

DublinR - Machine Learning 101

Introduction with Examples

Eoin Brazil - <https://github.com/braz/DublinR-ML-treesandforests>

Machine Learning Techniques in R

A bit of context around ML

How can you interpret their results?

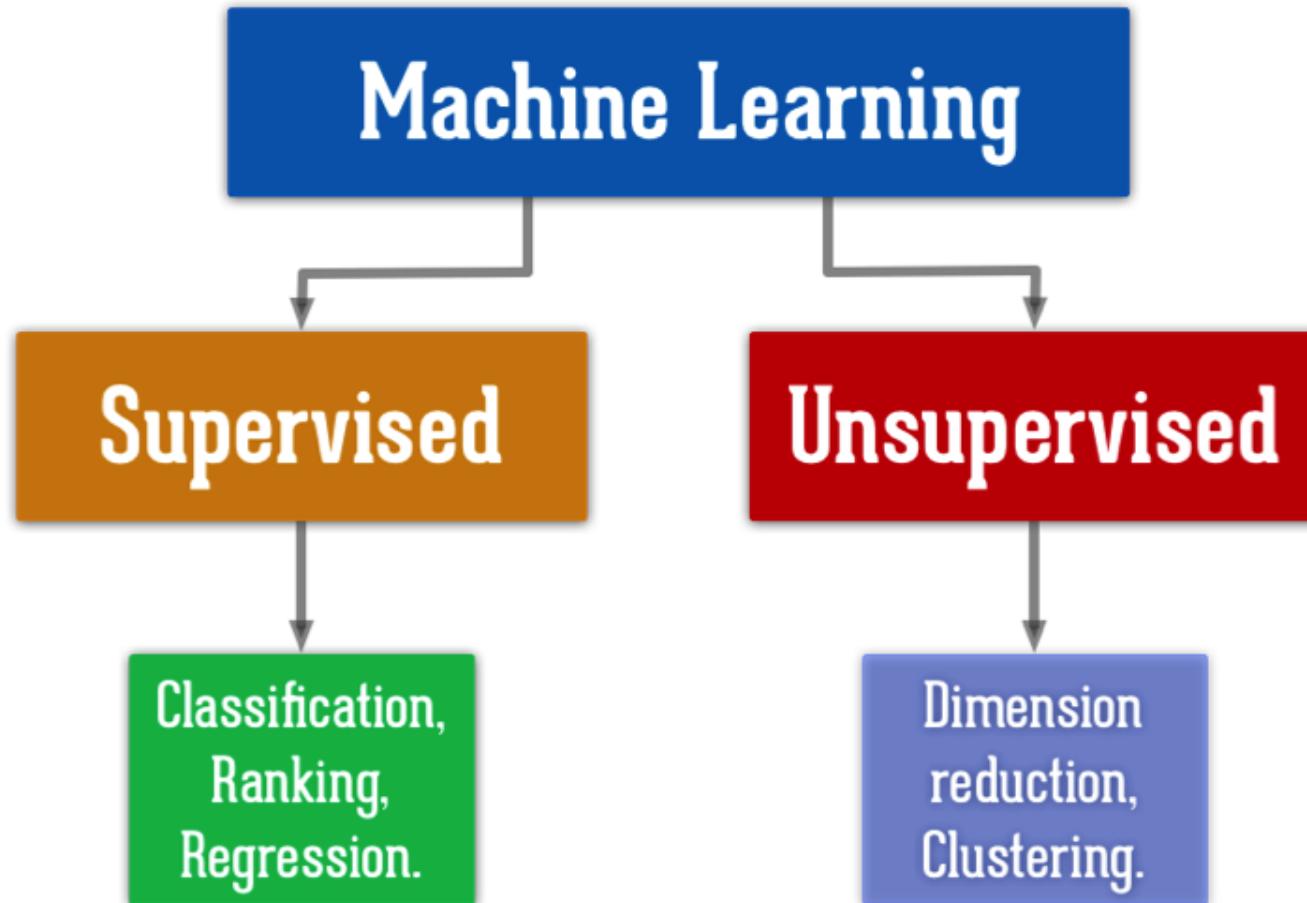
A few techniques to improve prediction / reduce over-fitting

Kaggle & similar competitions - using ML for fun & profit

Nuts & Bolts - 4 data sets and 6 techniques

A brief tour of some useful data handling / formatting tools

Types of Learning



A bit of context around ML

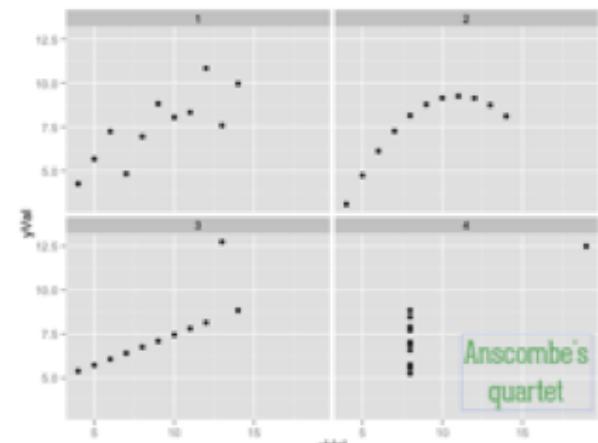
Clean



Simple

var1
var2
var3 } Pareto
Principle

Visual



```
> aggregate(.~group,data=mydata,mean)
  group  xVal  yVal
1     1    9 7.500000
2     2    9 7.500000
3     3    9 7.500000
4     4    9 7.500000
> aggregate(.~group,data=mydata,sd)
  group  xVal  yVal
1     1 3.316625 2.031568
2     2 3.316625 2.031857
3     3 3.316625 2.030424
4     4 3.316625 2.030579
```

