

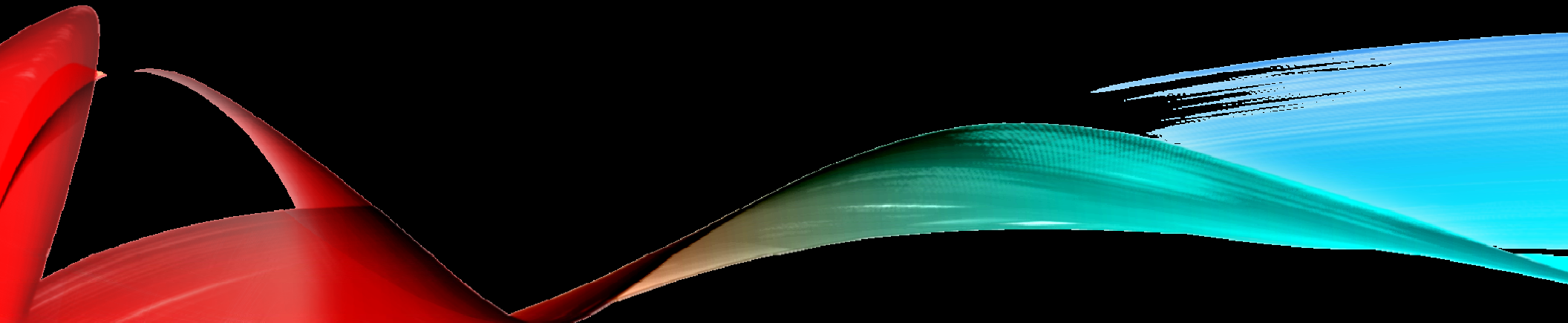


# ADVANCED DECISION TREES

JMP Partition Platform

Advanced Trees | Jim Grayson, PhD

# ADVANCED TREES



# Partition Modeling

Y	X	Objective	Predictive Accuracy	Statistical Significance Measure	Model Fit
Categorical	Continuous or Categorical	Classification	Misclassification Rate   Confusion Matrix	LogWorth	RMSE, MAD, ROC Curve

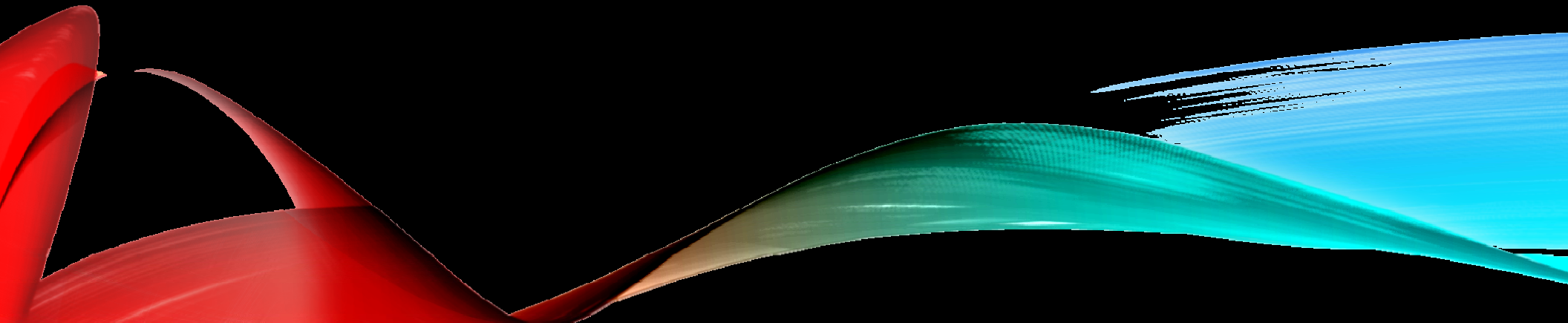
Y	X	Objective	Predictive Accuracy	Statistical Significance Measure	Model Fit
Continuous	Continuous or Categorical	Prediction	Rsquare, RMSE	LogWorth	RMSE, MAD, Lift Curve

## Advantages and Disadvantages of Trees

- ▲ Trees are very easy to explain to people. In fact, they are even easier to explain than linear regression!
- ▲ Some people believe that decision trees more closely mirror human decision-making than do the regression and classification approaches seen in previous chapters.
- ▲ Trees can be displayed graphically, and are easily interpreted even by a non-expert (especially if they are small).
- ▲ Trees can easily handle qualitative predictors without the need to create dummy variables.
- ▼ Unfortunately, trees generally do not have the same level of predictive accuracy as some of the other regression and classification approaches seen in this book.

However, by aggregating many decision trees, the predictive performance of trees can be substantially improved. We introduce these concepts next.

# BOOTSTRAP FOREST



# Bagging

6

- *Bootstrap aggregation*, or *bagging*, is a general-purpose procedure for reducing the variance of a statistical learning method; we introduce it here because it is particularly useful and frequently used in the context of decision trees.
- Recall that given a set of  $n$  independent observations  $Z_1, \dots, Z_n$ , each with variance  $\sigma^2$ , the variance of the mean  $\bar{Z}$  of the observations is given by  $\sigma^2/n$ .
- In other words, *averaging a set of observations reduces variance*. Of course, this is not practical because we generally do not have access to multiple training sets.



## Bagging— continued

- Instead, we can bootstrap, by taking repeated samples from the (single) training data set.
- In this approach we generate  $B$  different bootstrapped training data sets. We then train our method on the  $b$ th bootstrapped training set in order to get  $\hat{f}^{*b}(x)$ , the prediction at a point  $x$ . We then average all the predictions to obtain

$$\hat{f}_{\text{bag}}(x) = \frac{1}{B} \sum_{b=1}^B \hat{f}^{*b}(x).$$

This is called *bagging*.

## Bagging classification trees

- The above prescription applied to regression trees
- For classification trees: for each test observation, we record the class predicted by each of the  $B$  trees, and take a *majority vote*: the overall prediction is the most commonly occurring class among the  $B$  predictions.

Stanford Online Course: Course Notes, Statistical Learning by Hastie and Tibshirani, Winter 2014 | Based on [An Introduction to Statistical Learning](#) by James, Witten, Hastie and Tibshirani, Springer (2013)



## Random Forests

- *Random forests* provide an improvement over bagged trees by way of a small tweak that *decorrelates* the trees. This reduces the variance when we average the trees.
- As in bagging, we build a number of decision trees on bootstrapped training samples.
- But when building these decision trees, each time a split in a tree is considered, *a random selection of  $m$  predictors* is chosen as split candidates from the full set of  $p$  predictors. The split is allowed to use only one of those  $m$  predictors.
- A fresh selection of  $m$  predictors is taken at each split, and typically we choose  $m \approx \sqrt{p}$  — that is, the number of predictors considered at each split is approximately equal to the square root of the total number of predictors (4 out of the 13 for the Heart data).

Stanford Online Course: Course Notes, Statistical Learning by Hastie and Tibshirani, Winter 2014 | Based on [An Introduction to Statistical Learning](#) by James, Witten, Hastie and Tibshirani, Springer (2013)

February 28–March 1, 2012  
Santa Clara, California

O'REILLY\*

**Strata**  
**CONFERENCE**  
Making Data Work

## The Two Most Important Algorithms in Predictive Modeling Today

*Jeremy Howard (Kaggle), Mike Bowles (Biomatica)*

1:30pm Tuesday, 02/28/2012

Data Science, Ballroom CD

When doing predictive modelling, there are two situations in which you might find yourself:

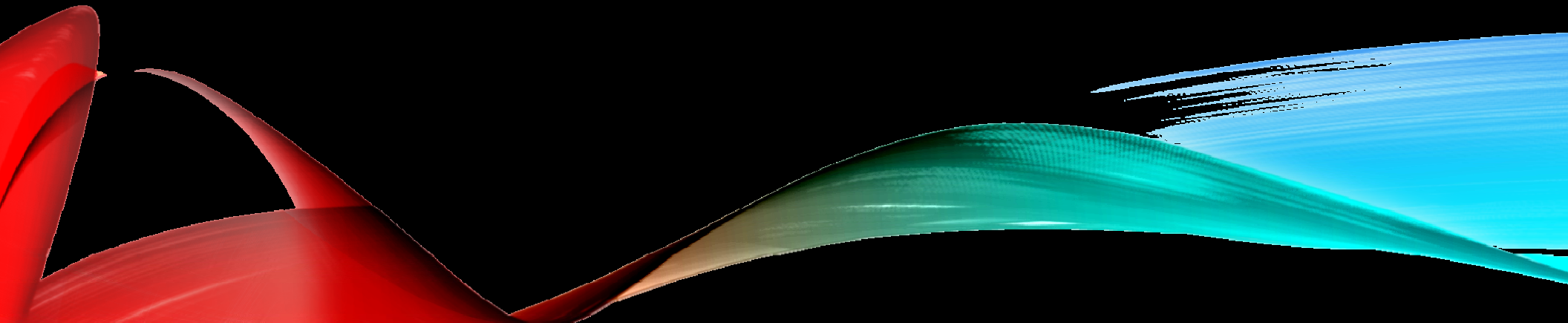
- ▶ You need to fit a well-defined parameterised model to your data, so you require a learning algorithm which can find those parameters on a large data set without over-fitting
- ▶ You just need a “black box” which can predict your dependent variable as accurately as possible, so you need a learning algorithm which can automatically identify the structure, interactions, and relationships in the data

For case (1), lasso and elastic-net regularized generalized linear models are a set of modern algorithms which meet all these needs. They are fast, work on huge data sets, and avoid over-fitting automatically. They are available in the “glmnet” package in R.

For case (2), ensembles of decision trees (often known as “Random Forests”) have been the most successful general-purpose algorithm in modern times. For instance, most Kaggle competitions have at least one top entry that heavily uses this approach. This algorithm is very simple to understand, and is fast and easy to apply. It is available in the “randomForest” package in R.

