

# Big Data: Data Analysis Boot Camp

## Titanic Dataset

Chuck Cartledge, PhD

22 September 2017

# Table of contents (1 of 1)

1 Introduction

2 Background

3 Classification problem

4 Techniques

5 Hands-on

6 Q & A

7 Conclusion

8 References

9 Files

# What are we going to cover?

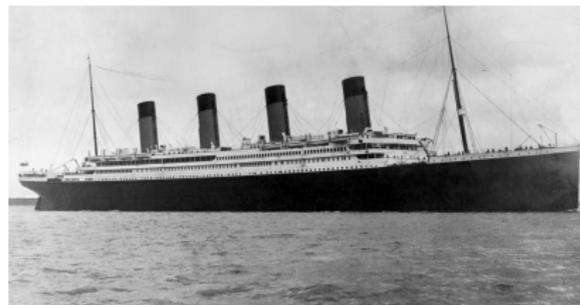
We're going to talk about:

- R's RMS Titanic dataset.
- Other Titanic datasets that contain different data.
- Modeling the datasets to see who will live and who will die.



# Basic information

- Ordered: 17 Sep. 1908
- Completed: 2 Apr. 1912
- Maiden voyage: 10 Apr. 1912
- Sank: 14 Apr. 1912



“Well” settled data

## Where she was damaged

- Red are water tight bulkheads
- Green is where the iceberg hit

As the bow settled, water overflowed the bulkheads

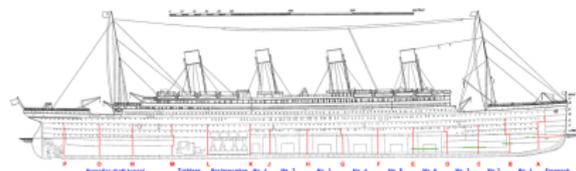
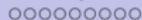


Image from [11].



“Well” settled data

# Same image.

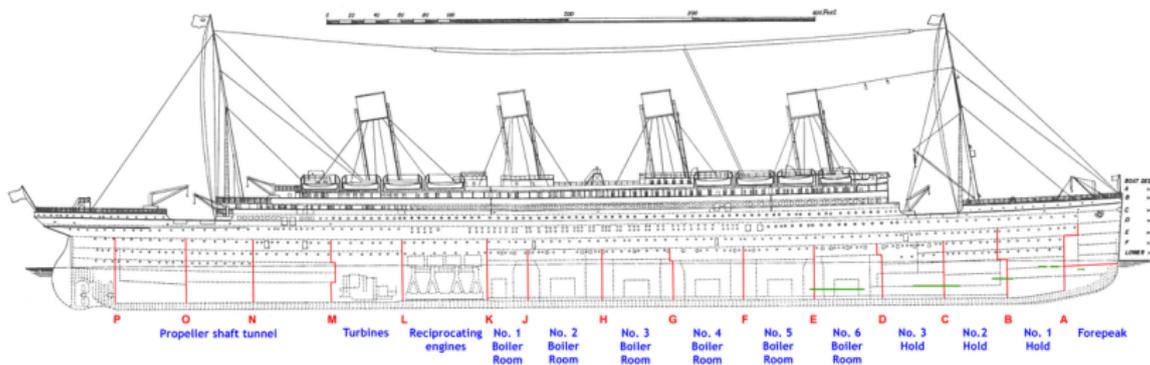


Image from [11].

“Well” settled data

## How many died and why?

- Sailing capacity (passengers and crew): 3,372
- Lifeboat capacity: 1,178
- Number of people on board (accounts vary): 2,201
- Number of people who survived: ~706 - 712 (R thinks 711)

Passengers, crew, builder's men, and others.



Image from [8].

Data from diverse places

## Expected first class passengers

Lots of lists of 1<sup>st</sup> class passengers. Even, some of 2<sup>nd</sup>, and 3<sup>rd</sup> class passengers[10]. Lists of non-passengers (ship's crew, and builder's technicians) are more challenging[6]. R has a built-in Titanic dataset: `Titanic`



Image from [7].

# A crew list

A reasonable collection of crew and builder's representatives is available.

- Name, job, status (lost or not)
- Age, place of birth

## A

**Abbott, Ernest Owen** Pantryman *Lost*  
98, Northumberland Road, Nicholstown, Southampton  
Age: 21 Place of Birth: Hants, Southampton

**Abraham, C**  
*see Abrams, William*

**Abrams, H**  
*see Abrams, William*

**Abrams, William** Fireman *Lost*  
3 or 11, Charles Street, Chapel, Southampton  
Age: 35 Place of Birth: Northwich

**Adams, R** Fireman *Lost*  
168, Romsey Road, Shirley, Southampton  
Age: 26 Place of Birth: Hants  
*Crew Agreement has 168 Pound Tree Road*

**Ahier, Percy Snowden** Steward *Lost*  
136, Northumberland Road, Nicholstown, Southampton  
Age: 20 Place of Birth: Jersey

**Akerman, Albert** Steward *Lost*  
25, Rochester Street, Northam, Southampton  
Age: 28 Place of Birth: Salisbury, Wiltshire

**Akerman, Joseph Francis** Assistant Pantryman *Lost*  
25, Rochester Street, Northam, Southampton  
Age: 35 Place of Birth: Salisbury, Wiltshire

Image from [9].



