

# PM 591 Final project

## Introduction

The National Inpatient Sample (NIS) data, collected by the Healthcare Cost and Utilization Project (HCUP), is the largest publicly available dataset that contains information on inpatient healthcare in hospitals throughout the United States. The NIS is used by policymakers and health officials to make national estimates of healthcare utilization, and observe key features of inpatient care. The NIS was first started in 1998 by the Healthcare Cost and Utilization Project, and contains information such as patient demographics, classification of diseases, total hospital bill, length of stay, and many other features that characterize hospital care. The goal of this assignment will be to build a model to predict inpatient mortality and determine what factors contribute to an increased risk of death during hospitalization.

The data that will be used in this assignment consists of a random subset of 200,000 patients from the 2012 National Inpatient Sample. The data was taken from the Healthcare Cost and Utilization Project (HCUP), which is the largest collection of hospital care data in the United States. The data was taken from discharge records from all hospitals that are participating with the HCUP, and use state guidelines to help identify the hospitals that qualify for the data collection process. 47 states and the District of Columbia participate in the NIS, and data is available for hospitals in those states. The outcome of interest is the inpatient mortality, of whether the patient died during the period of hospitalization. Features such as patient demographic, severity of disease, risk of mortality, and comorbidities were incorporated to determine if a patient was likely to die during hospitalization. This can be used to identify features that increase the risk of patient mortality in hospitals and seek to prevent such deaths in the future.

## Methods

First the relevant features to the outcome of interest were sorted out from the 175 original features that were present.

Then the data was then reevaluated and factors were added when necessary.

Out of the 175 possible features that were present in the original dataset, only 44 variables were selected to be included in analysis and model building. These 44 include data regarding patient demographics (age, race, gender), comorbidities (such as alcohol abuse and COPD), and the risks of patient mortality. Each variable was examined and was made into factor variables as was appropriate. A majority of the features were converted into dummy variables, however some remained as strings and integers. In examining the missing data, there was less than 1% of the total sample size that was missing from the target variable, whether the patient died. Because the sample was small compared to the dataset, the missing values of the target variable were removed before the analysis.

## Logistic Regression

I will be comparing 3 different methods to build a predictive model for patient mortality. The first will be logistic regression model. The logistic regression model is one of the most commonly used and basic binary classifiers. Because the desired goal is to determine if a patient died during their hospitalization, the outcome is a binary outcome. Given the extremely large sample size of the data with around 200,000 observations, both the training and testing sets will be large enough to ensure an accurate prediction model.

Forward selection was used to determine the features that will be included in the logistic regression model. According to the forward selection process, only the **APRDRG\_Risk\_Mortality**, a factor variable that characterizes the risk of patient mortality, was determined to be significant in the data. However, the race variable was also included to determine the effect of patient demographics on mortality. There will only be a couple of features included in the actual logistic prediction model, therefore the model will be a simpler one indicating that the model will have a higher bias. However, the large sample size of the data and the use of cross validation will be used to determine the accuracy of the results. I will be using the misclassification error and the AUC as performance metrics to determine the effectiveness of the logistic model and compare it to other models that I will be using.

K-fold cross validation will be used to reduce the error that comes from different training/testing splits. The resulting misclassification error from the k-fold cross validation will be compared with the misclassification error from the initial training/testing split.

### **Balanced Random Forests**

The data itself is very unbalanced, with 184,598 patients that were successfully discharged compared to the 3,412 that died in the hospital. This could lead to an optimistically low misclassification error. Therefore, balanced random forests will be used to help balance the two binary outcomes and correct over optimistic misclassification errors.

I will be using all of the 44 variables that were selected from the original dataset in the balanced random forest model. All of the variables are included because the random forest model will tune the parameters and adjust the model according to which variables are considered important. The variable importance plot generated from the balanced random forest model will be compared to the variables that were considered important in the forward selection algorithm used in the logistic regression. The AUC will be used as a performance metric as the goal of the model is to correctly predict patient mortality.

### **Lasso Regression**

Lasso regression will also be used to build a predictive model for determining patient mortality. Lasso regression was chosen over ridge regression because it is likely that only a small number of predictors will be significant in determining patient mortality than the all of the features that we have available. The regression model itself will chose which variables are important and act similarly to a feature selection algorithm. Therefore, we will continue to increase the tuning parameter to determine which parameters are important in determining patient mortality.