Mike and Jeremy will explain in simple terms, using no complex math, how these algorithms work, and will also explain using numerous examples how to apply them using R. They will also provide advice on how to select from these algorithms, and will show how to prepare the data, and how to use the trained models in practice.

# TITANTIC DATA



"The sinking of the RMS Titanic is one of the most infamous shipwrecks in history. On April 15, 1912, during her maiden voyage, the Titanic sank after colliding with an iceberg, killing 1502 out of 2224 passengers and crew. This sensational tragedy shocked the international community and led to better safety regulations for ships.

One of the reasons that the shipwreck led to such loss of life was that there were not enough lifeboats for the passengers and crew. Although there was some element of luck involved in surviving the sinking, some groups of people were more likely to survive than others."

From a Kaggle competition: "Titanic: Machine Learning from Disaster", <http://bit.ly/1f2crzi>, data accessed 08/2014.



This data table describes the survival status of 1309 of the 1324 individual passengers on the Titanic. Information on the 899 crew members is not included.

Name: Passenger Name

Survived: Yes or No

Passenger Class: 1, 2, or 3 corresponding to 1st, 2nd or 3rd class

Sex: Passenger sex

Age: Passenger age

Siblings and Spouses: The number of siblings and spouses aboard

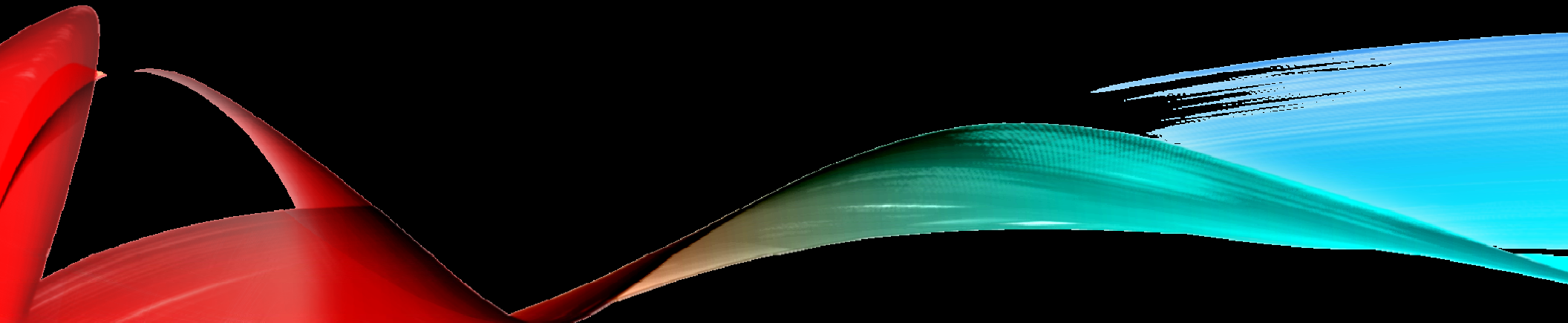
Parents and Children: The number of parents and children aboard

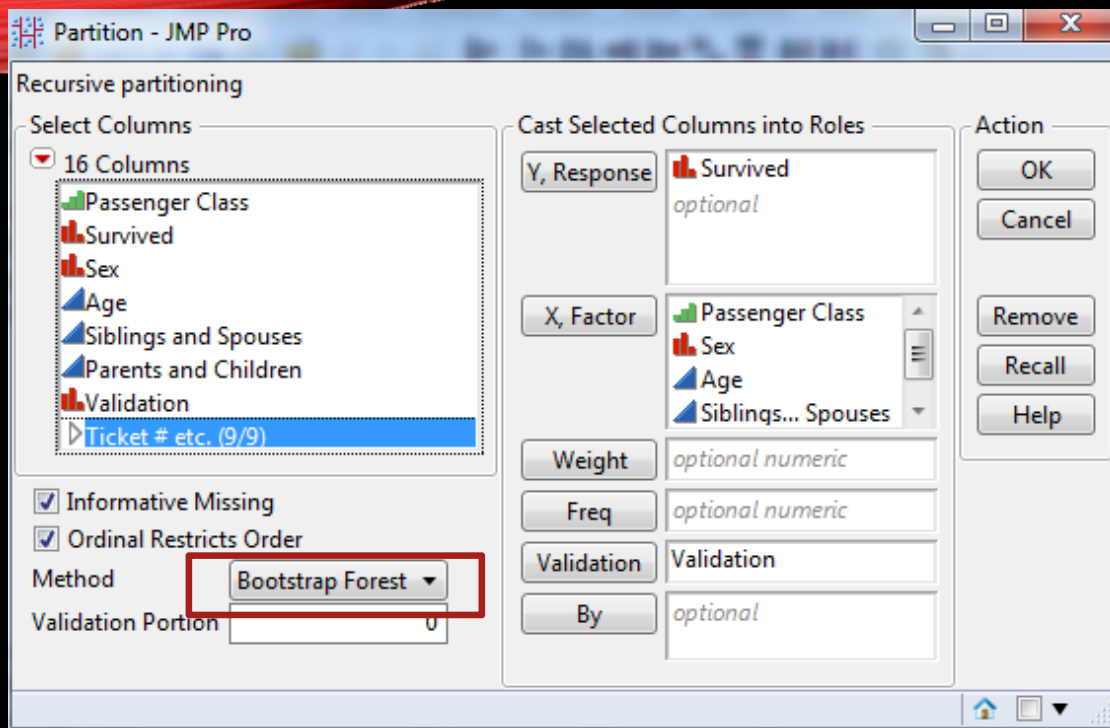
Fare: The passenger fare

Port: Port of embarkment (C = Cherbourg; Q = Queenstown; S = Southampton)

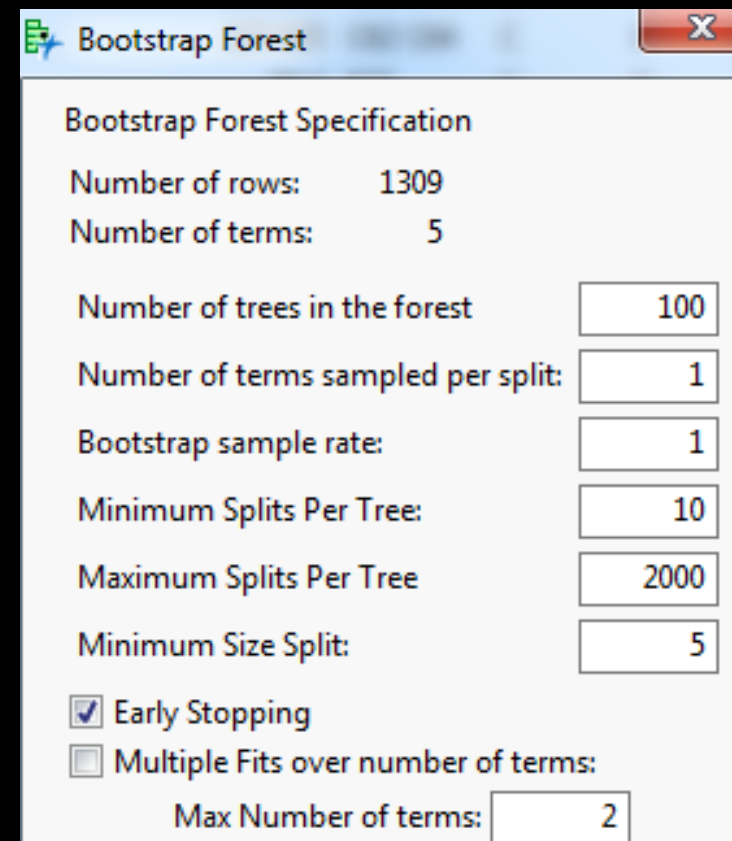
Home/Destination: The home or final intended destination of the passenger

# BOOTSTRAP FOREST CATEGORICAL RESPONSE





## Use Defaults

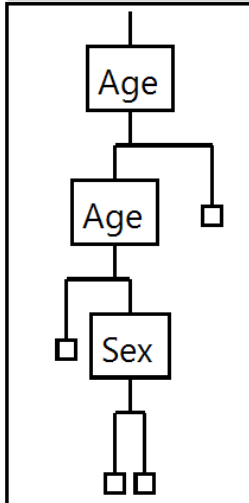


## Tree Views

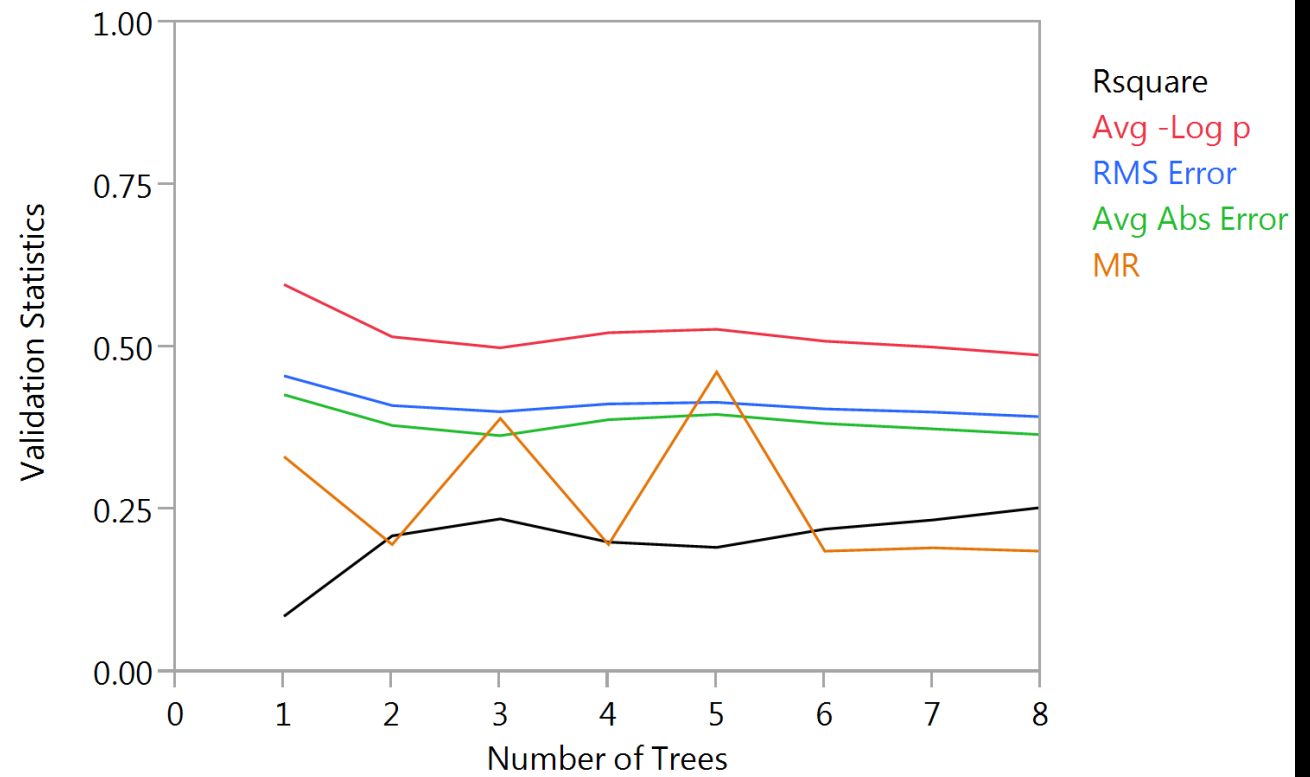
- Tree1
- Tree2
- Tree3
- Tree4
- Tree5
- Tree6
- Tree7
- Tree8

