

Auto Insurance prediction

Sri Seshadri

10/20/2017

Contents

1. Introduction	2
1.1 Analysis Process	2
1.2 Executive summary	2
2. Data	2
2.1 Exploratory Data Analysis (EDA)	3
2.1.1 Customer profile	3
2.1.2 Attribute variables	7
2.1.3 Handling missing data	7
3. Feature Selection	7
3.1 Training and Test data partition	7
3.2 Decision Tree	7
3.3 Penalized model - Lasso	8
3.3.1 Logistic regression as penalized models - Lasso	10
3.4. Features selection conclusion.	11
4. Modeling	12
4.1. Logistic regression - use all features selected	12

1. Introduction

An insurance company is interested in predicting which customers are likely to be in an accident and what would be the likely payout. The company requires this prediction to price the insurance policy. A predictive model is required to be deployed at point of request for quote or sale. The insurance company has been collecting data on which a predictive model would be trained and tested.

1.1 Analysis Process

The following process steps were used for building a predictive models:

- Exploratory Data Analysis
 - Perform data quality checks, quantify missing data.
 - Check for systemic loss in data
 - Understand relationships amongst predictors and between target variables and predictors.
 - Create attribute or indicator variables to aid data cleaning.
 - Filter out clean data for feature selection and model building.
- Feature Selection
 - Subset complete records to model wins in season
 - Use different modeling techniques to select candidate predictors.
 - If data is missing for candidate predictors, identify imputing methods.
- Model Building
 - Test models that were build using complete records on the entire data set with imputed data.
 - Compare models based on Deviance, ROC and MAE
 - Check if models make physical sense.
- Initial model deployment
 - Deploy model to predict wins on out of sample data.
 - Discuss models and results with subject matter experts.
 - Fine tune model and re-test
- Final model deployment

1.2 Executive summary

2. Data

The insurance company has data collected from almost 8200 customers. The dictionary of the data is provided in the appendix A.1. Tables 1 and 2 show the summary statistics of numeric and non-numeric features of the data. It is seen that some features have missing values. The missing values may need to be imputed if the features are deemed important predictors of likelihood of customer involving in a crash or payout.

It is seen that the minimum age of car is -3. Which is not rational, the data is filled in with +3 assuming it is a typographical error. Also it is seen that there is white space in the JOB column of the data.

Table 1: Summary statistics

	min	Q1	median	Q3	max	mean	sd	n	missing
INDEX	1	2559.00	5133.00	7745.00	10302.00	5151.87	2978.89	8161	0
TARGET_FLAG	0	0.00	0.00	1.00	1.00	0.26	0.44	8161	0
TARGET_AMT	0	0.00	0.00	1036.00	107586.14	1504.32	4704.03	8161	0
KIDSDRIV	0	0.00	0.00	0.00	4.00	0.17	0.51	8161	0
AGE	16	39.00	45.00	51.00	81.00	44.79	8.63	8155	6

	min	Q1	median	Q3	max	mean	sd	n	missing
HOMEKIDS	0	0.00	0.00	1.00	5.00	0.72	1.12	8161	0
YOJ	0	9.00	11.00	13.00	23.00	10.50	4.09	7707	454
INCOME	0	28096.97	54028.17	85986.21	367030.26	61898.10	47572.69	7716	445
HOME_VAL	0	0.00	161159.53	238724.45	885282.34	154867.29	129123.78	7697	464
TRAVTIME	5	22.45	32.87	43.81	142.12	33.49	15.90	8161	0
BLUEBOOK	1500	9280.00	14440.00	20850.00	69740.00	15709.90	8419.73	8161	0
TIF	1	1.00	4.00	7.00	25.00	5.35	4.15	8161	0
OLDCLAIM	0	0.00	0.00	4636.00	57037.00	4037.08	8777.14	8161	0
CLM_FREQ	0	0.00	0.00	2.00	5.00	0.80	1.16	8161	0
MVR_PTS	0	0.00	1.00	3.00	13.00	1.70	2.15	8161	0
CAR_AGE	-3	1.00	8.00	12.00	28.00	8.33	5.70	7651	510

Table 2: Sanity check of non numeric variables

	# Unique	n	missing	Blanks
PARENT1	2	8161	0	0
MSTATUS	2	8161	0	0
SEX	2	8161	0	0
EDUCATION	5	8161	0	0
JOB	9	8161	0	526
CAR_USE	2	8161	0	0
CAR_TYPE	6	8161	0	0
RED_CAR	2	8161	0	0
REVOKED	2	8161	0	0
URBANICITY	2	8161	0	0

2.1 Exploratory Data Analysis (EDA)

In this section interesting observations in the data are noted and used to characterize the population of customers.

2.1.1 Customer profile

Figure 1 shows the customer portfolio of the insurance company. It is seen from figure 1 that majority of the customers hold blue collar and clerical roles. The income and home values are zero inflated, likely contributed by students and home makers. Nearly 25% of the customers have been involved in a car crash. Most of the customers are with the policy less than 2 years.

It is seen that majority of the customers own cars that are new; less than 2 years old. The older cars are owned by professionals and by the group who has their nature of job missing (potentially not disclosed). Also interestingly that cars that are between 14 and 15 years of age are missing from the population. As seen in histogram in figure 2.

Figure 3 explores the missing JOB data. The missing JOB values lines up with the population of white collar jobs. From the income, age and year on the job perspective. Also the zero inflation in income is contributed by students and home makers. The jobs above the red dashed lines the income boxplot is defined as white collar.