Cross validation will be used to tune the parameters in the LASSO regression and the misclassification error and AUC will be calculated to compare the performance o the LASSO regression model with the logistic regression and balanced random forest model.

## **Results**

### **Logistic Regression:**

Below is the AUC and associated 95% Confidence Interval for the logistic regression on the training data.

```
## Area under the curve: 0.9362
```

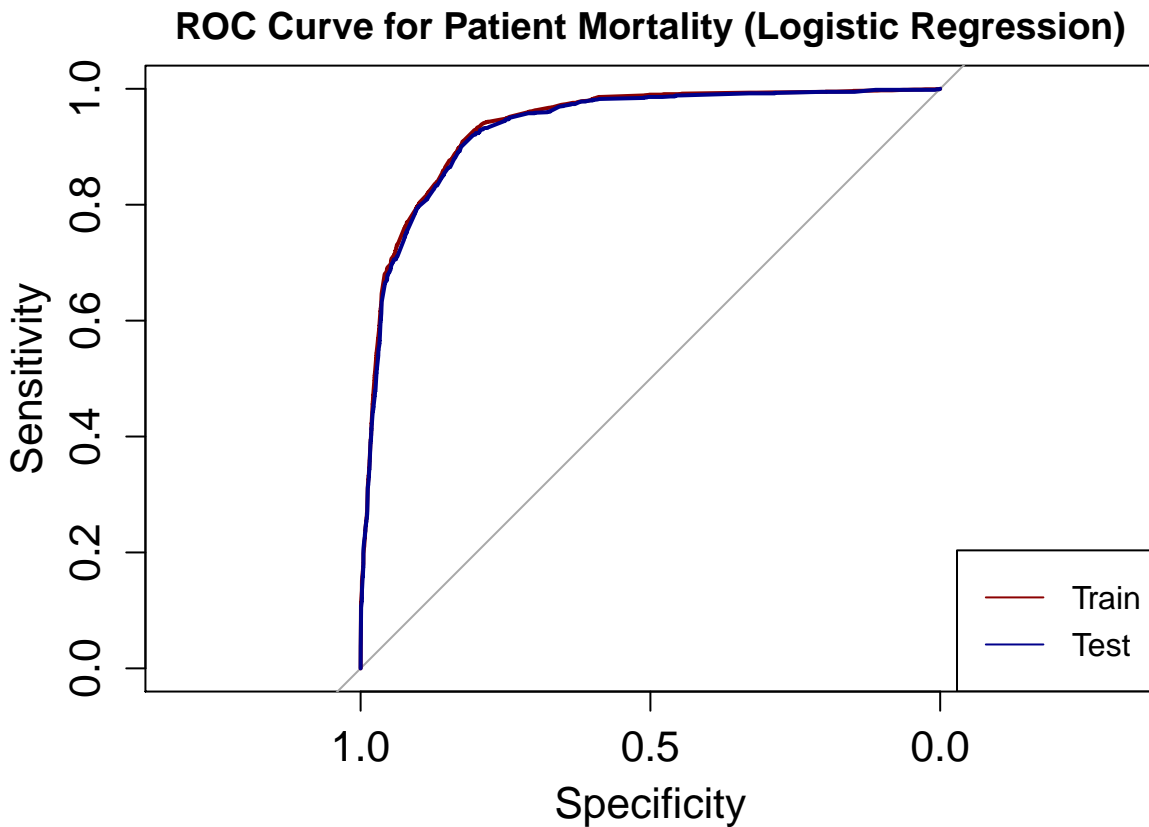
```
## 95% CI: 0.9317-0.9408 (DeLong)
```

Below is the AUC and associated 95% confidence interval for the logistic regression on the testing data.

## Area under the curve: 0.9324

## 95% CI: 0.9257-0.9391 (DeLong)

The plot below shows the ROC curve for both the training and testing data for the logistic regression model on patient mortality.



