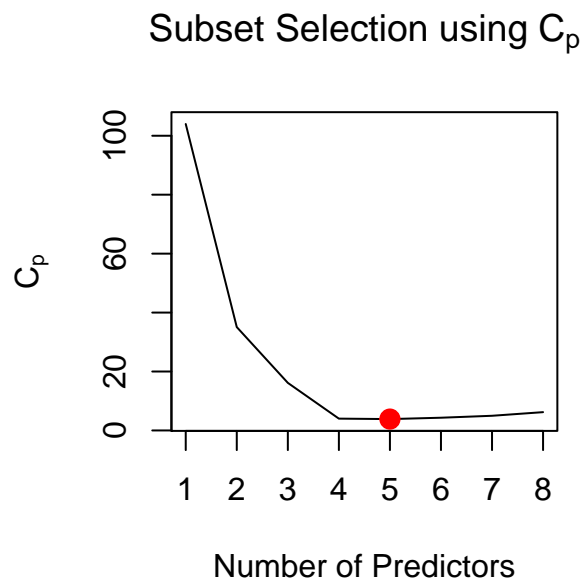


	0	1
0	214	23
1	37	192

This model provides an area under the curve of 0.9570 and an accuracy of 0.8691.

### Model 2: Mallows's $C_p$

The second model created utilizes Mallows's  $C_p$  to determine the number of predictors to use and which predictors should be used.



The left plot above shows that  $C_p$  is minimized using 5 predictors. The plot on the right shows that the 5 predictors with the lowest  $C_p$  are `nox`, `age`, `rad`, `ptratio`, and `medv`. As such, a model is created using these predictors; this model is presented below.

	Estimate	Std. Error	z value	Pr(> z )
<b>nox</b>	25.33	4.084	6.203	5.53e-10
<b>age</b>	0.0194	0.009308	2.085	0.03711
<b>rad</b>	0.5126	0.1148	4.464	8.027e-06
<b>ptratio</b>	0.2742	0.09874	2.777	0.005486
<b>medv</b>	0.08544	0.02798	3.054	0.002259
<b>(Intercept)</b>	-24.94	3.683	-6.77	1.289e-11

(Dispersion parameter for binomial family taken to be 1 )

Null deviance:	645.9 on 465 degrees of freedom
Residual deviance:	224.7 on 460 degrees of freedom

As in model 1, the coefficients in this model indicate that nitrogen oxides concentration has the strongest, as well as the most statistically significant, effect on on the target variable. This may be due to the fact that a single part per 10 million in concentration merits a unit increase in this variable. Each of the estimated coefficients are statistically significant.

The ROC curve and confusion matrix returned by this model are shown below:

This model provides an area under the curve of 0.9605 and an accuracy of 0.8691.

### Model 3: Transformed BIC

Model 3 is created using the same BIC selection from Model 1, with modification from transformation. Based on the distributions of the **nox** and **rad** variables, log transformations of these variables are used. The model using these transformed variables, the model is presented below.

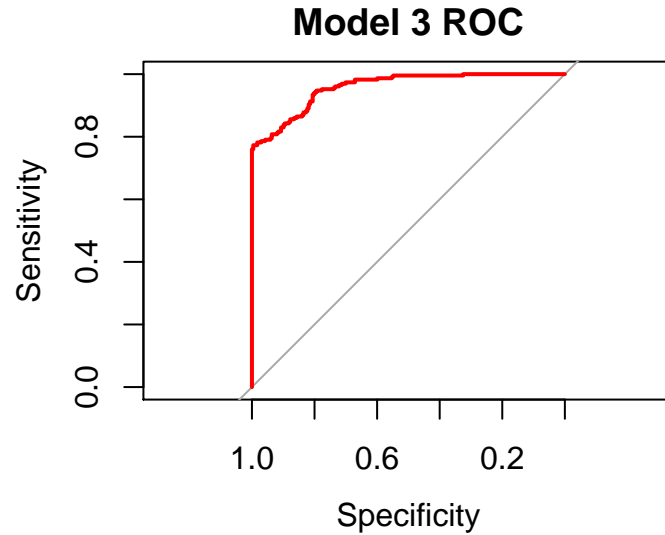
	Estimate	Std. Error	z value	Pr(> z )
<b>log(nox)</b>	13	2.097	6.199	5.681e-10
<b>age</b>	0.01696	0.009147	1.855	0.06366
<b>log(rad)</b>	2.269	0.4494	5.05	4.428e-07
<b>medv</b>	0.0489	0.02354	2.078	0.03774
<b>(Intercept)</b>	1.96	1.925	1.018	0.3086