## Tree Views

## Tree1



## Cumulative Validation



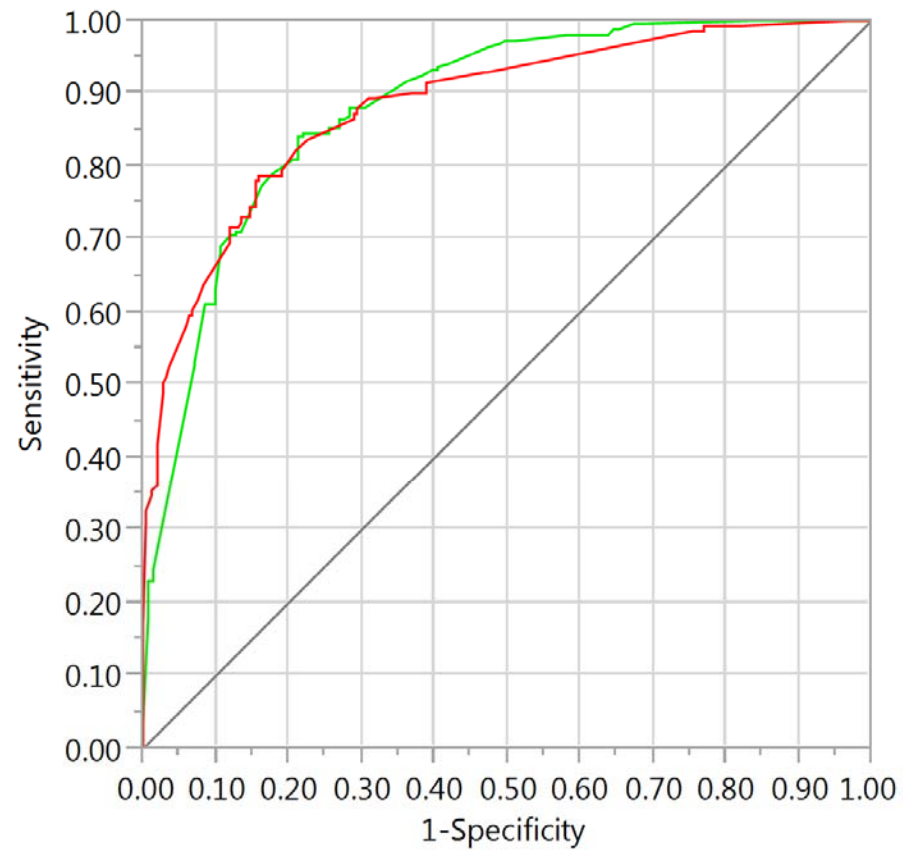
## Overall Statistics

Measure	Training	Validation	Definition
Entropy RSquare	0.2188	0.2527	$1 - \text{Loglike}(\text{model}) / \text{Loglike}(0)$
Generalized RSquare	0.3442	0.3854	$(1 - (L(0)/L(\text{model}))^{2/n}) / (1 - L(0)^{2/n})$
Mean -Log p	0.5231	0.4878	$\sum -\text{Log}(p[j]) / n$
RMSE	0.4116	0.3918	$\sqrt{\sum (y[j] - p[j])^2 / n}$
Mean Abs Dev	0.3803	0.3644	$\sum  y[j] - p[j]  / n$
Misclassification Rate	0.2063	0.1858	$\sum (p[j] \neq p_{\text{Max}}) / n$
N	916	393	n

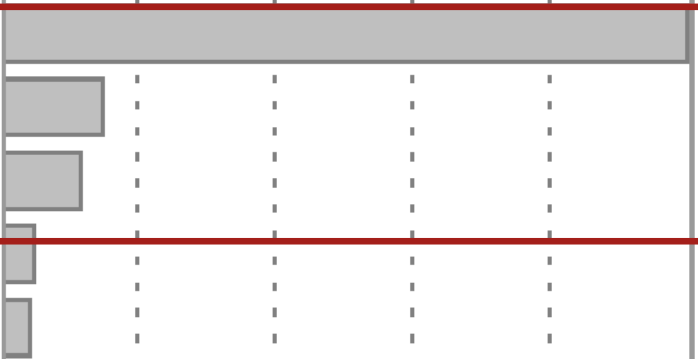
## Confusion Matrix

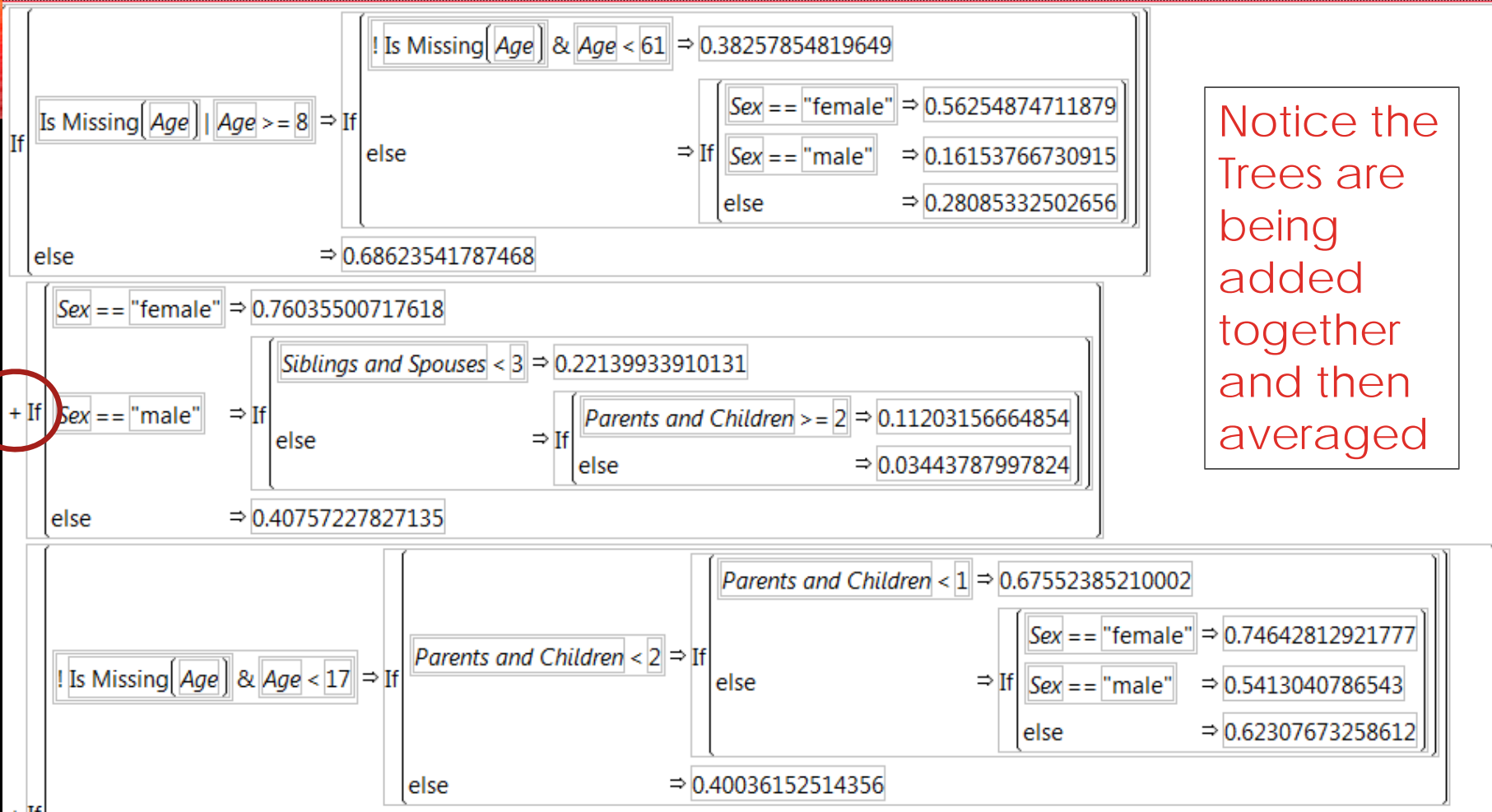
Actual	Predicted		Actual	Predicted	
Training	No	Yes	Validation	No	Yes
No	482	75	No	215	37
Yes	114	245	Yes	36	105