Data from diverse places

# titanic3 dataset from PASWR[14]

- Part of the PASWR library
- Thomas Cason of UVA has greatly updated and improved the titanic data frame using the *Encyclopedia Titanic*.
- Focuses and expands the passenger data.

A1	id	name	sex	age	sibsp	parch	ticket	fare	cabin	status
2	1	Allen, Mrs. Elizabeth Walter	female	29	0	0	24108	## 85		
3	1	Alison, Mavis Helen Traver	female	18	1	2	11781	## C22 C26		
4	1	Alison, Miss Helen Louisa	female	2	1	2	11781	## C22 C26		
5	1	Alison, Miss Helen Louisa	female	38	1	2	11781	## C22 C26		
6	1	Alison, Mrs. Hudson FC (Helen Mable Denzil)	female	25	1	2	11781	## C22 C26		
7	1	Andrew, Mr Henry	male	48	0	0	19952	28 5508 812		
8	1	Andrew, Miss Elizabeth Theodosia	female	82	0	0	12862	71 8088 187		
9	1	Andrew, Mr Thomas D	male	39	0	0	13208	8 8088 436		
10	1	Appleton, Mrs Edward Dale (Frances Louisa)	female	52	2	0	11781	16 8700 C101		
11	1	Argueyrie, Mr Rouen	male	71	0	0	1JC 17689	48 5042		
12	1	Aster, Col John Jacob	male	47	1	0	1JC 17787	## C82 C84		
13	1	Aster, Mrs John Jacob (Muriel Rose Yolande Fox)	female	18	1	0	1JC 17787	## C82 C84		
14	1	Baker, Mrs Lavonia Pauline	female	24	0	0	1JC 17477	68 3088 835		
15	1	Baker, Mrs Ellen "Nellie"	female	26	0	0	19871	18 8088		
16	1	Barkworth, Mr Algeman Henry Wilson	male	89	0	0	27042	18 8088 423		
17	1	Barnes, Mr John D	male	61	0	0	1JC 17588	18 8088		
18	1	Barnes, Mr Oleg Edward	male	24	0	0	1JC 17588	## 858 868		
19	1	Barnes, Mrs James (Helen DeLondres Chapin)	female	50	0	0	1JC 17588	## 858 868		
20	1	Barron, Mrs John	female	22	0	0	10811	16 2917 003		
21	1	Baxter, Mr Thomas	male	36	0	0	13058	75 2417 C8		
22	1	Berkich, Mr Richard Leonard	male	37	1	1	18751	16 1948 003		
23	1	Berkich, Mrs Richard Leonard (Julie Margaret)	female	47	1	1	18751	12 5542 C05		
24	1	Berg, Mr Carl Frederick	male	26	0	0	10188	18 8088 C46		
25	1	Berg, Mrs Rosalie	female	42	0	0	1JC 17483	## C87		
26	1	Bird, Mrs Ellen	female	29	0	0	1JC 17483	## C87		
27	1	Birkbeck, Mr John	male	23	0	0	19883	28 8088		
28	1	Bishop, Mr Dickson H	male	25	1	0	10567	16 8782 849		
29	1	Bishop, Mrs Dickson H (Julia Helen)	female	18	1	0	10567	16 8782 849		
30	1	Bishop, Miss Amelia	female	35	0	0	1JC 17369	## C39		
31	1	Blackburne, Mrs Minnie Helen	female	28	0	0	10984	26 5508 C52		
32	1	Blackburne, Mr Stephen Lloyd	male	41	0	0	11781	16 8088 F		
33	1	Blank, Mr Henry	male	49	0	0	13277	18 8088 431		

Image from [2].



# titanic3 attributes/variables

Name	Explanation
Pclass	Passenger Class (1 = 1 <sup>st</sup> ; 2 = 2 <sup>nd</sup> ; 3 = 3 <sup>rd</sup> )
survival	Survival (0 = No; 1 = Yes)
name	Name
sex	Sex
age	Age
sibsp	Number of Siblings/Spouses Aboard
parch	Number of Parents/Children Aboard
ticket	Ticket Number
fare	Passenger Fare (British pound)
cabin	Cabin
embarked	Port of Embarkation (C = Cherbourg; Q = Queenstown; S = Southampton)
boat	Lifeboat
body	Body Identification Number
home.dest	Home/Destination

# Bringing the pieces together

Combining:

- Passenger data from `titanic3`
- Crew data from Southampton
- Not all data in both datasets

Get a reasonable estimation of who survived, or not when the RMS Titanic went down.