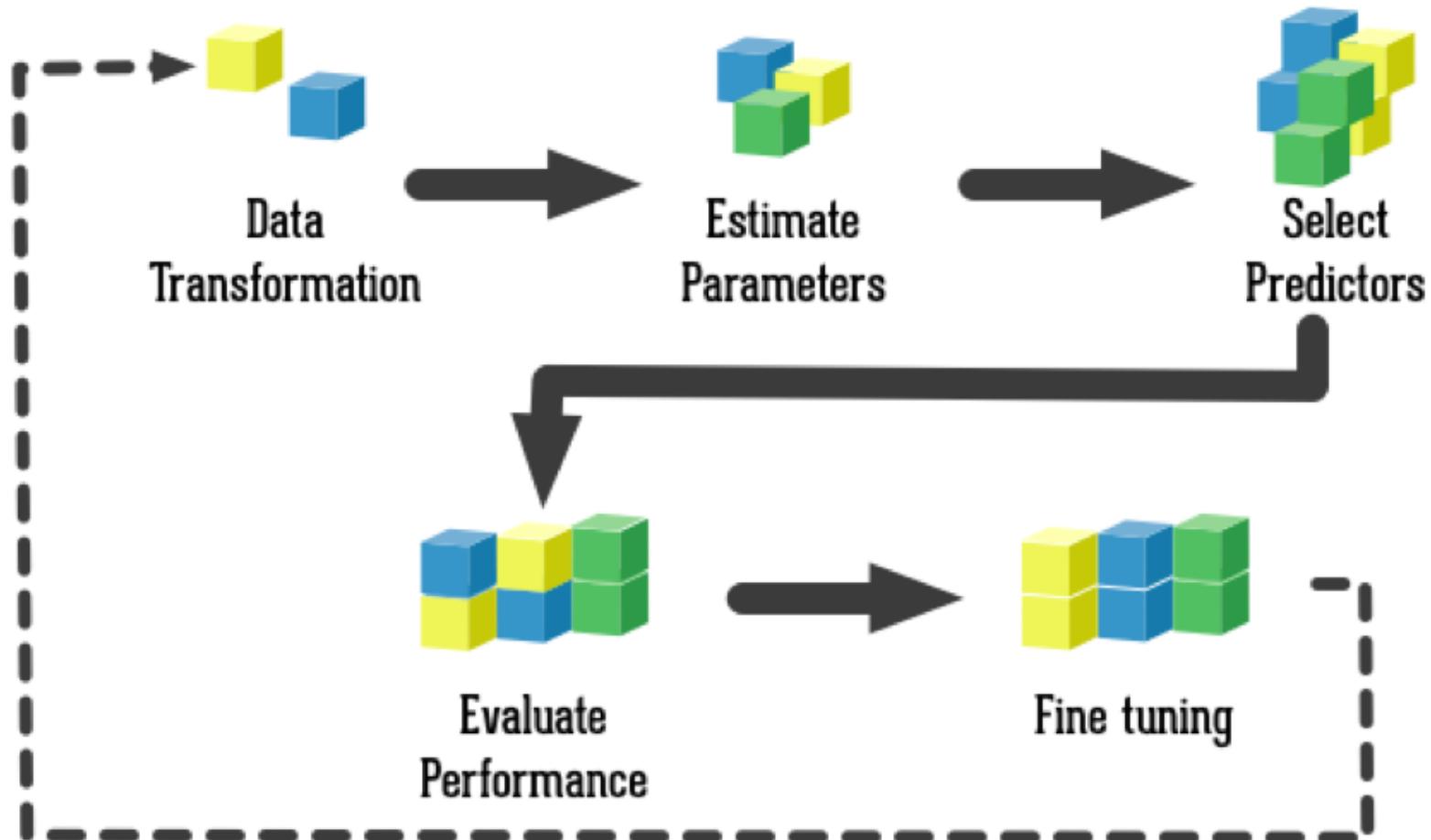


Data Transformations

$$\left. \begin{array}{l} C \text{ Zero Mean} \\ S \text{ One SD} \\ T \sqrt{\log \text{INV}} \end{array} \right\} \text{1 - Predictors - M} \left\{ \begin{array}{l} \text{SS M-Dim Sphere} \\ \text{PCA} \\ \text{PLS} \end{array} \right.$$

Correlation, Dummy Variables, Filtering

Model Building Process



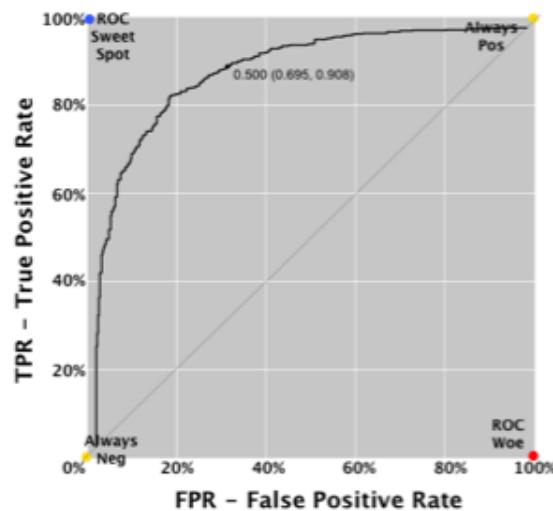
Model Selection and Model Assessment

Assessment [1 & 2]

1

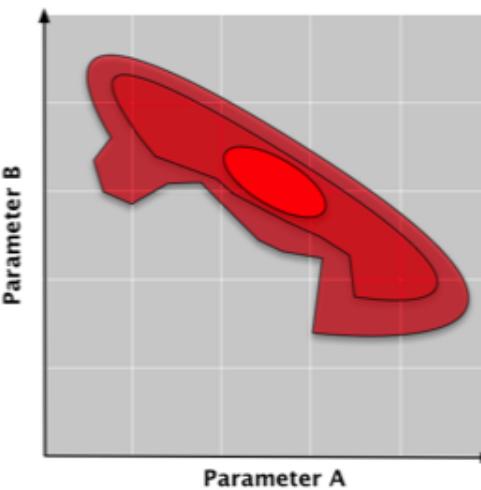
Confusion Matrix		REFERENCE (ACTUAL)
		HEALTHY (normal)
PREDICTED	HEALTHY (normal)	Actual Healthy Patient and Predicted as Healthy Patient
	DISEASED (patient)	Actual Diseased Patient but Predicted as Healthy Patient
		Actual Diseased Patient and Predicted as Diseased Patient

2

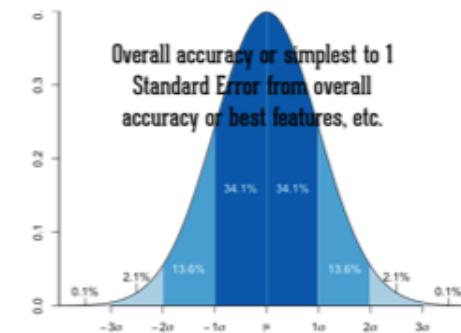


Selection [3 & 4]