**Receiver Operating Characteristic on Validation Data**

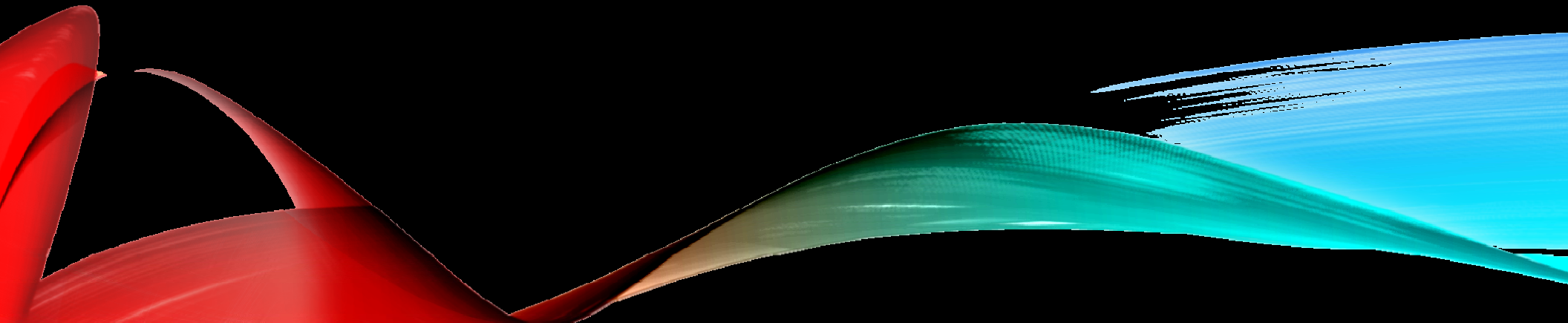
## Column Contributions

Term	Number of Splits	G <sup>2</sup>		Portion
Sex	7	747.889083		0.7371
Passenger Class	3	111.684997		0.1101
Age	8	86.7514698		0.0855
Siblings and Spouses	7	37.0069685		0.0365
Parents and Children	9	31.2527691		0.0308



Notice the  
Trees are  
being  
added  
together  
and then  
averaged

# BOOSTED TREE



# BOOSTED TREE

23

Boosting is the process of **building a large, additive decision tree** by fitting a sequence of smaller trees. **Each of the smaller trees is fit on the scaled residuals of the previous tree**. The trees are combined to form the larger final tree. The process can use validation to assess how many stages to fit, not to exceed the specified number of stages.

The tree at each stage is short, typically 1-5 splits. **After the initial tree, each stage fits the residuals from the previous stage**. The process continues until the specified number of stages is reached, or, if validation is used, until fitting an additional stage no longer improves the validation statistic. **The final prediction is the sum of the estimates for each terminal node over all the stages**.

If the response is categorical, the residuals fit at each stage are offsets of linear logits. The final prediction is a logistic transformation of the sum of the linear logits over all the stages. For categorical responses, only those with two levels are supported.

[JMP, Chapter 3, Specialized Models, Partition Models]



# Boosting

24

- Like bagging, boosting is a general approach that can be applied to many statistical learning methods for regression or classification. We only discuss boosting for decision trees.
- Recall that bagging involves creating multiple copies of the original training data set using the bootstrap, fitting a separate decision tree to each copy, and then combining all of the trees in order to create a single predictive model.
- Notably, each tree is built on a bootstrap data set, independent of the other trees.
- Boosting works in a similar way, except that the trees are grown *sequentially*: each tree is grown using information from previously grown trees.

Stanford Online Course: Course Notes, Statistical Learning by Hastie and Tibshirani, Winter 2014 | Based on [An Introduction to Statistical Learning](#) by James, Witten, Hastie and Tibshirani, Springer (2013)

## What is the idea behind this procedure?

- Unlike fitting a single large decision tree to the data, which amounts to *fitting the data hard* and potentially overfitting, the boosting approach instead *learns slowly*.
- Given the current model, we fit a decision tree to the residuals from the model. We then add this new decision tree into the fitted function in order to update the residuals.
- Each of these trees can be rather small, with just a few terminal nodes, determined by the parameter  $d$  in the algorithm.
- By fitting small trees to the residuals, we slowly improve  $\hat{f}$  in areas where it does not perform well. The shrinkage parameter  $\lambda$  slows the process down even further, allowing more and different shaped trees to attack the residuals.

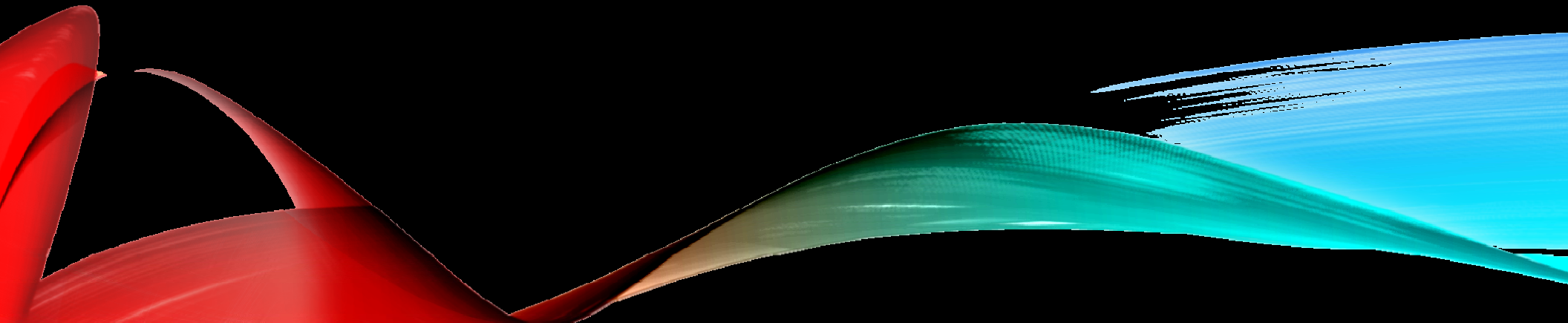
Stanford Online Course: Course Notes, Statistical Learning by Hastie and Tibshirani, Winter 2014 | Based on [An Introduction to Statistical Learning](#) by James, Witten, Hastie and Tibshirani, Springer (2013)

## Tuning parameters for boosting

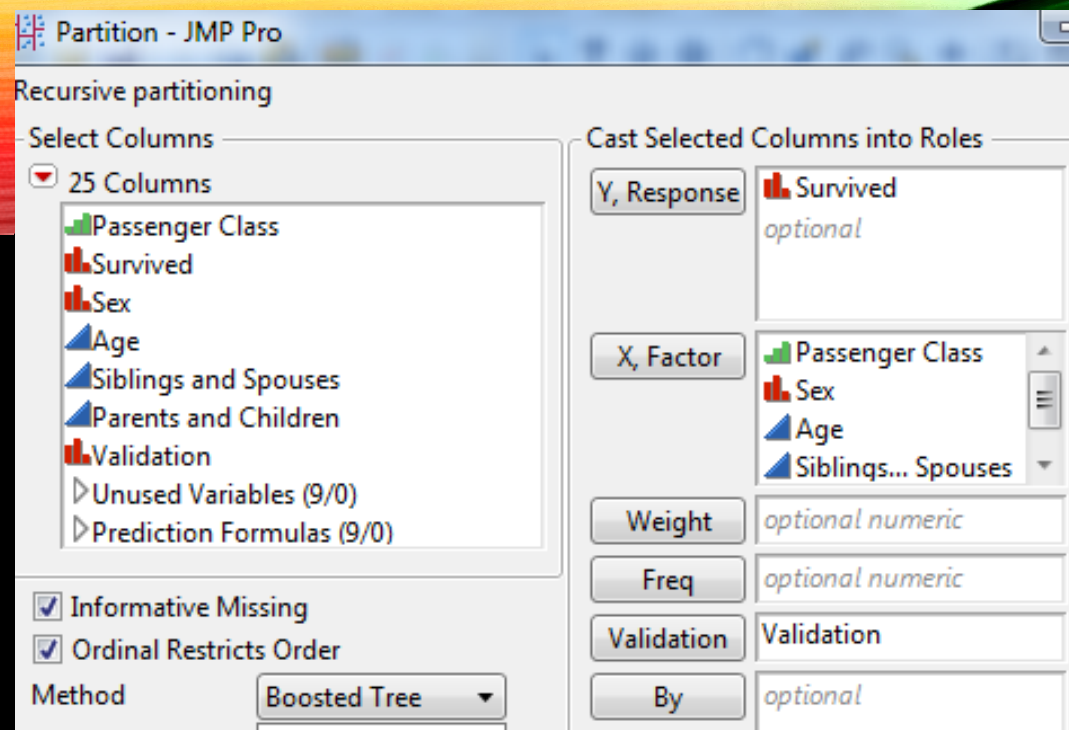
1. The *number of trees*  $B$ . Unlike bagging and random forests, boosting can overfit if  $B$  is too large, although this overfitting tends to occur slowly if at all. We use cross-validation to select  $B$ .
2. The *shrinkage parameter*  $\lambda$ , a small positive number. This controls the rate at which boosting learns. Typical values are 0.01 or 0.001, and the right choice can depend on the problem. Very small  $\lambda$  can require using a very large value of  $B$  in order to achieve good performance.
3. The *number of splits*  $d$  in each tree, which controls the complexity of the boosted ensemble. Often  $d = 1$  works well, in which case each tree is a *stump*, consisting of a single split and resulting in an additive model. More generally  $d$  is the *interaction depth*, and controls the interaction order of the boosted model, since  $d$  splits can involve at most  $d$  variables.