# A definition

*“Classification is the task of learning a **target function**  $f$  that maps each attribute set  $\mathbf{x}$  to one of the predefined class labels  $\mathbf{y}$ .”*

*Tan, et al. [12]*



## What is it?

# As a picture

- 1 A collection of correctly labeled data (training data) is available.
- 2 The supervised data is processed by some sort of machine learning algorithm (there are many) to create a model (or classifier).
- 3 Unlabeled (test or new) data, is processed by the model and predictions are made.

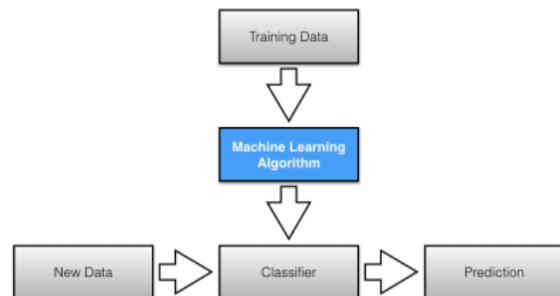


Image from [5].



What is it?

# Same image.

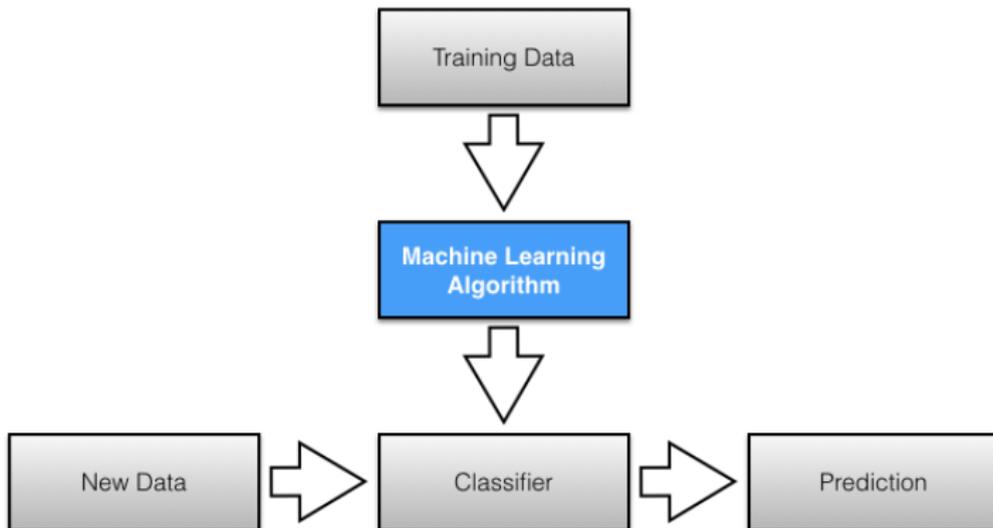


Image from [5].

# Supervised vs. Unsupervised learning

- **Supervised learning**

A training dataset with correct answers (labels) is “mined” to create a model

- **Unsupervised learning**

Data are provided with no apriori knowledge of labels or patterns. The goal is to discover labels and patterns.

- **Semi-supervised learning**

Knowledge from one dataset is applied to another dataset to help with mining, analysis, classification, and interpretation.





What is it?

## Supervised vs. Unsupervised learning techniques

	<i>Supervised Learning</i>	<i>Unsupervised Learning</i>
<i>Discrete</i>	classification or categorization	clustering
<i>Continuous</i>	regression	dimensionality reduction

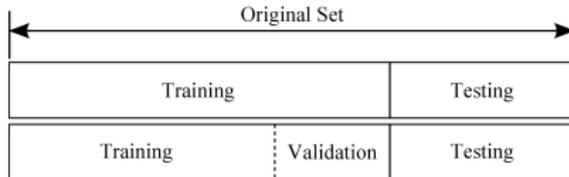
With the Titanic dataset, we will be focusing on classification.

# Working with data

Supervised learning requires:

- Training data – usually about 70% of available data
- Testing data – usually about 30% of available data