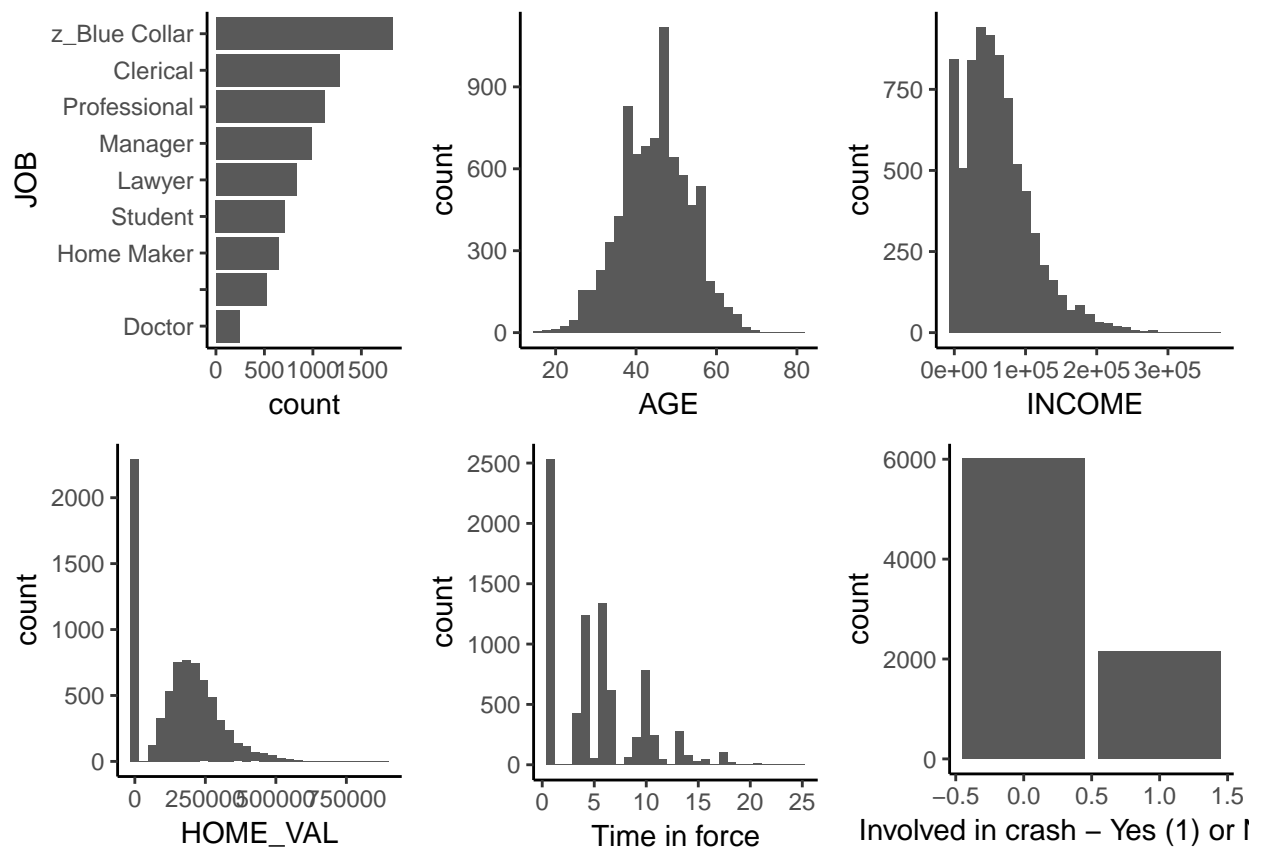


Figure 1: Customer portfolio

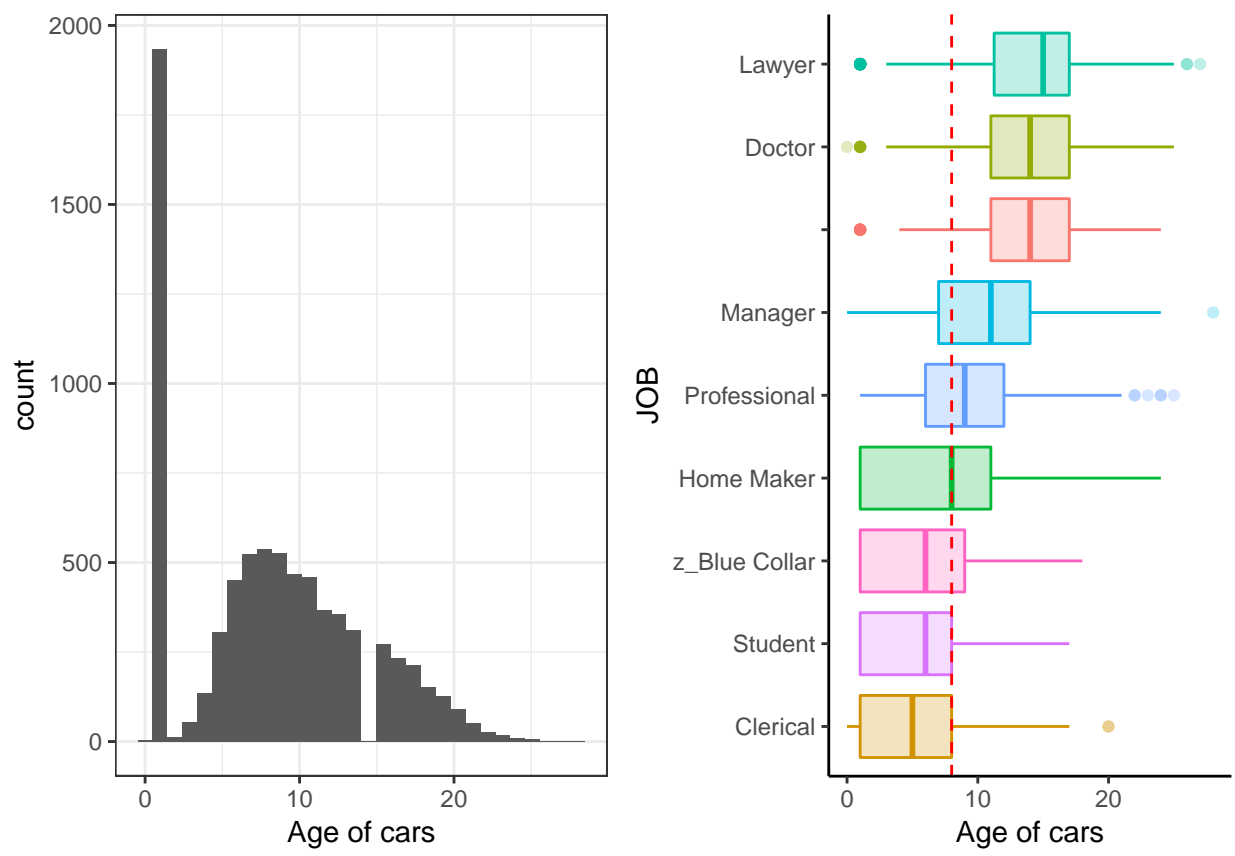


Figure 2: car age

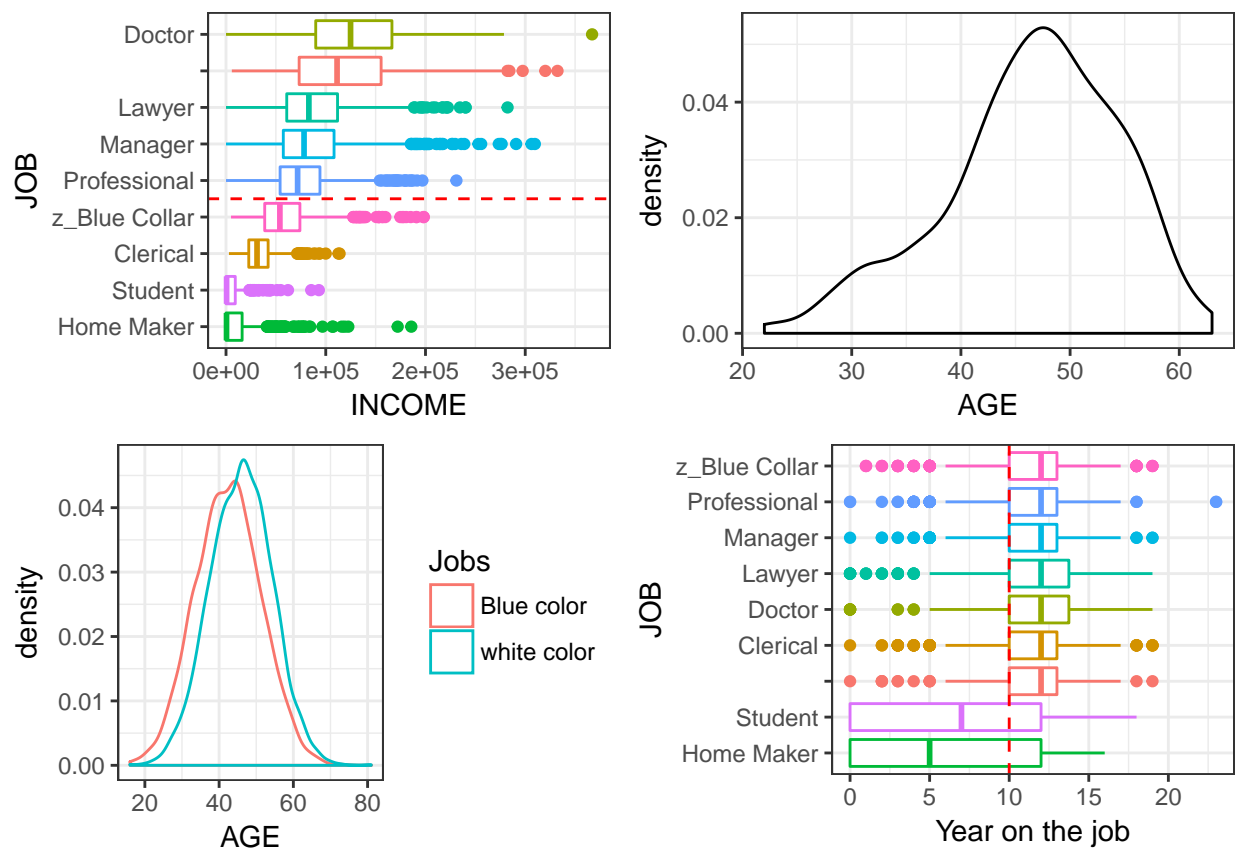


Figure 3: Missing JOB profile

The categorical variables are transformed into indicator variables to be used in modeling. The missing data points are manifested as indicator variables.

The complete cases would be used to determine key features required for modeling. If data for key features are missing, an imputation strategy would be determined.

In this section we consider various feature selection methodologies such as 1. Decision Trees, 2. Penalized model - Lasso

In order to test if the feature selection are really useful, feature selections need to be cross validated on a hold out test data set. 80% of the data is used as training and the rest is used as hold out for testing. A stratified sampling method is used to capture identical percentage of samples involved in crash.

Decision tree model is fitted on the training set to identify stand out splits in the data based on Gini index. The decision tree is shown in figure 4. Bagging technique is used to minimize variance in the model to ensure we have a reliable feature selection. The important features are shown in figure 5.



```