(Dispersion parameter for binomial family taken to be 1 )

Null deviance:	645.9 on 465 degrees of freedom
Residual deviance:	231.7 on 461 degrees of freedom

In model 3, the coefficient associated with nitrogen oxide concentration is once again the largest in magnitude, although it has decreased in magnitude. The estimated coefficient for home value is now statistically significant; however, the coefficient for age no longer is.

The ROC curve and confusion matrix returned by this model are shown below:



	0	1
0	213	24
1	37	192

This model has an area under the curve of 0.9584 and an accuracy of 0.8691.

### Model 4: Significant BIC

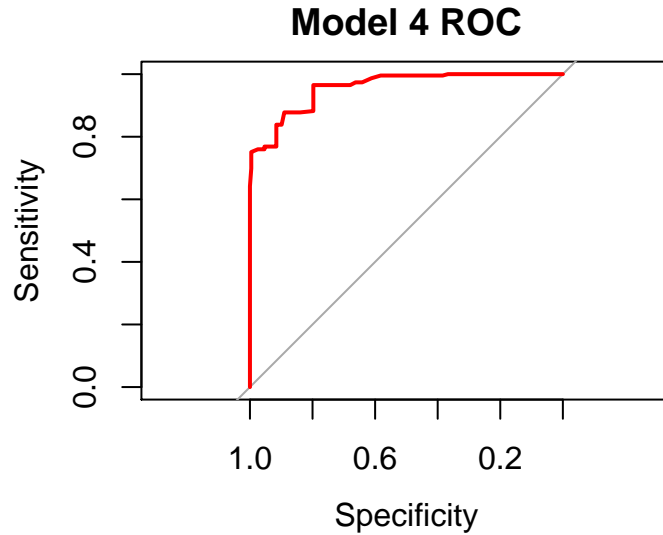
Like Model 1, this model is formed using best subsets regression with BIC as a criterion for choosing the number of predictors; however, with this model, predictors are removed sequentially until all predictors are statistically significant ( $p < 0.05$ ). This leads to the removal of age and median home values as predictors.

	Estimate	Std. Error	z value	Pr(> z )
<b>nox</b>	27.2	3.232	8.415	3.915e-17
<b>rad</b>	0.5139	0.1082	4.75	2.036e-06
<b>(Intercept)</b>	-17.45	1.949	-8.956	3.376e-19

(Dispersion parameter for binomial family taken to be 1 )

Null deviance:	645.9 on 465 degrees of freedom
Residual deviance:	239.5 on 463 degrees of freedom

The coefficients for this model are similar to the coefficients associated with the two variables in the larger BIC model (Model 1), but show far greater statistical significance with the additional variables removed.



	0	1
0	213	24
1	37	192

This model has an area under the curve of 0.9575 and an accuracy of 0.8691.

## Model 5: Best GLM

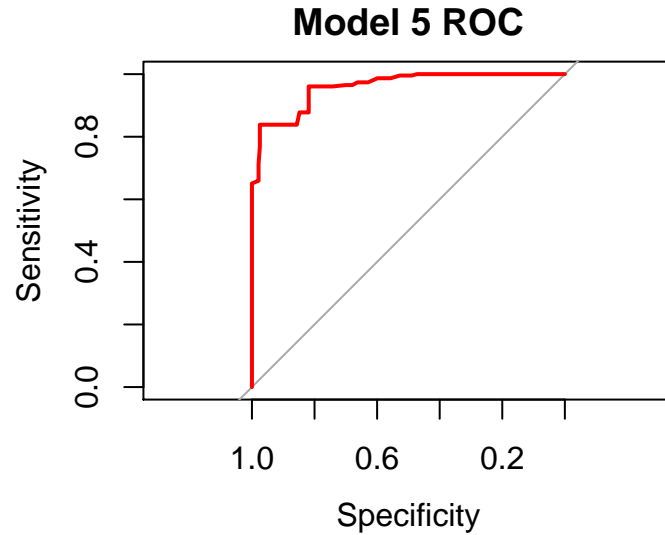
The `bestglm` package is implemented to “[select] the best subset of inputs for the GLM family.” The BIC is still used as the information criterion. The model generated by the package is presented below.

