

Eric_Hirsch_622_Final_Assignment

Predicting the Space Titanic Kaggle Competition

Eric Hirsch

12/3/2022

Contents

Discussion	2
Introduction	2
Prediction using the Kaggle Spaceship Titanic Data Set	3
The Business Problem	3
Data Summary	3
Distributions	3
Missing Values	4
Multicollinearity	4
Outliers	4
Data Preparation	4
Feature Engineering	4
Modelling	4
Choosing and Testing Models	4
Hyperparameter Tuning	4
Results	4
Discussion	5
The Importance of Knowing the Data	5
The Importance of Choosing the Right Model	5
The Importance of Hypertuning	5
Conclusion	6
Code	6
1. Data Exploration	6
A. Summary Statistics	6
B. Distributions: Skewedness and Outliers	7
C. Multicollinearity	15

D. Missing Values	16
E. First Pass Logistic Regression: 9th percentile	20
Data Preparation and Feature Engineering	22
1. Create groups based on the Passenger ID	22
2. Create Cabin Variables	22
3. Create Dummy Variables	23
4. Implement Interaction Features	23
5. Perform Logistic Regression With Engineered Features: 26th Percentile	23
More Complex Models	25
1. Perform SVM: 70th Percentile	25
2. Perform limited Neural Networks	29
3. Tree Algorithms 1: Perform Random Forest: 34th Percentile	30
4. Tree Algorithms 2: Perform XGBoost Untuned: 73rd Percentile	32
Hypertuning	32
1. Tune XGBoost: 78th Percentile	32
2. Tune SVM Radial: No Improvement	46
Discussion	47

Discussion

Introduction

In machine learning, we predict target variables based on input variables. For this final exercise, we will apply various machine learning algorithms to a Kaggle data set (Spaceship Titanic) in order to predict which passengers have been transported to another dimension.

While it's tempting to throw as many algorithms at the problem as possible to see what sticks, the statistical fact is that while it is rare that a poor model will perform well on a holdout set, the chances of making false conclusions based on performance increases if we simply try one model after another. Besides, if we don't understand our model and our data, and the model becomes much more difficult to troubleshoot and maintain.

When choosing models, we are balancing simplicity and complexity, and therefore tendencies to underfit or overfit. When the relationships in the data are simple and certain statistical conditions are met, parametric methods like OLS work well and have the advantage of being easily interpretable. If, for example, we are predicting height from weight, the relationship is simple enough that we can create a linear regression model and capture most of the variation that can be explained for these two variables.

When we increase our dimensions and/or complexity of relationships within the dataset, parametric methods are likely to underfit the data. Even in our simple height and weight example, if the relationship between height and weight varies considerably at lower weights, medium weights and higher weights, spline regression or another nonparametric technique will be necessary. As dimensions and complexity increases, we adopt techniques that are more powerful at morphing the data shape so that we can model the underlying structure, such as trees, SVM and neural nets.

Choosing the more complex algorithm will likely fit the training data better, but may be less interpretable and more subject to overfitting. With this in mind, each of these techniques has its advantages and disadvantages.

In my experience with earlier datasets in this class, trees will pick up autonomous clusters in the data set better than SVMs. For example, if there were a small but significant anomalous cluster of individuals for whom height and weight were inversely related, trees will incorporate the cluster while SVMs will ignore it. Of course, clusters like this might signal a missing variable, but not all of the necessary variables will be found in any given data set. Trees may be bagged (e.g., Random Forest) or boosted (e.g. xgBoost) - either will generally perform better than a single decision tree. Because xgBoost is not a lazy learner, it will often have the upper hand in fitting the training data. On the other hand, when the relationships are more systematic and class boundaries are clear, SVMs may perform better because the kernel trick allows SVMs to radically change the data shape in order to find the class boundary. SVMs can also perform better when there is less data.

One of the biggest advantages of neural networks is that they effectively do the feature engineering for you if you can apply enough layers. They are also subject to the “double descent” phenomenon, which helps with managing underfitting. However, for a student using a home computer like myself, it’s often impractical to take advantage of these facts as the algorithm would run too long. Neural networks, like SVMs, also powerfully change the data shape in order to find class boundaries.

Accurate prediction depends not only on algorithm choice. We also need to engineer features (except possibly in very large neural nets) and tune hyperparameters. We also need to choose metrics that tell us whether or not our model is effective.

Prediction using the Kaggle Spaceship Titanic Data Set

For this exercise I’ve chosen a Kaggle Competition – the Kaggle Spaceship data set. The advantages of using this a competition data set are that we can compare our performance those of others. Achieving 90% on a holdout set in and of itself tells us nothing - we don’t know if achieving 95% would have been easy or impossible. In this competition, the 2,000 or so submitted accuracies on the leaderboard range from about 76% to 82%, which gives us a good idea of how well our model is working.

The main disadvantages of this data set are that the data is made up and the scenario a bit far-fetched. However, I wanted a data set that had a simple class as a target, as opposed to an image example, and the standard Titanic data set has been over analyzed, this was one of the few good choices.

The Business Problem

In the year 2912, the Spaceship Titanic, an interstellar passenger liner with almost 13,000 passengers on board, collided with a spacetime anomaly hidden within a dust cloud. Though the ship stayed intact, almost half of the passengers were transported to an alternate dimension. Our job is to predict which passengers were transported by the anomaly using records recovered from the spaceship’s damaged computer system.

Data Summary

The data set consists of 8693 records and 13 variables, including spending on the ship’s various amenities (VR Deck, Spa, Room Service, Food Court, Shopping Mall, cabin number, whether the individual was traveling with the group, whether the individual was a VIP, planet of origin and destination, and so on. These columns map to some degree with the original Titanic database. The target variable, Transported, is roughly equally distributed between false (4315) and true (4378).

Distributions

Missing Values 1073, or 12%, of records have missing values. The vast majority of missing values are found in the amenity expenditure columns. Oddly, the amenity expenditure rows with missing values are completely independent of each other - there are no records where more than one of these values is missing. This may be an artifact of the fact that the data is manufactured. In order to confirm that there is no systematic relationship between missing data and the target variable, we look at the Chi square between the target and a flag designating missing data. We do this for each amenity expenditure column and find no relationship between missing data and the target variable. We therefore eliminate rows with missing values for the training set. The test set, we impute the median.

Multicollinearity There is very little, even surprisingly little, multicollinearity in the database. In the case of variables that track spending on amenities this is most surprising, and may suggest that passengers were working within a budget and only spent money on the activities they liked most.

Outliers All of the spending variables are highly skewed, with very large ending occurring at the very end of the distribution. However, as most of our techniques are robust for outliers, records with extreme values remain in the database, as there is no reason to think that the spending is a data entry error or an anomalous occurrence.

Data Preparation

Feature Engineering The data set holds a number of opportunities for feature engineering. Through testing, it was found that the following new features were significant in predicting transportation. They are:

Modelling

Choosing and Testing Models

Hyperparameter Tuning

Results

Accuracy_Holdout	Accuracy_Kaggle	Position_Kaggle	Model	Kaggle.Percentile
0.8000	0.80383	529	XGBoost Tuned (14 Rounds)	78.2%
0.7900	0.80243	635	XGBoost Untuned (55 Rounds)	73.8%
0.7715	0.80149	726	SVM Rad	70.0%
0.7741	0.80032	943	SVM Rad Lower Sigma	61.1%
0.7669	0.79682	1035	SVM Rad Higher Sigma	57.3%
0.7800	0.79775	1047	RVM Linear	56.8%
0.7800	0.79611	1084	SVM poly	55.3%
0.7900	0.78793	1586	Random Forest	34.5%
0.7800	0.77975	1789	Logistic Regression with features	26.2%
0.7600	0.77881	1805	Neural Net (severely limited due to computer power)	25.5%
0.7800	0.69227	2194	Logistic Regression First Pass	9.5%

Discussion

For this exercise, we entered a Kaggle competition (Space Titanic) in order to make predictions and test their accuracy against other competition participants. The target variable was whether or not a passenger was transported to another dimension while the ship was moving through a dangerous and destructive space anomaly. We reached the 78th percentile, with an idea of further modifications which might help us do better.

In a sense, measuring performance in a Kaggle is a cheat - after a time the “unknown” data in the Kaggle test set becomes increasingly “known” as you test more and more models against it. However, it also provides a real-world check on how well your model is actually performing.

After an initial logistic regression, we improved performance in three ways:

The Importance of Knowing the Data

Understanding the data is a key element in solid prediction. First, we needed to clean the data of anomalies (missing values, outliers, etc.). In our case, missing values appear to be at random and outliers didn’t present a problem.

Second, an understanding of the data highly conditioned our ability to perform feature engineering. Transported passengers tended to spend their money at the food court rather than get room service, they occupied lower-class cabins, and many of them were in cryosleep during the journey. Knowing this helped us parse out cabin locations and discover interaction terms within the data.

The Importance of Choosing the Right Model

Understanding what models are doing and how is a key part of prediction. We compared tree models, svm, neural net, and linear regression.

The first task is to understand the requirements of the business problem. In this case, we have a partial roster of passengers where we know who was transported and who wasn’t. As for the rest of the passengers, it’s our job to predict which of them were transported as well. Insight into the data is only necessary insofar as it helps us make better predictions. Accuracy is our metric, as this is the metric used in the Kaggle (there is no penalty or false positives or false negatives which might suggest a different metric). This suggests we should use the most complex model we can that has the highest accuracy.

We used tenfold cross validation and compared models on the metric of accuracy – our results varied widely. There was some, but not perfect, match between the accuracy predicted by the cross validation, and the accuracy achieved in the Kaggle. For example, going on cross validation alone, the svm with the radial kernel underperformed compared to a linear and polynomial kernel - however, in the Kaggle it significantly outperformed both kernels. Of course, without the Kaggle we would not have known this and would not have chosen the radial kernel to put into production.

On the other hand, our best model (xgboost) also performed best on the Kaggle. We suspect the neural net would have performed well also, but we did not have the computer power to hypertune it properly.

The Importance of Hypertuning

The caret package in R automatically hypertunes models on basic parameters and chooses the most optimal based on the metric determined by the user. With the exception of xgBoost, all of our models were tuned that way. This left two tasks – tune xgBoost, and experiment with manual tuning of other models.

The tuned xgBoost model (number of rounds was reduce from 55 to 14) increased accuracy from 79% to 80%, and improved our percentile position in the Kaggle from the 73rd to the 78th percentile.

The fact that our radial-kernal SVM had lower cross validation accuracy but higher Kaggle accuracy might have a number of different causes - but one is the possibility that there are idiosyncrasies in the training set that are not found in the Kaggle set. In this case, we can increase sigma to decrease the risk of overfitting.

As we experimented by increasing sigma and decreasing sigma on the radial SVM model. Reducing sigma (tighter fit) improved accuracy during cross validation, but not when applied to the Kaggle. Because the differences were very small (.0026 improvement in accuracy), the models probably only varied by a handful of predictions. Increasing sigma (looser fit) reduced accuracy both in cross validation and on the Kaggle set.

Conclusion

All of our models came in between 77% and 80% accuracy (based on a no information rate of 51%) after tenfold cross validation and optimal hypertuning. This was also the range for individuals who submitted predictions in the Kaggle competition - and so a score of 77% achieved a much lower percentile in the competition than a score a few percentage points higher. Without the Kaggle, we still would have chosen the xgBoost model, but if we were going to choose an SVM we would likely have chosen a less than optimal kernel.

While it makes sense that the boosted tree model outperforms the bagged tree model, it more difficult to speculate on why tree model forms better than the Support Vector Machine. Given the small difference in accuracy between them (.003), it is possible that there is no significant difference, and that with slightly different feature engineering the SVM would be the better model.

There would be a number of ways to improve accuracy even further:

1. Search for more features, particularly interaction features between cryosleep and spending, and understanding better how groupings and cabins work.
2. Borrow a more powerful computer to perform a properly tuned neural net, and then create a prediction set based on neural net, xgBoost and radial SVM.
3. Experiment more with manual tuning of models.

In the meantime, however, we are satisfied with the nearly 80th percentile compared to the 26th percentile using logistical regression which we would have achieved prior to this class.

Code

1. Data Exploration

A. Summary Statistics

The data consists of 8693 records and 14 variables (6 numeric and 8 character). There are a number of missing values and what appear to be skewed distributions among the numeric variables.

```
## PassengerId      HomePlanet      CryoSleep      Cabin
## Length:8693      Length:8693      Length:8693      Length:8693
## Class :character  Class :character  Class :character  Class :character
## Mode  :character  Mode  :character  Mode  :character  Mode  :character
##
##
##
##
## Destination      Age      VIP      RoomService
```

```

## Length:8693      Min.   : 0.00      Length:8693      Min.   :    0.0
## Class :character  1st Qu.:19.00      Class :character  1st Qu.:    0.0
## Mode  :character  Median :27.00      Mode  :character  Median :    0.0
##                                     Mean  :28.83      Mean   :   224.7
##                                     3rd Qu.:38.00      3rd Qu.:   47.0
##                                     Max.   :79.00      Max.    :14327.0
##                                     NA's   :179       NA's    :181
##      FoodCourt      ShoppingMall      Spa      VRDeck
## Min.   :    0.0      Min.   :    0.0      Min.   :    0.0      Min.   :    0.0
## 1st Qu.:    0.0      1st Qu.:    0.0      1st Qu.:    0.0      1st Qu.:    0.0
## Median :    0.0      Median :    0.0      Median :    0.0      Median :    0.0
## Mean   :   458.1      Mean   :   173.7      Mean   :   311.1      Mean   :   304.9
## 3rd Qu.:    76.0      3rd Qu.:    27.0      3rd Qu.:    59.0      3rd Qu.:    46.0
## Max.   :29813.0      Max.   :23492.0      Max.   :22408.0      Max.   :24133.0
## NA's   :183         NA's   :208         NA's   :183         NA's   :188
##      Name      Transported
## Length:8693      Min.   :0.0000
## Class :character  1st Qu.:0.0000
## Mode  :character  Median :1.0000
##                                     Mean  :0.5036
##                                     3rd Qu.:1.0000
##                                     Max.   :1.0000
##
## 'data.frame':    8693 obs. of  14 variables:
## $ PassengerId : chr  "0001_01" "0002_01" "0003_01" "0003_02" ...
## $ HomePlanet  : chr  "Europa" "Earth" "Europa" "Europa" ...
## $ CryoSleep   : chr  "False" "False" "False" "False" ...
## $ Cabin       : chr  "B/O/P" "F/O/S" "A/O/S" "A/O/S" ...
## $ Destination : chr  "TRAPPIST-1e" "TRAPPIST-1e" "TRAPPIST-1e" "TRAPPIST-1e" ...
## $ Age         : num  39 24 58 33 16 44 26 28 35 14 ...
## $ VIP         : chr  "False" "False" "True" "False" ...
## $ RoomService : num  0 109 43 0 303 0 42 0 0 0 ...
## $ FoodCourt   : num  0 9 3576 1283 70 ...
## $ ShoppingMall: num  0 25 0 371 151 0 3 0 17 0 ...
## $ Spa         : num  0 549 6715 3329 565 ...
## $ VRDeck      : num  0 44 49 193 2 0 0 NA 0 0 ...
## $ Name        : chr  "Maham Ofracculy" "Juanna Vines" "Altark Susent" "Solam Susent" ...
## $ Transported : num  0 1 0 0 1 1 1 1 1 1 ...

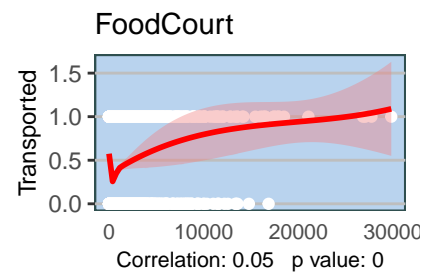
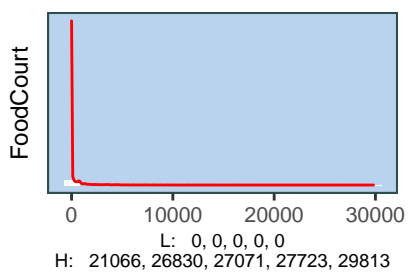
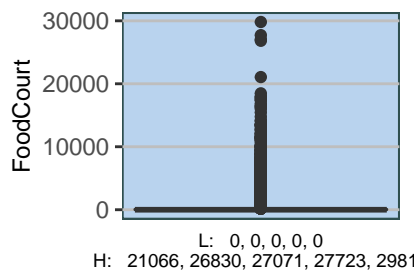
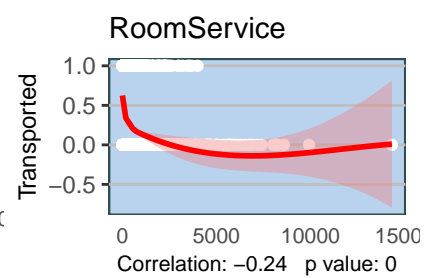
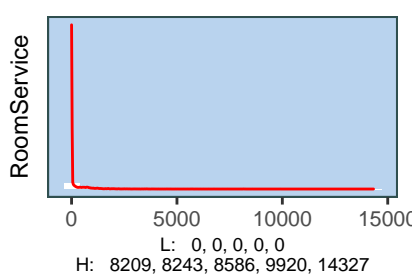
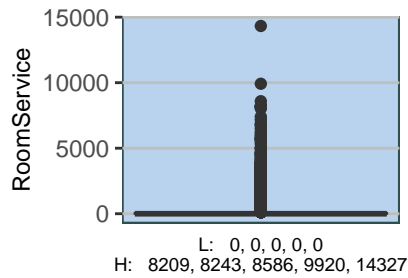
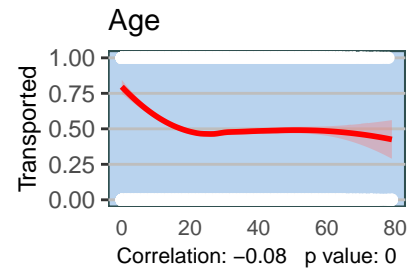
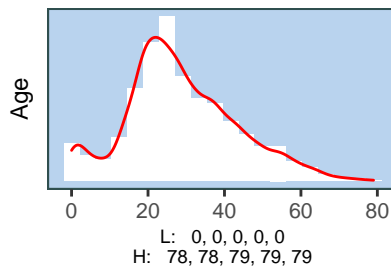
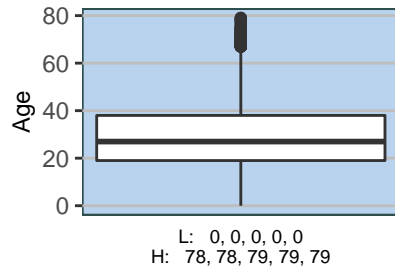
## dfTrain$Transported      n
## 1                        0 4315
## 2                        1 4378

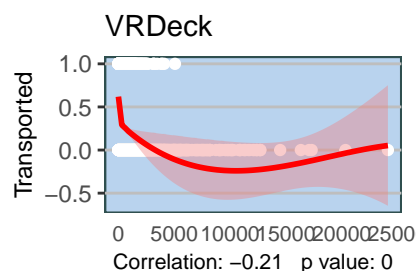
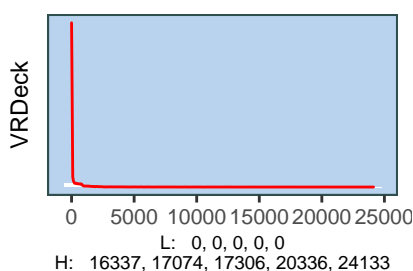
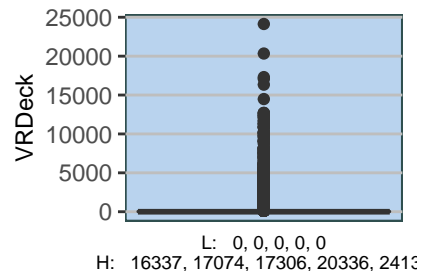
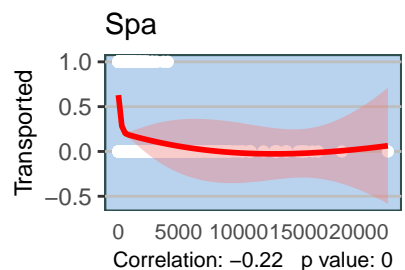
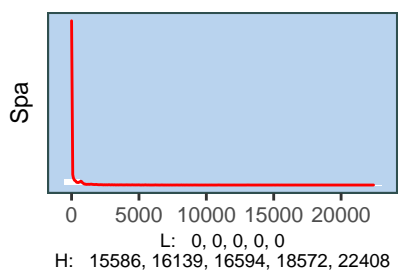
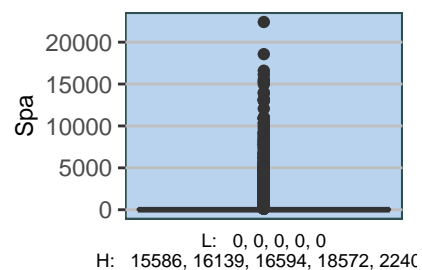
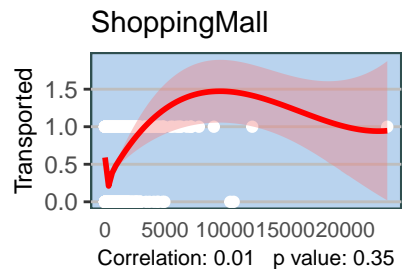
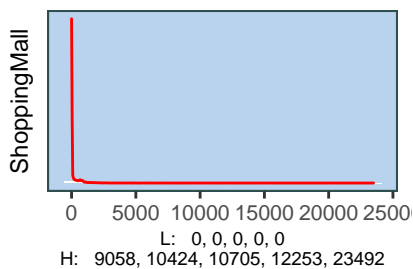
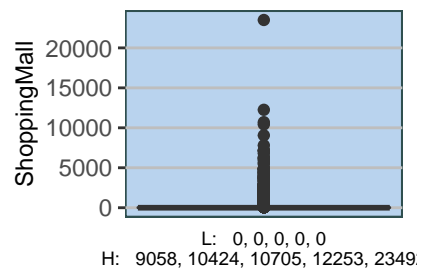
```