3

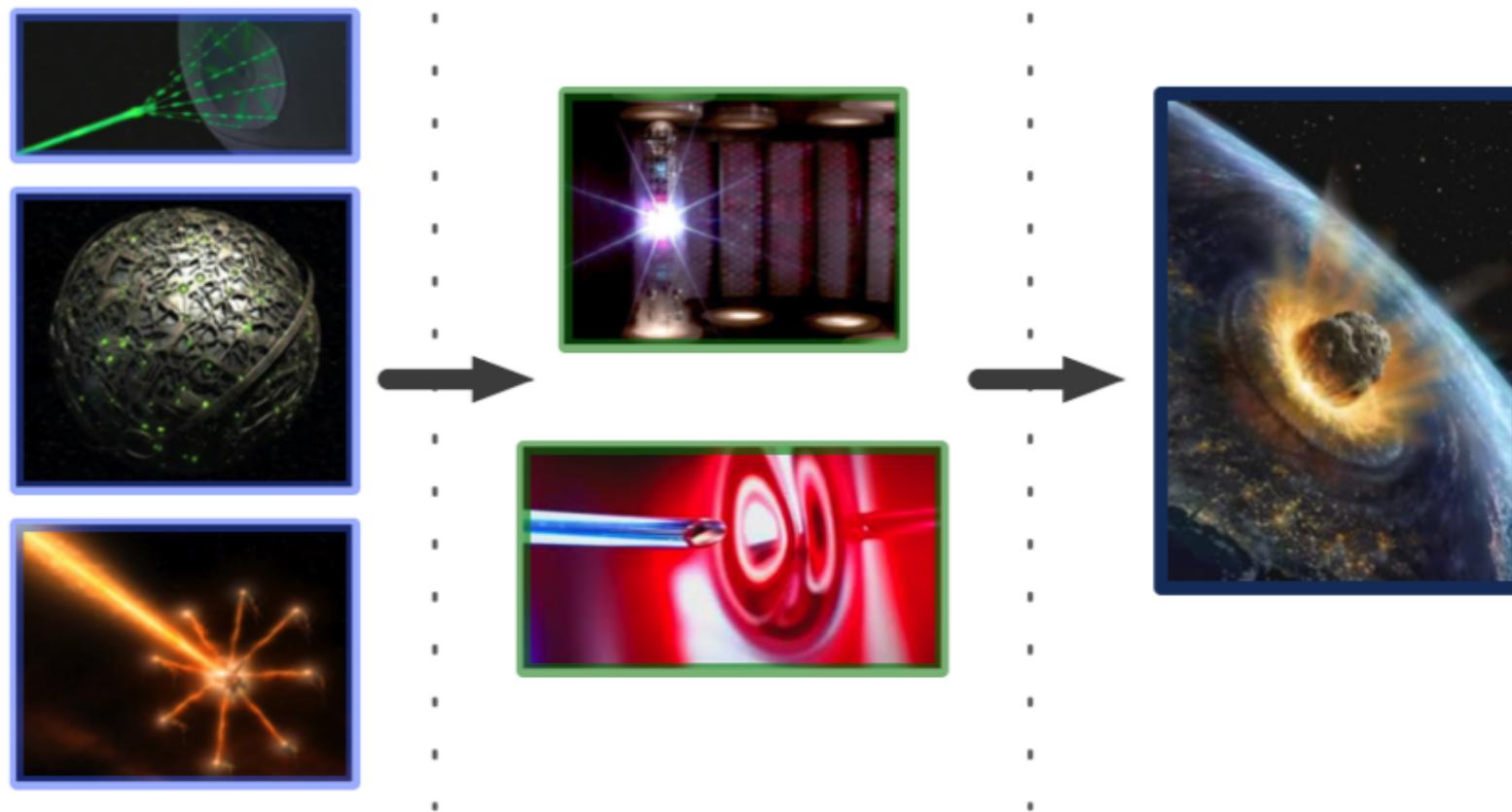


4



Model Choice - Move from Adaptability to Simplicity

Steps for between model choice



Interpreting A Confusion Matrix

Confusion Matrix		REFERENCE (ACTUAL)	
PREDICTED	Predicted Positive		
	Predicted Positive	Predicted Negative	
Positive Examples	TP true positive – a hit	FN false negative – positive but classified as negative, Type II error	✓ ✗
Negative Examples	FP false positive – negative but classified as positive, Type I error	TN true negative – a correct rejection	✗ ✓

Interpreting A Confusion Matrix Example

Confusion Matrix		REFERENCE (ACTUAL)	
		HEALTHY (normal)	DISEASED (patient)
PREDICTED	HEALTHY (normal)	TP Actual Healthy Patient and Predicted as Healthy Patient	FN Actual Diseased Patient but Predicted as Healthy Patient
	DISEASED (patient)	FP Actual Healthy Patient but Predicted as Diseased Patient	TN Actual Diseased Patient and Predicted as Diseased Patient