The table below shows the parameter estimates and the associated p-values for the logistic regression. The risk of mortality variable is statistically significant across almost all of its factor levels, and the length of stay variable is also statistically significant. The race variable is statistically significant only if the individual is Black. However, the parameter estimates for Race and Length of Stay are not very large compared to that of the Risk Mortality, indicating that the Risk Mortality has the highest influence in determining the probability of patient mortality.

##	Estimate	Pr(> z )
## (Intercept)	-2.715474760	9.069181e-11
## APRDRG_Risk_MortalityMinor Likelihood	-4.212201722	5.005929e-21
## APRDRG_Risk_MortalityModerate Likelihood	-1.602078731	1.622046e-04
## APRDRG_Risk_MortalityMajor Likelihood	0.209901893	6.175816e-01
## APRDRG_Risk_MortalityExtreme Likelihood	2.347800631	2.188817e-08
## RACEBlack	-0.266228485	1.814111e-04
## RACEHispanic	-0.087103289	3.181046e-01
## RACEAsian	0.006590652	9.662414e-01
## RACENative American	-0.269350057	4.107313e-01
## RACEOther	-0.030889503	8.115125e-01
## LOS	-0.009580439	2.065340e-79

### Balanced Random Forests:

```
## Area under the curve: 0.9492
```

```
## 95% CI: 0.9457-0.9527 (DeLong)
```

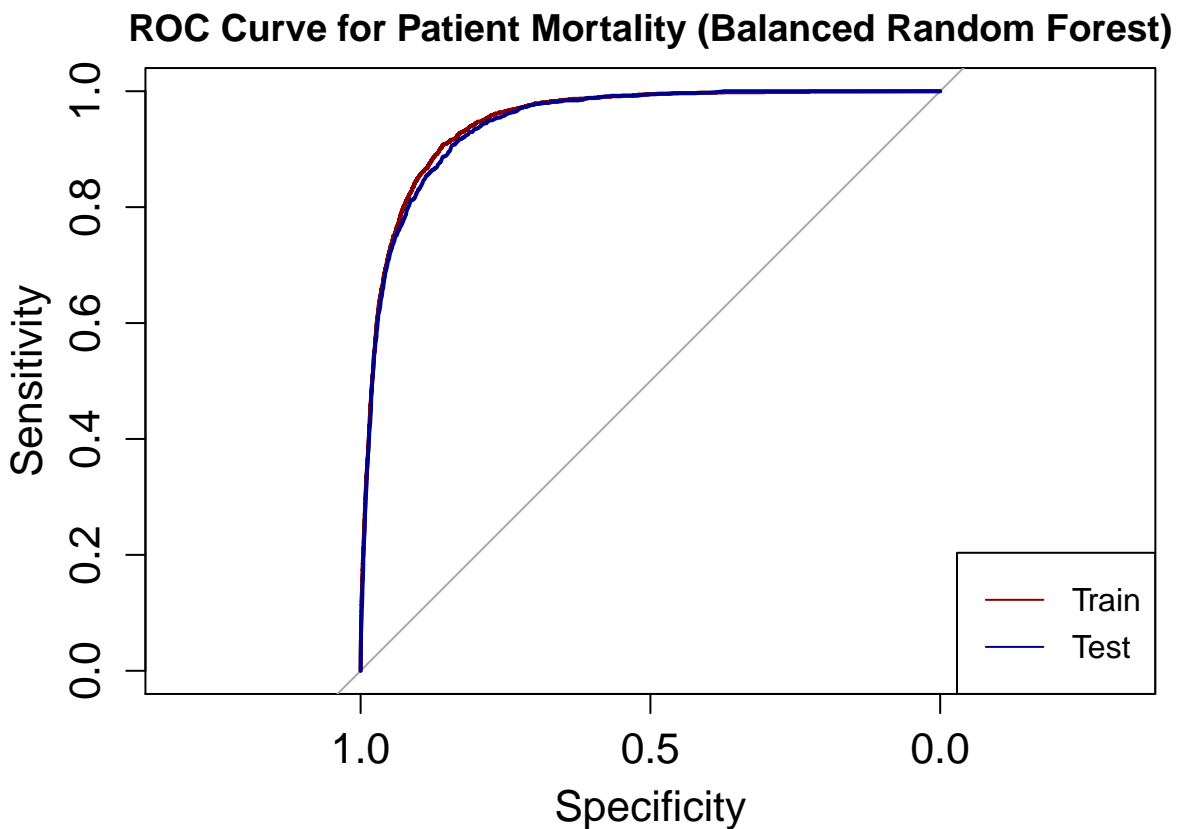
```
auc(rf_roc_test)
```

```
## Area under the curve: 0.9459
```