Stanford Online Course: Course Notes, Statistical Learning by Hastie and Tibshirani, Winter 2014 | Based on [An Introduction to Statistical Learning](#) by James, Witten, Hastie and Tibshirani, Springer (2013)

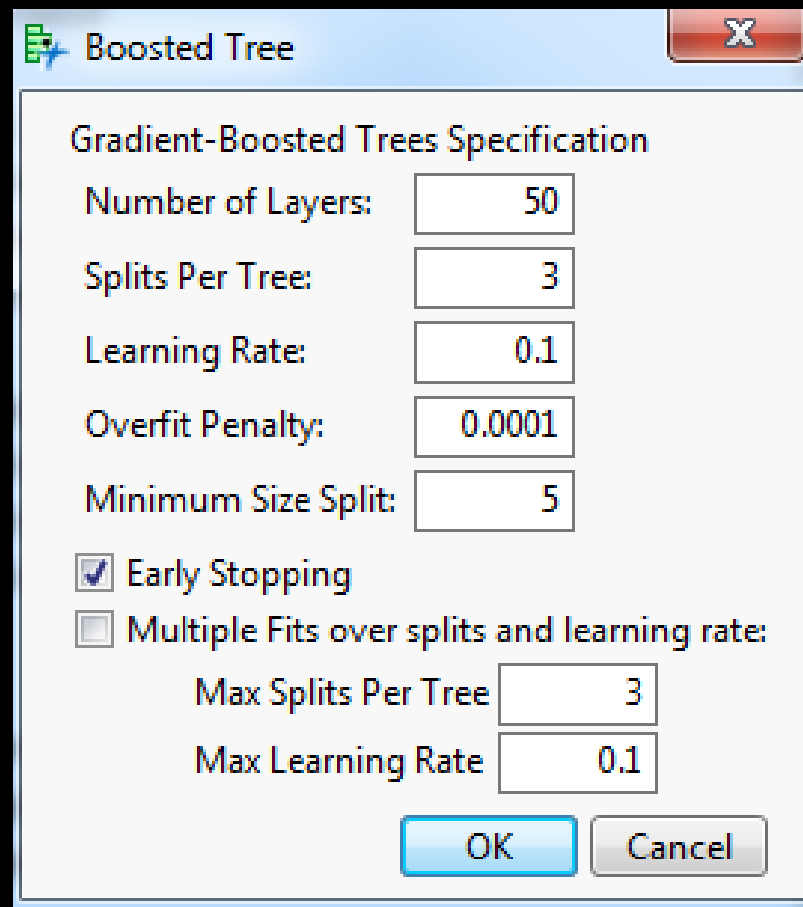
# BOOSTED TREE CATEGORICAL RESPONSE







Use Defaults





## Overall Statistics

Measure	Training	Validation	Definition
Entropy RSquare	0.3513	0.3788	$1 - \text{Loglike}(\text{model}) / \text{Loglike}(0)$
Generalized RSquare	0.5085	0.5352	$(1 - (L(0)/L(\text{model}))^{2/n}) / (1 - L(0)^{2/n})$
Mean -Log p	0.4344	0.4055	$\sum -\text{Log}(p[j]) / n$
RMSE	0.3716	0.3537	$\sqrt{\sum (y[j] - p[j])^2 / n}$
Mean Abs Dev	0.2881	0.2791	$\sum  y[j] - p[j]  / n$
Misclassification Rate	0.1932	0.1756	$\sum (p[j] \neq p_{\text{Max}}) / n$
N	916	393	n






## Confusion Matrix

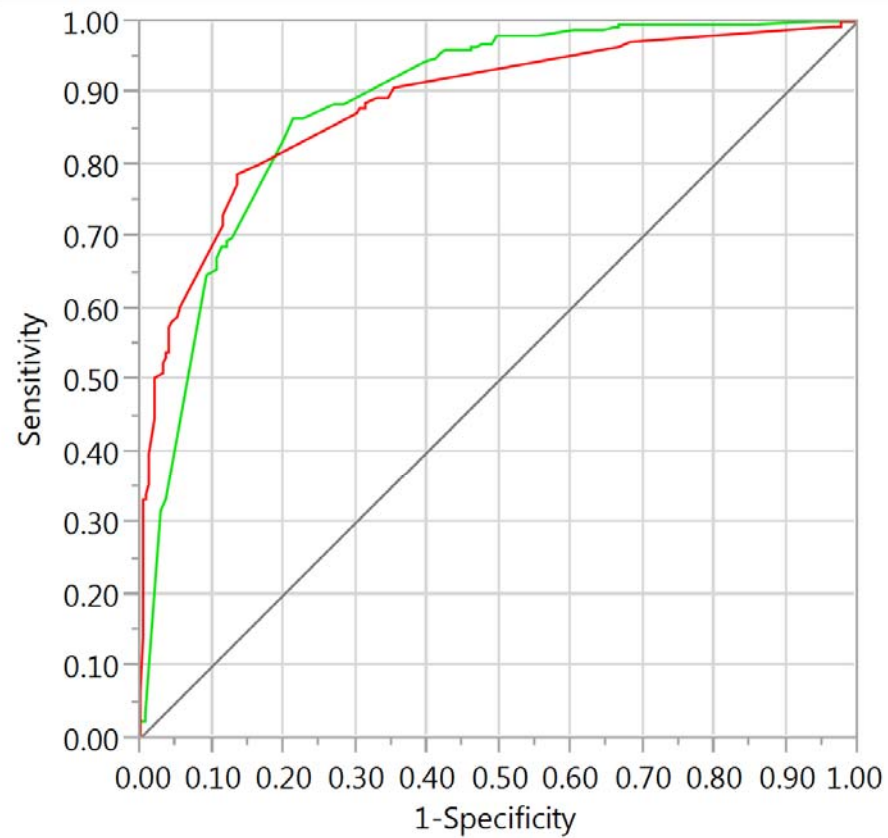
Actual	Predicted	
Training	No	Yes
No	509	48
Yes	129	230

Actual	Predicted	
Validation	No	Yes
No	223	29
Yes	40	101

## Column Contributions

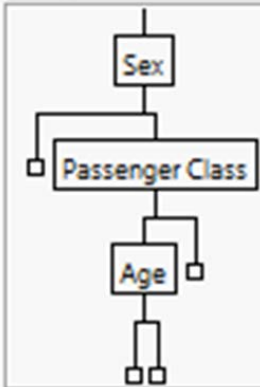
Term	Number of Splits	G <sup>2</sup>		Portion
Sex	57	54527.9454		0.5299
Age	23	20263.8193		0.1969
Passenger Class	39	19984.6722		0.1942
Siblings and Spouses	13	4208.00667		0.0409
Parents and Children	18	3922.45484		0.0381

**Receiver Operating Characteristic on Validation Data**

Survived	Area
No	0.8855
Yes	0.8855

## Tree Views

### Layer1



### Layer2

### Layer3

### Layer4

### Layer5

### Layer6

### Layer7

Logist

-0.439242851439

Sex == "male" ⇒ -0.0893205561555

+ If Sex == "female" ⇒ If

    Passenger Class == 1 | Passenger Class == 2 ⇒ If

        ! Is Missing Age & Age < 60 ⇒ 0.32458983863027

    else ⇒ 0.21130867824872

    Passenger Class == 3 ⇒ 0.03702503610437

    else ⇒ 0

else ⇒ 0

Sex == "male" ⇒ -0.0803885005401

+ If Sex == "female" ⇒ If

    Passenger Class == 1 | Passenger Class == 2 ⇒ If

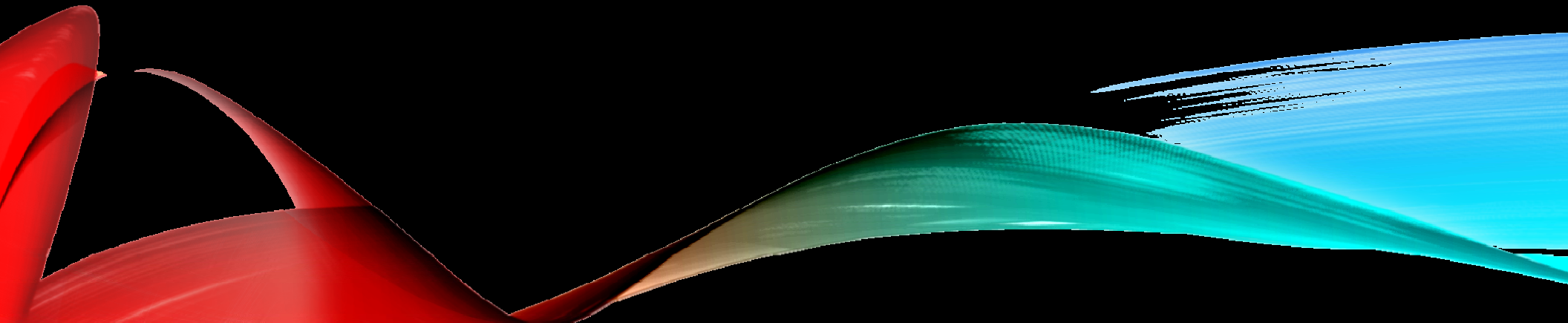
        ! Is Missing Age & Age < 60 ⇒ 0.29213084003711