	Estimate	Std. Error	z value	Pr(> z )
<b>nox</b>	35.63	4.524	7.877	3.35e-15
<b>rad</b>	0.6376	0.1194	5.338	9.377e-08
<b>tax</b>	-0.008146	0.002332	-3.493	0.0004776
<b>(Intercept)</b>	-19.87	2.368	-8.389	4.911e-17

(Dispersion parameter for binomial family taken to be 1 )

Null deviance:	645.9 on 465 degrees of freedom
Residual deviance:	224.5 on 462 degrees of freedom

The coefficients for this model again indicate the strong influence of a unit increase in nitrogen oxide concentration. The inclusion of property tax rate in the model lead to the first negative coefficient seen. Each of the coefficient estimates is of very high statistical significance.



	0	1
0	222	15
1	37	192

This model has an area under the curve of 0.8876 and an accuracy of 0.8691.

## Model 6: Transformed Best GLM

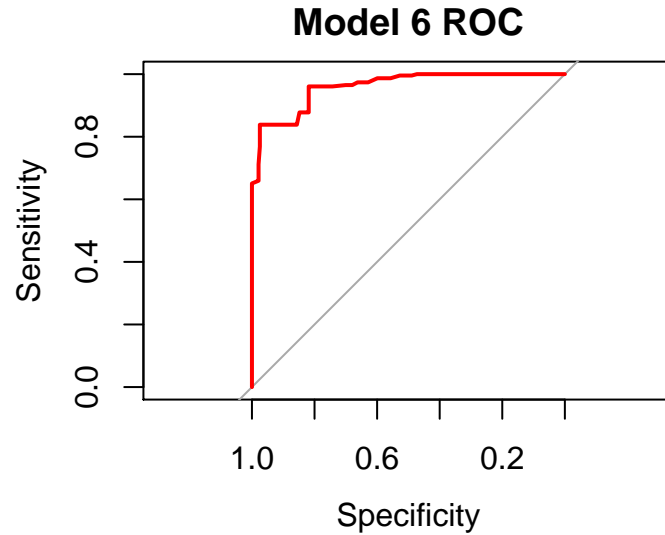
Using the three variables identified by the `bestglm` function for Model 5, the `nox` and `rad` predictors are transformed using logarithms, as in Model 3. The model using these transformations is presented below.

	Estimate	Std. Error	z value	Pr(> z )
<b>log(nox)</b>	19.35	2.519	7.682	1.564e-14
<b>log(rad)</b>	3.356	0.5447	6.162	7.196e-10
<b>tax</b>	-0.008214	0.002335	-3.518	0.0004345
<b>(Intercept)</b>	9.385	1.829	5.132	2.866e-07

(Dispersion parameter for binomial family taken to be 1 )

Null deviance:	645.9 on 465 degrees of freedom
Residual deviance:	223.5 on 462 degrees of freedom

The coefficient for nitrogen oxide concentration decreased following the transformation, as did the standard error of the estimate of the coefficient. While the p-value related to the coefficients changed following the transformation, all estimates remain statistically significant.



	0	1
0	222	15
1	37	192

This model has an area under the curve of 0.9594 and an accuracy of 0.8884.

## Model 7: Full Model

For completeness, a full generalized linear model is created using all explanatory variables.