Training data can also be partitioned into validation data





# Lots of different things can be done with training data

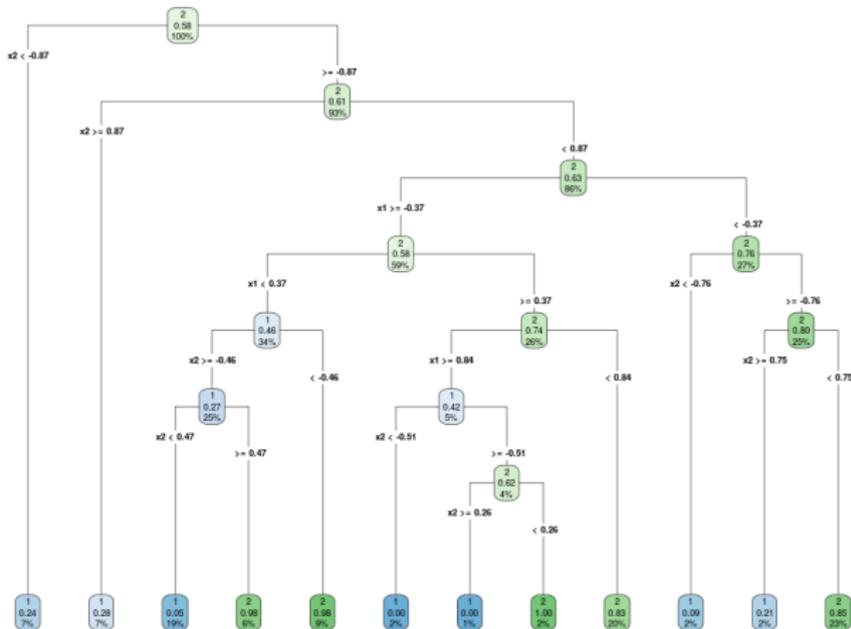
- Use as one monolithic entity
- Randomly sample data (with and without replacement)
- Divide original training data into training and validation subsets to create multiple models
- With multiple models:
  - Choose best one,
  - Use all and vote on the outcome





## Types of errors

# Same image.

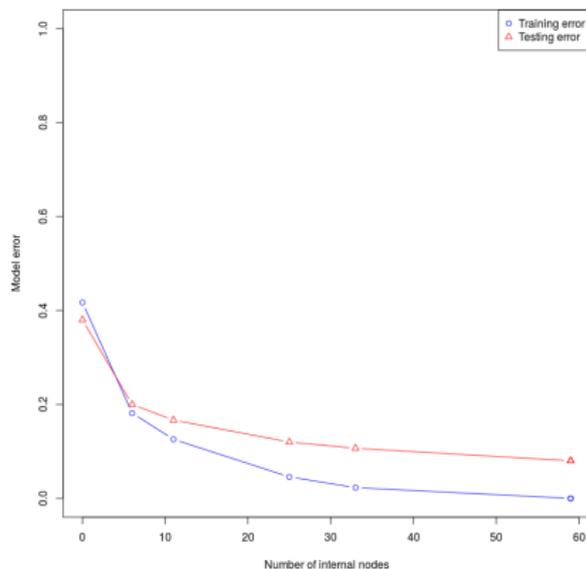


Attached file.



# Errors in machine learning

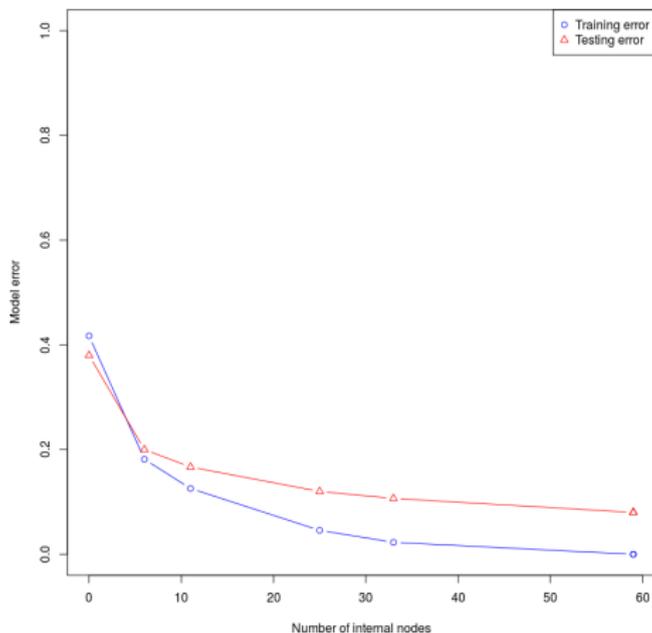
- Total sample divided into training (70%) and testing (30%) datasets
- Training dataset was partitioned into different sized decision trees (models)
- Training and testing datasets were classified using each model
- Results were compared to the original data
- Initially models **under-fitted** until around 6 nodes
- Finally models **over-fitted** beyond 25 nodes



Training and testing errors

## Types of errors

# Same image.



## Training and testing errors

# A collection of decision tree techniques

**rpart** from the **rpart** library. "Recursive partitioning for classification, regression and survival trees. An implementation of most of the functionality of the 1984 book by Breiman, Friedman, Olshen and Stone." [13]

**C50** from the **C50** library. "C5.0 decision trees and rule-based models for pattern recognition." [4]

**Random Forest** from the **randomForest** library. "Classification and regression based on a forest of trees using random inputs." [1]