Confusion Matrix - Calculations

REFERENCE (ACTUAL)		
PREDICTED	TP	FN
FP	TP	TN
	FN	TN

REFERENCE (ACTUAL)		
PREDICTED	TP	FN
FP	TP	TN
	FN	TN

ACCURACY

REFERENCE (ACTUAL)		
PREDICTED	TP	FN
FP	TP	TN
	FN	TN

Neg Pred Val: NPV

REFERENCE (ACTUAL)		
PREDICTED	TP	FN
FP	TP	TN
	FN	TN

Specificity: SPC

REFERENCE (ACTUAL)		
PREDICTED	TP	FN
FP	TP	TN
	FN	TN

False Discover R: FDR

REFERENCE (ACTUAL)		
PREDICTED	TP	FN
FP	TP	TN
	FN	TN

False Pos R: FPR

REFERENCE (ACTUAL)		
PREDICTED	TP	FN
FP	TP	TN
	FN	TN

True Pos R: TPR

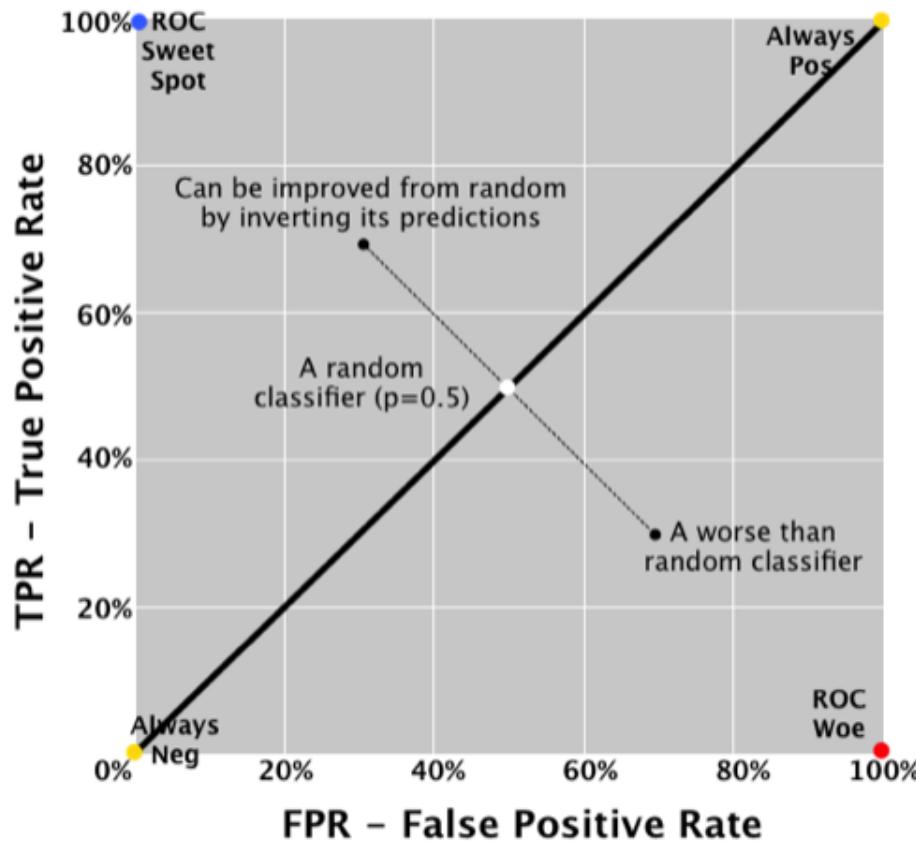
REFERENCE (ACTUAL)		
PREDICTED	TP	FN
FP	TP	TN
	FN	TN

Pos Pred Val: PPV

ROC Curve

Precision Recall Curve

Interpreting A ROC Plot



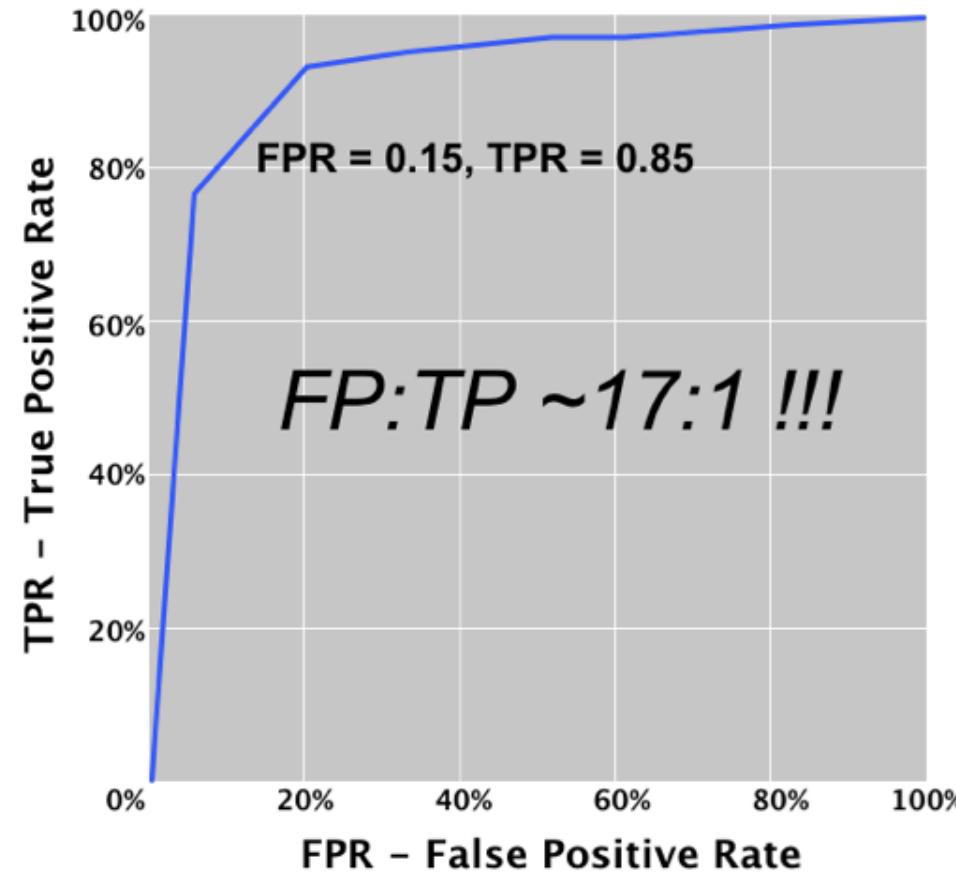
- A point in this plot is better than another if it is to the northwest (TPR higher / FPR lower / or both)
- ``Conservatives'' - on LHS and near the X-axis - only make positive classification with strong evidence and making few FP errors but low TP rates
- ``Liberals'' - on upper RHS - make positive classifications with weak evidence so nearly all positives identified however high FP rates

ROC Dangers

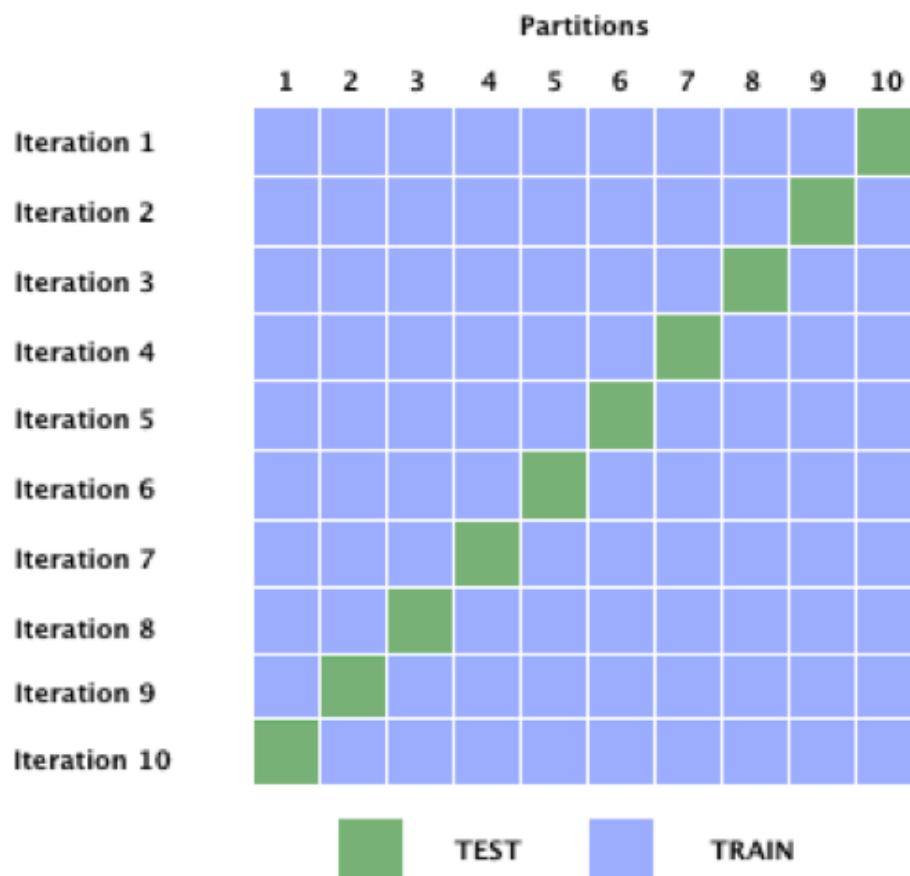
REFERENCE (ACTUAL)		
PREDICTED	TP	FN
FP	TP	TN