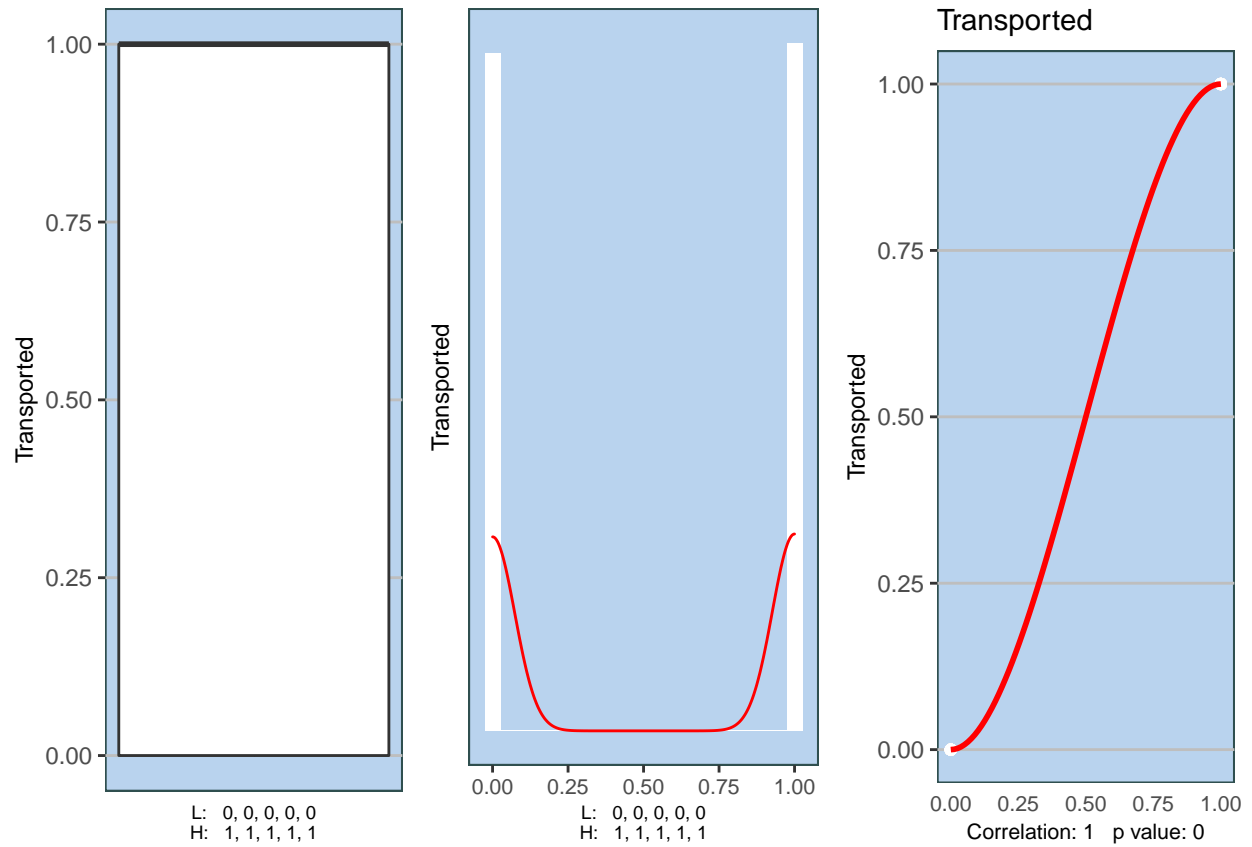
B. Distributions: Skewedness and Outliers

We at the distribution of all numeric variables. Spending variables are highly skewed - most passengers spend no money while a few spend a great deal. We can see that spending on luxuries (the spa, room service, etc.) is strongly negatively correlated with being transported - this supports the supposition that the rich were spared. Spending on more popular amenities like the food court and shopping mall are also negatively correlated but less so. Age has a small negative correlation as well.

We decide not to log transform the numeric variables as normal distributions for predictors are not required by our models and interpretability suffers.

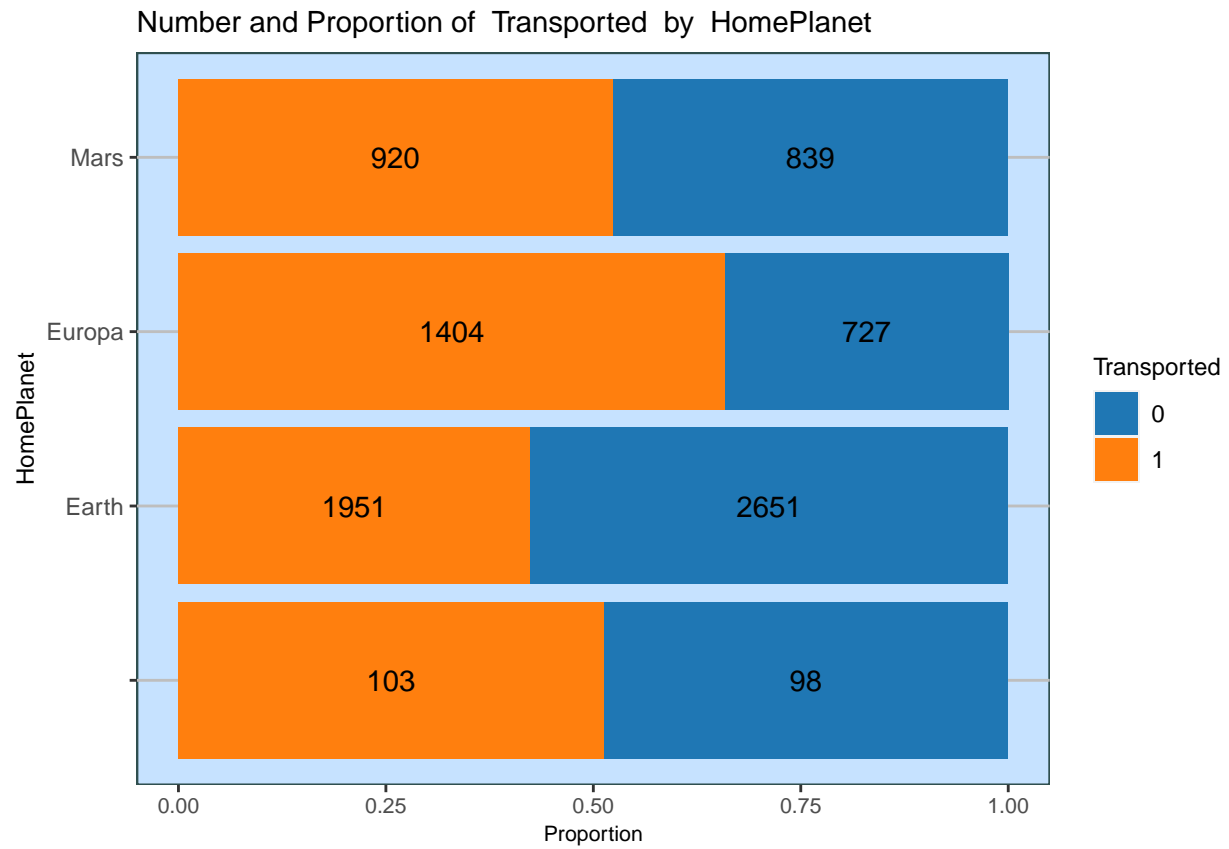




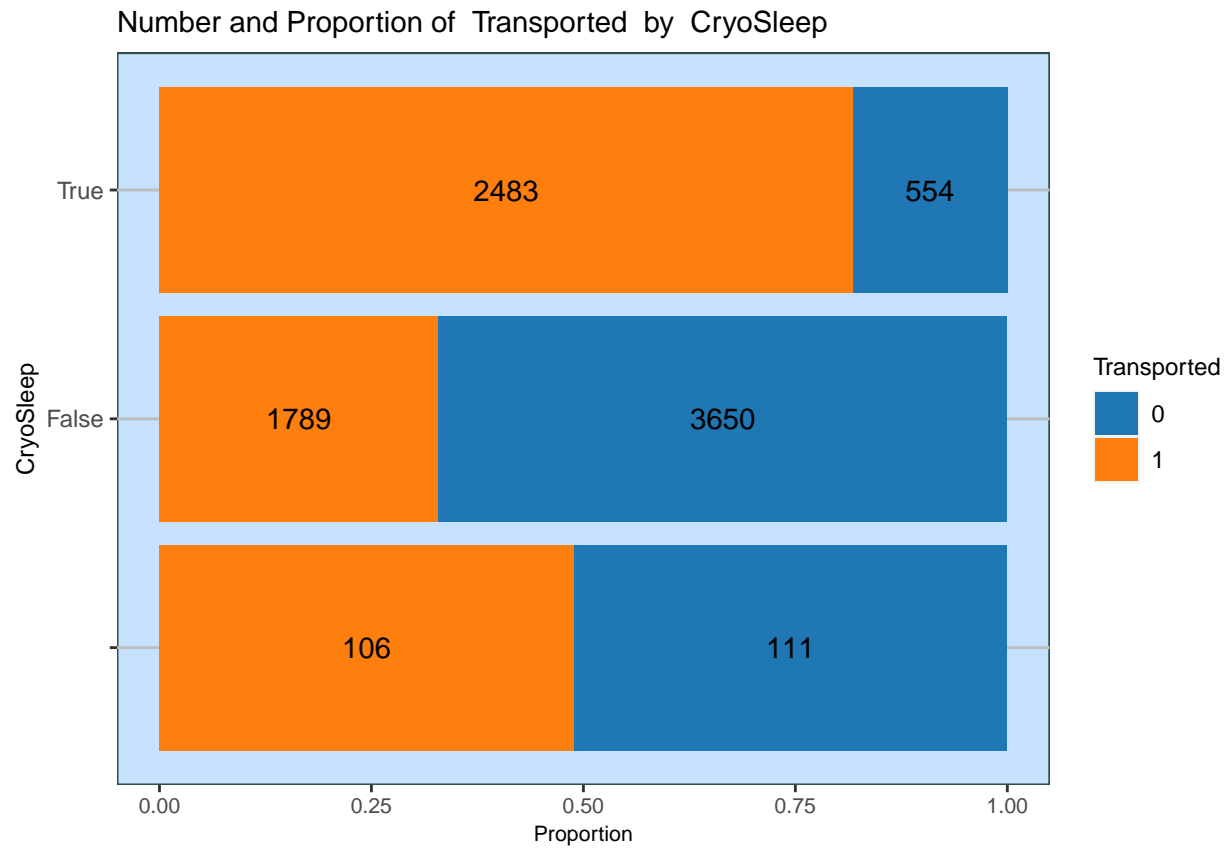


Here we look at count plots for character variables. Home and destination have an association with transported, but cryoSleep is especially important - over 75% of those in cryosleep were transported.