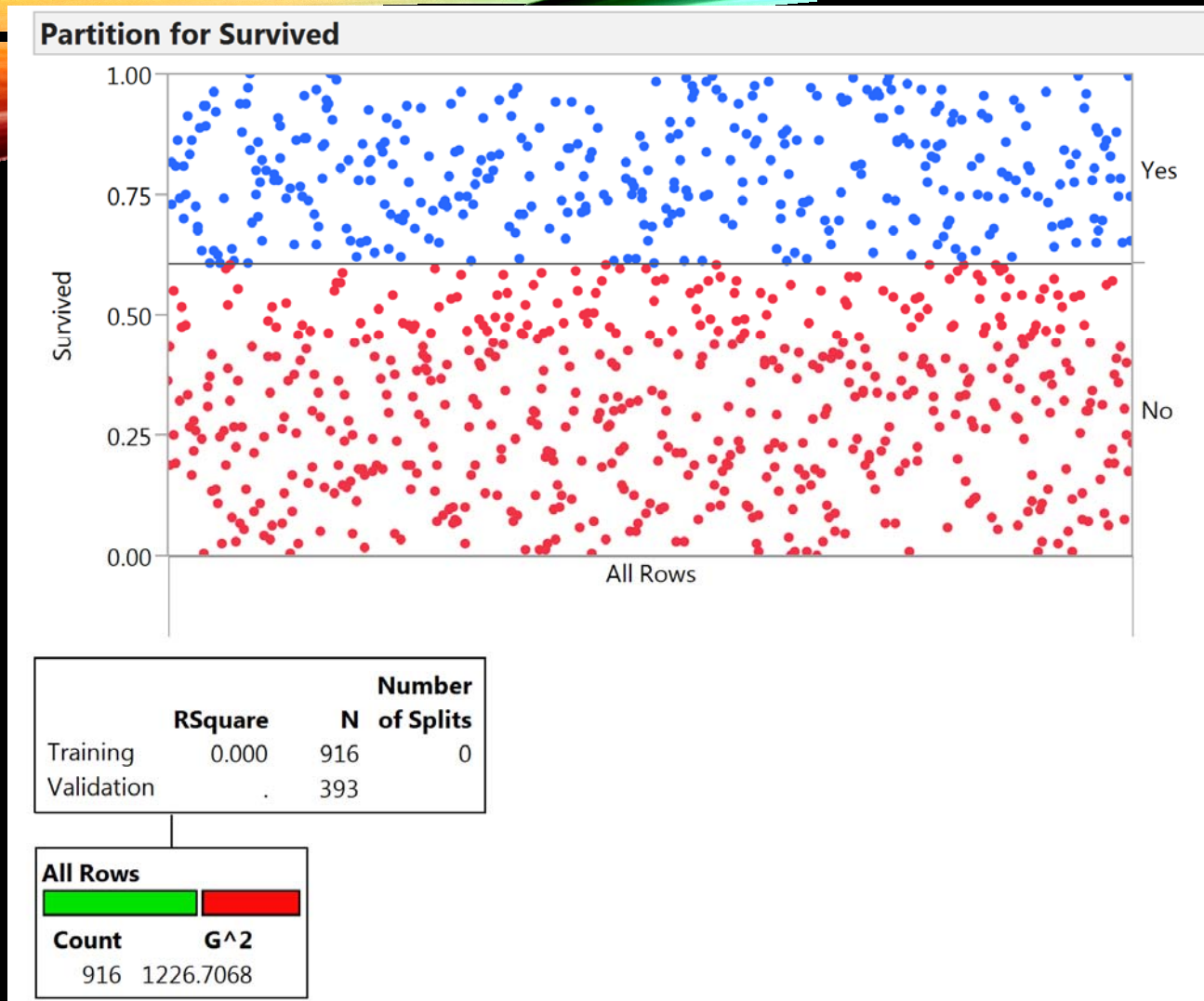
    else ⇒ 0.19017781025151

    Passenger Class == 3 ⇒ 0.03332263313794

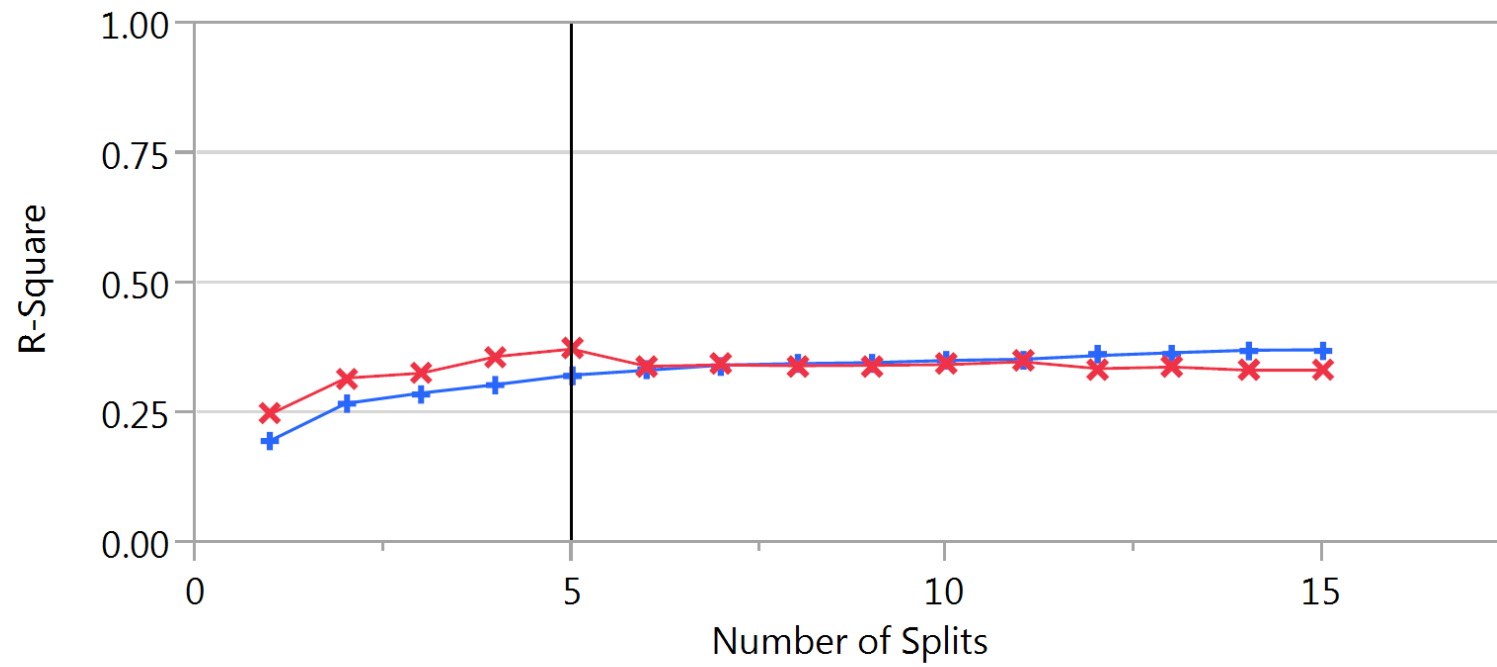
    else ⇒ 0

# DECISION TREE



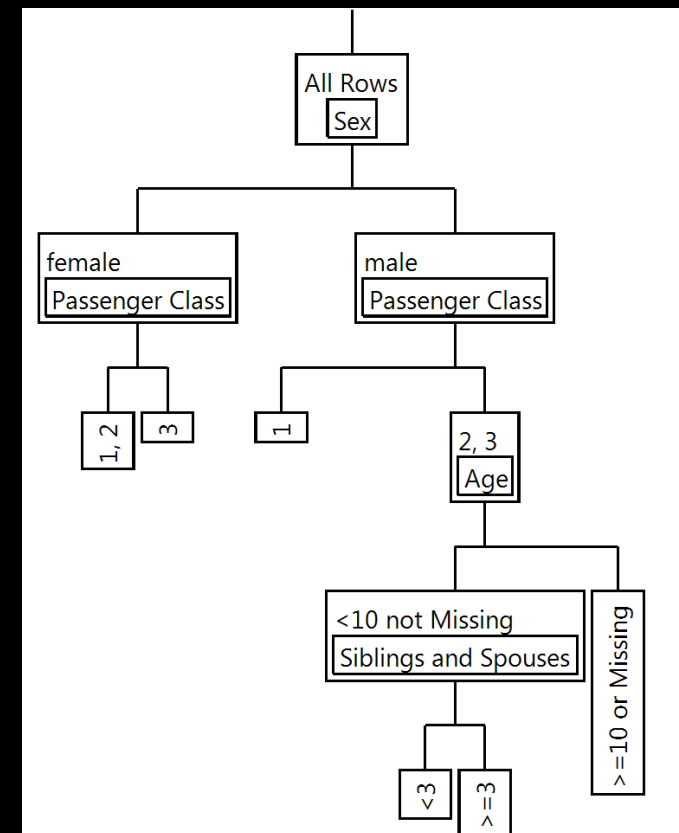
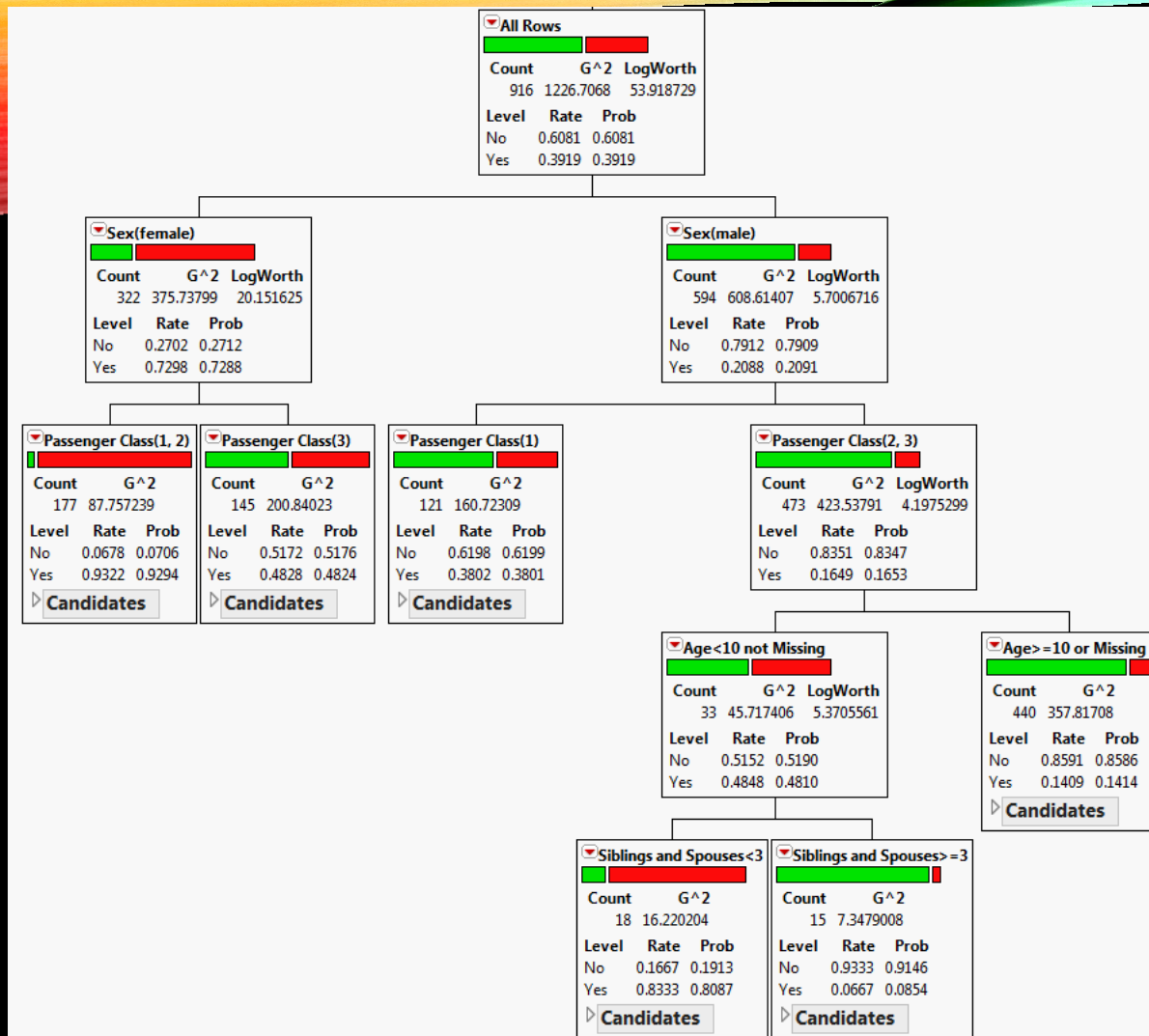