REFERENCE (ACTUAL)		
PREDICTED	85	1485
15	8415	

REFERENCE (ACTUAL)		
PREDICTED	0.85	14.85
0.15	%	84.15

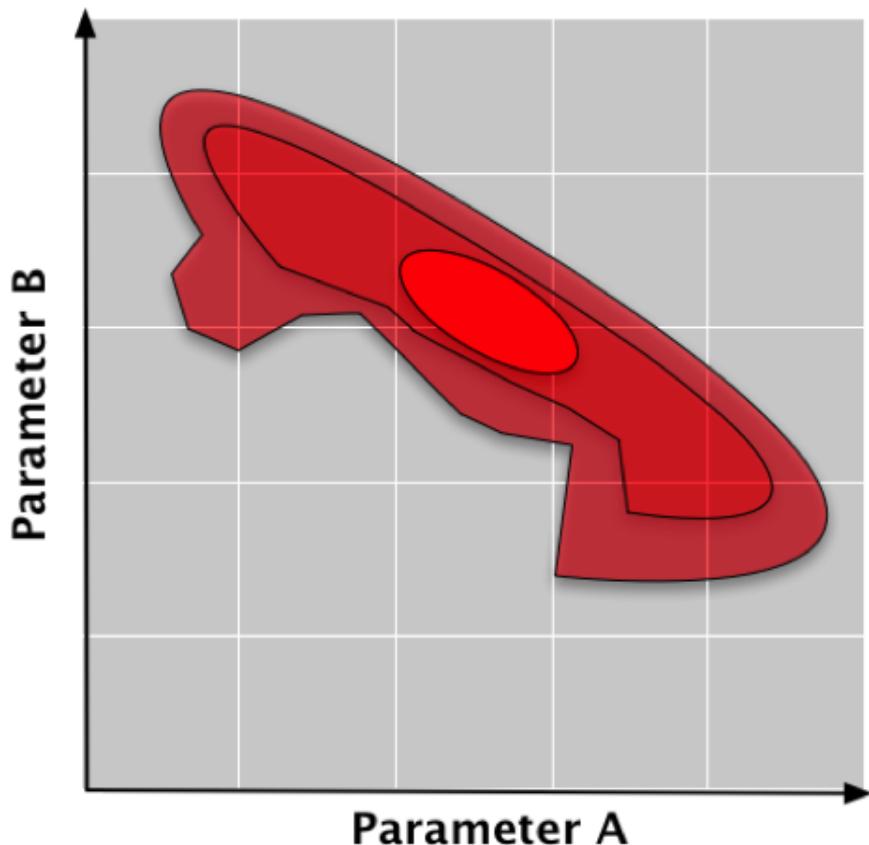


Addressing Prediction Error



- K-fold Cross-Validation (e.g. 10-fold)
 - Allows for averaging the error across the models
- Bootstrapping, draw B random samples with replacement from data set to create B bootstrapped data sets with same size as original. These are used as training sets with the original used as the test set.
- Other variations on above:
 - Repeated cross validation
 - The '.632' bootstrap

Addressing Feature Selection



CARET and Feature Selection

- Recursive Feature Elimination Algorithm

```
> cvCtrl <- trainControl(method = "repeatedcv", repeats = 3,  
+ summaryFunction = twoClassSummary, classProbs = TRUE)  
> rpartTune <- train(hadaffair ~ ., data = affairs.df.train,  
+ method = "rpart", tuneLength = 10, metric = "ROC", trControl = cvCtrl)  
> rpartTune  
481 samples  
9 predictors  
2 classes: 'No', 'Yes'
```

No pre-processing
Resampling: Cross-Validation (10 fold, repeated 3 times)

Summary of sample sizes: 432, 433, 433, 433, 433, 433, ...

Resampling results across tuning parameters:

cp	ROC	Sens	Spec	ROC SD	Sens SD	Spec SD
0	0.654	0.911	0.225	0.0793	0.0597	0.11
0.00417	0.652	0.911	0.222	0.0775	0.0597	0.106
0.00833	0.594	0.932	0.161	0.121	0.0563	0.105
0.0125	0.572	0.949	0.156	0.133	0.045	0.0972
0.0167	0.562	0.958	0.119	0.133	0.0383	0.0839
0.0208	0.548	0.955	0.119	0.134	0.044	0.092
0.025	0.554	0.963	0.103	0.124	0.0474	0.092
0.0292	0.559	0.976	0.0806	0.118	0.0417	0.0861
0.0333	0.542	0.988	0.0528	0.101	0.037	0.0804
0.0375	0.516	1	0	0.0692	0	0

ROC was used to select the optimal model using the largest value.
The final value used for the model was cp = 0.

Kaggle - using ML for fun & profit

The image is a collage of screenshots from Kaggle and Synapse websites, illustrating various machine learning competitions:

- Synapse Beta:** A screenshot of the Synapse Beta interface showing the "HPN-DREAM breast cancer network inference challenge". It includes a sidebar with options like Dashboard, Home, Data, Information, Description, Evaluation, Rules, Dos and Don'ts, FAQ, Milestone Winners, and Timeline.
- kaggle:** A screenshot of the Kaggle homepage featuring the "Titanic: Machine Learning from Disaster" competition. It shows a score of 6.58, a belkin logo, and a progress bar indicating 12 months to go.
- Heritage Provider Network Health Prize:** A screenshot of the competition page for "Improve Healthcare, Win \$3,000,000." It features a large blue banner with the text "Examples of Competitions" and "Improve Healthcare, Win \$3,000,000."
- Marinexplore Whale Detection Challenge:** A screenshot of the Marinexplore Whale Detection Challenge page. It features an illustration of a right whale and text about the competition's goal: "recognition solutions to detect and classify right whale sightings using BIG data mining and exploration studies".
- Modeling Hospitalization Outcomes with Random Decision Trees and Bayesian Feature Selection:** A screenshot of a research paper titled "De-identification Methods for Open Health Data: The Case of the Heritage Health Prize Claims Dataset". It includes author information for Thomson Van Nguyen and Bhushanewar Mishra.