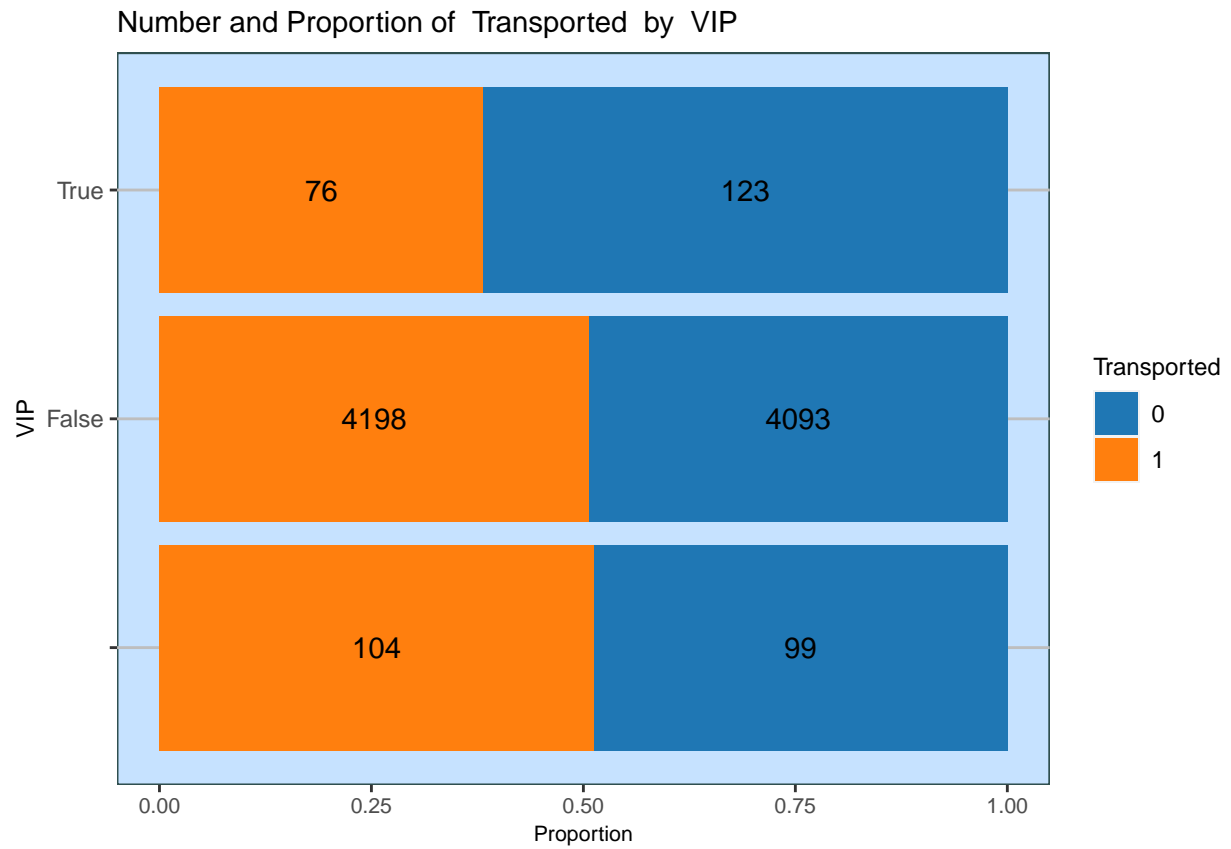
```
## [[1]]
```



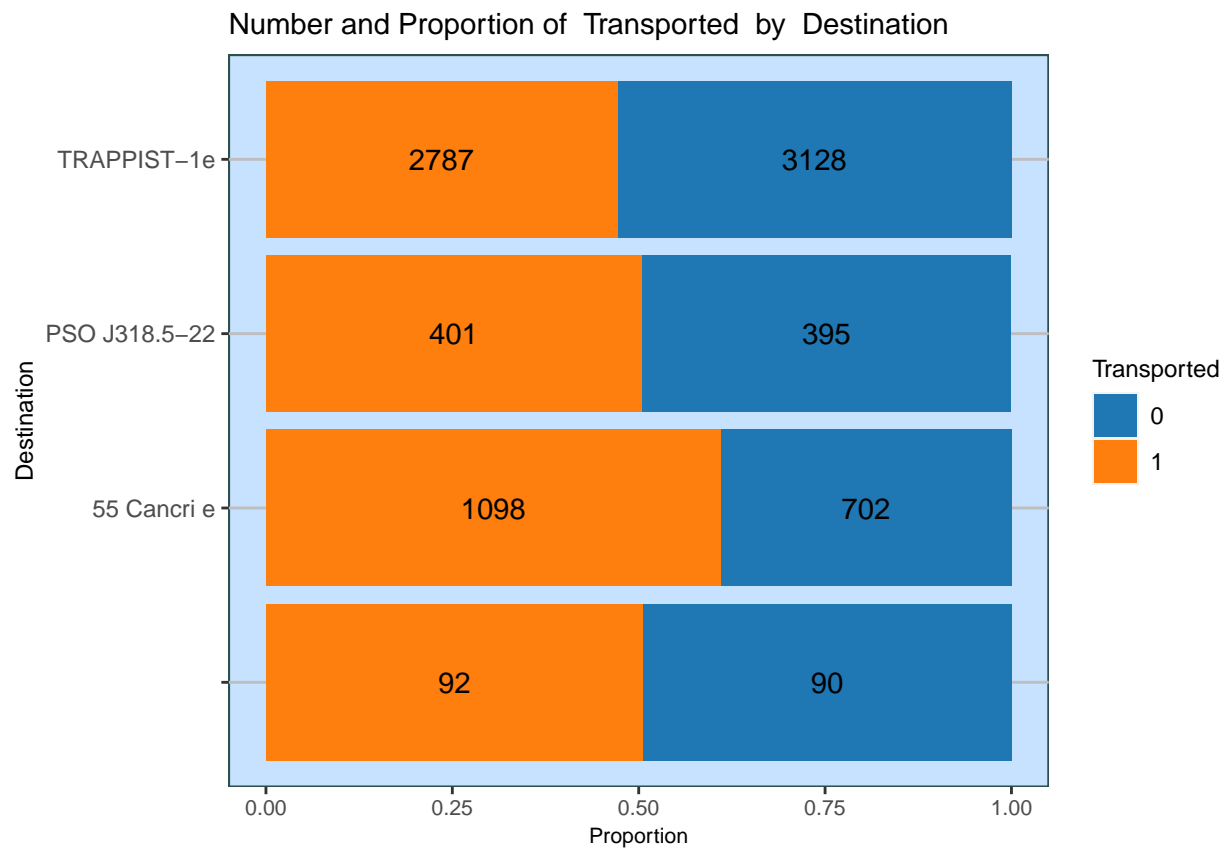
```
##  
## [[2]]
```



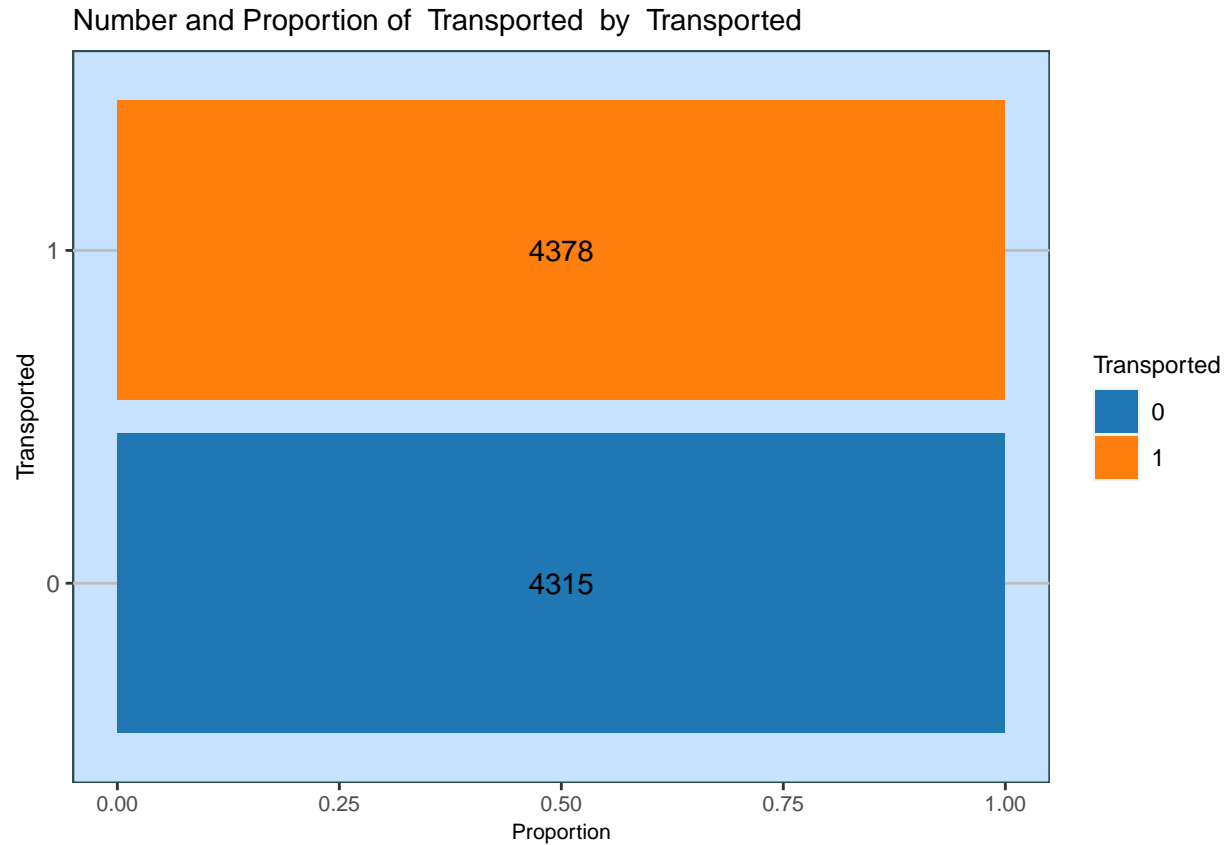
```
##  
## [[3]]
```



```
##  
## [[4]]
```



```
##  
## [[5]]
```

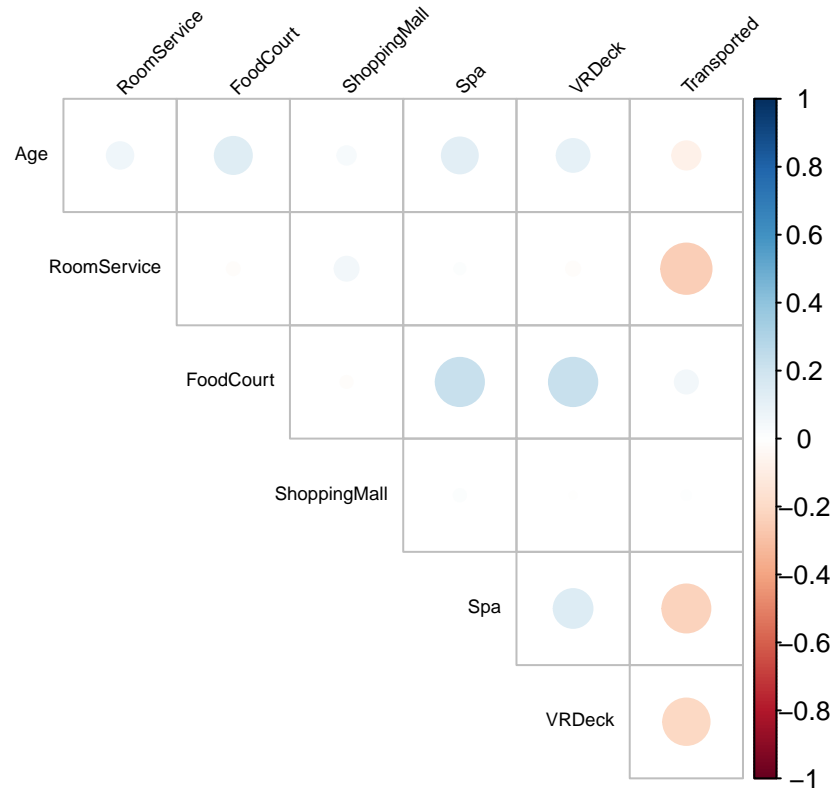


Because we have no reason to think that outliers are errors in data entry or significant data anomalies, and because our algorithms are relatively resistant to outliers, we do not remove outliers from the dataset.

C. Multicollinearity

While we were aware of the correlations with Transported, it is interesting to note that the correlations among various forms of spending are actually quite mild. We do not need to address multicollinearity in any systematic way.

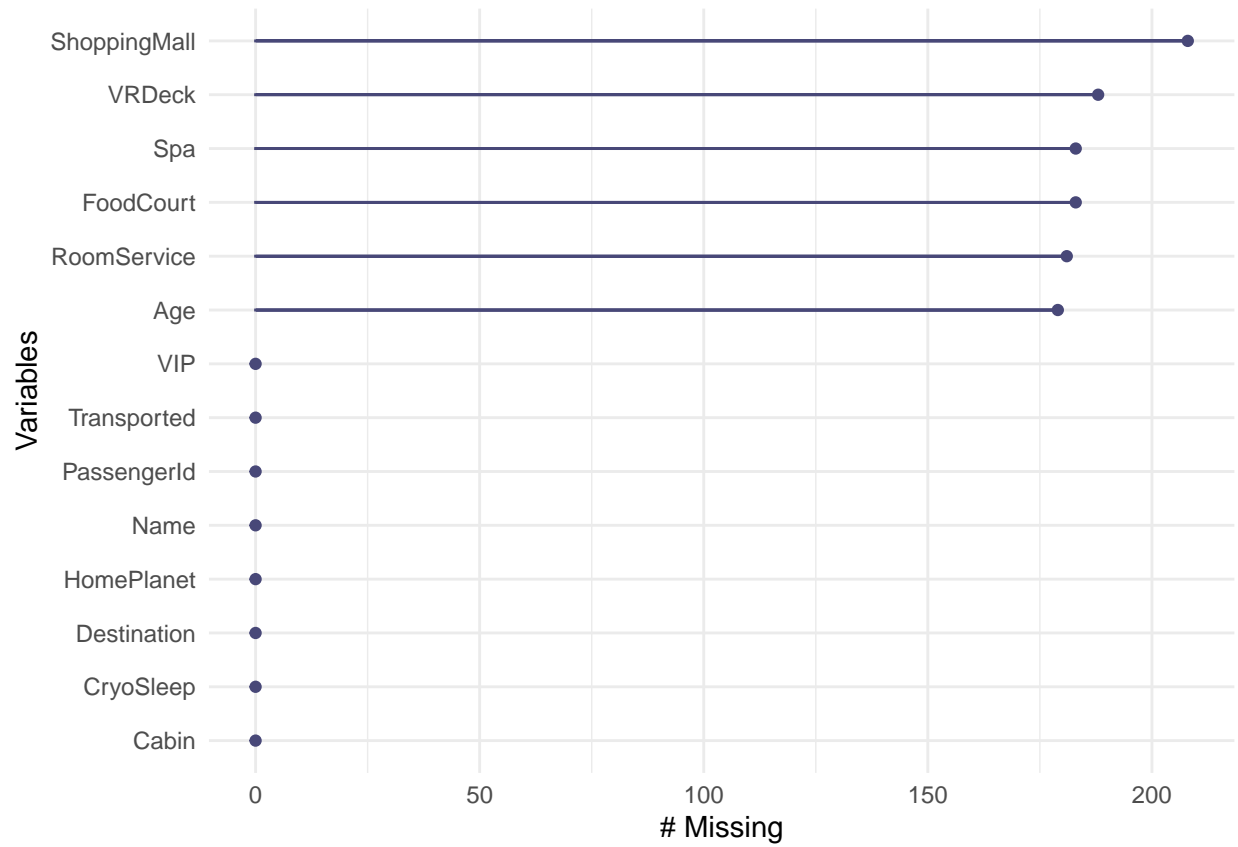
Heatmap for Multicollinearity Analysis



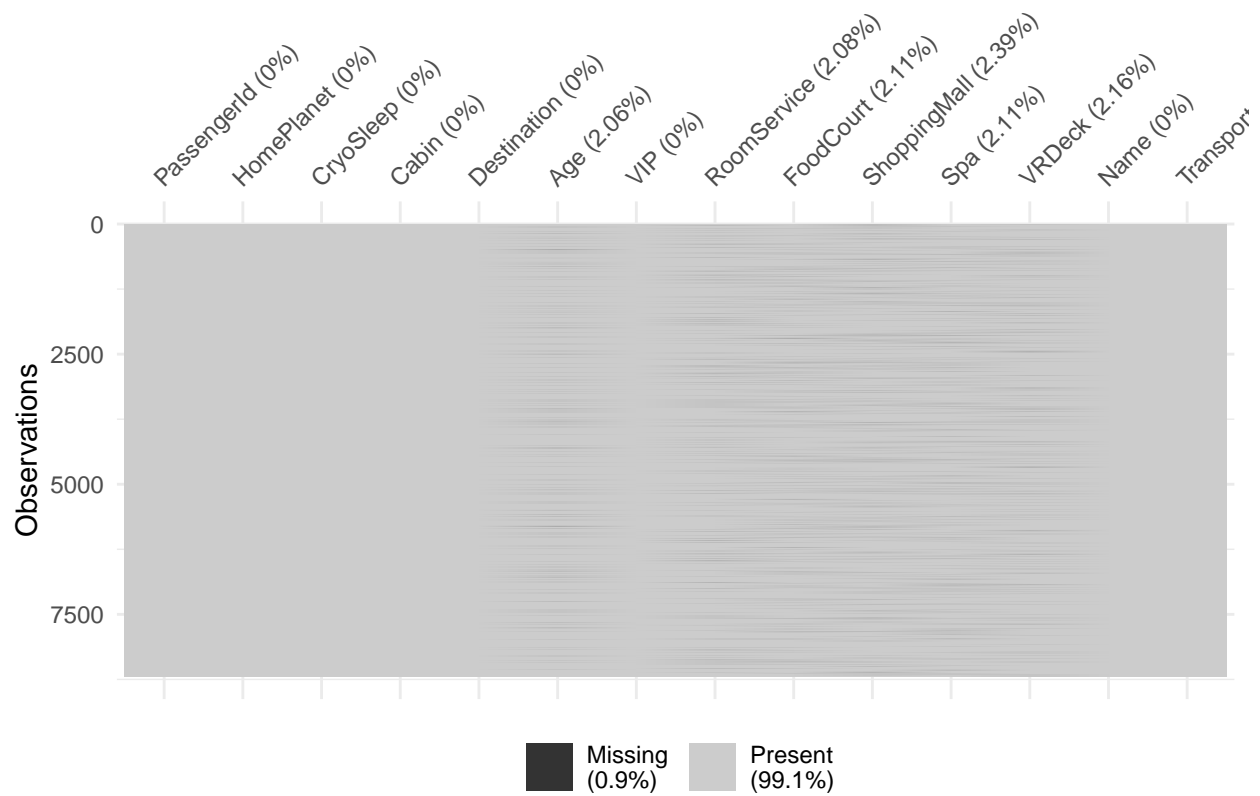
D. Missing Values

Missing values mainly appear for the amenities spending variables in the dataset. There are over 1,000 (12% of the database).

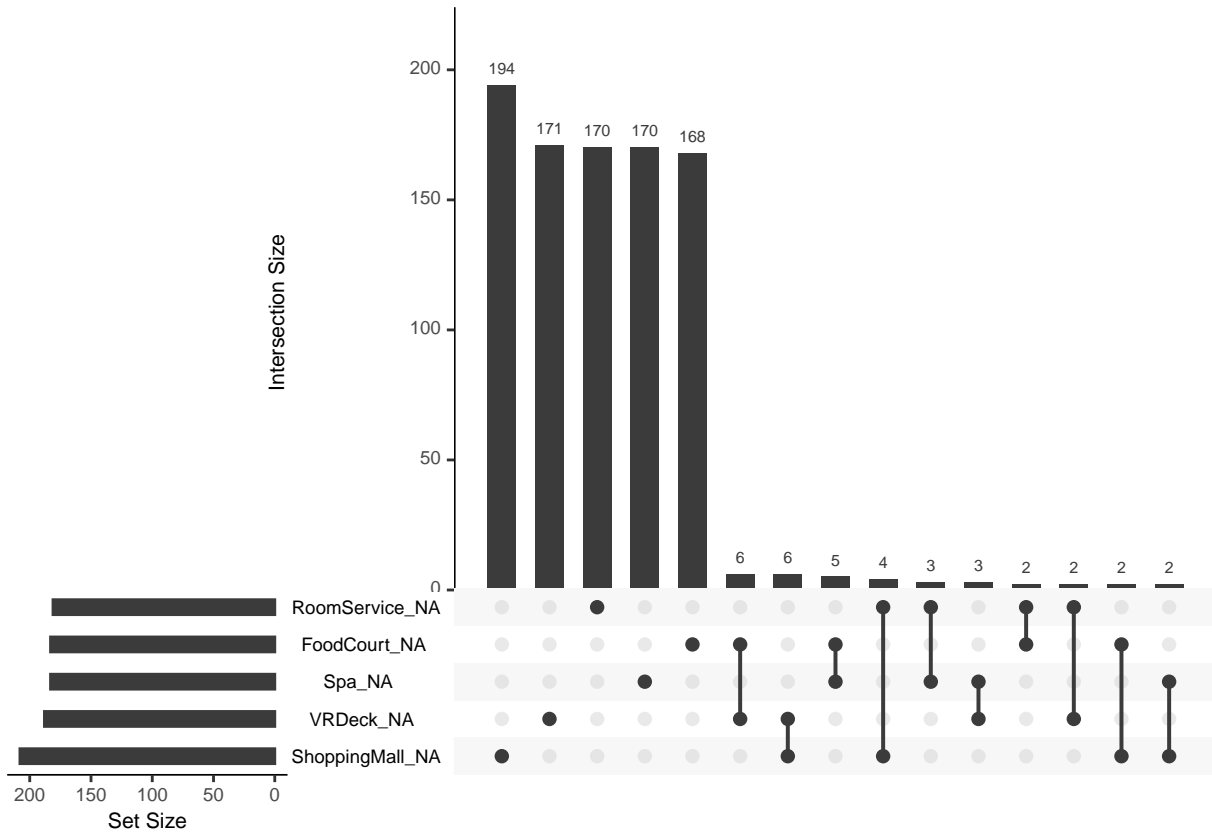
```
## [[1]]
```

```
##  
## [[2]]
```



```
##  
## [[3]]
```