##
## Call:
## randomForest(formula = as.factor(TARGET_FLAG2) ~ . - TARGET_FLAG,      data = training, mtry = 28,
##               Type of random forest: classification
##               Number of trees: 500
## No. of variables tried at each split: 28
##
## OOB estimate of error rate: 21.1%
## Confusion matrix:
##      No Yes class.error
## No  3495 312  0.08195429
## Yes   777 578  0.57343173
```

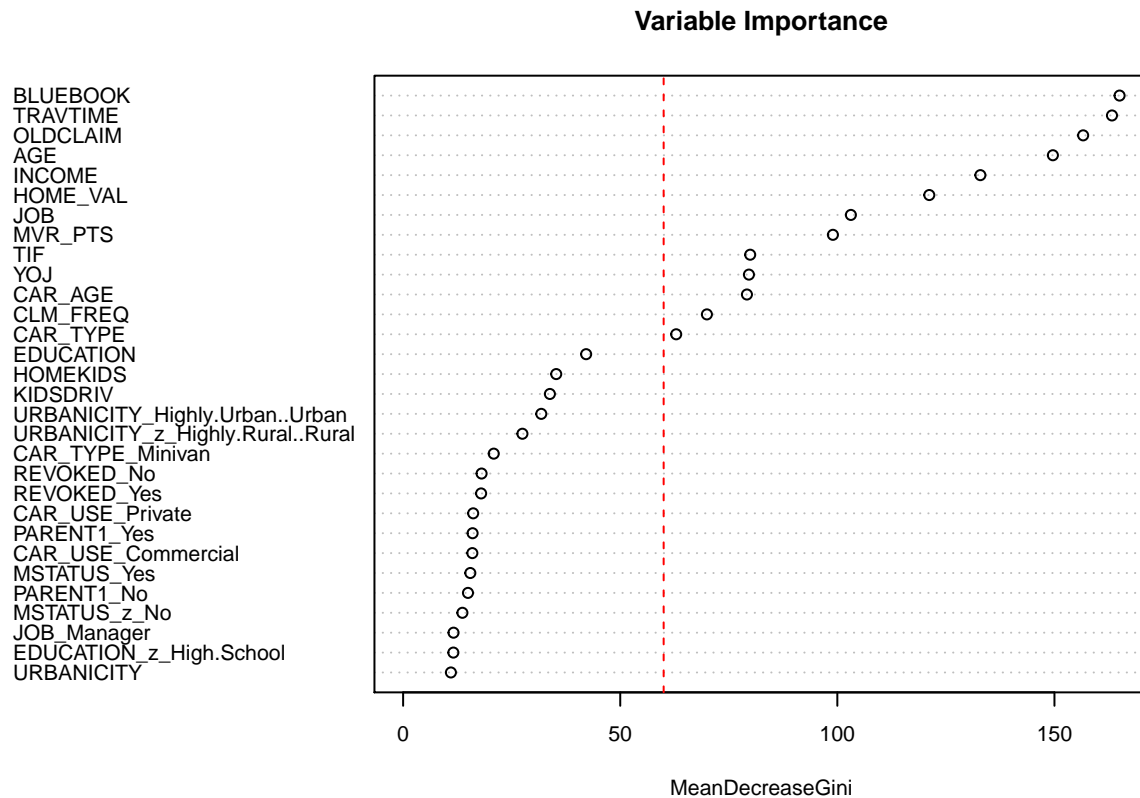


Figure 5: Variable Importance Plot and test prediction

The decision tree was used to predict the hold out test data and the AUC was found to be 80 % as shown in figure 6. Therefore we'll use the top 13 predictors as shown in figure 4. In the next section we will explore penalized models to gather important predictors.

3.3 Penalized model - Lasso

The variable selection property of Lasso is used to aid with automated variable selection. Again the model is trained on training data set and tested on a hold out dataset, as in the decision tree method. However a 10 fold cross validation method was used to identify the optimal penalization parameter - lambda. Figure 7 shows the coefficients that are not zero in the decreasing order of absolute value of coefficients. The predictors falling above the red dashed line, drawn at the elbow in figure 7 are deemed good predictors.