Nuts & Bolts - Data sets and Techniques



- 1-Marketing-Survey-Dataset.R**
Associative Rules
- 2-Car-Dataset.R**
Decision Trees [Simple]
- 3-Image-Segmentation-Dataset.R**
Decision Trees
- 4-Wine-Portugal-Dataset.R**
Random Forests
k Nearest Neighbors
Neural Networks
Support Vector Machines
- 5-Affairs-Dataset.R**
Random Forests
Naive Bayesian

<https://github.com/braz/DublinR-ML-treesandforests>

Model comparison

Data Partitioning
Tuning [e.g. cross-fold validation]
Confusion Matrices & ROCs



Associate Rule Learning & Apriori Algorithm

- ▶ Relationships
- ▶ Item-sets / Rules
- ▶ Temporal
- ▶ Cross-Selling



Aside - How does associative analysis work?

1

Customer	Items
1	Orange juice, Apples
2	Milk, orange juice, bread
3	Washing powder, milk
4	Orange juice, Milk
5	Milk, bread
6	Washing powder, orange juice
⋮	⋮

ITEMSET1 IF a CUST buys MILK then the CUST also buys OJ

ITEMSET2 IF a CUST buys APPLES then the CUST also buys OJ

2

	OJ	Apples	Milk	Bread	Washing Powder
OJ	4	1	2	1	1
Apples	1	1	0	0	0
Milk	2	0	4	1	0
Bread	1	0	2	2	0
Washing Powder	1	0	1	0	2

3 If customer buys bread, then the customer also buys milk

Support: Num of trans containing all items in the rule

Confidence: Measure of predicting the RHS (after the THEN) of a rule

Lift: Measures the power of the rule to the full by randomly guessing RHS

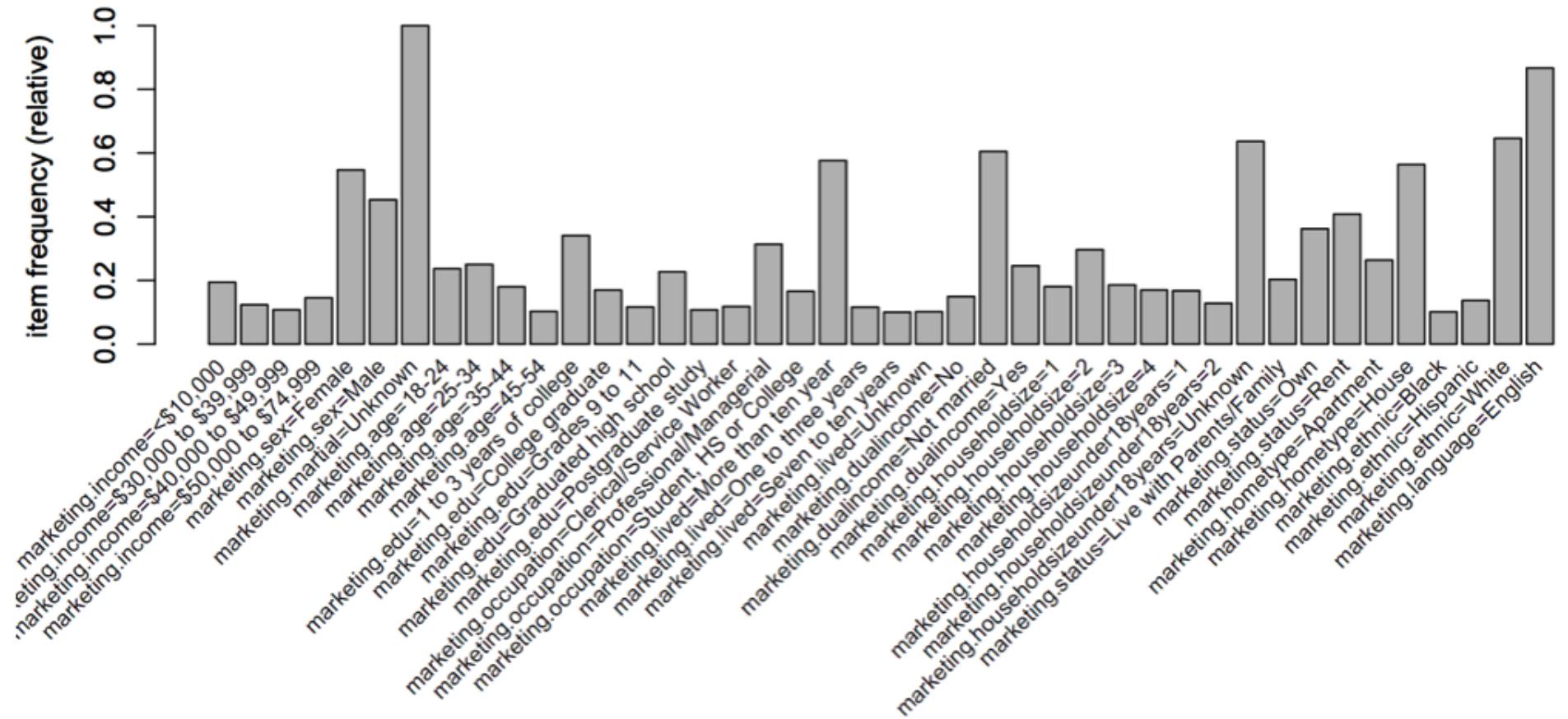
Chi-Square: Measures the probability that the contingency table behind the rule could be produced randomly