The missing values are not correlated with each other, suggesting they are probably missing at random. To further support this hypothesis we create flags for missing values and perform Chi Square tests against the target variable. None of the flags are significant. We will therefore remove the records with missing values from the training set (we will do this after some feature engineering), and impute the median for the test set.

```
##
## Pearson's Chi-squared test with Yates' continuity correction
##
## data:  table(dfMissingRecordsFlagAny)
## X-squared = 0.15887, df = 1, p-value = 0.6902

##
## Pearson's Chi-squared test with Yates' continuity correction
##
## data:  table(dfMissingRecordsFlag_SPA)
## X-squared = 0.0098187, df = 1, p-value = 0.9211

##
## Pearson's Chi-squared test with Yates' continuity correction
##
## data:  table(dfMissingRecordsFlag_FoodCourt)
## X-squared = 0.89665, df = 1, p-value = 0.3437

##
```

```
## Pearson's Chi-squared test with Yates' continuity correction
##
## data:  table(dfMissingRecordsFlag_VRDeck)
## X-squared = 0.1728, df = 1, p-value = 0.6776

##
## Pearson's Chi-squared test with Yates' continuity correction
##
## data:  table(dfMissingRecordsFlag_ShoppingMall)
## X-squared = 1.5072, df = 1, p-value = 0.2196

##
## Pearson's Chi-squared test with Yates' continuity correction
##
## data:  table(dfMissingRecordsFlag_RoomService)
## X-squared = 1.3229, df = 1, p-value = 0.2501
```

E. First Pass Logistic Regression: 9th percentile

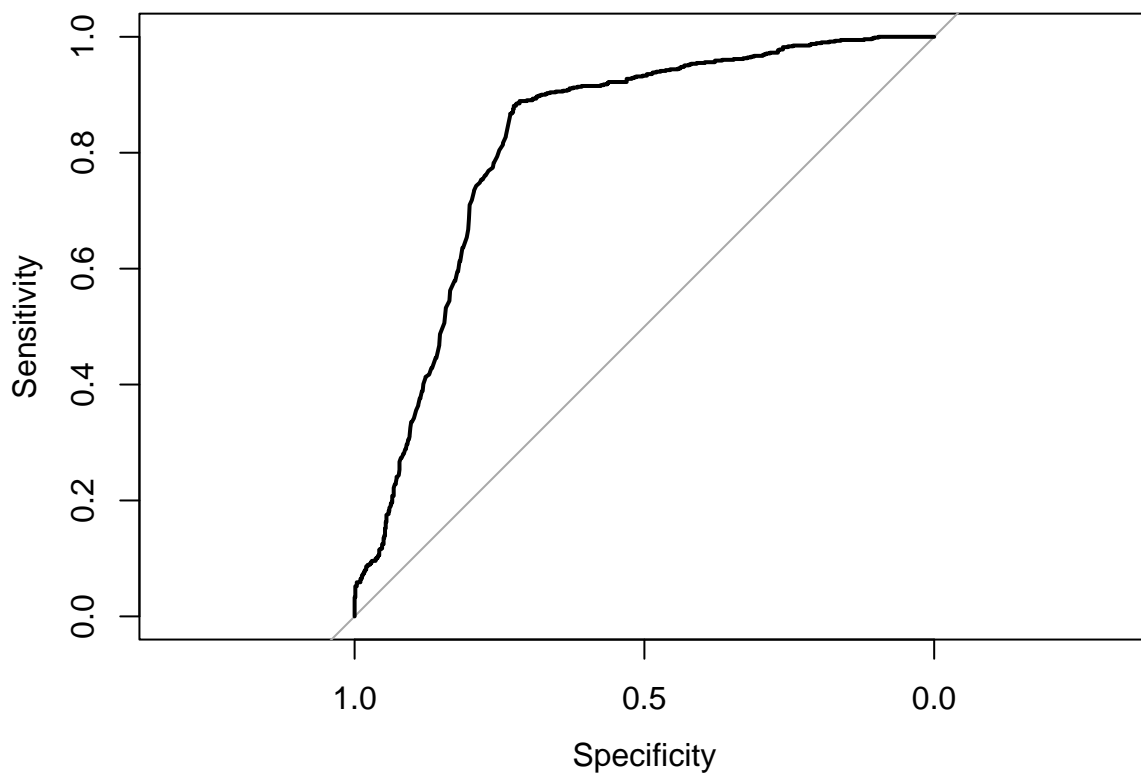
We perform a logistic regression with what we have and post to Kaggle just to get a baseline. Accuracy on training is 77%, significantly better than the 51% no information rate, but gives us only 69% on the Kaggle set which puts us at the 9th percentile.

```
##
## Call:
## glm(formula = fla, family = "binomial", data = train_reg)
##
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
## -2.2014  -0.8333   0.1278   0.8613   4.8998
##
## Coefficients:
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept)  7.174e-01  6.616e-02  10.844 < 2e-16 ***
## Age          3.329e-03  2.128e-03   1.564  0.11776
## RoomService -2.204e-03  1.095e-04 -20.123 < 2e-16 ***
## FoodCourt    8.032e-04  4.826e-05  16.645 < 2e-16 ***
## ShoppingMall 1.678e-04  6.178e-05   2.716  0.00662 **
## Spa         -2.371e-03  1.287e-04 -18.424 < 2e-16 ***
## VRDeck       -2.277e-03  1.225e-04 -18.587 < 2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
##      Null deviance: 8437.7  on 6087  degrees of freedom
## Residual deviance: 6226.1  on 6081  degrees of freedom
## (867 observations deleted due to missingness)
## AIC: 6240.1
##
## Number of Fisher Scoring iterations: 7
##
## Confusion Matrix and Statistics
```

```

##
##           Reference
## Prediction  0   1
##           0 536  72
##           1 265 659
##
##           Accuracy : 0.78
##           95% CI : (0.7584, 0.8005)
##           No Information Rate : 0.5228
##           P-Value [Acc > NIR] : < 2.2e-16
##
##           Kappa : 0.5642
##
## Mcnemar's Test P-Value : < 2.2e-16
##
##           Sensitivity : 0.6692
##           Specificity : 0.9015
##           Pos Pred Value : 0.8816
##           Neg Pred Value : 0.7132
##           Prevalence : 0.5228
##           Detection Rate : 0.3499
##           Detection Prevalence : 0.3969
##           Balanced Accuracy : 0.7853
##
##           'Positive' Class : 0
##

```



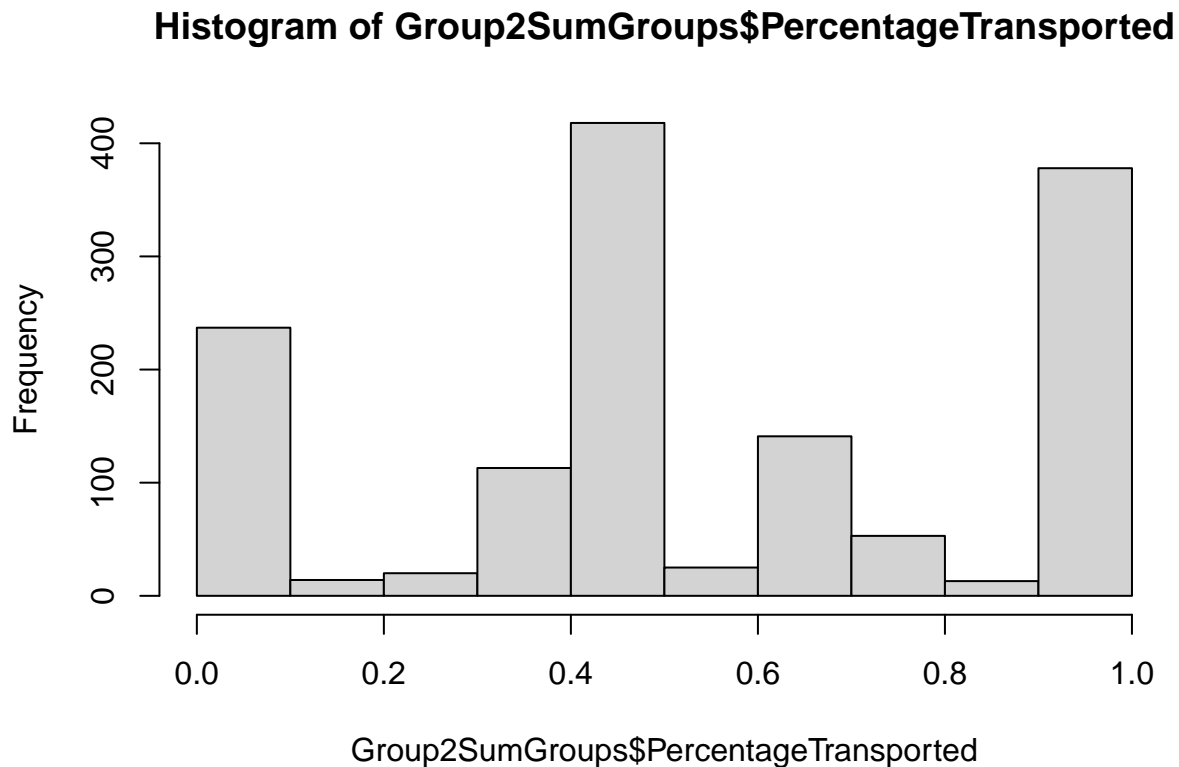
```
##
## Call:
## roc.default(response = dfPred_raw$class, predictor = dfPred_raw$predict_reg,      plot = TRUE)
##
## Data: dfPred_raw$predict_reg in 801 controls (dfPred_raw$class 0) < 731 cases (dfPred_raw$class 1).
## Area under the curve: 0.8183
```

Data Preparation and Feature Engineering

1. Create groups based on the Passenger ID

Passenger IDs are constructed to identify passengers travelling in groups. We create groupings from the ID.

How likely is it that if the majority of members of a group transported, then they all transported? Only somewhat likely. A histogram shows the distribution of percentages of transported within groups. Most often, half the members transported and half did not.



2. Create Cabin Variables

Cabin variables consist of 3 parts in the form of a/b/c which indicate the location of the cabin on the ship. Here we extract out parts “a” and “c” - b appears to have no influence on the target.

3. Create Dummy Variables

Now that we have engineered Cabin, we create dummy variables to handle category variables throughout the dataset.

4. Implement Interaction Features

5. Perform Logistic Regression With Engineered Features: 26th Percentile

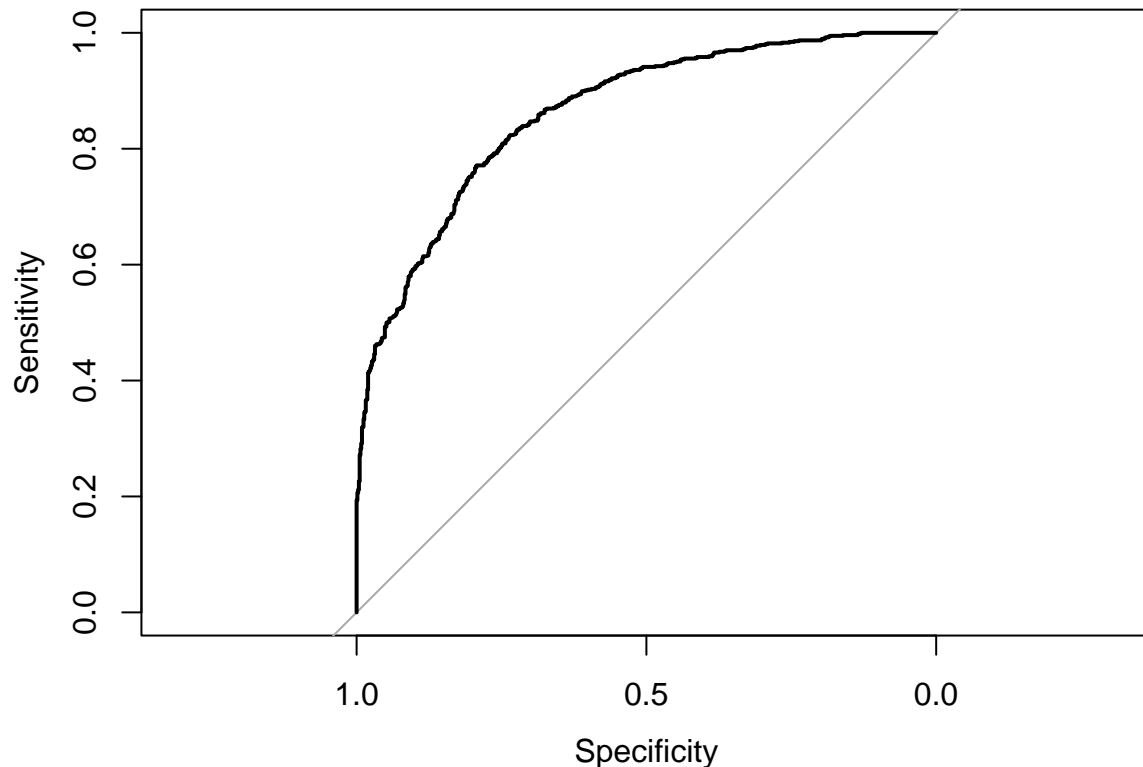
We perform a logistic regression with what we have and post to Kaggle just to get a baseline. Accuracy on training is 79% (compared to 77% on the untransformed training set), but more importantly, this gives us only 78% on the Kaggle set which puts us at the 926th percentile.

```
##
## Call:
## glm(formula = fla, family = "binomial", data = train_reg)
##
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
## -3.0545  -0.6575   0.0253   0.6986   3.3416
##
## Coefficients:
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept)    4.001e-01  1.200e-01   3.335 0.000853 ***
## Age           -7.428e-03  2.533e-03  -2.932 0.003363 **
## RoomService   -1.547e-03  1.096e-04 -14.113 < 2e-16 ***
## FoodCourt      5.075e-04  4.767e-05  10.647 < 2e-16 ***
## ShoppingMall    6.959e-04  1.057e-04   6.583 4.62e-11 ***
## Spa           -2.153e-03  1.262e-04 -17.069 < 2e-16 ***
## VRDeck        -2.000e-03  1.227e-04 -16.296 < 2e-16 ***
## InAGroup       1.462e-01  8.036e-02   1.819 0.068919 .
## HomePlanet_    5.212e-01  2.244e-01   2.323 0.020191 *
## HomePlanet_Europa 1.731e+00  2.658e-01   6.513 7.37e-11 ***
## HomePlanet_Mars 4.984e-01  1.142e-01   4.366 1.27e-05 ***
## CryoSleep_     2.776e-01  2.038e-01   1.363 0.172999
## CryoSleep_True 1.324e+00  9.803e-02  13.503 < 2e-16 ***
## Destination_    4.001e-01  2.372e-01   1.687 0.091640 .
## Destination_55.Cancr i.e 6.230e-01  9.908e-02   6.287 3.23e-10 ***
## Destination_PSO.J318.5.22 6.041e-02  1.108e-01   0.545 0.585512
## VIP_           7.915e-02  2.246e-01   0.352 0.724522
## VIP_True       -1.836e-01  3.104e-01  -0.591 0.554329
## Cabin1_        -4.619e-01  2.386e-01  -1.936 0.052873 .
## Cabin1_A       -9.247e-01  3.528e-01  -2.621 0.008759 **
## Cabin1_B        3.802e-01  3.234e-01   1.176 0.239685
## Cabin1_C        1.735e+00  3.588e-01   4.835 1.33e-06 ***
## Cabin1_D       -1.181e-01  2.000e-01  -0.591 0.554704
## Cabin1_E       -7.058e-01  1.256e-01  -5.621 1.90e-08 ***
## Cabin1_G       -5.125e-01  1.000e-01  -5.124 2.99e-07 ***
## Cabin1_T       -1.295e+00  1.797e+00  -0.721 0.470937
## Cabin2_P       -6.652e-01  7.072e-02  -9.406 < 2e-16 ***
## Inter_CountShop -4.067e-04  1.520e-04  -2.677 0.007434 **
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

```

##
## (Dispersion parameter for binomial family taken to be 1)
##
##      Null deviance: 8450.6   on 6095   degrees of freedom
## Residual deviance: 5125.6   on 6068   degrees of freedom
## AIC: 5181.6
##
## Number of Fisher Scoring iterations: 7
##
## Confusion Matrix and Statistics
##
##           Reference
## Prediction    0    1
##           0 564 144
##           1 195 621
##
##           Accuracy : 0.7776
##           95% CI : (0.7558, 0.7982)
##      No Information Rate : 0.502
##      P-Value [Acc > NIR] : < 2.2e-16
##
##           Kappa : 0.555
##
## Mcnemar's Test P-Value : 0.006615
##
##           Sensitivity : 0.7431
##           Specificity : 0.8118
##           Pos Pred Value : 0.7966
##           Neg Pred Value : 0.7610
##           Prevalence : 0.4980
##           Detection Rate : 0.3701
##           Detection Prevalence : 0.4646
##           Balanced Accuracy : 0.7774
##
##           'Positive' Class : 0
##

```

```
##
## Call:
## roc.default(response = dfPred_raw$class, predictor = dfPred_raw$predict_reg,      plot = TRUE)
##
## Data: dfPred_raw$predict_reg in 759 controls (dfPred_raw$class 0) < 765 cases (dfPred_raw$class 1).
## Area under the curve: 0.8634
```

More Complex Models

Given the apparent complexity of the data shape, we turn to more complex nonparametric models to improve our predictions.

1. Perform SVM: 70th Percentile

We begin with Support Vector Machines and try three kernels - linear, poly and radial. We perform ten-fold cross validation and optimal hypertuning based on the caret package. Radial performs the best (accuracy=80.1%) and boosts us to the 70th percentile.