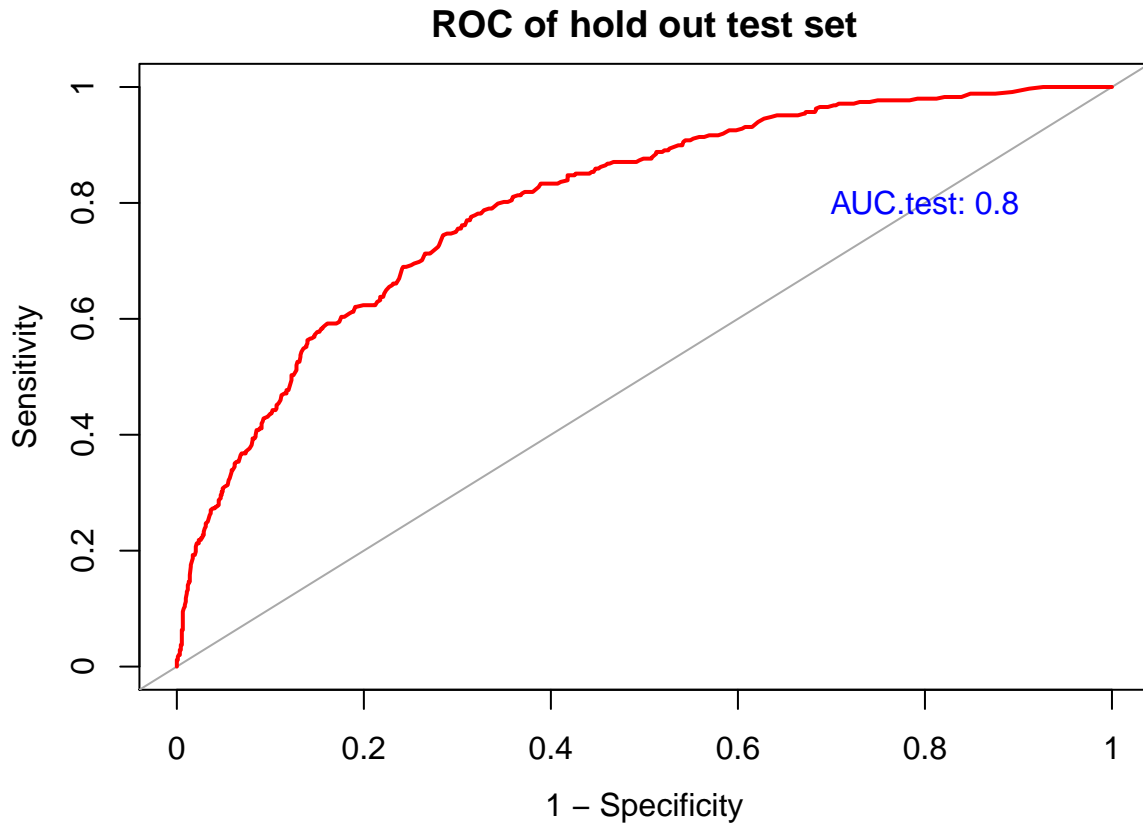


Figure 6: Prediction of bagged decision tree

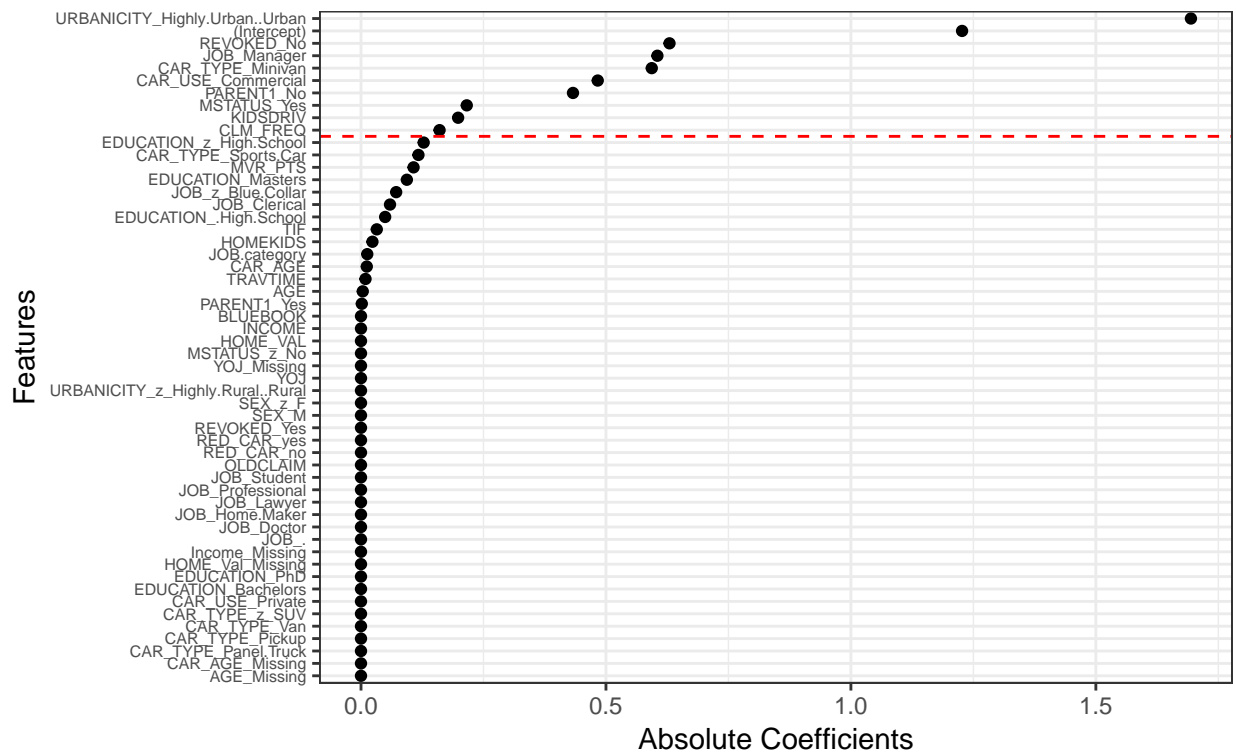


Figure 7: Feature selection - Lasso

3.3.1 Logistic regression as penalized models - Lasso

While the penalized logistic regression model is used for feature selection, it may be used for prediction as well. All the predictors corresponding to non zero coefficients are considered as predictors. The ROC curves for the training and test samples are shown in figure 8 and the Gain chart is shown in figure 9.