4 SKU, EAN, Hierarchies

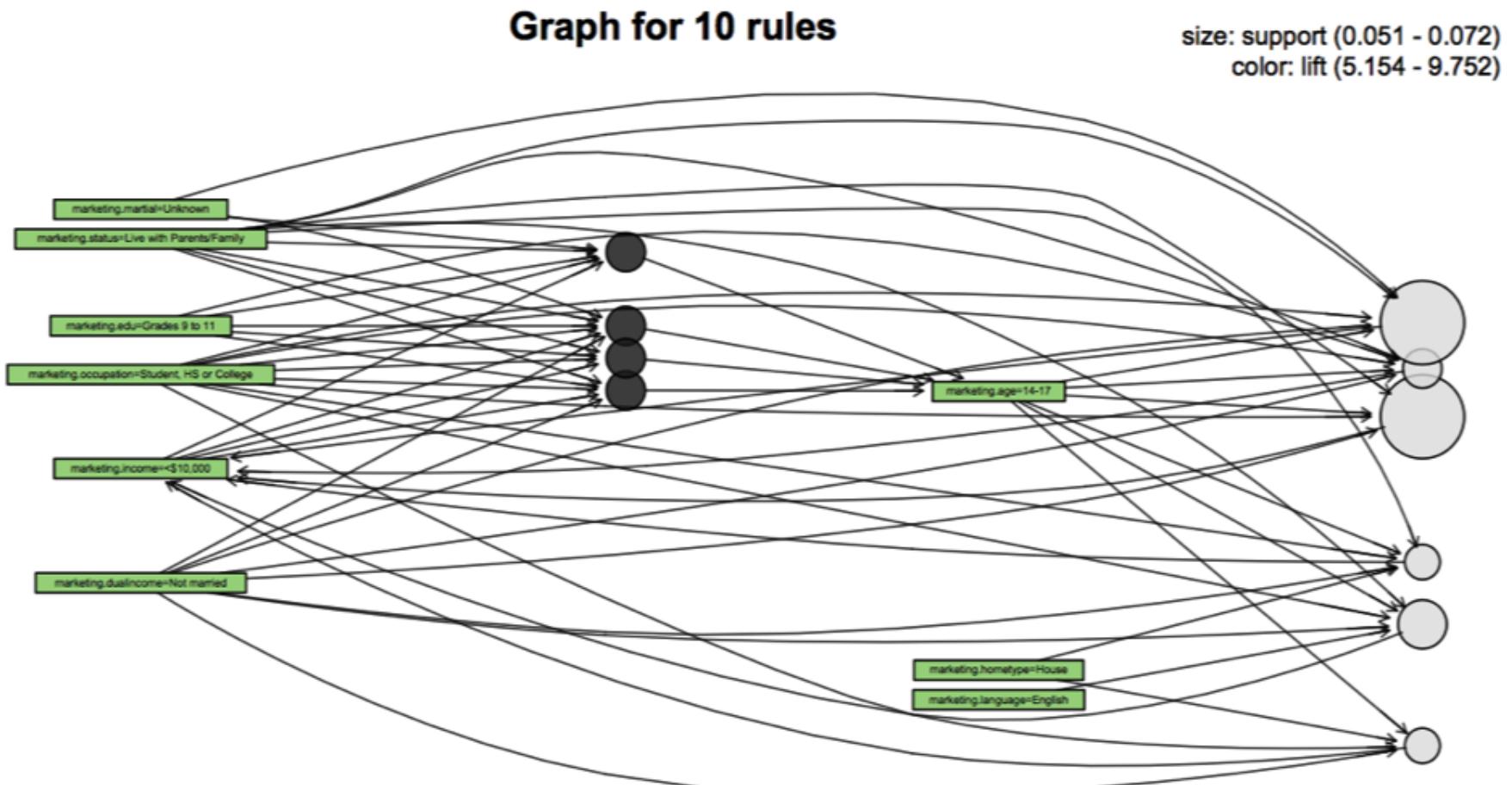


What are they good for?

Marketing Survey Data - Part 1



Marketing Survey Data - Part 2

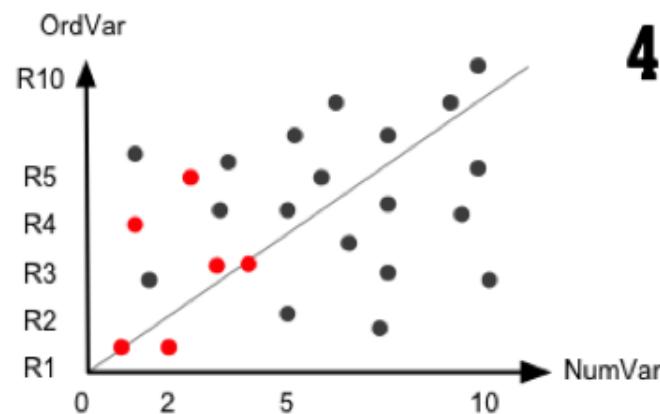


Aside - How do decision trees work?

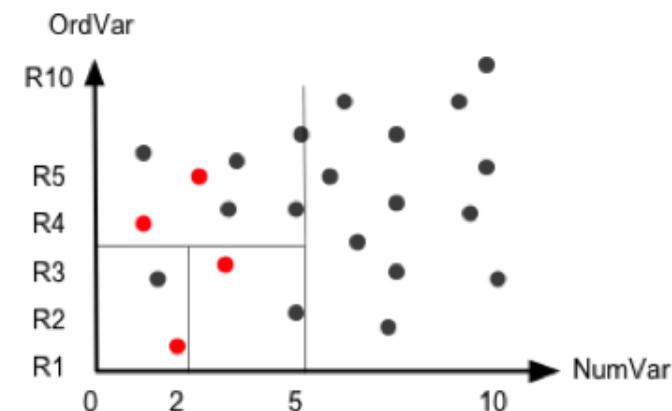
1

NumVar	OrdVar	CatVar
1	R2	Good
1.5	R4	Bad
2	R1	Bad
2.5	R5	Good
3	R3	Good
3.5	R3	Bad
⋮	⋮	⋮
⋮	⋮	⋮