```
## [1] "Linear: -----"

## Support Vector Machines with Linear Kernel
##
## 6097 samples
## 27 predictor
```

```

##      2 classes: '0', '1'
##
## Pre-processing: centered (27), scaled (27)
## Resampling: Cross-Validated (10 fold, repeated 3 times)
## Summary of sample sizes: 5487, 5487, 5487, 5488, 5488, 5488, ...
## Resampling results across tuning parameters:
##
##      C          Accuracy   Kappa
##  0.0100000  0.7760621  0.5524086
##  0.1147368  0.7987524  0.5974004
##  0.2194737  0.8002288  0.6003060
##  0.3242105  0.8000653  0.5999630
##  0.4289474  0.8004484  0.6007181
##  0.5336842  0.8001749  0.6001643
##  0.6384211  0.8001748  0.6001610
##  0.7431579  0.8002845  0.6003763
##  0.8478947  0.8003938  0.6005952
##  0.9526316  0.8004481  0.6007033
##  1.0573684  0.8001202  0.6000454
##  1.1621053  0.8004484  0.6007013
##  1.2668421  0.8002297  0.6002636
##  1.3715789  0.8002845  0.6003721
##  1.4763158  0.8002845  0.6003719
##  1.5810526  0.8002844  0.6003714
##  1.6857895  0.8002845  0.6003698
##  1.7905263  0.8005032  0.6008077
##  1.8952632  0.8003390  0.6004790
##  2.0000000  0.8003937  0.6005881
##
## Accuracy was used to select the optimal model using the largest value.
## The final value used for the model was C = 1.790526.
## Confusion Matrix and Statistics
##
##              Reference
## Prediction    0    1
##           0 562 133
##           1 195 633
##
##              Accuracy : 0.7846
##              95% CI : (0.7631, 0.805)
##      No Information Rate : 0.503
##      P-Value [Acc > NIR] : < 2.2e-16
##
##              Kappa : 0.569
##
## Mcnemar's Test P-Value : 0.0007567
##
##              Sensitivity : 0.7424
##              Specificity : 0.8264
##      Pos Pred Value : 0.8086
##      Neg Pred Value : 0.7645
##              Prevalence : 0.4970
##      Detection Rate : 0.3690
##      Detection Prevalence : 0.4563

```

```

##          Balanced Accuracy : 0.7844
##
##          'Positive' Class : 0
##

## [1] "Poly: -----"

## Support Vector Machines with Polynomial Kernel
##
## 6097 samples
## 27 predictor
## 2 classes: '0', '1'
##
## Pre-processing: centered (27), scaled (27)
## Resampling: Cross-Validated (10 fold, repeated 3 times)
## Summary of sample sizes: 5487, 5487, 5487, 5488, 5488, ...
## Resampling results across tuning parameters:
##
##  degree  scale  C      Accuracy  Kappa
##  1        0.001 0.25  0.7551225 0.5108272
##  1        0.001 0.50  0.7498187 0.5002867
##  1        0.001 1.00  0.7554498 0.5114740
##  1        0.010 0.25  0.7633219 0.5271118
##  1        0.010 0.50  0.7703213 0.5410180
##  1        0.010 1.00  0.7759529 0.5521904
##  1        0.100 0.25  0.7869427 0.5740405
##  1        0.100 0.50  0.7943245 0.5886851
##  1        0.100 1.00  0.7984247 0.5967629
##  2        0.001 0.25  0.7497093 0.5000669
##  2        0.001 0.50  0.7552857 0.5111456
##  2        0.001 1.00  0.7615722 0.5236397
##  2        0.010 0.25  0.7717418 0.5438533
##  2        0.010 0.50  0.7814206 0.5631022
##  2        0.010 1.00  0.7927376 0.5855954
##  2        0.100 0.25  0.8023051 0.6045642
##  2        0.100 0.50  0.8050391 0.6099812
##  2        0.100 1.00  0.8044376 0.6087612
##  3        0.001 0.25  0.7527161 0.5060429
##  3        0.001 0.50  0.7588931 0.5183145
##  3        0.001 1.00  0.7650716 0.5305831
##  3        0.010 0.25  0.7802723 0.5608270
##  3        0.010 0.50  0.7912077 0.5825486
##  3        0.010 1.00  0.7976053 0.5952331
##  3        0.100 0.25  0.7935057 0.5869583
##  3        0.100 0.50  0.7950370 0.5900214
##  3        0.100 1.00  0.7952015 0.5903619
##
## Accuracy was used to select the optimal model using the largest value.
## The final values used for the model were degree = 2, scale = 0.1 and C = 0.5.
## Confusion Matrix and Statistics
##
##          Reference
## Prediction  0    1
##          0 570 148

```

```

##          1 187 618
##
##          Accuracy : 0.78
##          95% CI : (0.7584, 0.8006)
##    No Information Rate : 0.503
##    P-Value [Acc > NIR] : < 2e-16
##
##          Kappa : 0.5599
##
##    McNemar's Test P-Value : 0.03788
##
##          Sensitivity : 0.7530
##          Specificity : 0.8068
##    Pos Pred Value : 0.7939
##    Neg Pred Value : 0.7677
##          Prevalence : 0.4970
##    Detection Rate : 0.3743
##    Detection Prevalence : 0.4714
##    Balanced Accuracy : 0.7799
##
##    'Positive' Class : 0
##

## [1] "Radial: -----"

## Support Vector Machines with Radial Basis Function Kernel
##
## 6097 samples
## 27 predictor
## 2 classes: '0', '1'
##
## Pre-processing: centered (27), scaled (27)
## Resampling: Cross-Validated (10 fold, repeated 3 times)
## Summary of sample sizes: 5487, 5487, 5488, 5488, 5488, ...
## Resampling results across tuning parameters:
##
##    C      Accuracy   Kappa
##    0.25  0.7841545  0.5685024
##    0.50  0.7899499  0.5799809
##    1.00  0.7951430  0.5902743
##
## Tuning parameter 'sigma' was held constant at a value of 0.04367297
## Accuracy was used to select the optimal model using the largest value.
## The final values used for the model were sigma = 0.04367297 and C = 1.
## Confusion Matrix and Statistics
##
##          Reference
## Prediction  0    1
##          0 574 165
##          1 183 601
##
##          Accuracy : 0.7715
##          95% CI : (0.7496, 0.7924)
##    No Information Rate : 0.503

```

```

##      P-Value [Acc > NIR] : <2e-16
##
##              Kappa : 0.5429
##
## Mcnemar's Test P-Value : 0.3621
##
##      Sensitivity : 0.7583
##      Specificity : 0.7846
##      Pos Pred Value : 0.7767
##      Neg Pred Value : 0.7666
##      Prevalence : 0.4970
##      Detection Rate : 0.3769
##      Detection Prevalence : 0.4852
##      Balanced Accuracy : 0.7714
##
##      'Positive' Class : 0
##

```

2. Perform limited Neural Networks

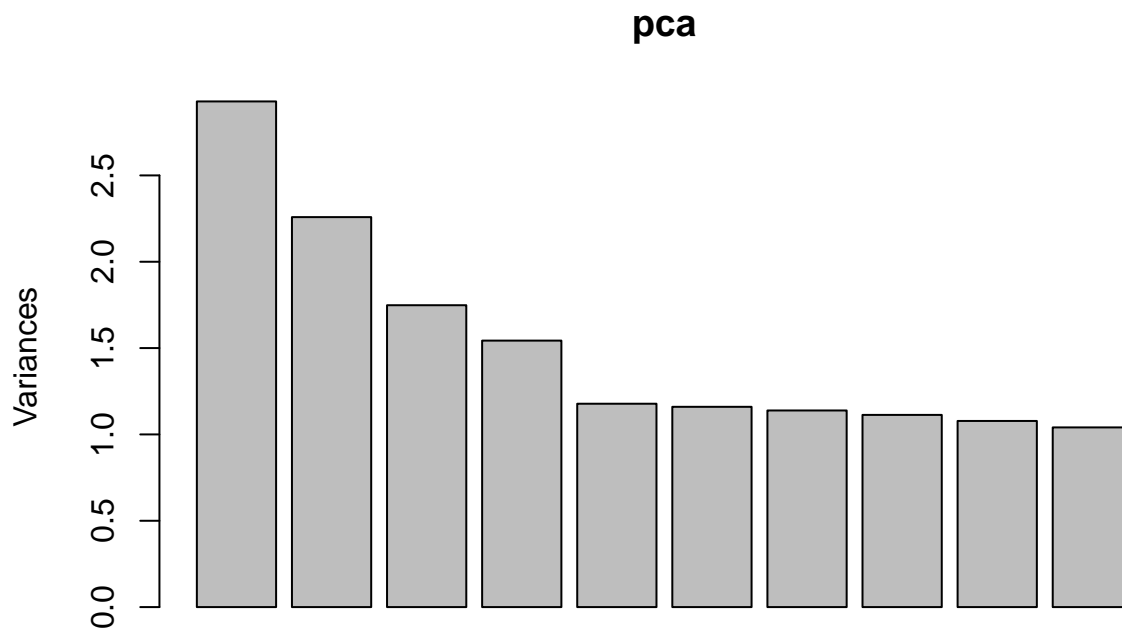
Neural Networks require a great deal of computer resources and time. A simple first pass with two hidden layers took a great deal of time, needed a high stepmax to converge and provided poor results (15th percentile). The algorithm was therefore difficult to hypertune.

In order to address the long time until convergence (many hours), we experimented with dimensionality reduction, but this was ineffective. PCA, e.g., did not result in a small number of components taking the largest share of variance. Taking a sample of records or manually eliminating columns allowed for faster run times but hurt performance. Below is the result of PCA analysis:

```

## Importance of components:
##
##      PC1      PC2      PC3      PC4      PC5      PC6      PC7
## Standard deviation  1.7113 1.50279 1.32219 1.24229 1.08522 1.07691 1.06709
## Proportion of Variance 0.1046 0.08066 0.06243 0.05512 0.04206 0.04142 0.04067
## Cumulative Proportion 0.1046 0.18524 0.24768 0.30279 0.34486 0.38627 0.42694
##
##      PC8      PC9      PC10     PC11     PC12     PC13     PC14
## Standard deviation  1.05502 1.0382 1.02008 1.01195 1.00481 0.99789 0.98678
## Proportion of Variance 0.03975 0.0385 0.03716 0.03657 0.03606 0.03556 0.03478
## Cumulative Proportion 0.46669 0.5052 0.54236 0.57893 0.61499 0.65055 0.68533
##
##      PC15     PC16     PC17     PC18     PC19     PC20     PC21
## Standard deviation  0.98295 0.93825 0.93370 0.9150 0.90835 0.88785 0.87515
## Proportion of Variance 0.03451 0.03144 0.03114 0.0299 0.02947 0.02815 0.02735
## Cumulative Proportion 0.71983 0.75127 0.78241 0.8123 0.84178 0.86993 0.89729
##
##      PC22     PC23     PC24     PC25     PC26     PC27     PC28
## Standard deviation  0.81848 0.8111 0.76818 0.64984 0.50335 0.45711 0.27096
## Proportion of Variance 0.02393 0.0235 0.02108 0.01508 0.00905 0.00746 0.00262
## Cumulative Proportion 0.92121 0.9447 0.96578 0.98087 0.98992 0.99738 1.00000

```



3. Tree Algorithms 1: Perform Random Forest: 34th Percentile

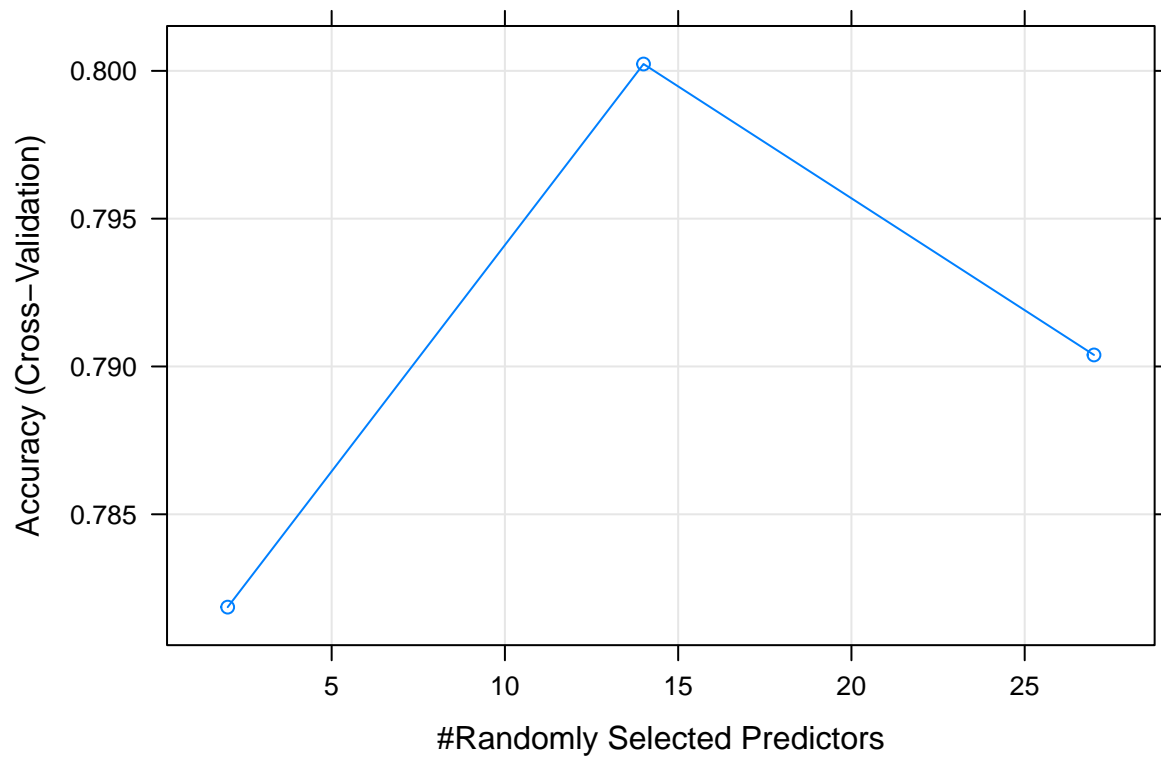
We begin with tree models. Random forest is our first. We use the parallel library to run the model on multiple cores.

We perform ten-fold cross validation and optimal hypertuning based on the caret package.

This improves accuracy on the test set to 78.7% which puts us in the 34th percentile

```
## Random Forest
##
## 6097 samples
## 27 predictor
## 2 classes: '0', '1'
##
## No pre-processing
## Resampling: Cross-Validated (10 fold)
## Summary of sample sizes: 5487, 5487, 5487, 5488, 5488, 5488, ...
## Resampling results across tuning parameters:
##
## mtry Accuracy Kappa
## 2 0.7818587 0.5639198
## 14 0.8002320 0.6005062
## 27 0.7903895 0.5809306
##
## Accuracy was used to select the optimal model using the largest value.
```

```
## The final value used for the model was mtry = 14.
```



```
## rf variable importance
##
##   only 20 most important variables shown (out of 27)
##
##                                     Overall
## CryoSleep_True                    100.000
## Spa                                85.968
## Age                                85.931
## VRDeck                             80.377
## RoomService                        76.410
## FoodCourt                          67.359
## ShoppingMall                       45.336
## Cabin1_G                           20.058
## Cabin2_P                           15.122
## Cabin1_E                           13.814
## HomePlanet_Europa                  12.745
## Inter_CountShop                    10.717
## Destination_55.Cancr.i.e           9.720
## InAGroup                           9.616
## HomePlanet_Mars                     8.855
## Destination_PS0.J318.5.22          8.058
## HomePlanet_                         4.233
## CryoSleep_                         3.845
```

```

## VIP_                                3.684
## Cabin1_C                            3.332
## Confusion Matrix and Statistics
##
##           Reference
## Prediction  0    1
##           0 606 169
##           1 151 597
##
##           Accuracy : 0.7899
##           95% CI : (0.7686, 0.8101)
##           No Information Rate : 0.503
##           P-Value [Acc > NIR] : <2e-16
##
##           Kappa : 0.5798
##
## Mcnemar's Test P-Value : 0.3419
##
##           Sensitivity : 0.8005
##           Specificity : 0.7794
##           Pos Pred Value : 0.7819
##           Neg Pred Value : 0.7981
##           Prevalence : 0.4970
##           Detection Rate : 0.3979
##           Detection Prevalence : 0.5089
##           Balanced Accuracy : 0.7900
##
##           'Positive' Class : 0
##
## [1] "Parameters:  mtry = 14 , ntree = 500 , nrnodes = 1605"

```

4. Tree Algorithms 2: Perform XGBoost Untuned: 73rd Percentile

At the 34th percentile we need a more powerful model. As an active learner, XGBoost is likely to fit our training model better than random forest, though it may overfit the data.

Our untuned model, with 55 rounds, achieves 80.2% accuracy and reaches the 74th percentile.

Hypertuning

1. Tune XGBoost: 78th Percentile

We perform ten-fold cross validation on 100 rounds and find the optimal rounds for accuracy (14). Our model, with 14 rounds, achieves 80.3% accuracy and reaches the 74th percentile.