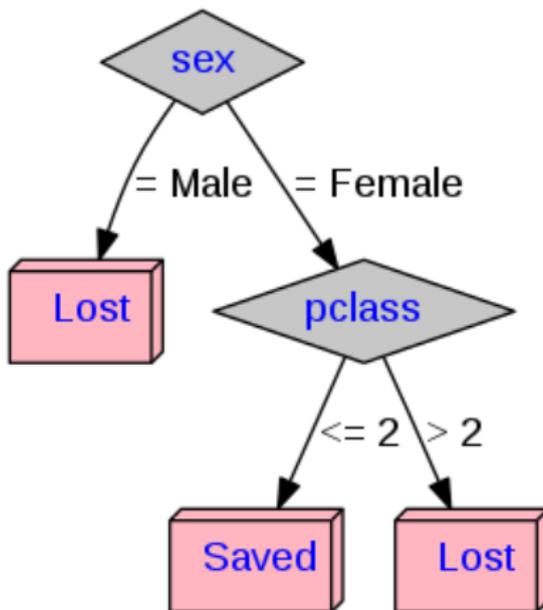
**J48** from the **RWeka** library. "An R interface to Weka (Version 3.9.1). Weka is a collection of machine learning algorithms for data mining tasks written in Java, containing tools for data pre-processing, classification, regression, clustering, association rules, and visualization." [3] The J48 algorithm is run in a pruned and unpruned mode.

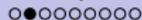
These and additional techniques to be covered in detail later.