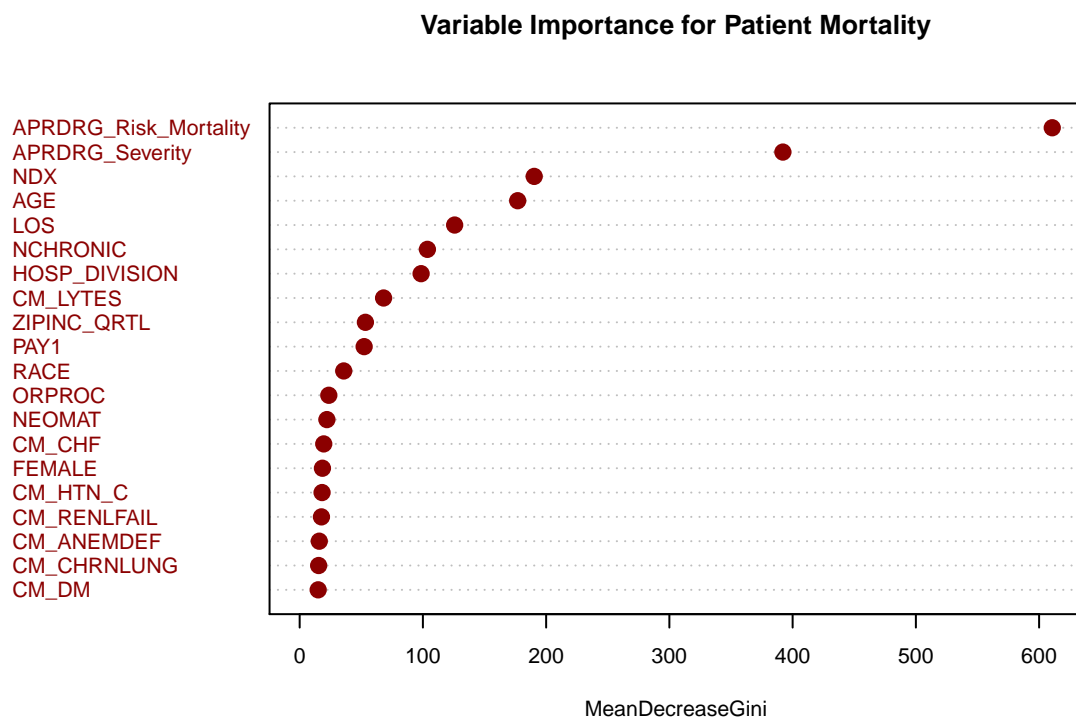
```
ci(rf_roc_test)
```

```
## 95% CI: 0.9409-0.9509 (DeLong)
```

The plot below shows the ROC curve for both the training and testing data for the balanced random forest method in predicting patient mortality.



Below is the variable importance in predicting Patient Mortality using the balanced random forests method. This indicates that the Risk Mortality and the Severity of the disease are the most important variables in predicting patient mortality. Other significant variables include age and the number of diagnoses coded on the patient's health record.