```

## [1] train-rmse:0.415047+0.000841    train-auc:0.813023+0.001987 train-error:0.282094+0.001992 tes
## [2] train-rmse:0.385845+0.001185    train-auc:0.860516+0.002116 train-error:0.222164+0.007927 tes
## [3] train-rmse:0.378852+0.001398    train-auc:0.870728+0.001990 train-error:0.209274+0.003398 tes
## [4] train-rmse:0.372729+0.001637    train-auc:0.880756+0.002273 train-error:0.203893+0.003009 tes
## [5] train-rmse:0.368863+0.001861    train-auc:0.886316+0.002073 train-error:0.201575+0.002140 tes
## [6] train-rmse:0.366441+0.001607    train-auc:0.889350+0.001927 train-error:0.199446+0.002347 tes
## [7] train-rmse:0.364508+0.001586    train-auc:0.891987+0.001954 train-error:0.196544+0.003051 tes
## [8] train-rmse:0.361741+0.002500    train-auc:0.895200+0.002572 train-error:0.192797+0.004358 tes

```


## [9]	train-rmse:0.359367+0.002674	train-auc:0.897881+0.002820	train-error:0.190580+0.005296	tes
## [10]	train-rmse:0.357505+0.002598	train-auc:0.900123+0.002623	train-error:0.188655+0.005272	tes
## [11]	train-rmse:0.356330+0.002516	train-auc:0.901504+0.002411	train-error:0.187474+0.005112	tes
## [12]	train-rmse:0.354930+0.002560	train-auc:0.903052+0.002540	train-error:0.185681+0.004593	tes
## [13]	train-rmse:0.353317+0.002560	train-auc:0.905031+0.002491	train-error:0.183129+0.003636	tes
## [14]	train-rmse:0.352609+0.002423	train-auc:0.905840+0.002382	train-error:0.181846+0.003486	tes
## [15]	train-rmse:0.351316+0.002204	train-auc:0.907134+0.002231	train-error:0.180782+0.003400	tes
## [16]	train-rmse:0.350333+0.002146	train-auc:0.908162+0.002161	train-error:0.179280+0.002695	tes
## [17]	train-rmse:0.349126+0.002166	train-auc:0.909447+0.002254	train-error:0.177924+0.003161	tes
## [18]	train-rmse:0.348295+0.002153	train-auc:0.910399+0.002297	train-error:0.177705+0.002517	tes
## [19]	train-rmse:0.347184+0.002109	train-auc:0.911515+0.002141	train-error:0.176363+0.002503	tes
## [20]	train-rmse:0.346005+0.001498	train-auc:0.912769+0.001518	train-error:0.175168+0.002447	tes
## [21]	train-rmse:0.345138+0.001769	train-auc:0.913644+0.001712	train-error:0.174161+0.003187	tes
## [22]	train-rmse:0.344329+0.001754	train-auc:0.914464+0.001689	train-error:0.173637+0.003224	tes
## [23]	train-rmse:0.343350+0.001915	train-auc:0.915459+0.001870	train-error:0.172907+0.003426	tes
## [24]	train-rmse:0.342475+0.001964	train-auc:0.916324+0.001901	train-error:0.171945+0.003266	tes
## [25]	train-rmse:0.341509+0.002072	train-auc:0.917216+0.002013	train-error:0.171289+0.003922	tes
## [26]	train-rmse:0.340689+0.001933	train-auc:0.918042+0.001888	train-error:0.170516+0.003508	tes
## [27]	train-rmse:0.339792+0.002083	train-auc:0.918946+0.002021	train-error:0.169714+0.004183	tes
## [28]	train-rmse:0.339044+0.002041	train-auc:0.919682+0.001946	train-error:0.169175+0.003845	tes
## [29]	train-rmse:0.338046+0.001954	train-auc:0.920647+0.001815	train-error:0.167760+0.003082	tes
## [30]	train-rmse:0.337128+0.001795	train-auc:0.921598+0.001656	train-error:0.166900+0.002514	tes
## [31]	train-rmse:0.336358+0.001755	train-auc:0.922325+0.001597	train-error:0.166229+0.002539	tes
## [32]	train-rmse:0.335645+0.001621	train-auc:0.922911+0.001497	train-error:0.165150+0.002449	tes
## [33]	train-rmse:0.334922+0.001529	train-auc:0.923596+0.001316	train-error:0.164071+0.002469	tes
## [34]	train-rmse:0.333932+0.001624	train-auc:0.924532+0.001408	train-error:0.163459+0.002936	tes
## [35]	train-rmse:0.332912+0.001554	train-auc:0.925443+0.001347	train-error:0.162321+0.002618	tes
## [36]	train-rmse:0.332262+0.001438	train-auc:0.926035+0.001194	train-error:0.161738+0.002974	tes
## [37]	train-rmse:0.331246+0.001364	train-auc:0.926959+0.001162	train-error:0.161534+0.002805	tes
## [38]	train-rmse:0.330345+0.001437	train-auc:0.927820+0.001221	train-error:0.159740+0.003067	tes
## [39]	train-rmse:0.329495+0.001612	train-auc:0.928612+0.001336	train-error:0.158953+0.003136	tes
## [40]	train-rmse:0.328807+0.001581	train-auc:0.929218+0.001320	train-error:0.158093+0.003108	tes
## [41]	train-rmse:0.327937+0.001694	train-auc:0.930023+0.001379	train-error:0.157480+0.003106	tes
## [42]	train-rmse:0.327184+0.001656	train-auc:0.930684+0.001377	train-error:0.156416+0.003253	tes
## [43]	train-rmse:0.326425+0.001521	train-auc:0.931303+0.001260	train-error:0.155497+0.002907	tes
## [44]	train-rmse:0.325603+0.001620	train-auc:0.931990+0.001330	train-error:0.154666+0.002618	tes
## [45]	train-rmse:0.324805+0.001692	train-auc:0.932658+0.001348	train-error:0.153733+0.002719	tes
## [46]	train-rmse:0.324063+0.001575	train-auc:0.933289+0.001216	train-error:0.153223+0.003100	tes
## [47]	train-rmse:0.323439+0.001762	train-auc:0.933793+0.001388	train-error:0.153135+0.003544	tes
## [48]	train-rmse:0.322854+0.001788	train-auc:0.934262+0.001377	train-error:0.152391+0.003528	tes
## [49]	train-rmse:0.322252+0.001745	train-auc:0.934807+0.001355	train-error:0.151750+0.003334	tes
## [50]	train-rmse:0.321600+0.001629	train-auc:0.935366+0.001298	train-error:0.151035+0.002842	tes
## [51]	train-rmse:0.320875+0.001747	train-auc:0.935970+0.001400	train-error:0.150423+0.003509	tes
## [52]	train-rmse:0.320162+0.001893	train-auc:0.936542+0.001521	train-error:0.149329+0.003485	tes
## [53]	train-rmse:0.319647+0.001986	train-auc:0.936932+0.001590	train-error:0.149081+0.003185	tes
## [54]	train-rmse:0.319180+0.002051	train-auc:0.937360+0.001665	train-error:0.148921+0.003133	tes
## [55]	train-rmse:0.318495+0.002023	train-auc:0.937930+0.001619	train-error:0.148046+0.002944	tes
## [56]	train-rmse:0.317847+0.001816	train-auc:0.938420+0.001408	train-error:0.147171+0.002895	tes
## [57]	train-rmse:0.317274+0.001728	train-auc:0.938898+0.001332	train-error:0.146530+0.002161	tes
## [58]	train-rmse:0.316655+0.001743	train-auc:0.939373+0.001312	train-error:0.145669+0.002031	tes
## [59]	train-rmse:0.316001+0.001710	train-auc:0.939908+0.001291	train-error:0.145217+0.002054	tes
## [60]	train-rmse:0.315394+0.001733	train-auc:0.940395+0.001311	train-error:0.144736+0.002320	tes
## [61]	train-rmse:0.314790+0.001697	train-auc:0.940824+0.001289	train-error:0.144430+0.002118	tes
## [62]	train-rmse:0.314159+0.001587	train-auc:0.941314+0.001225	train-error:0.143992+0.001886	tes

## [63]	train-rmse:0.313643+0.001755	train-auc:0.941697+0.001368	train-error:0.143336+0.002088	tes
## [64]	train-rmse:0.313051+0.001694	train-auc:0.942164+0.001256	train-error:0.142695+0.002148	tes
## [65]	train-rmse:0.312558+0.001828	train-auc:0.942522+0.001354	train-error:0.142695+0.002514	tes
## [66]	train-rmse:0.312013+0.001876	train-auc:0.942916+0.001425	train-error:0.141951+0.002486	tes
## [67]	train-rmse:0.311434+0.001698	train-auc:0.943358+0.001266	train-error:0.141441+0.002384	tes
## [68]	train-rmse:0.310875+0.001640	train-auc:0.943777+0.001186	train-error:0.141134+0.002037	tes
## [69]	train-rmse:0.310278+0.001617	train-auc:0.944223+0.001195	train-error:0.140245+0.001969	tes
## [70]	train-rmse:0.309696+0.001594	train-auc:0.944673+0.001148	train-error:0.139531+0.002061	tes
## [71]	train-rmse:0.309181+0.001532	train-auc:0.945069+0.001093	train-error:0.138947+0.001563	tes
## [72]	train-rmse:0.308746+0.001576	train-auc:0.945353+0.001151	train-error:0.138320+0.001831	tes
## [73]	train-rmse:0.308254+0.001624	train-auc:0.945715+0.001192	train-error:0.138087+0.001517	tes
## [74]	train-rmse:0.307611+0.001518	train-auc:0.946229+0.001126	train-error:0.137270+0.001555	tes
## [75]	train-rmse:0.307122+0.001547	train-auc:0.946569+0.001098	train-error:0.136468+0.001889	tes
## [76]	train-rmse:0.306571+0.001594	train-auc:0.946961+0.001127	train-error:0.135944+0.001965	tes
## [77]	train-rmse:0.305946+0.001564	train-auc:0.947416+0.001071	train-error:0.135185+0.001789	tes
## [78]	train-rmse:0.305473+0.001577	train-auc:0.947754+0.001094	train-error:0.134835+0.002163	tes
## [79]	train-rmse:0.305042+0.001672	train-auc:0.948028+0.001165	train-error:0.134354+0.002092	tes
## [80]	train-rmse:0.304559+0.001754	train-auc:0.948360+0.001245	train-error:0.133596+0.002231	tes
## [81]	train-rmse:0.304144+0.001801	train-auc:0.948627+0.001221	train-error:0.133231+0.002122	tes
## [82]	train-rmse:0.303739+0.001734	train-auc:0.948914+0.001159	train-error:0.133348+0.002223	tes
## [83]	train-rmse:0.303312+0.001882	train-auc:0.949232+0.001253	train-error:0.132458+0.002218	tes
## [84]	train-rmse:0.302778+0.001926	train-auc:0.949588+0.001289	train-error:0.131715+0.002569	tes
## [85]	train-rmse:0.302109+0.001988	train-auc:0.950044+0.001358	train-error:0.131030+0.002782	tes
## [86]	train-rmse:0.301626+0.002012	train-auc:0.950331+0.001356	train-error:0.130767+0.002845	tes
## [87]	train-rmse:0.300992+0.001877	train-auc:0.950830+0.001264	train-error:0.130300+0.002366	tes
## [88]	train-rmse:0.300468+0.001871	train-auc:0.951155+0.001258	train-error:0.129863+0.002225	tes
## [89]	train-rmse:0.299921+0.001991	train-auc:0.951500+0.001340	train-error:0.129251+0.002430	tes
## [90]	train-rmse:0.299382+0.002110	train-auc:0.951856+0.001388	train-error:0.128653+0.002764	tes
## [91]	train-rmse:0.298813+0.002120	train-auc:0.952226+0.001396	train-error:0.128055+0.002812	tes
## [92]	train-rmse:0.298246+0.002226	train-auc:0.952622+0.001474	train-error:0.127647+0.002876	tes
## [93]	train-rmse:0.297793+0.002160	train-auc:0.952924+0.001428	train-error:0.127238+0.002568	tes
## [94]	train-rmse:0.297216+0.002172	train-auc:0.953288+0.001439	train-error:0.126947+0.002558	tes
## [95]	train-rmse:0.296801+0.002170	train-auc:0.953534+0.001385	train-error:0.126465+0.002678	tes
## [96]	train-rmse:0.296367+0.002323	train-auc:0.953836+0.001481	train-error:0.125722+0.002673	tes
## [97]	train-rmse:0.295864+0.002369	train-auc:0.954187+0.001540	train-error:0.125882+0.002883	tes
## [98]	train-rmse:0.295462+0.002359	train-auc:0.954423+0.001510	train-error:0.125299+0.002879	tes
## [99]	train-rmse:0.295035+0.002460	train-auc:0.954663+0.001593	train-error:0.124993+0.002777	tes
## [100]	train-rmse:0.294489+0.002418	train-auc:0.955032+0.001500	train-error:0.124555+0.002784	

xgb.cv 10-folds

##	iter	train_rmse_mean	train_rmse_std	train_auc_mean	train_auc_std
##	1	0.4150472	0.0008407872	0.8130229	0.001987129
##	2	0.3858453	0.0011852734	0.8605165	0.002116012
##	3	0.3788522	0.0013976987	0.8707282	0.001989507
##	4	0.3727295	0.0016374695	0.8807559	0.002273087
##	5	0.3688626	0.0018613585	0.8863161	0.002073216
##	6	0.3664408	0.0016074771	0.8893501	0.001927122
##	7	0.3645081	0.0015860282	0.8919871	0.001953559
##	8	0.3617411	0.0025001798	0.8952004	0.002572208
##	9	0.3593669	0.0026737110	0.8978812	0.002819932
##	10	0.3575046	0.0025977840	0.9001226	0.002623137
##	11	0.3563304	0.0025156403	0.9015038	0.002411376
##	12	0.3549304	0.0025602207	0.9030519	0.002540477
##	13	0.3533172	0.0025595232	0.9050305	0.002490587

##	14	0.3526090	0.0024227740	0.9058402	0.002381847
##	15	0.3513156	0.0022040046	0.9071336	0.002231236
##	16	0.3503330	0.0021458483	0.9081623	0.002161142
##	17	0.3491261	0.0021664241	0.9094473	0.002254360
##	18	0.3482954	0.0021525416	0.9103991	0.002296694
##	19	0.3471842	0.0021090532	0.9115153	0.002141252
##	20	0.3460052	0.0014983255	0.9127690	0.001518343
##	21	0.3451379	0.0017689474	0.9136437	0.001711919
##	22	0.3443292	0.0017538449	0.9144641	0.001688819
##	23	0.3433504	0.0019147572	0.9154587	0.001869824
##	24	0.3424754	0.0019642783	0.9163236	0.001900750
##	25	0.3415089	0.0020722321	0.9172161	0.002012797
##	26	0.3406887	0.0019334262	0.9180421	0.001887678
##	27	0.3397922	0.0020830397	0.9189460	0.002021011
##	28	0.3390438	0.0020406827	0.9196823	0.001946301
##	29	0.3380463	0.0019536782	0.9206469	0.001815359
##	30	0.3371280	0.0017953080	0.9215984	0.001656312
##	31	0.3363576	0.0017554653	0.9223249	0.001597113
##	32	0.3356455	0.0016210755	0.9229110	0.001497178
##	33	0.3349221	0.0015292546	0.9235958	0.001315839
##	34	0.3339316	0.0016236697	0.9245321	0.001407679
##	35	0.3329116	0.0015539443	0.9254432	0.001346647
##	36	0.3322624	0.0014381169	0.9260354	0.001193753
##	37	0.3312456	0.0013643091	0.9269594	0.001161588
##	38	0.3303449	0.0014368349	0.9278201	0.001220801
##	39	0.3294946	