○○○○  
○○○○○○○○○○  
○○  
○○○○○

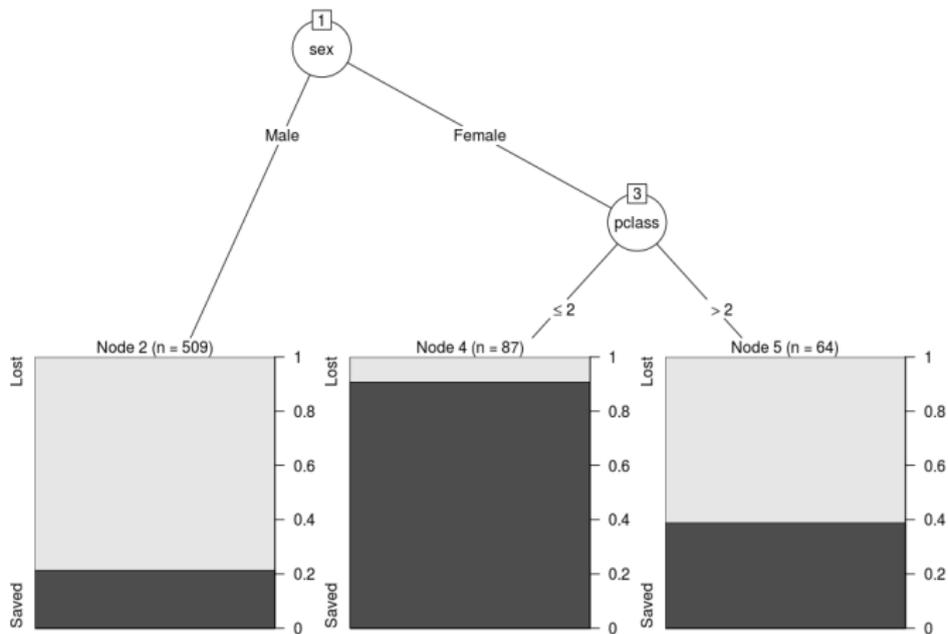
●○○○○○○○

# C50 decision tree





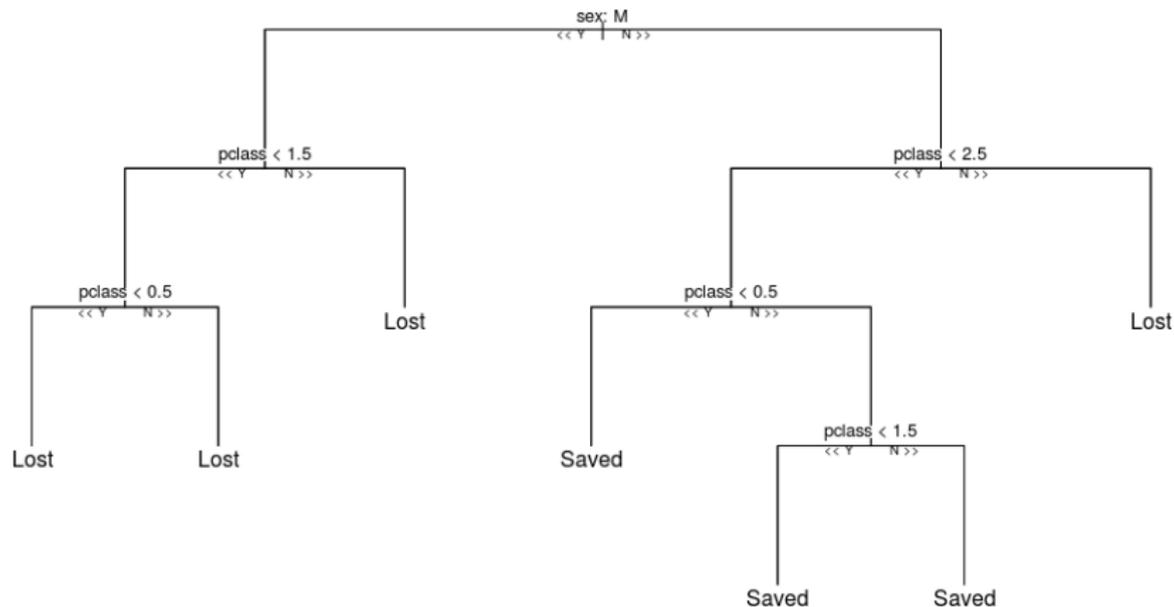
# Random Forest decision tree



○○○○  
○○○○○○○○○○  
○○  
○○○○○

○○●○○○○○

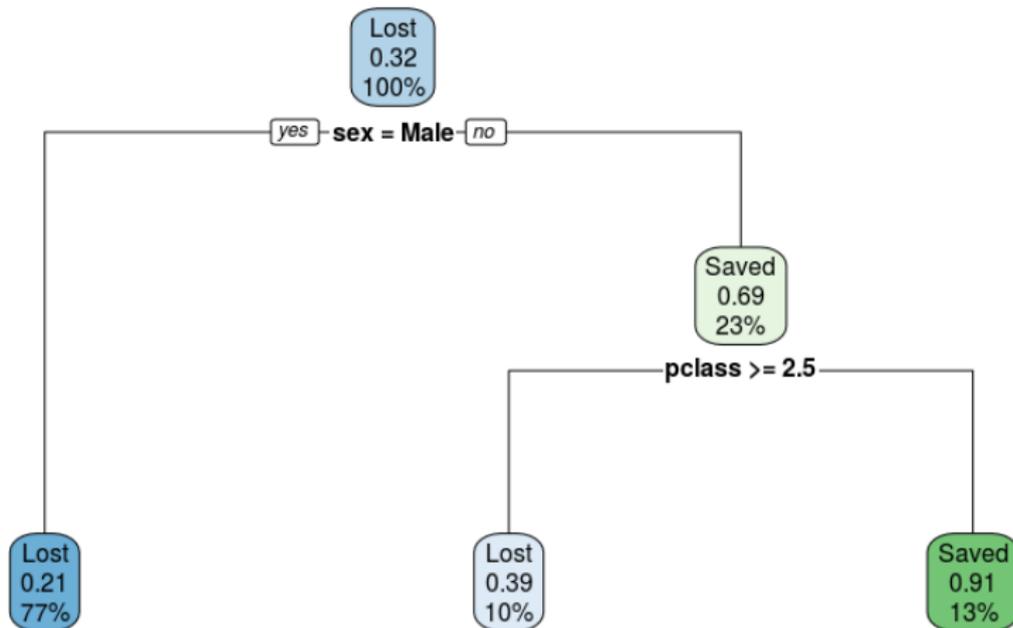
## rpart decision tree



○○○○  
○○○○○○○○○○  
○○  
○○○○○

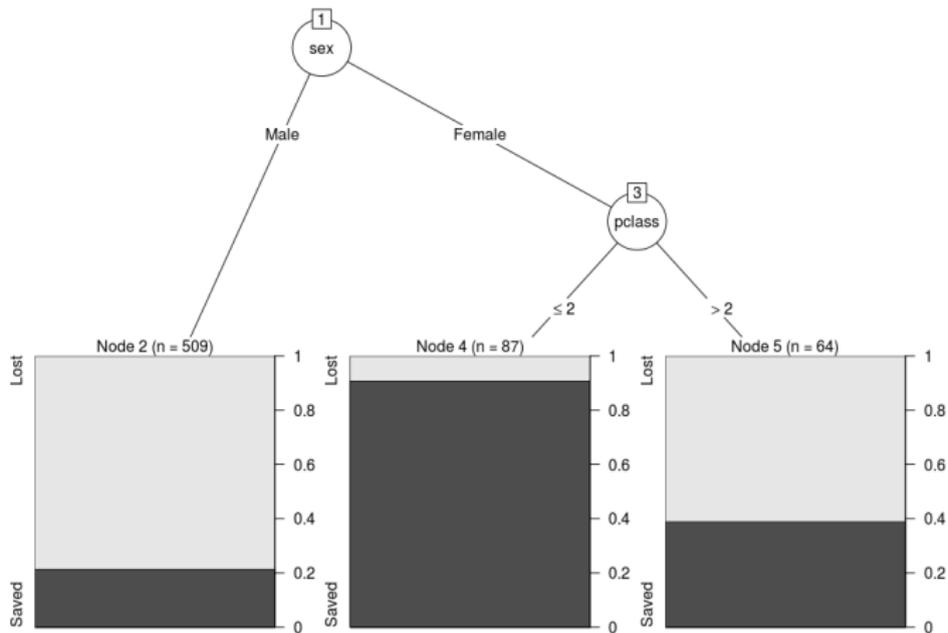
○○●○○○○○

# J48 (unpruned) decision tree



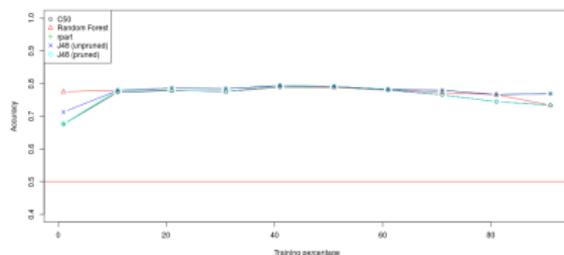


# J48 (pruned) decision tree



## Accuracy based on training percentage

The horizontal line at 50% represents the accuracy that would be achieved based on using an unbiased coin to decide the likelihood of survival.



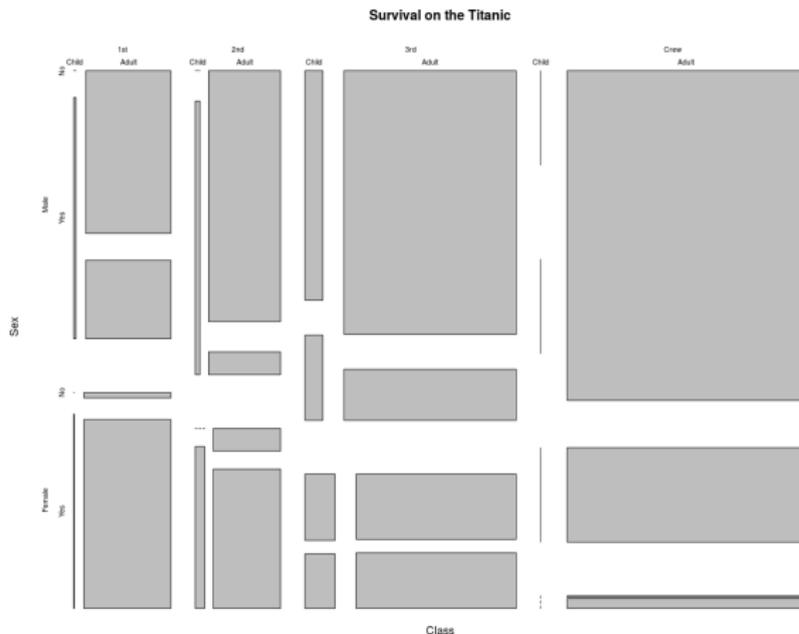
Using training percentages from about 10 to 60 result in all algorithms having nearly identical accuracies. Below 10%, the Random Forest approach appears best.





## Results

## Same image.



A summary of all personnel on the RMS Titanic broken down by gender, by survival or not, and class. It is interesting to look at the data and consider the adage: "women and children first."

# Some simple exercises to get familiar with data analysis

- 1 Using the Titanic report as a guide, create a recursive partition decision tree modeling survival based on sex and number of siblings
- 2 Create a recursive partition decision tree modeling survival based on all available data

○○○○  
○○○○○○○○○○  
○○  
○○○○○

○○○○○○○○○

## Q & A time.

Q: How many Harvard MBA's does it take to screw in a light bulb?