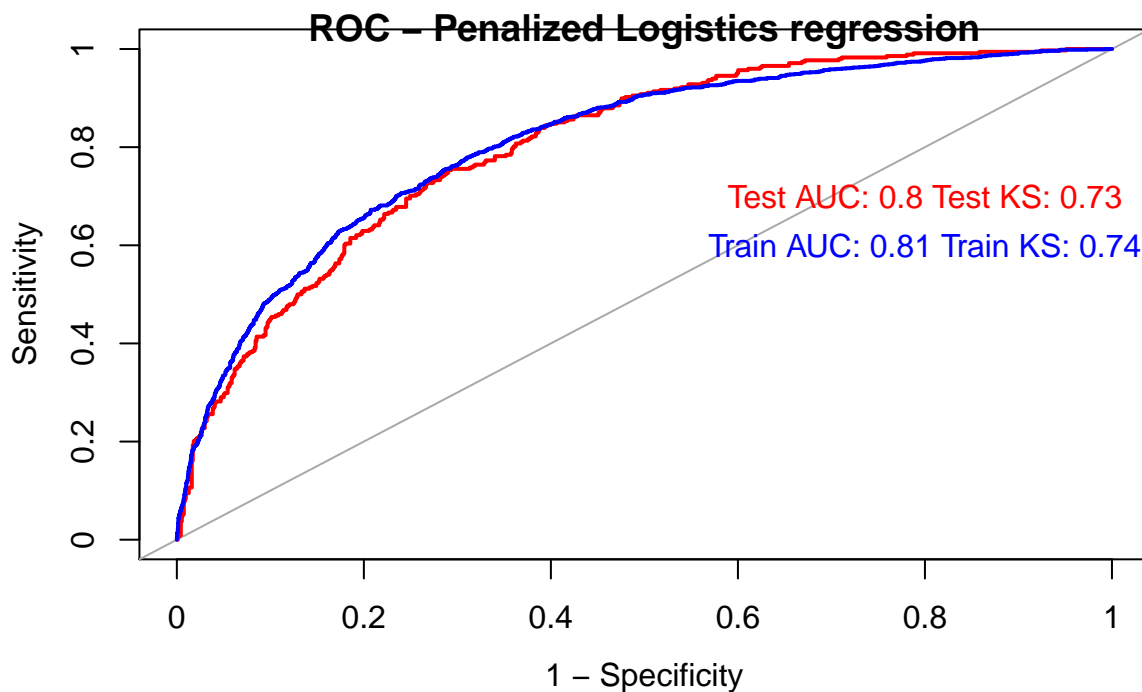
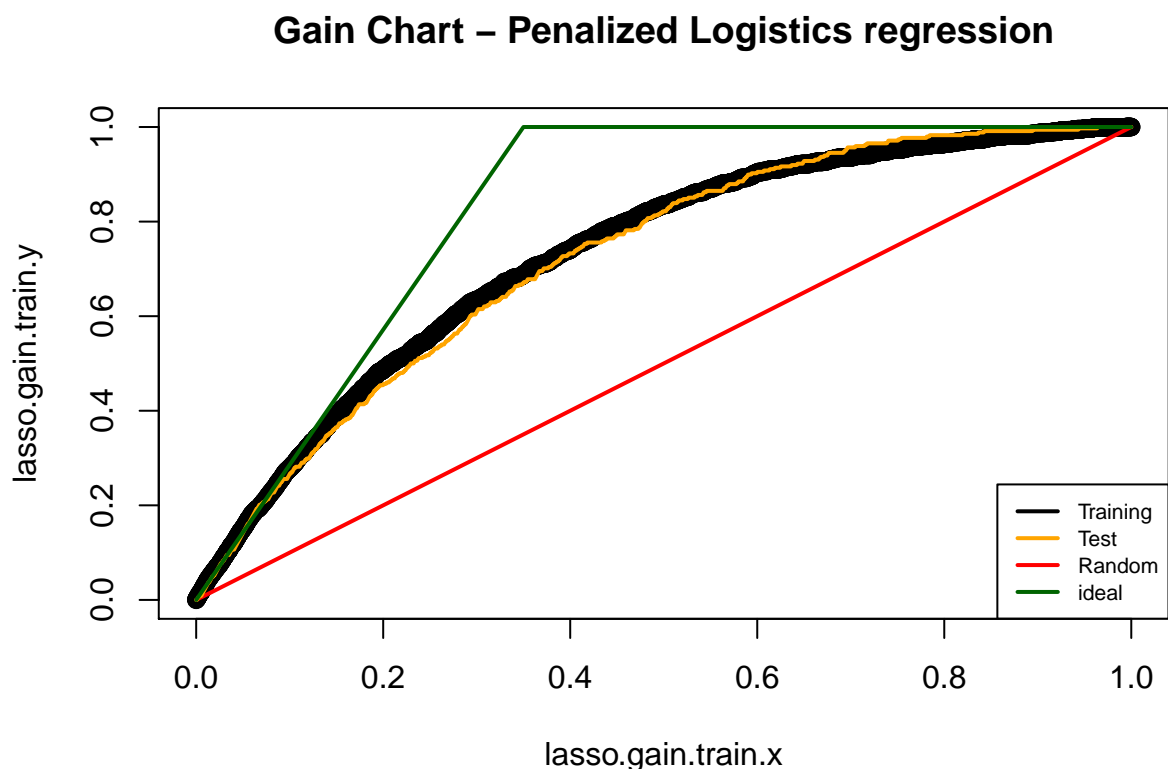


Figure 8: Lasso ROC curves for test and training data



It can

be seen from figure 8 and 9, lasso regression has 74% AUC and is better than a random model. The test performance is identical to the training performance. The misclassification is about 23%. However performance on the records with missing values is yet evaluated and will be discussed in the modeling section.

```
##           Reference
## Prediction    0    1
##           0 3616  902
##           1  191  453
```

Table 3: confusion matrix statistics - lasso logistic train

Accuracy	0.79
Kappa	0.34
AccuracyLower	0.78
AccuracyUpper	0.80
AccuracyNull	0.74
AccuracyPValue	0.00
McnemarPValue	0.00

```
##           Reference
## Prediction    0    1
##           0  886  240
##           1   52  108
```

Table 4: confusion matrix statistics - lasso logistic test

Accuracy	0.77
Kappa	0.31
AccuracyLower	0.75
AccuracyUpper	0.80
AccuracyNull	0.73
AccuracyPValue	0.00
McnemarPValue	0.00

3.4. Features selection conclusion.

The following predictors are deemed useful for modeling purposes after the above feature selection process.

```
BLUEBOOK
TRAVTIME
OLDCLAIM
AGE
INCOME
HOME_VAL
JOB
MVR_PTS
TIF
YOJ
CAR_AGE
CLM_FREQ
URBANICITY_Highly.Urban..Urban
JOB_Manager
CAR_TYPE_Minivan
```