0.0016117879	0.9286122	0.001335569
##	40	0.3288075	0.0015811839	0.9292182	0.001319932
##	41	0.3279367	0.0016937697	0.9300226	0.001378730
##	42	0.3271845	0.0016557488	0.9306839	0.001376820
##	43	0.3264252	0.0015210003	0.9313029	0.001260430
##	44	0.3256029	0.0016199637	0.9319899	0.001330492
##	45	0.3248051	0.0016917590	0.9326583	0.001347714
##	46	0.3240628	0.0015749109	0.9332886	0.001216220
##	47	0.3234387	0.0017624129	0.9337930	0.001388292
##	48	0.3228536	0.0017878012	0.9342619	0.001376842
##	49	0.3222521	0.0017451265	0.9348072	0.001354840
##	50	0.3216003	0.0016293940	0.9353660	0.001298031
##	51	0.3208746	0.0017466966	0.9359704	0.001399812
##	52	0.3201616	0.0018929619	0.9365422	0.001520536
##	53	0.3196474	0.0019863154	0.9369315	0.001590095
##	54	0.3191798	0.0020510904	0.9373596	0.001664705
##	55	0.3184953	0.0020233171	0.9379301	0.001619470
##	56	0.3178470	0.0018159202	0.9384201	0.001408072
##	57	0.3172745	0.0017284655	0.9388979	0.001331981
##	58	0.3166554	0.0017430991	0.9393732	0.001312450
##	59	0.3160015	0.0017098937	0.9399080	0.001291092
##	60	0.3153941	0.0017327089	0.9403952	0.001311240
##	61	0.3147899	0.0016971185	0.9408238	0.001289445
##	62	0.3141589	0.0015871036	0.9413142	0.001224871
##	63	0.3136426	0.0017545167	0.9416971	0.001367983
##	64	0.3130513	0.0016935131	0.9421638	0.001256322
##	65	0.3125581	0.0018283325	0.9425217	0.001354058
##	66	0.3120128	0.0018758543	0.9429164	0.001425453
##	67	0.3114345	0.0016978604	0.9433579	0.001266203

##	68	0.3108745	0.0016403151	0.9437774	0.001186087
##	69	0.3102778	0.0016174699	0.9442233	0.001194706
##	70	0.3096960	0.0015938708	0.9446733	0.001147959
##	71	0.3091811	0.0015323005	0.9450690	0.001093086
##	72	0.3087464	0.0015762548	0.9453532	0.001151400
##	73	0.3082539	0.0016238172	0.9457153	0.001191572
##	74	0.3076112	0.0015179670	0.9462288	0.001126127
##	75	0.3071218	0.0015471330	0.9465694	0.001098388
##	76	0.3065711	0.0015937457	0.9469606	0.001127085
##	77	0.3059460	0.0015644019	0.9474155	0.001071474
##	78	0.3054734	0.0015765148	0.9477536	0.001094148
##	79	0.3050415	0.0016721615	0.9480278	0.001164874
##	80	0.3045590	0.0017535599	0.9483599	0.001245139
##	81	0.3041436	0.0018011682	0.9486268	0.001220946
##	82	0.3037391	0.0017341519	0.9489142	0.001158834
##	83	0.3033118	0.0018817398	0.9492318	0.001253110
##	84	0.3027782	0.0019258601	0.9495880	0.001288597
##	85	0.3021095	0.0019879170	0.9500436	0.001358454
##	86	0.3016264	0.0020116503	0.9503314	0.001355558
##	87	0.3009923	0.0018768301	0.9508298	0.001264026
##	88	0.3004677	0.0018706476	0.9511545	0.001258138
##	89	0.2999206	0.0019907126	0.9515000	0.001339754
##	90	0.2993822	0.0021101515	0.9518561	0.001388026
##	91	0.2988134	0.0021202069	0.9522256	0.001395612
##	92	0.2982457	0.0022264526	0.9526222	0.001474256
##	93	0.2977932	0.0021597408	0.9529243	0.001427907
##	94	0.2972163	0.0021721367	0.9532876	0.001439034
##	95	0.2968009	0.0021696078	0.9535342	0.001384682
##	96	0.2963672	0.0023232846	0.9538355	0.001481315
##	97	0.2958643	0.0023686242	0.9541874	0.001539564
##	98	0.2954618	0.0023590862	0.9544232	0.001509885
##	99	0.2950347	0.0024596957	0.9546631	0.001592549
##	100	0.2944892	0.0024184142	0.9550320	0.001500150
##	iter train_rmse_mean train_rmse_std train_auc_mean train_auc_std				
##	train_error_mean train_error_std test_rmse_mean test_rmse_std test_auc_mean				
##	0.2820939	0.001992099	0.4163444	0.005940971	0.8103009
##	0.2221641	0.007926785	0.3905971	0.007574149	0.8523477
##	0.2092737	0.003397503	0.3860413	0.008923041	0.8609303
##	0.2038932	0.003008937	0.3814645	0.009476891	0.8682670
##	0.2015747	0.002139571	0.3787298	0.009713631	0.8723888
##	0.1994459	0.002346903	0.3769712	0.010198804	0.8756100
##	0.1965442	0.003050911	0.3757651	0.010404697	0.8769801
##	0.1927967	0.004358211	0.3746049	0.011306811	0.8789893
##	0.1905803	0.005296333	0.3731903	0.010896145	0.8807695
##	0.1886555	0.005272195	0.3727517	0.010769920	0.8816784
##	0.1874744	0.005112225	0.3724146	0.010664375	0.8823298
##	0.1856809	0.004592811	0.3718681	0.010613391	0.8831247
##	0.1831291	0.003636310	0.3710748	0.011118640	0.8842339
##	0.1818460	0.003486183	0.3705184	0.011414561	0.8849245
##	0.1807815	0.003400295	0.3708519	0.011694643	0.8841869
##	0.1792796	0.002695176	0.3707264	0.011704054	0.8845269
##	0.1779235	0.003161466	0.3712928	0.011491559	0.8838648
##	0.1777048	0.002517482	0.3718740	0.011599427	0.8833306
##	0.1763633	0.002503280	0.3717453	0.011720242	0.8834212

##	0.1751676	0.002446691	0.3713145	0.012274063	0.8843503
##	0.1741615	0.003186965	0.3716268	0.012131086	0.8838853
##	0.1736365	0.003223589	0.3712662	0.012412699	0.8844709
##	0.1729075	0.003425837	0.3717313	0.013099009	0.8839264
##	0.1719451	0.003266054	0.3719164	0.012877542	0.8839594
##	0.1712889	0.003921531	0.3718655	0.013102286	0.8840172
##	0.1705161	0.003508495	0.3721261	0.012656256	0.8836352
##	0.1697141	0.004182673	0.3729312	0.012649515	0.8829951
##	0.1691746	0.003845228	0.3734957	0.012742949	0.8824653
##	0.1677602	0.003081937	0.3737875	0.012831275	0.8821490
##	0.1668999	0.002513610	0.3738110	0.013433636	0.8820928
##	0.1662292	0.002538956	0.3737938	0.013033989	0.8822264
##	0.1651502	0.002449222	0.3737518	0.012882601	0.8821574
##	0.1640712	0.002469257	0.3745227	0.012648728	0.8813211
##	0.1634588	0.002935600	0.3748835	0.013186861	0.8810441
##	0.1623214	0.002618447	0.3747488	0.013704860	0.8812463
##	0.1617381	0.002974293	0.3751152	0.013389884	0.8811790
##	0.1615340	0.002805438	0.3748097	0.013509896	0.8814768
##	0.1597405	0.003067497	0.3751836	0.013597112	0.8809343
##	0.1589530	0.003136424	0.3755461	0.014105492	0.8804968
##	0.1580928	0.003107901	0.3757008	0.013685368	0.8804653
##	0.1574803	0.003105683	0.3758811	0.013404811	0.8801763
##	0.1564159	0.003252921	0.3756843	0.013582188	0.8803998
##	0.1554972	0.002906947	0.3755217	0.013491695	0.8807061
##	0.1546661	0.002617860	0.3758504	0.013920480	0.8806825
##	0.1537329	0.002718963	0.3761600	0.013886949	0.8804458
##	0.1532225	0.003099631	0.3761690	0.013544243	0.8805973
##	0.1531350	0.003543852	0.3759833	0.013561906	0.8808755
##	0.1523913	0.003528368	0.3761295	0.013322055	0.8807837
##	0.1517497	0.003334097	0.3763730	0.013047615	0.8805153
##	0.1510352	0.002841739	0.3765837	0.012697463	0.8802232
##	0.1504228	0.003509029	0.3768346	0.012753244	0.8800933
##	0.1493292	0.003484773	0.3771535	0.012650149	0.8797425
##	0.1490814	0.003184958	0.3772315	0.012662881	0.8796460
##	0.1489210	0.003132930	0.3778558	0.012504540	0.8790490
##	0.1480461	0.002944064	0.3776862	0.012396650	0.8791791
##	0.1471712	0.002895272	0.3780824	0.012285204	0.8788216
##	0.1465296	0.002161138	0.3783945	0.012455634	0.8784282
##	0.1456693	0.002031106	0.3787235	0.012390253	0.8783465
##	0.1452173	0.002053897	0.3793426	0.012629892	0.8775908
##	0.1447361	0.002320276	0.3794049	0.012468829	0.8776270
##	0.1444298	0.002118282	0.3791911	0.012942039	0.8779519
##	0.1439924	0.001885917	0.3792233	0.012687439	0.8778298
##	0.1433363	0.002087925	0.3788662	0.012774251	0.8783715
##	0.1426947	0.002148479	0.3784912	0.012298131	0.8788565
##	0.1426947	0.002514097	0.3783710	0.012427593	0.8790959
##	0.1419510	0.002486093	0.3785079	0.012471644	0.8789987
##	0.1414407	0.002383815	0.3788439	0.012370071	0.8786583
##	0.1411345	0.002036997	0.3789487	0.012547271	0.8785655
##	0.1402450	0.001968808	0.3792602	0.012676776	0.8782160
##	0.1395305	0.002060956	0.3795440	0.012525504	0.8778857
##	0.1389472	0.001562765	0.3797347	0.012657868	0.8776412
##	0.1383202	0.001831433	0.3796498	0.012405659	0.8777717
##	0.1380869	0.001517111	0.3794444	0.012358788	0.8779069

##	0.1372704	0.001555493	0.3796003	0.012349827	0.8778526
##	0.1364684	0.001889184	0.3800292	0.012010576	0.8774506
##	0.1359435	0.001965240	0.3807024	0.012300158	0.8767426
##	0.1351853	0.001789483	0.3806841	0.012528177	0.8767919
##	0.1348353	0.002163130	0.3806865	0.012479562	0.8767785
##	0.1343541	0.002092018	0.3808431	0.012557640	0.8766432
##	0.1335959	0.002230848	0.3807394	0.012845058	0.8767672
##	0.1332313	0.002121707	0.3810106	0.013132345	0.8765135
##	0.1333480	0.002223465	0.3812867	0.012874867	0.8763985
##	0.1324585	0.002217818	0.3814256	0.013057811	0.8762678
##	0.1317148	0.002569301	0.3818120	0.012985211	0.8759474
##	0.1310295	0.002782151	0.3819941	0.012992894	0.8758744
##	0.1307670	0.002845036	0.3822309	0.012889672	0.8759527
##	0.1303004	0.002365917	0.3831844	0.013099716	0.8749906
##	0.1298630	0.002224556	0.3832133	0.013012175	0.8751051
##	0.1292505	0.002429759	0.3832353	0.012890424	0.8750778
##	0.1286527	0.002764170	0.3838789	0.012765165	0.8746764
##	0.1280548	0.002811584	0.3839534	0.012664215	0.8746476
##	0.1276466	0.002875570	0.3837816	0.012788417	0.8748136
##	0.1272383	0.002568386	0.3837339	0.013001662	0.8748475
##	0.1269467	0.002557710	0.3839414	0.012872987	0.8745956
##	0.1264655	0.002677593	0.3840430	0.012998917	0.8745925
##	0.1257218	0.002672741	0.3843602	0.012980077	0.8742902
##	0.1258822	0.002882943	0.3841278	0.012834125	0.8746246
##	0.1252990	0.002879313	0.3842362	0.012830571	0.8746688
##	0.1249928	0.002777072	0.3843985	0.012750189	0.8746084
##	0.1245553	0.002784244	0.3844450	0.012598636	0.8744459
##	train_error_mean	train_error_std	test_rmse_mean	test_rmse_std	test_auc_mean
##	test_auc_std	test_error_mean	test_error_std		
##	0.01108728	0.2905516	0.01223988		
##	0.01083094	0.2332003	0.01163912		
##	0.01309413	0.2206039	0.01678296		
##	0.01300392	0.2181086	0.01183428		
##	0.01294677	0.2154841	0.01394771		
##	0.01269970	0.2104967	0.01333915		
##	0.01277184	0.2074795	0.01159622		
##	0.01406098	0.2087951	0.01698399		
##	0.01321164	0.2061703	0.01379299		
##	0.01274347	0.2055146	0.01479766		
##	0.01248041	0.2045943	0.01288102		
##	0.01235360	0.2040695	0.01321176		
##	0.01307136	0.2014428	0.01347696		
##	0.01331109	0.1999994	0.01399277		
##	0.01342856	0.2009180	0.01419129		
##	0.01362023	0.2014428	0.01602215		
##	0.01332994	0.2021000	0.01483841		
##	0.01341941	0.2022307	0.01518117		
##	0.01330782	0.2018368	0.01482716		
##	0.01402412	0.2005243	0.01586857		
##	0.01388941	0.2026237	0.01627588		
##	0.01416349	0.2026235	0.01594229		
##	0.01474153	0.2019677	0.01561097		
##	0.01453009	0.2031490	0.01636020		
##	0.01482702	0.2027555	0.01678262		

##	0.01439909	0.2028865	0.01669788
##	0.01430383	0.2036731	0.01754895
##	0.01434710	0.2049856	0.01713376
##	0.01457266	0.2043289	0.01730056
##	0.01534312	0.2030160	0.01849838
##	0.01483236	0.2034089	0.01685060
##	0.01454883	0.2035408	0.01537859
##	0.01419541	0.2048547	0.01429535
##	0.01491990	0.2057725	0.01662786
##	0.01542082	0.2047228	0.01744095
##	0.01514652	0.2055105	0.01761378
##	0.01519824	0.2028867	0.01773024
##	0.01535407	0.2043289	0.01988539
##	0.01564721	0.2038039	0.01967663
##	0.01510704	0.2047238	0.01835648
##	0.01473146	0.2064293	0.01814378
##	0.01490411	0.2062984	0.01835259
##	0.01467348	0.2069542	0.01839799
##	0.01507226	0.2077420	0.01974727
##	0.01502624	0.2090543	0.01761777
##	0.01477085	0.2078731	0.01792716
##	0.01478820	0.2070862	0.01778459
##	0.01463435	0.2074795	0.01688302
##	0.01456083	0.2065614	0.01687622
##	0.01438638	0.2072172	0.01589736
##	0.01436449	0.2073488	0.01707709
##	0.01412330	0.2077423	0.01629112
##	0.01413643	0.2069542	0.01653529
##	0.01408184	0.2069546	0.01713234
##	0.01404165	0.2072164	0.01813468
##	0.01387336	0.2085294	0.01612108
##	0.01426612	0.2086603	0.01723773
##	0.01405308	0.2093168	0.01808271
##	0.01445949	0.2107607	0.01966456
##	0.01435124	0.2104977	0.01776324
##	0.01472845	0.2103658	0.01771531
##	0.01449095	0.2119404	0.01844850
##	0.01473005	0.2106278	0.01975248
##	0.01396577	0.2107586	0.01901973
##	0.01437535	0.2103660	0.01939672
##	0.01418317	0.2082664	0.01732583
##	0.01413923	0.2093163	0.01784670
##	0.01430384	0.2087920	0.01781780
##	0.01456424	0.2102353	0.01621831
##	0.01457869	0.2101049	0.01569640
##	0.01484526	0.2095793	0.01589035
##	0.01448711	0.2099742	0.01452432
##	0.01428116	0.2099740	0.01458114
##	0.01435274	0.2111541	0.01579642
##	0.01404611	0.2108918	0.01582710
##	0.01462188	0.2119410	0.01692485
##	0.01483964	0.2110227	0.01583146
##	0.01477349	0.2116780	0.01500679
##	0.01495237	0.2120715	0.01597333