## Split History



Validation Data in Red





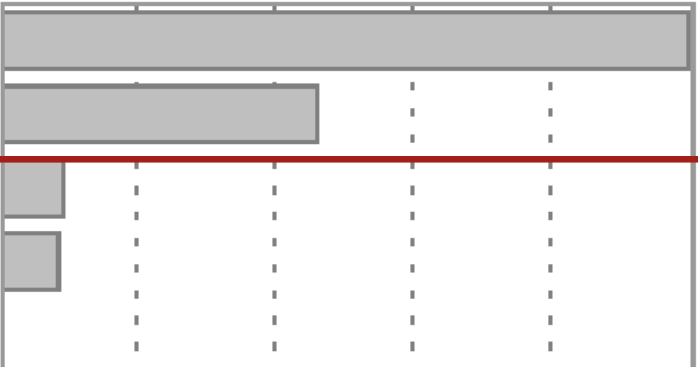
## Fit Details

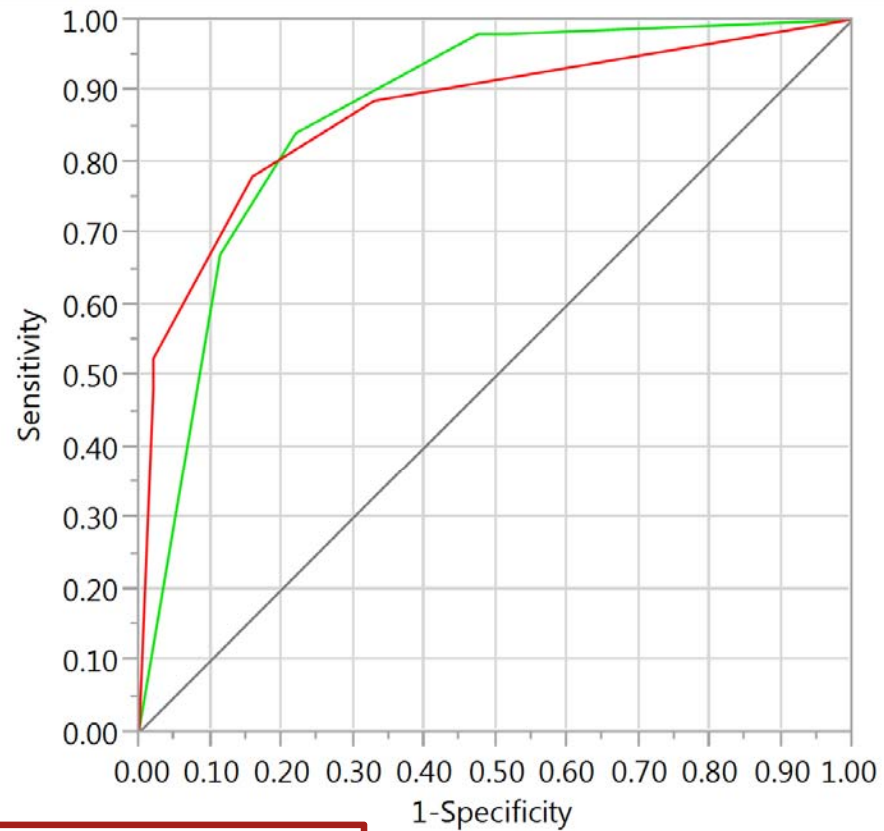
Measure	Training	Validation	Definition
Entropy RSquare	0.3227	0.3729	$1 - \text{Loglike}(\text{model}) / \text{Loglike}(0)$
Generalized RSquare	0.4755	0.5287	$(1 - (L(0)/L(\text{model}))^{2/n}) / (1 - L(0)^{2/n})$
Mean -Log p	0.4535	0.4093	$\sum -\text{Log}(p[j]) / n$
RMSE	0.3805	0.3560	$\sqrt{\sum (y[j] - p[j])^2 / n}$
Mean Abs Dev	0.2908	0.2782	$\sum  y[j] - p[j]  / n$
Misclassification Rate	0.2118	0.1832	$\sum (p[j] \neq p_{\text{Max}}) / n$
N	916	393	n

## Confusion Matrix

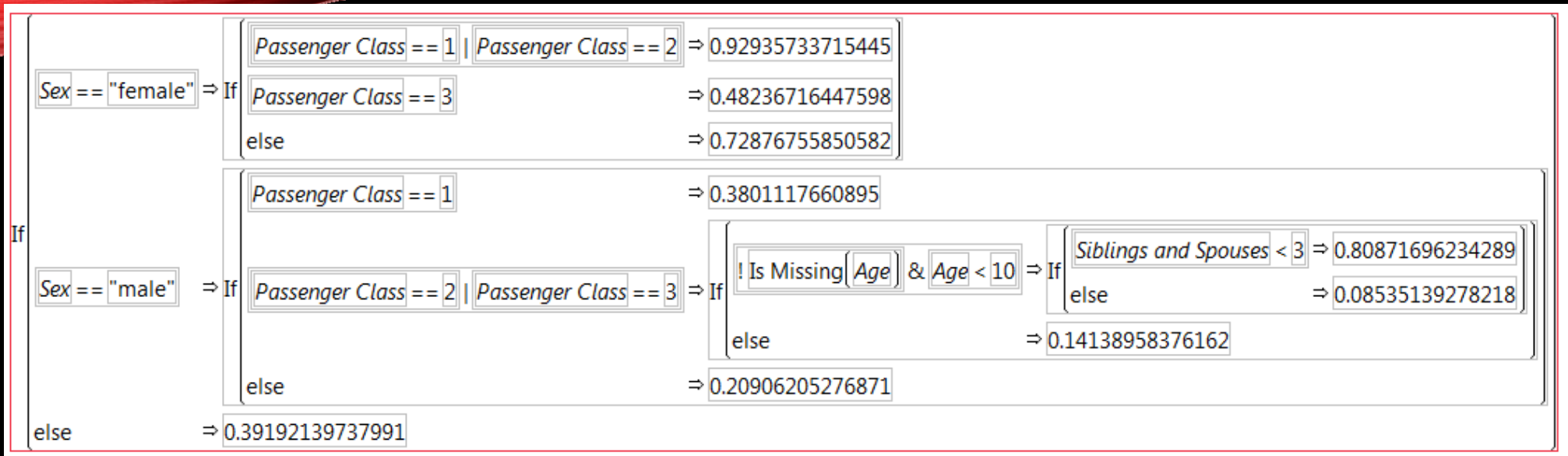
Actual	Predicted		Actual	Predicted	
Training	No	Yes	Validation	No	Yes
No	542	15	No	247	5
Yes	179	180	Yes	67	74

## Column Contributions

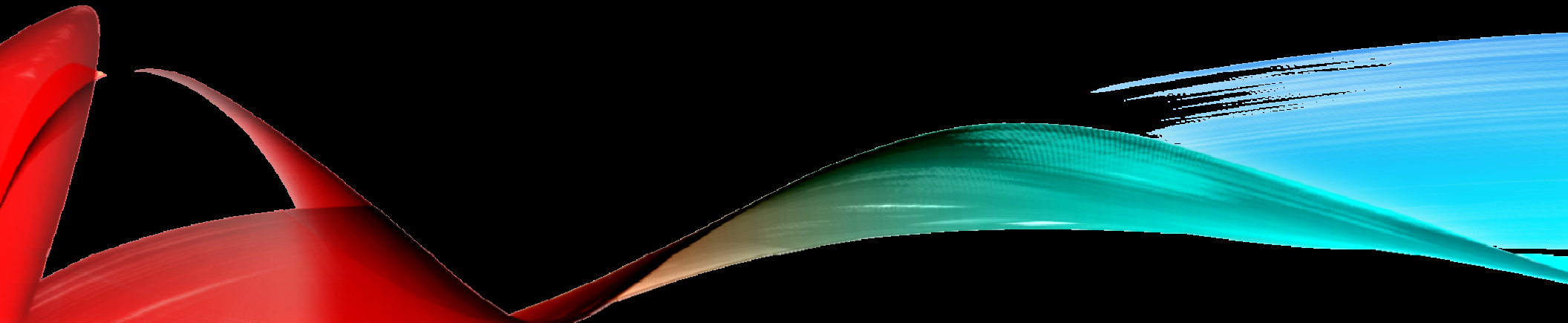
Term	Number of Splits	G <sup>2</sup>		Portion
Sex	1	242.354767		0.6120
Passenger Class	2	111.493595		0.2815
Siblings and Spouses	1	22.1493019		0.0559
Age	1	20.003416		0.0505
Parents and Children	0	0		0.0000