CAR_USE_Commercial
PARENT1_No
MSTATUS_Yes
KIDSDRIV

4. Modeling

4.1. Logistic regression - use all features selected

A logistic regression with the above selected predictors (section 3.4) is fitted. The summary of the fit is shown below. It can be seen from the Chi-squared goodness of fit that the model is adequate. There are some predictors whose slope is not significantly different from 0.

```
##
## Call:
## glm(formula = TARGET_FLAG ~ ., family = binomial, data = training.logistic)
##
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
## -2.5611  -0.7159  -0.4098   0.6132   3.1197
##
## Coefficients: (2 not defined because of singularities)
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept)   -9.598e-01  3.063e-01  -3.134 0.001724 **
## BLUEBOOK      -2.573e-05  5.136e-06  -5.009 5.48e-07 ***
## TRAVTIME       1.531e-02  2.362e-03   6.482 9.03e-11 ***
## OLDCLAIM      -5.777e-06  4.936e-06  -1.170 0.241828
## AGE           -6.339e-03  4.584e-03  -1.383 0.166708
## INCOME        -2.250e-06  1.367e-06  -1.646 0.099758 .
## HOME_VAL      -1.215e-06  4.352e-07  -2.792 0.005240 **
## MVR_PTS       1.233e-01  1.706e-02   7.227 4.93e-13 ***
## TIF           -5.167e-02  9.171e-03  -5.635 1.75e-08 ***
## YOJ           -6.790e-03  1.040e-02  -0.653 0.513890
## CAR_AGE       -1.992e-02  8.108e-03  -2.456 0.014034 *
## CLM_FREQ      2.023e-01  3.526e-02   5.736 9.68e-09 ***
## URBANICITY_Highly.Urban..Urban 2.212e+00  1.380e-01  16.035 < 2e-16 ***
## REVOKED_No    -8.805e-01  1.170e-01  -7.526 5.22e-14 ***
## CAR_TYPE_Minivan -8.280e-01  9.421e-02  -8.789 < 2e-16 ***
## CAR_USE_Commercial 6.378e-01  9.771e-02   6.527 6.69e-11 ***
## PARENT1_No    -4.979e-01  1.231e-01  -4.044 5.25e-05 ***
## MSTATUS_Yes   -3.514e-01  1.045e-01  -3.361 0.000776 ***
## KIDSDRIV      3.512e-01  7.006e-02   5.012 5.37e-07 ***
## JOB_         -4.831e-01  1.810e-01  -2.669 0.007611 **
## JOB_Clerical   1.437e-01  1.287e-01   1.116 0.264352
## JOB_Doctor    -4.345e-01  2.850e-01  -1.525 0.127318
## JOB_Home.Maker -2.067e-01  1.866e-01  -1.108 0.267922
## JOB_Lawyer    -2.361e-01  1.783e-01  -1.324 0.185534
## JOB_Manager   -1.046e+00  1.608e-01  -6.506 7.73e-11 ***
## JOB_Professional -2.567e-01  1.368e-01  -1.876 0.060607 .
## JOB_Student   -3.779e-02  1.648e-01  -0.229 0.818615
## JOB_z_Blue.Collar NA          NA          NA          NA
## JOB.category   NA          NA          NA          NA
```

```

## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
##      Null deviance: 5943.0  on 5161  degrees of freedom
## Residual deviance: 4628.8  on 5135  degrees of freedom
## AIC: 4682.8
##
## Number of Fisher Scoring iterations: 5
##
## Analysis of Deviance Table
##
## Model: binomial, link: logit
##
## Response: TARGET_FLAG
##
## Terms added sequentially (first to last)
##
##
##              Df Deviance Resid. Df Resid. Dev  Pr(>Chi)
## NULL                                5161      5943.0
## BLUEBOOK                1    65.153      5160      5877.9 6.930e-16
## TRAVTIME                 1    10.877      5159      5867.0 0.0009739
## OLDCLAIM                 1   129.809      5158      5737.2 < 2.2e-16
## AGE                     1    51.392      5157      5685.8 7.563e-13
## INCOME                  1    31.756      5156      5654.0 1.748e-08
## HOME_VAL                1    57.773      5155      5596.3 2.941e-14
## MVR_PTS                 1   160.041      5154      5436.2 < 2.2e-16
## TIF                     1    26.291      5153      5409.9 2.937e-07
## YOJ                     1     0.194      5152      5409.7 0.6595193
## CAR_AGE                 1    14.879      5151      5394.9 0.0001146
## CLM_FREQ                1    68.838      5150      5326.0 < 2.2e-16
## URBANICITY_Highly.Urban..Urban 1   295.112      5149      5030.9 < 2.2e-16
## REVOKED_No              1    60.882      5148      4970.0 6.060e-15
## CAR_TYPE_Minivan        1   110.763      5147      4859.3 < 2.2e-16
## CAR_USE_Commercial       1    81.130      5146      4778.1 < 2.2e-16
## PARENT1_No              1    54.336      5145      4723.8 1.690e-13
## MSTATUS_Yes             1     7.242      5144      4716.5 0.0071215
## KIDSDRIV                1    26.154      5143      4690.4 3.152e-07
## JOB_                    1     0.748      5142      4689.6 0.3870173
## JOB_Clerical            1    10.875      5141      4678.8 0.0009745
## JOB_Doctor              1     0.015      5140      4678.8 0.9033436
## JOB_Home.Maker          1     0.006      5139      4678.8 0.9382485
## JOB_Lawyer              1     2.567      5138      4676.2 0.1090987
## JOB_Manager             1    43.842      5137      4632.3 3.559e-11
## JOB_Professional        1     3.519      5136      4628.8 0.0606790
## JOB_Student             1     0.053      5135      4628.8 0.8185791
## JOB_z_Blue.Collar       0     0.000      5135      4628.8
## JOB.category            0     0.000      5135      4628.8
##
## NULL
## BLUEBOOK                ***
## TRAVTIME                 ***

```

```

## OLDCLAIM ***
## AGE ***
## INCOME ***
## HOME_VAL ***
## MVR_PTS ***
## TIF ***
## YOJ
## CAR_AGE ***
## CLM_FREQ ***
## URBANICITY_Highly.Urban..Urban ***
## REVOKED_No ***
## CAR_TYPE_Minivan ***
## CAR_USE_Commercial ***
## PARENT1_No ***
## MSTATUS_Yes **
## KIDSDRIV ***
## JOB_
## JOB_Clerical ***
## JOB_Doctor
## JOB_Home.Maker
## JOB_Lawyer
## JOB_Manager ***
## JOB_Professional .
## JOB_Student
## JOB_z_Blue.Collar
## JOB.category
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

## [1] "Goodness of fit test... P-Value 1"

## Confusion Matrix and Statistics
##
##           Reference
## Prediction  0    1
##           0 869 208
##           1  69 140
##
##           Accuracy : 0.7846
##           95% CI : (0.7611, 0.8068)
##           No Information Rate : 0.7294
##           P-Value [Acc > NIR] : 3.01e-06
##
##           Kappa : 0.376
##           McNemar's Test P-Value : < 2.2e-16
##
##           Sensitivity : 0.4023
##           Specificity : 0.9264
##           Pos Pred Value : 0.6699
##           Neg Pred Value : 0.8069
##           Prevalence : 0.2706
##           Detection Rate : 0.1089
##           Detection Prevalence : 0.1625
##           Balanced Accuracy : 0.6644
##

```

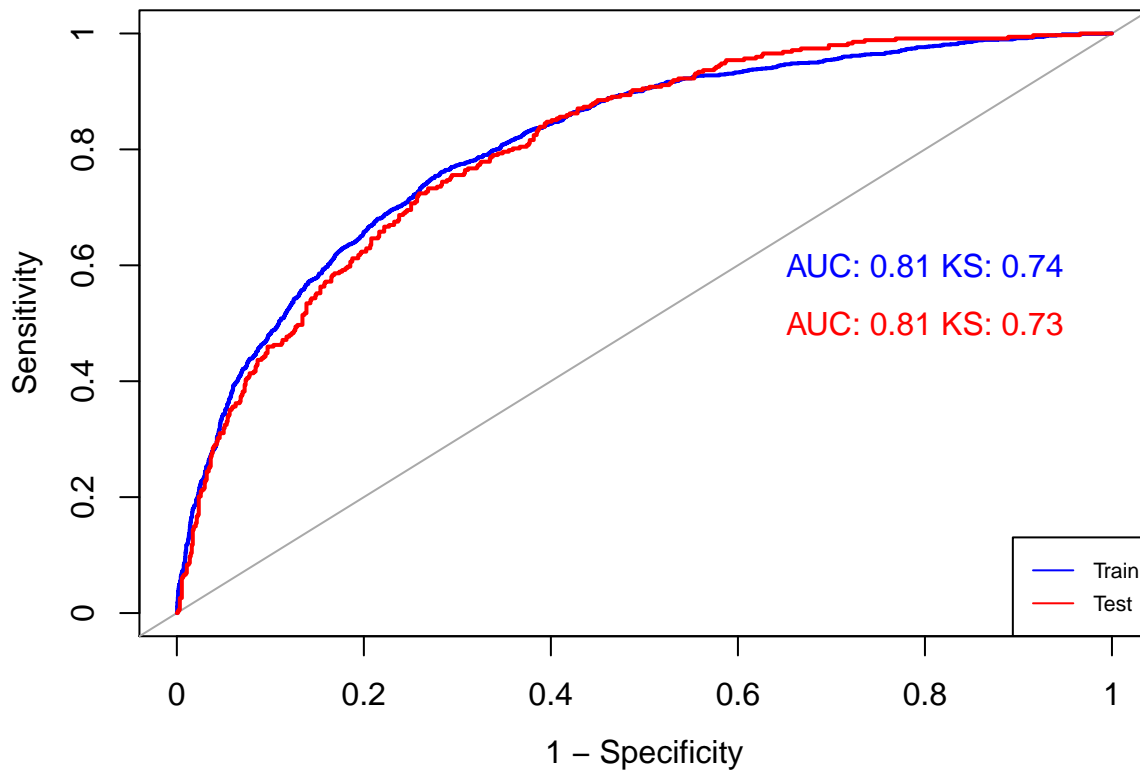


Figure 9: ROC curve for Logistic regression - full chosen predictors

```
##      'Positive' Class : 1
##
```

The performance of the model on complete cases seem to be good and identical to the Lasso logistic regression model. However, the model does require features that has data missing. We'll use the KNN imputation method to impute data.