```

##      0.01530188      0.2118092      0.01744904
##      0.01549512      0.2119406      0.01787080
##      0.01525579      0.2120720      0.01902123
##      0.01553935      0.2120730      0.01820921
##      0.01547304      0.2120730      0.01795197
##      0.01544001      0.2124654      0.01862590
##      0.01527018      0.2128591      0.01839634
##      0.01538256      0.2133843      0.01814966
##      0.01529104      0.2137779      0.01737426
##      0.01501692      0.2144335      0.01739449
##      0.01476730      0.2148272      0.01683999
##      0.01482685      0.2154831      0.01693515
##      0.01505861      0.2161391      0.01703639
##      0.01525643      0.2150892      0.01704109
##      0.01495951      0.2150894      0.01705315
##      0.01518626      0.2150894      0.01738361
##      0.01510705      0.2149588      0.01698411
##      0.01493328      0.2152211      0.01654235
##      0.01479555      0.2141716      0.01684381
##      0.01490047      0.2140405      0.01618553
##      0.01492134      0.2144342      0.01497779
## test_auc_std test_error_mean test_error_std

## ##### xgb.cv 10-folds
## call:
##   xgb.cv(data = xgb_train, nrounds = 100, nfold = 10, metrics = list("rmse",
##     "auc", "error"), nthread = 2, max_depth = 3, eta = 1, objective = "binary:logistic")
## params (as set within xgb.cv):
##   nthread = "2", max_depth = "3", eta = "1", objective = "binary:logistic", eval_metric = "rmse", ev
## callbacks:
##   cb.print.evaluation(period = print_every_n, showsd = showsd)
##   cb.evaluation.log()
## niter: 100
## evaluation_log:
##   iter train_rmse_mean train_rmse_std train_auc_mean train_auc_std
##      1      0.4150472    0.0008407872    0.8130229    0.001987129
##      2      0.3858453    0.0011852734    0.8605165    0.002116012
##      3      0.3788522    0.0013976987    0.8707282    0.001989507
##      4      0.3727295    0.0016374695    0.8807559    0.002273087
##      5      0.3688626    0.0018613585    0.8863161    0.002073216
##      6      0.3664408    0.0016074771    0.8893501    0.001927122
##      7      0.3645081    0.0015860282    0.8919871    0.001953559
##      8      0.3617411    0.0025001798    0.8952004    0.002572208
##      9      0.3593669    0.0026737110    0.8978812    0.002819932
##     10      0.3575046    0.0025977840    0.9001226    0.002623137
##     11      0.3563304    0.0025156403    0.9015038    0.002411376
##     12      0.3549304    0.0025602207    0.9030519    0.002540477
##     13      0.3533172    0.0025595232    0.9050305    0.002490587
##     14      0.3526090    0.0024227740    0.9058402    0.002381847
##     15      0.3513156    0.0022040046    0.9071336    0.002231236
##     16      0.3503330    0.0021458483    0.9081623    0.002161142
##     17      0.3491261    0.0021664241    0.9094473    0.002254360
##     18      0.3482954    0.0021525416    0.9103991    0.002296694
##     19      0.3471842    0.0021090532    0.9115153    0.002141252

```


##	20	0.3460052	0.0014983255	0.9127690	0.001518343
##	21	0.3451379	0.0017689474	0.9136437	0.001711919
##	22	0.3443292	0.0017538449	0.9144641	0.001688819
##	23	0.3433504	0.0019147572	0.9154587	0.001869824
##	24	0.3424754	0.0019642783	0.9163236	0.001900750
##	25	0.3415089	0.0020722321	0.9172161	0.002012797
##	26	0.3406887	0.0019334262	0.9180421	0.001887678
##	27	0.3397922	0.0020830397	0.9189460	0.002021011
##	28	0.3390438	0.0020406827	0.9196823	0.001946301
##	29	0.3380463	0.0019536782	0.9206469	0.001815359
##	30	0.3371280	0.0017953080	0.9215984	0.001656312
##	31	0.3363576	0.0017554653	0.9223249	0.001597113
##	32	0.3356455	0.0016210755	0.9229110	0.001497178
##	33	0.3349221	0.0015292546	0.9235958	0.001315839
##	34	0.3339316	0.0016236697	0.9245321	0.001407679
##	35	0.3329116	0.0015539443	0.9254432	0.001346647
##	36	0.3322624	0.0014381169	0.9260354	0.001193753
##	37	0.3312456	0.0013643091	0.9269594	0.001161588
##	38	0.3303449	0.0014368349	0.9278201	0.001220801
##	39	0.3294946	0.0016117879	0.9286122	0.001335569
##	40	0.3288075	0.0015811839	0.9292182	0.001319932
##	41	0.3279367	0.0016937697	0.9300226	0.001378730
##	42	0.3271845	0.0016557488	0.9306839	0.001376820
##	43	0.3264252	0.0015210003	0.9313029	0.001260430
##	44	0.3256029	0.0016199637	0.9319899	0.001330492
##	45	0.3248051	0.0016917590	0.9326583	0.001347714
##	46	0.3240628	0.0015749109	0.9332886	0.001216220
##	47	0.3234387	0.0017624129	0.9337930	0.001388292
##	48	0.3228536	0.0017878012	0.9342619	0.001376842
##	49	0.3222521	0.0017451265	0.9348072	0.001354840
##	50	0.3216003	0.0016293940	0.9353660	0.001298031
##	51	0.3208746	0.0017466966	0.9359704	0.001399812
##	52	0.3201616	0.0018929619	0.9365422	0.001520536
##	53	0.3196474	0.0019863154	0.9369315	0.001590095
##	54	0.3191798	0.0020510904	0.9373596	0.001664705
##	55	0.3184953	0.0020233171	0.9379301	0.001619470
##	56	0.3178470	0.0018159202	0.9384201	0.001408072
##	57	0.3172745	0.0017284655	0.9388979	0.001331981
##	58	0.3166554	0.0017430991	0.9393732	0.001312450
##	59	0.3160015	0.0017098937	0.9399080	0.001291092
##	60	0.3153941	0.0017327089	0.9403952	0.001311240
##	61	0.3147899	0.0016971185	0.9408238	0.001289445
##	62	0.3141589	0.0015871036	0.9413142	0.001224871
##	63	0.3136426	0.0017545167	0.9416971	0.001367983
##	64	0.3130513	0.0016935131	0.9421638	0.001256322
##	65	0.3125581	0.0018283325	0.9425217	0.001354058
##	66	0.3120128	0.0018758543	0.9429164	0.001425453
##	67	0.3114345	0.0016978604	0.9433579	0.001266203
##	68	0.3108745	0.0016403151	0.9437774	0.001186087
##	69	0.3102778	0.0016174699	0.9442233	0.001194706
##	70	0.3096960	0.0015938708	0.9446733	0.001147959
##	71	0.3091811	0.0015323005	0.9450690	0.001093086
##	72	0.3087464	0.0015762548	0.9453532	0.001151400
##	73	0.3082539	0.0016238172	0.9457153	0.001191572

##	74	0.3076112	0.0015179670	0.9462288	0.001126127
##	75	0.3071218	0.0015471330	0.9465694	0.001098388
##	76	0.3065711	0.0015937457	0.9469606	0.001127085
##	77	0.3059460	0.0015644019	0.9474155	0.001071474
##	78	0.3054734	0.0015765148	0.9477536	0.001094148
##	79	0.3050415	0.0016721615	0.9480278	0.001164874
##	80	0.3045590	0.0017535599	0.9483599	0.001245139
##	81	0.3041436	0.0018011682	0.9486268	0.001220946
##	82	0.3037391	0.0017341519	0.9489142	0.001158834
##	83	0.3033118	0.0018817398	0.9492318	0.001253110
##	84	0.3027782	0.0019258601	0.9495880	0.001288597
##	85	0.3021095	0.0019879170	0.9500436	0.001358454
##	86	0.3016264	0.0020116503	0.9503314	0.001355558
##	87	0.3009923	0.0018768301	0.9508298	0.001264026
##	88	0.3004677	0.0018706476	0.9511545	0.001258138
##	89	0.2999206	0.0019907126	0.9515000	0.001339754
##	90	0.2993822	0.0021101515	0.9518561	0.001388026
##	91	0.2988134	0.0021202069	0.9522256	0.001395612
##	92	0.2982457	0.0022264526	0.9526222	0.001474256
##	93	0.2977932	0.0021597408	0.9529243	0.001427907
##	94	0.2972163	0.0021721367	0.9532876	0.001439034
##	95	0.2968009	0.0021696078	0.9535342	0.001384682
##	96	0.2963672	0.0023232846	0.9538355	0.001481315
##	97	0.2958643	0.0023686242	0.9541874	0.001539564
##	98	0.2954618	0.0023590862	0.9544232	0.001509885
##	99	0.2950347	0.0024596957	0.9546631	0.001592549
##	100	0.2944892	0.0024184142	0.9550320	0.001500150
##	iter train_rmse_mean train_rmse_std train_auc_mean train_auc_std				
##	train_error_mean train_error_std test_rmse_mean test_rmse_std test_auc_mean				
##	0.2820939	0.001992099	0.4163444	0.005940971	0.8103009
##	0.2221641	0.007926785	0.3905971	0.007574149	0.8523477
##	0.2092737	0.003397503	0.3860413	0.008923041	0.8609303
##	0.2038932	0.003008937	0.3814645	0.009476891	0.8682670
##	0.2015747	0.002139571	0.3787298	0.009713631	0.8723888
##	0.1994459	0.002346903	0.3769712	0.010198804	0.8756100
##	0.1965442	0.003050911	0.3757651	0.010404697	0.8769801
##	0.1927967	0.004358211	0.3746049	0.011306811	0.8789893
##	0.1905803	0.005296333	0.3731903	0.010896145	0.8807695
##	0.1886555	0.005272195	0.3727517	0.010769920	0.8816784
##	0.1874744	0.005112225	0.3724146	0.010664375	0.8823298
##	0.1856809	0.004592811	0.3718681	0.010613391	0.8831247
##	0.1831291	0.003636310	0.3710748	0.011118640	0.8842339
##	0.1818460	0.003486183	0.3705184	0.011414561	0.8849245
##	0.1807815	0.003400295	0.3708519	0.011694643	0.8841869
##	0.1792796	0.002695176	0.3707264	0.011704054	0.8845269
##	0.1779235	0.003161466	0.3712928	0.011491559	0.8838648
##	0.1777048	0.002517482	0.3718740	0.011599427	0.8833306
##	0.1763633	0.002503280	0.3717453	0.011720242	0.8834212
##	0.1751676	0.002446691	0.3713145	0.012274063	0.8843503
##	0.1741615	0.003186965	0.3716268	0.012131086	0.8838853
##	0.1736365	0.003223589	0.3712662	0.012412699	0.8844709
##	0.1729075	0.003425837	0.3717313	0.013099009	0.8839264
##	0.1719451	0.003266054	0.3719164	0.012877542	0.8839594
##	0.1712889	0.003921531	0.3718655	0.013102286	0.8840172

##	0.1705161	0.003508495	0.3721261	0.012656256	0.8836352
##	0.1697141	0.004182673	0.3729312	0.012649515	0.8829951
##	0.1691746	0.003845228	0.3734957	0.012742949	0.8824653
##	0.1677602	0.003081937	0.3737875	0.012831275	0.8821490
##	0.1668999	0.002513610	0.3738110	0.013433636	0.8820928
##	0.1662292	0.002538956	0.3737938	0.013033989	0.8822264
##	0.1651502	0.002449222	0.3737518	0.012882601	0.8821574
##	0.1640712	0.002469257	0.3745227	0.012648728	0.8813211
##	0.1634588	0.002935600	0.3748835	0.013186861	0.8810441
##	0.1623214	0.002618447	0.3747488	0.013704860	0.8812463
##	0.1617381	0.002974293	0.3751152	0.013389884	0.8811790
##	0.1615340	0.002805438	0.3748097	0.013509896	0.8814768
##	0.1597405	0.003067497	0.3751836	0.013597112	0.8809343
##	0.1589530	0.003136424	0.3755461	0.014105492	0.8804968
##	0.1580928	0.003107901	0.3757008	0.013685368	0.8804653
##	0.1574803	0.003105683	0.3758811	0.013404811	0.8801763
##	0.1564159	0.003252921	0.3756843	0.013582188	0.8803998
##	0.1554972	0.002906947	0.3755217	0.013491695	0.8807061
##	0.1546661	0.002617860	0.3758504	0.013920480	0.8806825
##	0.1537329	0.002718963	0.3761600	0.013886949	0.8804458
##	0.1532225	0.003099631	0.3761690	0.013544243	0.8805973
##	0.1531350	0.003543852	0.3759833	0.013561906	0.8808755
##	0.1523913	0.003528368	0.3761295	0.013322055	0.8807837
##	0.1517497	0.003334097	0.3763730	0.013047615	0.8805153
##	0.1510352	0.002841739	0.3765837	0.012697463	0.8802232
##	0.1504228	0.003509029	0.3768346	0.012753244	0.8800933
##	0.1493292	0.003484773	0.3771535	0.012650149	0.8797425
##	0.1490814	0.003184958	0.3772315	0.012662881	0.8796460
##	0.1489210	0.003132930	0.3778558	0.012504540	0.8790490
##	0.1480461	0.002944064	0.3776862	0.012396650	0.8791791
##	0.1471712	0.002895272	0.3780824	0.012285204	0.8788216
##	0.1465296	0.002161138	0.3783945	0.012455634	0.8784282
##	0.1456693	0.002031106	0.3787235	0.012390253	0.8783465
##	0.1452173	0.002053897	0.3793426	0.012629892	0.8775908
##	0.1447361	0.002320276	0.3794049	0.012468829	0.8776270
##	0.1444298	0.002118282	0.3791911	0.012942039	0.8779519
##	0.1439924	0.001885917	0.3792233	0.012687439	0.8778298
##	0.1433363	0.002087925	0.3788662	0.012774251	0.8783715
##	0.1426947	0.002148479	0.3784912	0.012298131	0.8788565
##	0.1426947	0.002514097	0.3783710	0.012427593	0.8790959
##	0.1419510	0.002486093	0.3785079	0.012471644	0.8789987
##	0.1414407	0.002383815	0.3788439	0.012370071	0.8786583
##	0.1411345	0.002036997	0.3789487	0.012547271	0.8785655
##	0.1402450	0.001968808	0.3792602	0.012676776	0.8782160
##	0.1395305	0.002060956	0.3795440	0.012525504	0.8778857
##	0.1389472	0.001562765	0.3797347	0.012657868	0.8776412
##	0.1383202	0.001831433	0.3796498	0.012405659	0.8777717
##	0.1380869	0.001517111	0.3794444	0.012358788	0.8779069
##	0.1372704	0.001555493	0.3796003	0.012349827	0.8778526
##	0.1364684	0.001889184	0.3800292	0.012010576	0.8774506
##	0.1359435	0.001965240	0.3807024	0.012300158	0.8767426
##	0.1351853	0.001789483	0.3806841	0.012528177	0.8767919
##	0.1348353	0.002163130	0.3806865	0.012479562	0.8767785
##	0.1343541	0.002092018	0.3808431	0.012557640	0.8766432

##	0.1335959	0.002230848	0.3807394	0.012845058	0.8767672
##	0.1332313	0.002121707	0.3810106	0.013132345	0.8765135
##	0.1333480	0.002223465	0.3812867	0.012874867	0.8763985
##	0.1324585	0.002217818	0.3814256	0.013057811	0.8762678
##	0.1317148	0.002569301	0.3818120	0.012985211	0.8759474
##	0.1310295	0.002782151	0.3819941	0.012992894	0.8758744
##	0.1307670	0.002845036	0.3822309	0.012889672	0.8759527
##	0.1303004	0.002365917	0.3831844	0.013099716	0.8749906
##	0.1298630	0.002224556	0.3832133	0.013012175	0.8751051
##	0.1292505	0.002429759	0.3832353	0.012890424	0.8750778
##	0.1286527	0.002764170	0.3838789	0.012765165	0.8746764
##	0.1280548	0.002811584	0.3839534	0.012664215	0.8746476
##	0.1276466	0.002875570	0.3837816	0.012788417	0.8748136
##	0.1272383	0.002568386	0.3837339	0.013001662	0.8748475
##	0.1269467	0.002557710	0.3839414	0.012872987	0.8745956
##	0.1264655	0.002677593	0.3840430	0.012998917	0.8745925
##	0.1257218	0.002672741	0.3843602	0.012980077	0.8742902
##	0.1258822	0.002882943	0.3841278	0.012834125	0.8746246
##	0.1252990	0.002879313	0.3842362	0.012830571	0.8746688
##	0.1249928	0.002777072	0.3843985	0.012750189	0.8746084
##	0.1245553	0.002784244	0.3844450	0.012598636	0.8744459
##	train_error_mean	train_error_std	test_rmse_mean	test_rmse_std	test_auc_mean
##	test_auc_std	test_error_mean	test_error_std		
##	0.01108728	0.2905516	0.01223988		
##	0.01083094	0.2332003	0.01163912		
##	0.01309413	0.2206039	0.01678296		
##	0.01300392	0.2181086	0.01183428		
##	0.01294677	0.2154841	0.01394771		
##	0.01269970	0.2104967	0.01333915		
##	0.01277184	0.2074795	0.01159622		
##	0.01406098	0.2087951	0.01698399		
##	0.01321164	0.2061703	0.01379299		
##	0.01274347	0.2055146	0.01479766		
##	0.01248041	0.2045943	0.01288102		
##	0.01235360	0.2040695	0.01321176		
##	0.01307136	0.2014428	0.01347696		
##	0.01331109	0.1999994	0.01399277		
##	0.01342856	0.2009180	0.01419129		
##	0.01362023	0.2014428	0.01602215		
##	0.01332994	0.2021000	0.01483841		
##	0.01341941	0.2022307	0.01518117		
##	0.01330782	0.2018368	0.01482716		
##	0.01402412	0.2005243	0.01586857		
##	0.01388941	0.2026237	0.01627588		
##	0.01416349	0.2026235	0.01594229		
##	0.01474153	0.2019677	0.01561097		
##	0.01453009	0.2031490	0.01636020		
##	0.01482702	0.2027555	0.01678262		
##	0.01439909	0.2028865	0.01669788		
##	0.01430383	0.2036731	0.01754895		
##	0.01434710	0.2049856	0.01713376		
##	0.01457266	0.2043289	0.01730056		
##	0.01534312	0.2030160	0.01849838		
##	0.01483236	0.2034089	0.01685060		

##	0.01454883	0.2035408	0.01537859
##	0.01419541	0.2048547	0.01429535
##	0.01491990	0.2057725	0.01662786
##	0.01542082	0.2047228	0.01744095
##	0.01514652	0.2055105	0.01761378
##	0.01519824	0.2028867	0.01773024
##	0.01535407	0.2043289	0.01988539
##	0.01564721	0.2038039	0.01967663
##	0.01510704	0.2047238	0.01835648
##	0.01473146	0.2064293	0.01814378
##	0.01490411	0.2062984	0.01835259
##	0.01467348	0.2069542	0.01839799
##	0.01507226	0.2077420	0.01974727
##	0.01502624	0.2090543	0.01761777
##	0.01477085	0.2078731	0.01792716
##	0.01478820	0.2070862	0.01778459
##	0.01463435	0.2074795	0.01688302
##	0.01456083	0.2065614	0.01687622
##	0.01438638	0.2072172	0.01589736
##	0.01436449	0.2073488	0.01707709
##	0.01412330	0.2077423	0.01629112
##	0.01413643	0.2069542	0.01653529
##	0.01408184	0.2069546	0.01713234
##	0.01404165	0.2072164	0.01813468
##	0.01387336	0.2085294	0.01612108
##	0.01426612	0.2086603	0.01723773
##	0.01405308	0.2093168	0.01808271
##	0.01445949	0.2107607	0.01966456
##	0.01435124	0.2104977	0.01776324
##	0.01472845	0.2103658	0.01771531
##	0.01449095	0.2119404	0.01844850
##	0.01473005	0.2106278	0.01975248
##	0.01396577	0.2107586	0.01901973
##	0.01437535	0.2103660	0.01939672
##	0.01418317	0.2082664	0.01732583
##	0.01413923	0.2093163	0.01784670
##	0.01430384	0.2087920	0.01781780
##	0.01456424	0.2102353	0.01621831
##	0.01457869	0.2101049	0.01569640
##	0.01484526	0.2095793	0.01589035
##	0.01448711	0.2099742	0.01452432
##	0.01428116	0.2099740	0.01458114
##	0.01435274	0.2111541	0.01579642
##	0.01404611	0.2108918	0.01582710
##	0.01462188	0.2119410	0.01692485
##	0.01483964	0.2110227	0.01583146
##	0.01477349	0.2116780	0.01500679
##	0.01495237	0.2120715	0.01597333
##	0.01530188	0.2118092	0.01744904
##	0.01549512	0.2119406	0.01787080
##	0.01525579	0.2120720	0.01902123
##	0.01553935	0.2120730	0.01820921
##	0.01547304	0.2120730	0.01795197
##	0.01544001	0.2124654	0.01862590

```
##      0.01527018      0.2128591      0.01839634
##      0.01538256      0.2133843      0.01814966
##      0.01529104      0.2137779      0.01737426
##      0.01501692      0.2144335      0.01739449
##      0.01476730      0.2148272      0.01683999
##      0.01482685      0.2154831      0.01693515
##      0.01505861      0.2161391      0.01703639
##      0.01525643      0.2150892      0.01704109
##      0.01495951      0.2150894      0.01705315
##      0.01518626      0.2150894      0.01738361
##      0.01510705      0.2149588      0.01698411
##      0.01493328      0.2152211      0.01654235
##      0.01479555      0.2141716      0.01684381
##      0.01490047      0.2140405      0.01618553
##      0.01492134      0.2144342      0.01497779
## test_auc_std test_error_mean test_error_std
```

```
## [1] train-logloss:0.518785 eval-logloss:0.820399
## [2] train-logloss:0.460336 eval-logloss:0.881454
## [3] train-logloss:0.443383 eval-logloss:0.952815
## [4] train-logloss:0.429115 eval-logloss:1.006953
## [5] train-logloss:0.421462 eval-logloss:1.068660
## [6] train-logloss:0.417349 eval-logloss:1.082765
## [7] train-logloss:0.404430 eval-logloss:1.115529
## [8] train-logloss:0.401043 eval-logloss:1.130244
## [9] train-logloss:0.396576 eval-logloss:1.161850
## [10] train-logloss:0.393187 eval-logloss:1.185680
## [11] train-logloss:0.390200 eval-logloss:1.205786
## [12] train-logloss:0.388254 eval-logloss:1.226213
## [13] train-logloss:0.386829 eval-logloss:1.243061
## [14] train-logloss:0.383553 eval-logloss:1.241412
```

```
## [1] 0.850011170 0.199663937 0.003278726 0.004425377 0.187479198 0.122859113
```

2. Tune SVM Radial: No Improvement

We manually hypertune or svm radial model by increasing sigma (loosening the fit) and decreasing sigma (tightening the fit). A decreased sigma increases accuracy on corss validation but not on the Kaggle.

```
## * checking for file 'C:\Users\erico\AppData\Local\Temp\Rtmpg7iJQd\remotes2d583711944\ericonsi-EHData'
## * preparing 'EHData':
## * checking DESCRIPTION meta-information ... OK
## * checking for LF line-endings in source and make files and shell scripts
## * checking for empty or unneeded directories
## Omitted 'LazyData' from DESCRIPTION
## * creating default NAMESPACE file
## * building 'EHData_0.1.0.tar.gz'
##

## Support Vector Machines with Radial Basis Function Kernel
##
## 6097 samples
```

```

## 27 predictor
## 2 classes: '0', '1'
##
## Pre-processing: centered (27), scaled (27)
## Resampling: Cross-Validated (10 fold, repeated 3 times)
## Summary of sample sizes: 5487, 5487, 5487, 5488, 5488, 5488, ...
## Resampling results:
##
## Accuracy Kappa
## 0.7959085 0.5918032
##
## Tuning parameter 'sigma' was held constant at a value of 0.04
## Tuning
## parameter 'C' was held constant at a value of 1
## Confusion Matrix and Statistics
##
##           Reference
## Prediction  0    1
##           0 576 163
##           1 181 603
##
##           Accuracy : 0.7741
##           95% CI : (0.7523, 0.7949)
##           No Information Rate : 0.503
##           P-Value [Acc > NIR] : <2e-16
##
##           Kappa : 0.5482
##
## Mcnemar's Test P-Value : 0.3594
##
##           Sensitivity : 0.7609
##           Specificity : 0.7872
##           Pos Pred Value : 0.7794
##           Neg Pred Value : 0.7691
##           Prevalence : 0.4970
##           Detection Rate : 0.3782
##           Detection Prevalence : 0.4852
##           Balanced Accuracy : 0.7741
##
##           'Positive' Class : 0
##

```

Discussion

(see above)