A: Just one. He grasps it firmly and the universe revolves around him.



## What have we covered?

- All of the decision tree algorithms tested had comparable results (~76% accuracy) when the training dataset was between 10 and 60% of the entire dataset.
- Random forest performed most consistently over the widest range of training percentages of all tested algorithms.



Next: LPAR Chapter 2, basic data visualization

○○○○  
○○○○○○○○○○  
○○  
○○○○○

○○○○○○○○○

## References (1 of 5)

- [1] Leo Breiman, randomForest: Breiman and Cutlers random forests for classification  
<http://stat-www.berkeley.edu/users/breiman/RandomForests>, 2006.
- [2] Jr. Frank E. Harrell, Titanic Data,  
<http://biostat.mc.vanderbilt.edu/wiki/pub/Main/DataSets/titanic.html>, 2002.
- [3] K Hornik, A Zeileis, T Hothorn, and C Buchta, RWeka: an R interface to Weka, R package version 0.4-32 (2017).

## References (2 of 5)

- [4] M Kuhn, S Weston, N Coulter, M Culp, and R Quinlan, C5.0 Decision Trees and Rule-Based Models, R Package Version 0.1. 0 **24** (2015).
- [5] Sebastian Raschka, Predictive modeling, supervised machine learning, and pattern classification, [http://sebastianraschka.com/Articles/2014\\_intro\\_supervised\\_learning.html](http://sebastianraschka.com/Articles/2014_intro_supervised_learning.html), 2014.
- [6] Encyclopedia Titanica Staff, Encyclopedia Titanica, Titanic Facts, History and Biography, <https://www.encyclopedia-titanica.org/>, 2017.