Gain Chart – Lasso logistic and logistic

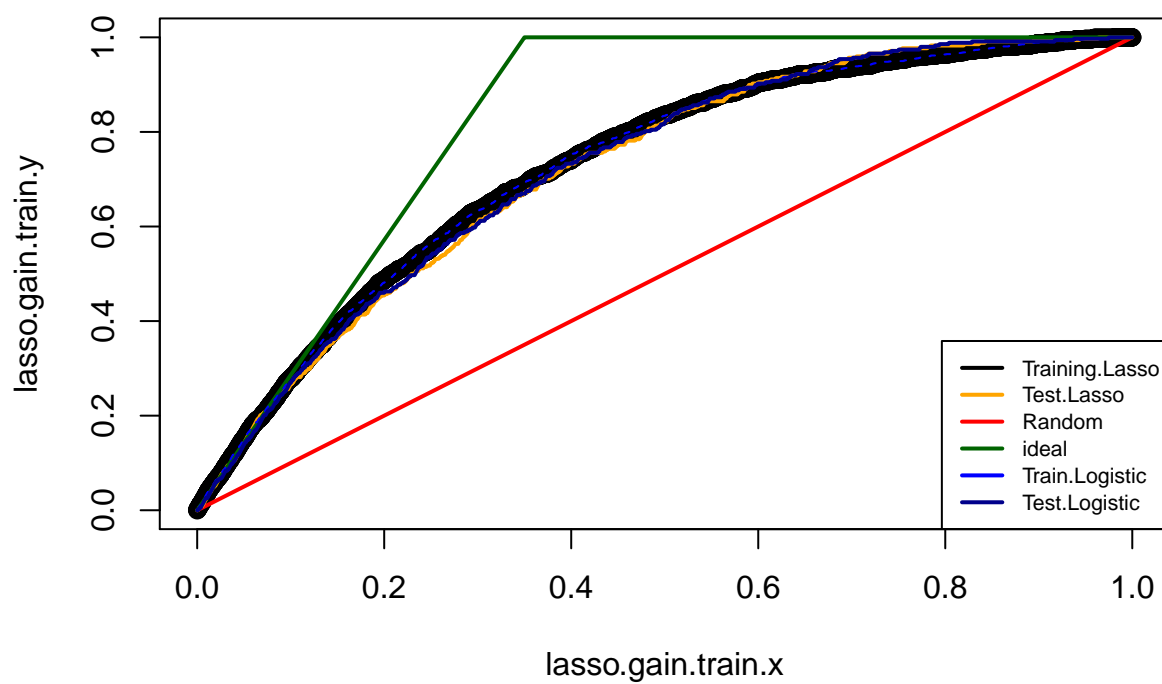
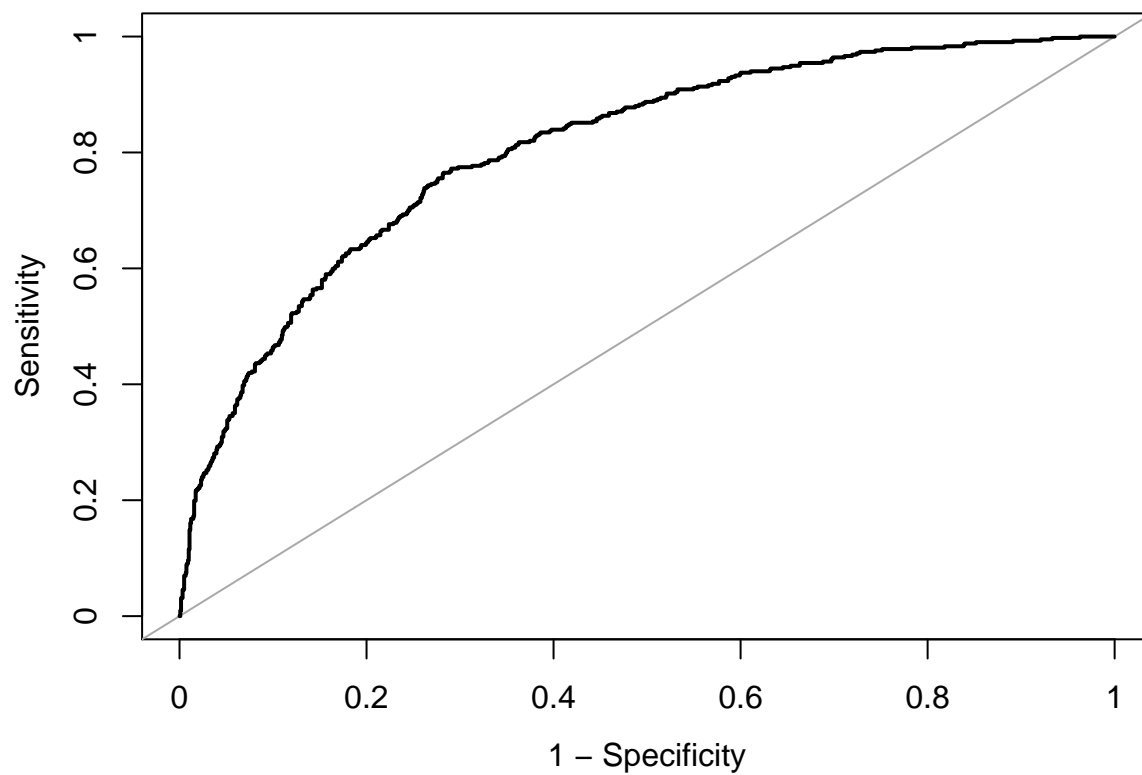
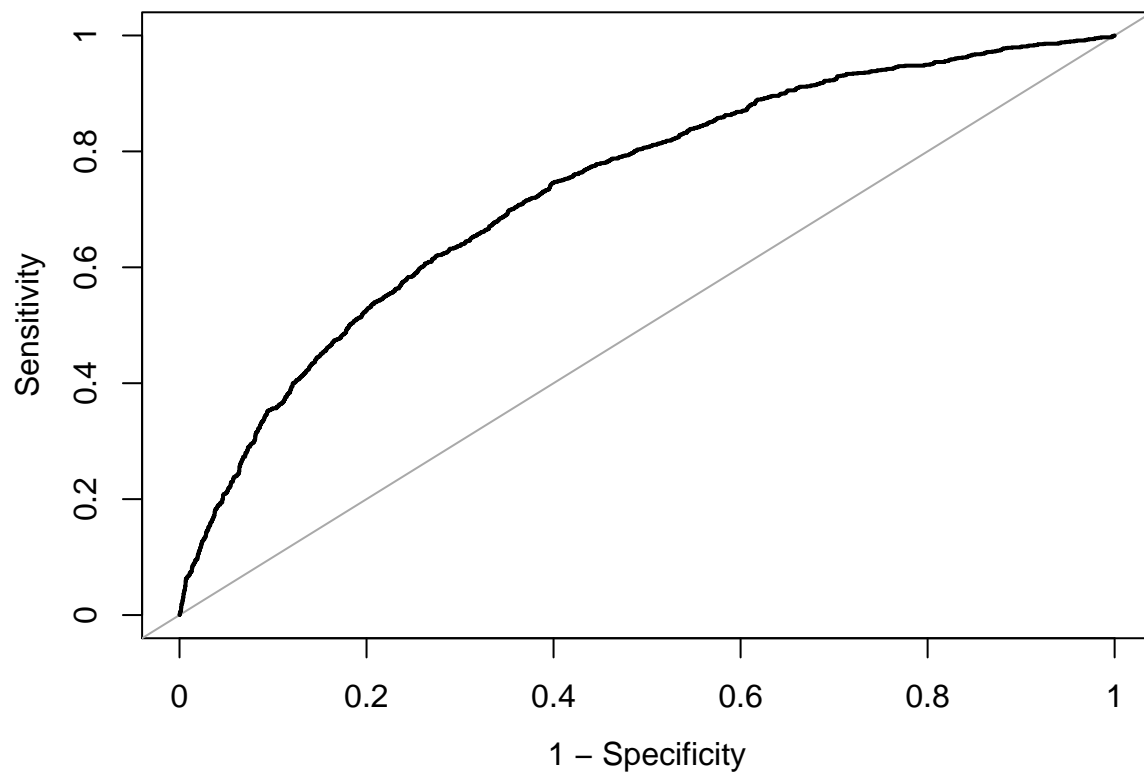