**Receiver Operating Characteristic on Validation Data**

Survived	Area
No	0.8704
Yes	0.8704



# MODEL COMPARISON



Analyze Graph Tools Add-Ins View Window Help

Distribution  
Fit Y by X  
Matched Pairs  
Tabulate  
Fit Model  
**Modeling**  
Multivariate Methods  
Quality and Process  
Reliability and Survival  
Consumer Research

Partition  
Neural  
**Model Comparison**  
Nonlinear  
Gaussian Process  
Time Series

Name	Prob(Survive d==No)	Prob(Survive d==Yes)	Mo
Elisabeth Walton	0.2607227958	0.7392772042	Yes
ster. Hudson Trevor	0.4844378517	0.5155621483	Yes
s. Helen Loraine	0.2607227958	0.7392772042	Yes
Hudson Joshua Creighton	0.8084776147	0.1915223853	No
	0.2607227958	0.7392772042	Yes
	0.8084776147	0.1915223853	No
	0.1554551458	0.8445448542	Yes
			No
			Yes
	0.8084776147	0.1915223853	No
	0.8084776147	0.1915223853	No
12 17.5 Astor, Mrs.	0.2607227958	0.7392772042	Yes

Find predictor columns for the same target response and compare how well they fit



## Cast Selected Columns into Roles

Y, Predictors *optional numeric*Group *optional*Weight *optional numeric*Freq *optional numeric*By Validation *optional*

If you choose no Predictor columns, it will find and analyze all predictors.

## Model Comparison Validation= Training

## Predictors

## Measures of Fit for Survived

Creator	.2 .4 .6 .8	Entropy RSquare	Generalized RSquare	Mean -Log p	RMSE	Mean Abs Dev	Misclassification Rate	N
Bootstrap Forest		0.2188	0.3442	0.5231	0.4116	0.3803	0.2063	916
Boosted Tree		0.3513	0.5085	0.4344	0.3716	0.2881	0.1932	916
Partition		0.3227	0.4755	0.4535	0.3805	0.2908	0.2118	916

## Model Comparison Validation= Validation

## Predictors

## Measures of Fit for Survived

Creator	.2 .4 .6 .8	Entropy RSquare	Generalized RSquare	Mean -Log p	RMSE	Mean Abs Dev	Misclassification Rate	N
Bootstrap Forest		0.2527	0.3854	0.4878	0.3918	0.3644	0.1858	393
Boosted Tree		0.3788	0.5352	0.4055	0.3537	0.2791	0.1756	393
Partition		0.3729	0.5287	0.4093	0.3560	0.2782	0.1832	393

## Model Comparison Validation=Validation

### Predictors

#### Target Column Predictors

Survived

#### Category Probability Column

No Prob(Survived==No) Bootstrap Forest

Yes Prob(Survived==Yes)

#### Category Probability Column

No Prob(Survived==No) 2 Boosted Tree




Yes Prob(Survived==Yes) 2

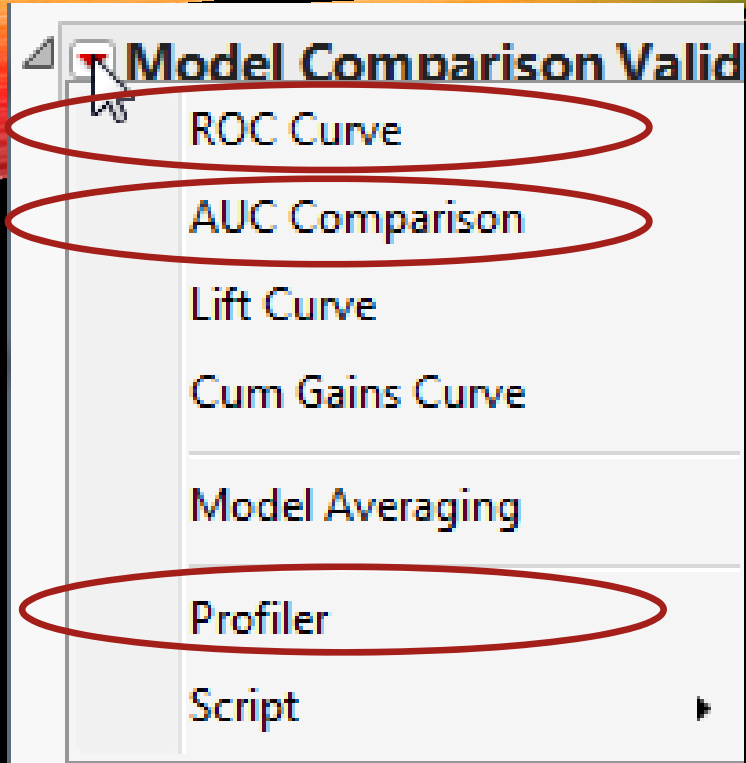
#### Category Probability Column

No Prob(Survived==No) 3 Partition

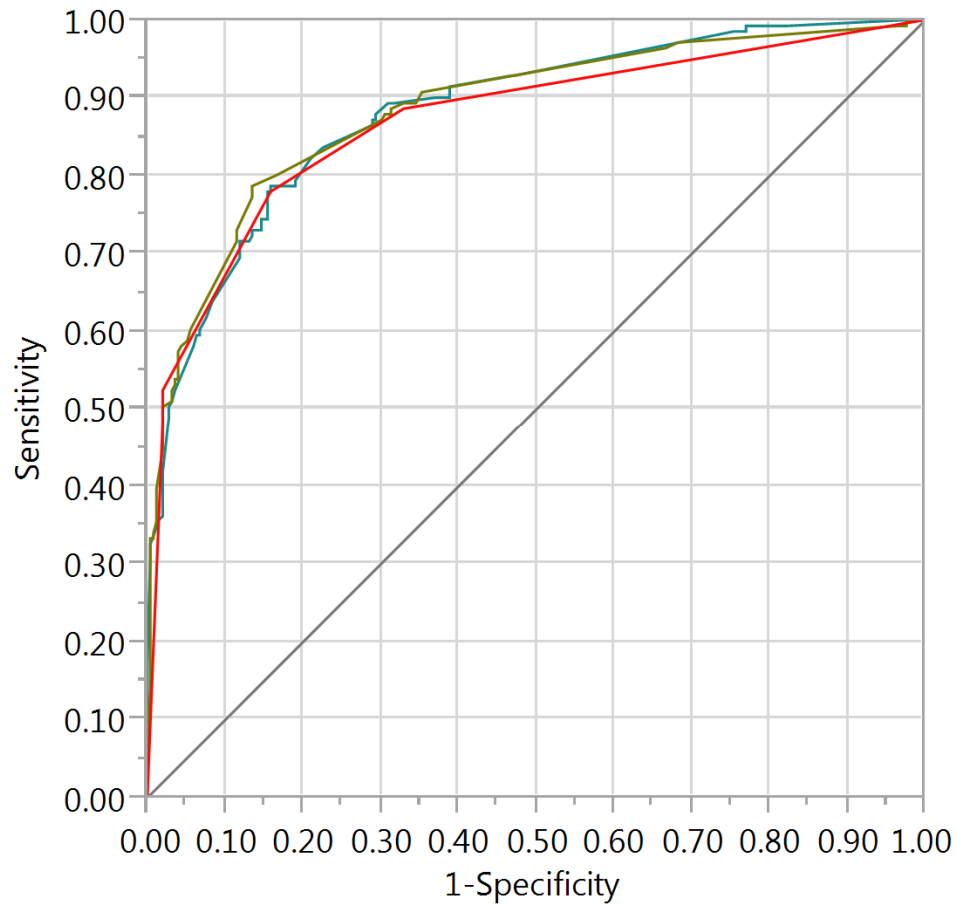
Yes Prob(Survived==Yes) 3

### Measures of Fit for Survived

Creator	.2 .4 .6 .8	Entropy	Generalized	Mean -Log p	RMSE	Mean	Misclassification	N
		RSquare	RSquare			Abs Dev	Rate	
Bootstrap Forest		0.2527	0.3854	0.4878	0.3918	0.3644	0.1858	393
Boosted Tree		0.3788	0.5352	0.4055	0.3537	0.2791	0.1756	393
Partition		0.3729	0.5287	0.4093	0.3560	0.2782	0.1832	393



### ROC Curve for Survived=Yes



Predictor	AUC
Prob(Survived==Yes)	0.8833
Prob(Survived==Yes) 2	0.8855
Prob(Survived==Yes) 3	0.8704

## AUC Comparison

### AUC Comparison for Survived=Yes

Predictor	AUC	Std Error	Lower 95%	Upper 95%
Prob(Survived==Yes)	0.8833	0.0178	0.8437	0.9139
Prob(Survived==Yes) 2	0.8855	0.0182	0.8446	0.9167
Prob(Survived==Yes) 3	0.8704	0.0192	0.8279	0.9037

		AUC					
Predictor	vs. Predictor	Difference	Std Error	Lower 95%	Upper 95%	ChiSquare	Prob>ChiSq
Prob(Survived==Yes)	Prob(Survived==Yes) 2	-0.002	0.0064	-0.015	0.0103	0.1176	0.7317
Prob(Survived==Yes)	Prob(Survived==Yes) 3	0.0129	0.0076	-0.002	0.0278	2.8789	0.0897
Prob(Survived==Yes) 2	Prob(Survived==Yes) 3	0.0151	0.0074	0.0007	0.0295	4.1977	0.0405*

Test	ChiSquare	DF	Prob>ChiSq
All AUCs equal	4.45261	2	0.1079

## Profiler Validation=Validation

### Prediction Profiler