○○○  
○○○○○○○○  
○○  
○○○○

○○○○○○○○

## References (3 of 5)

- [7] ISM Staff, [Titanic Survivor, Titanic Passenger List Booklet](http://www.phillyseaport.org/web_exhibits/mini_exhibits/titanic_passenger_list/titanic_passenger_list-object-passenger_list.html), [http://www.phillyseaport.org/web\\_exhibits/mini\\_exhibits/titanic\\_passenger\\_list/titanic\\_passenger\\_list-object-passenger\\_list.html](http://www.phillyseaport.org/web_exhibits/mini_exhibits/titanic_passenger_list/titanic_passenger_list-object-passenger_list.html), 2017.
- [8] OceanGate Staff, [Titanic Survey Expedition: 2018](http://www.oceangate.com/expeditions/titanic-survey-2018.html), <http://www.oceangate.com/expeditions/titanic-survey-2018.html>, 2017.
- [9] Southampton Staff, [Titanic crew list](http://www.plimsoll.org/Southampton/Titanic/titaniccrewlist/Default.asp), <http://www.plimsoll.org/Southampton/Titanic/titaniccrewlist/Default.asp>, 2017.

## References (4 of 5)

- [10] Titanic Facts Staff, [Titanic Passenger List](http://www.titanic-facts.com/titanic-passenger-list.html), <http://www.titanic-facts.com/titanic-passenger-list.html>, 2017.
- [11] Wikipedia Staff, [Sinking of the RMS Titanic](https://en.wikipedia.org/wiki/Sinking_of_the_RMS_Titanic), [https://en.wikipedia.org/wiki/Sinking\\_of\\_the\\_RMS\\_Titanic](https://en.wikipedia.org/wiki/Sinking_of_the_RMS_Titanic), 2017.
- [12] Pang-Ning Tan, Michael Steinbach, and Vipin Kumar, [Introduction to Data Mining](#), Pearson Education India, 2006.
- [13] Terry Therneau, Beth Atkinson, and Brian Ripley, [rpart](#), Available at CRAN. [R-project.org/package= rpart](http://R-project.org/package=rpart). Accessed May (2015).

○○○  
○○○○○○○○  
○○  
○○○○

○○○○○○○○

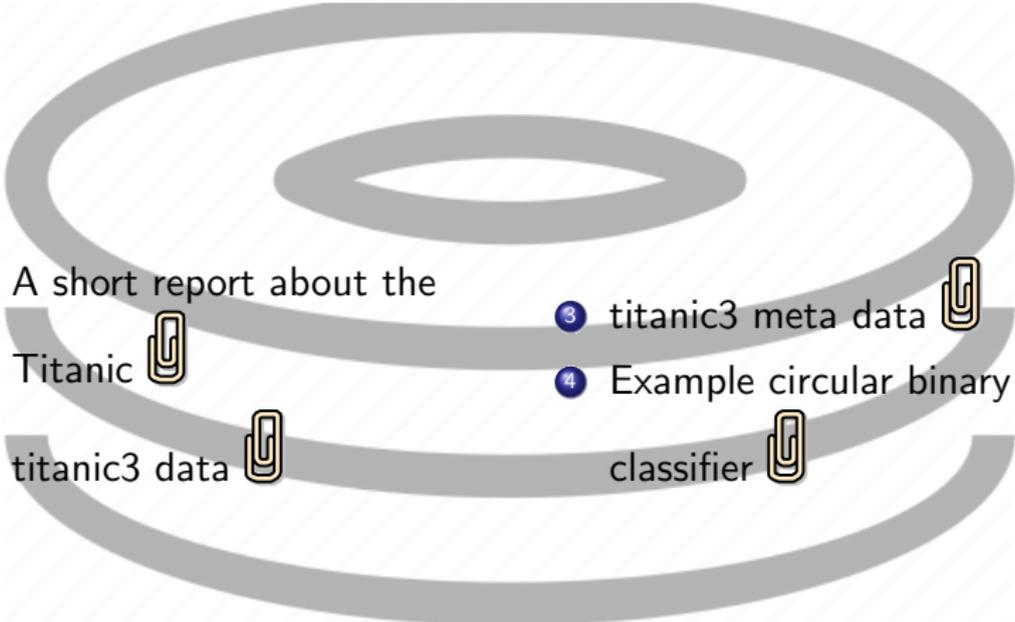
## References (5 of 5)

- [14] [Maria Dolores Ugarte, Ana F Militino, and Alan T Arnholt, Probability and Statistics with R, CRC Press, 2008.](#)

○○○○  
○○○○○○○○○○  
○○  
○○○○○

○○○○○○○○○

## Files of interest

- 
- 1 A short report about the Titanic 
  - 2 titanic3 data 
  - 3 titanic3 meta data 
  - 4 Example circular binary classifier 