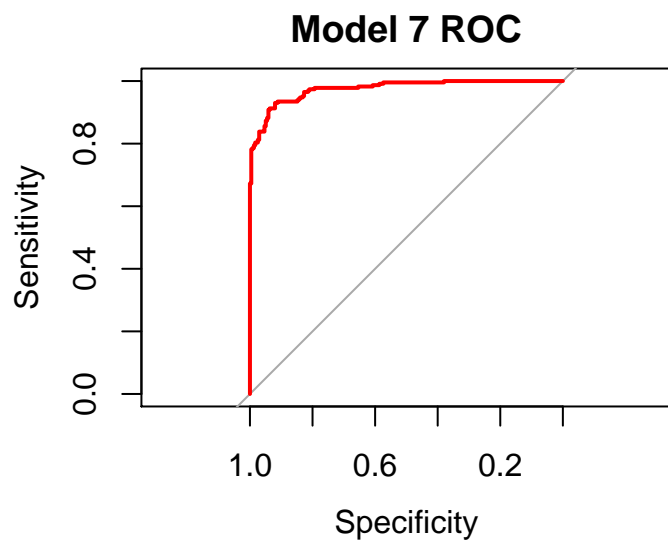
	Estimate	Std. Error	z value	Pr(> z )
<b>zn</b>	-0.06172	0.03441	-1.794	0.07287
<b>indus</b>	-0.07258	0.04855	-1.495	0.1349
<b>chas</b>	1.032	0.7596	1.359	0.1741
<b>nox</b>	50.16	8.05	6.231	4.623e-10
<b>rm</b>	-0.6921	0.7414	-0.9335	0.3505
<b>age</b>	0.03452	0.01388	2.487	0.01289
<b>dis</b>	0.7658	0.2344	3.267	0.001087
<b>rad</b>	0.663	0.1651	4.015	5.945e-05
<b>tax</b>	-0.006593	0.003064	-2.152	0.03142
<b>ptratio</b>	0.4422	0.1322	3.344	0.0008252
<b>black</b>	-0.01309	0.00668	-1.96	0.04997
<b>lstat</b>	0.04757	0.05451	0.8727	0.3828
<b>medv</b>	0.1997	0.07102	2.812	0.004919
<b>(Intercept)</b>	-36.84	7.029	-5.241	1.595e-07

(Dispersion parameter for binomial family taken to be 1 )

Null deviance:	645.9 on 465 degrees of freedom
Residual deviance:	186.1 on 452 degrees of freedom

The 13 coefficients in the full model vary very widely in magnitude, and five of the coefficients are negative. There is also a wide range of significances – 5 of the coefficients are not significant under any reasonable  $\alpha$ , and one is very nearly exactly 0.05. It is interesting to note that the coefficient for **indus** is negative, while the correlation of this variable with **target** is relatively strongly positive – this may be due to collinearity between **indus** and **nox** – highly industrialized neighborhoods

will likely have higher air pollution.



	0	1
0	222	15
1	20	209

This model has an area under the curve of 0.9753 and an accuracy of 0.9249.



# Model Selection and Prediction

## Model Comparison

The characteristics and performance of the seven models from the previous section are compared below:

Model #	# of Predictors	AUC	Accuracy
1	4	0.9570	0.8712
2	5	0.9605	0.8691
3	4	0.9584	0.8691
4	2	0.9575	0.8691
5	3	0.8876	0.8691
6	3	0.9594	0.8884
7	13	0.9753	0.9249

The model using only the significant predictors as determined by BIC (Model 3) was selected as the best model for prediction of **TARGET** in the crime data set. While the AUC value of this model was not the highest of the models tested, its mean cross-validation error indicates that it has the best predictive value for unseen data. Additionally, it is a parsimonious model, and the simplicity lends itself to easier understanding of the model by other users.

## 10-fold Cross Validation

### Mean CV Error

<b>Model1</b>	36.6
<b>Model2</b>	46.69
<b>Model3</b>	15.48
<b>Model4</b>	22.14
<b>Model5</b>	55.49
<b>Model6</b>	21.8
<b>Model7</b>	93.87

The linear model is applied to an evaluation dataset containing response variables for 259 cases. A table of the predicted team wins is presented below.

### Predicted Targets for Test Data

0	1
250	216

### Targets for Training Data

0	1
237	229

Similar to the training dataset, the predictions for the test data set predictions are weighted more toward crime being below the median. It is important to note that Model 3 did not achieve the highest AUC or accuracy on the training data but, given its cross-validation performance, we believe that Model 3 will produce the greatest generalizability of the models explored here.

A comparison of the full sets of predictions for the evaluation dataset is available in Appendix B.

## Appendix A

## Appendix B – Index-wise Results from Predictive Model

## Appendix C – R Code