Figure 10: Gain Chart for All models so far



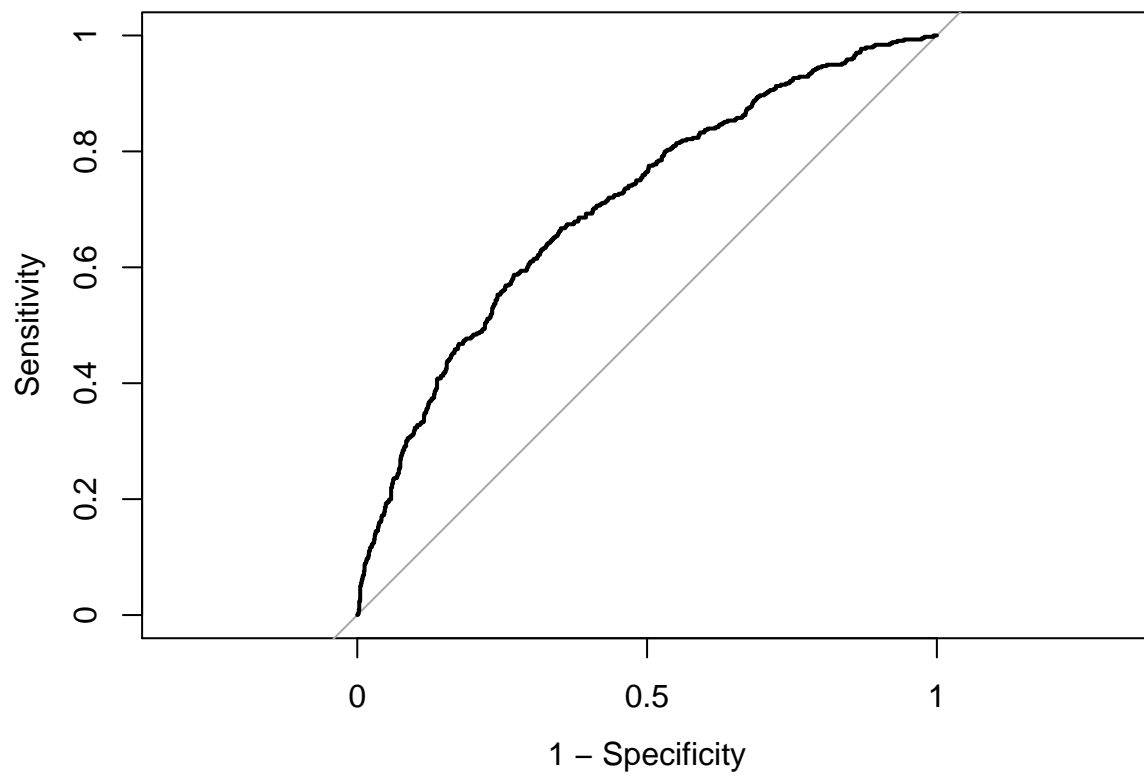
Area under the curve: 0.8066


```
##
## Call:
## glm(formula = y ~ ., family = binomial, data = dffit[complete.cases(dffit),
##      ])
##
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
## -1.9853  -0.7762  -0.5575   0.8636   2.9407
##
## Coefficients: (1 not defined because of singularities)
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept)    1.199e-01  2.438e-01   0.492 0.622809
## BLUEBOOK      -1.542e-05  4.884e-06  -3.158 0.001588 **
## TRAVTIME       7.834e-03  2.230e-03   3.513 0.000443 ***
## OLDCLAIM       4.994e-06  4.225e-06   1.182 0.237234
## AGE           -1.125e-02  4.232e-03  -2.659 0.007847 **
## INCOME         7.658e-08  1.288e-06   0.059 0.952573
## HOME_VAL      -2.841e-06  3.531e-07  -8.045 8.60e-16 ***
## MVR_PTS        1.468e-01  1.649e-02   8.900 < 2e-16 ***
## TIF           -5.450e-02  9.074e-03  -6.006 1.90e-09 ***
## YOJ           -8.364e-03  9.820e-03  -0.852 0.394394
## CAR_AGE       -9.947e-03  7.969e-03  -1.248 0.211919
## CLM_FREQ       2.768e-01  3.401e-02   8.138 4.03e-16 ***
## JOB_          -1.154e-01  1.838e-01  -0.628 0.530057
## JOB_Clerical   -3.188e-01  1.104e-01  -2.888 0.003883 **
## JOB_Doctor     -9.293e-01  2.836e-01  -3.277 0.001048 **
## JOB_Home.Maker -6.588e-01  1.626e-01  -4.051 5.11e-05 ***
## JOB_Lawyer     -5.203e-01  1.554e-01  -3.347 0.000816 ***
## JOB_Manager    -1.116e+00  1.533e-01  -7.275 3.45e-13 ***
## JOB_Professional -5.545e-01  1.259e-01  -4.403 1.07e-05 ***
## JOB_Student    -4.071e-01  1.520e-01  -2.678 0.007410 **
## JOB_z_Blue.Collar      NA         NA         NA         NA
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
##      Null deviance: 5548.7  on 4812  degrees of freedom
## Residual deviance: 4895.9  on 4793  degrees of freedom
## AIC: 4935.9
##
## Number of Fisher Scoring iterations: 4
```



```
## Area under the curve: 0.7338
```

```
##           test.y
## predictedClass  0   1
##           0 1123 334
##           1   76 102
```



Area under the curve: 0.7075