## Lasso Regression

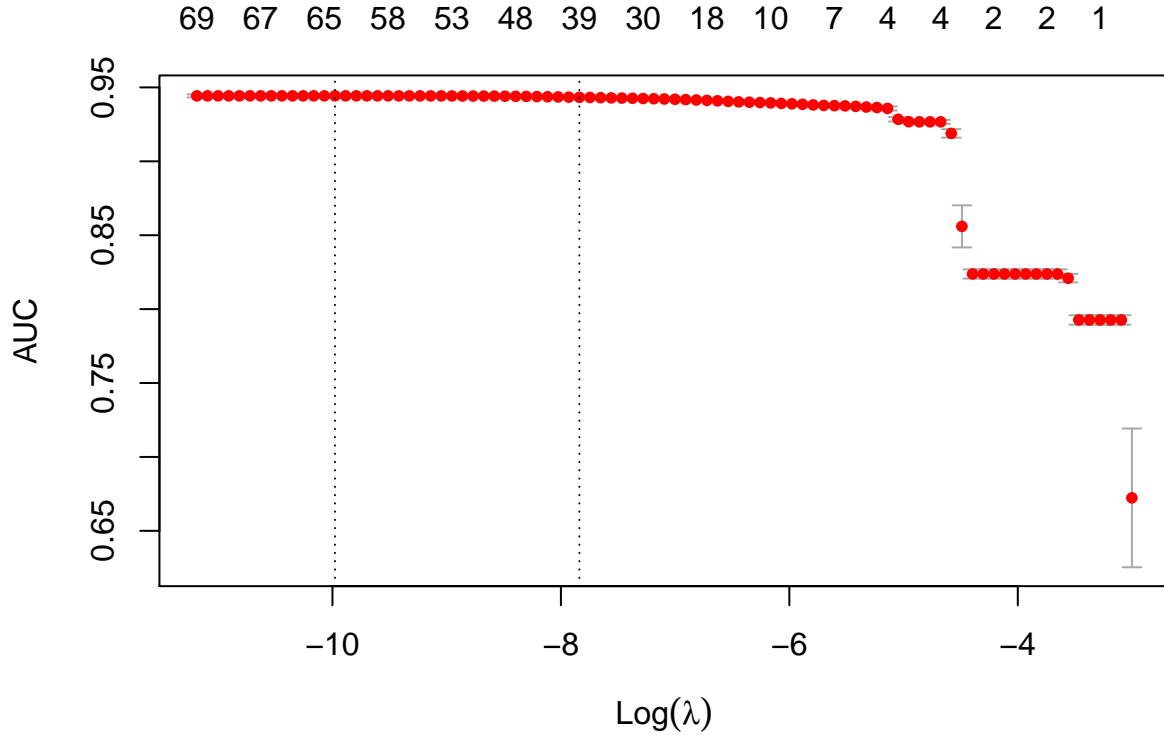
Below is the cross-validated AUC for the LASSO regression model in predicting patient mortality.

```
## [1] 0.9443041
```

Below is the cross-validated misclassification error for the LASSO regression model in predicting patient mortality.

```
##          mmce
## 0.01800968
```

Below is the plot of the cross-validated tuning process for the regression model. The numnber of coefficents in the top of the graph is greater than the 44 variables in the data because of factor variables that have multiple levels.



## Conclusions/discussion

Across all three prediction models, the AUC was above 0.90 indicated that the prediction models were more than adequate in correctly prediction patient mortality. Overall, all of the models indicated similar variables to be the most important in determining the probability of a patient dying in the hospital. These variables include the Risk of Mortality, the Severity of the Disease, the length of stay at the hospital, and age. Although there were concerns that the data was imbalanced, the results from the balanced random forests show an AUC similar to that of the logistic regression model and the LASSO regression.

The advantage of the logistic regression is that the parameter estimates can be used to determine how influential statistically significant variables are in determining the probability of death. It is much more detailed, however, the forward selection algorithm in determine which features to include severely limits the scope of variables that was examined. The balanced random forest is not as detailed in its description of the relationships between the parameters and the outcome like logistic regression, but it takes into account for all of the variables present. Because of this, the variable importance plot gives a fuller picture of how all of the variables affect the probability of patient mortality. All of the variables that were considered statistically significant in the logistic regression were considered important in the balanced random forest model, but there were some variables that the forward selection algorithm did not include. The variables that were considered important by the LASSO regression model aligned with what was considered important by the balanced random forest, although cross-validated tuning plot for LASSO regression also shows that the model performs the best with around 59 or 37 non-zero coefficient estimates. T

Overall, all three models show similar results in that the Risk Mortality variable was considered to be the most important in predicting whether the patient would die during hospitalization. Other key variables include the severity of the disease, length of stay at the hospital, and the number of diagnoses on the patient's record. As the AUC for all three models were very high and the misclassification error was low, I

am confident that we can use these models to predict patient mortality in hospitals.