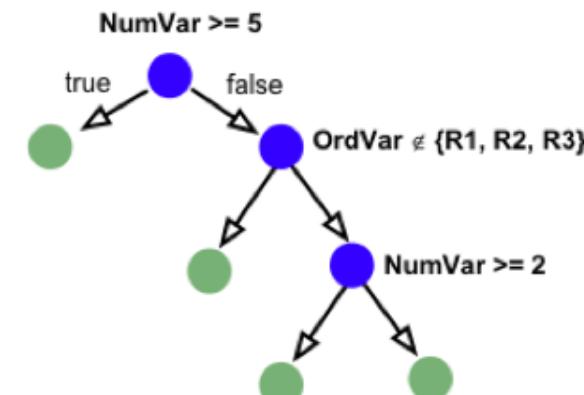
2



3

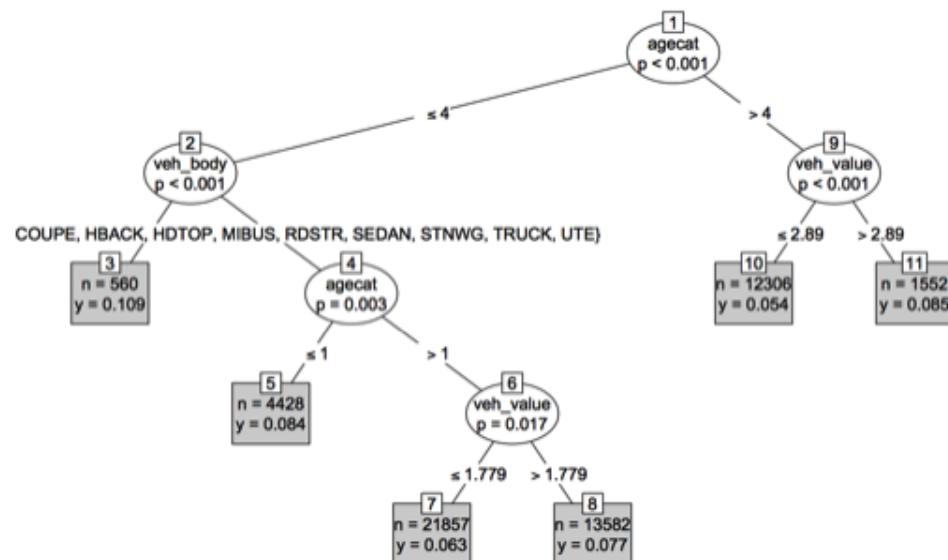


4



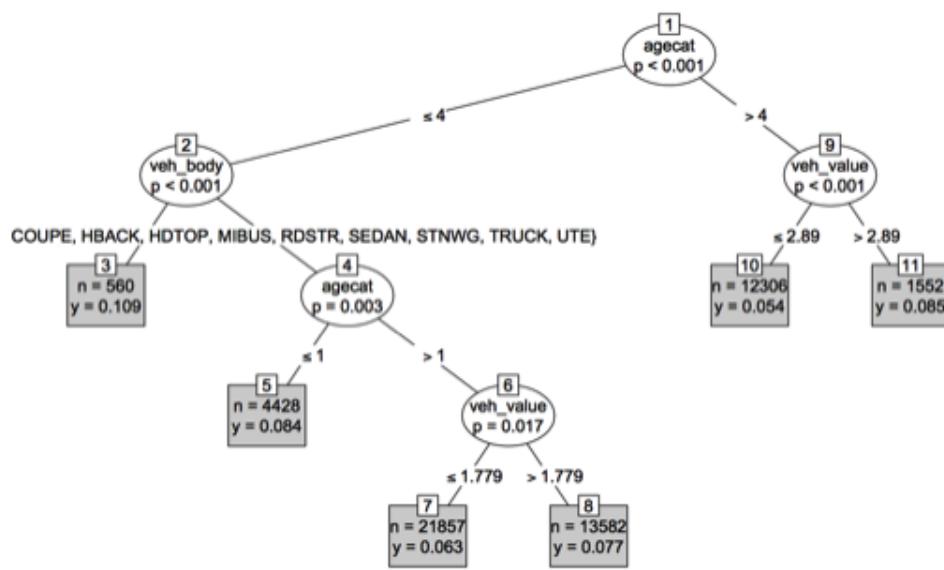
What are they good for?

Car Insurance Policy Exposure Management - Part 1



- Analysing insurance claim details of 67856 policies taken out in 2004 and 2005.
- The model maps each record into one of X mutually exclusive terminal nodes or groups.
- These groups are represented by their average response, where the node number is treated as the data group.
- The binary claim indicator uses 6 variables to determine a probability estimate for each terminal node determine if a insurance policyholder will claim on their policy.

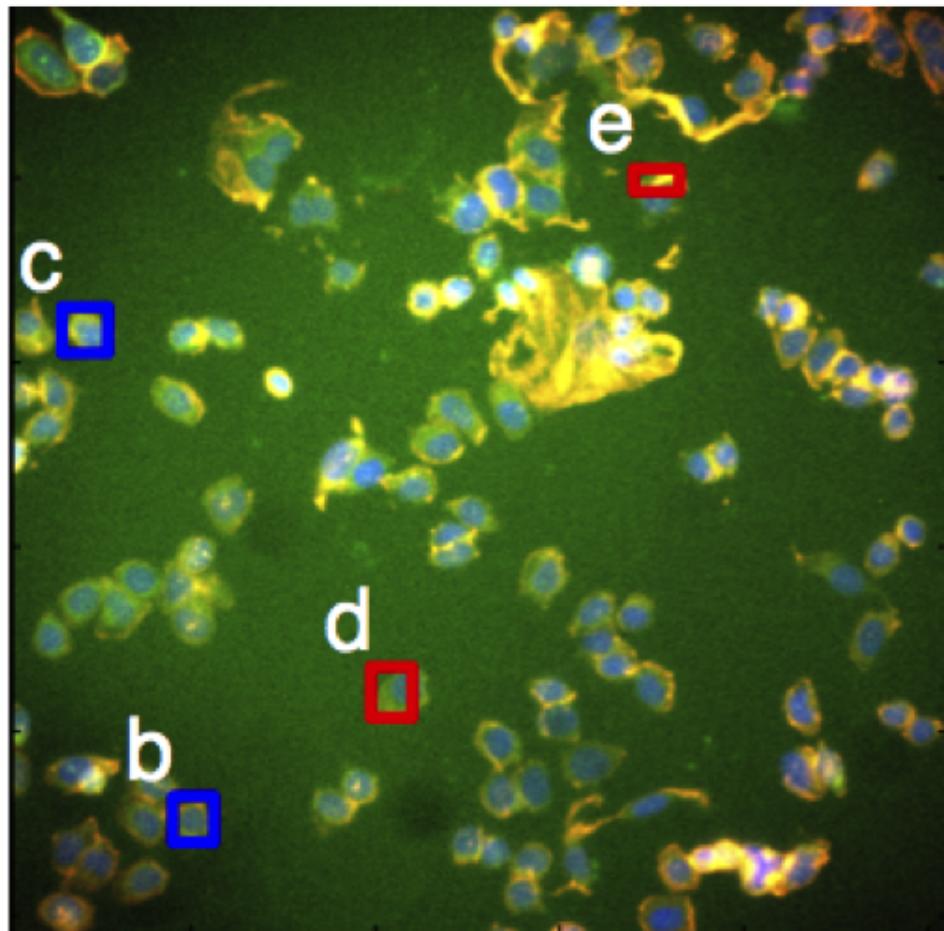
Car Insurance Policy Exposure Management - Part 2



- Root node, splits the data set on 'agecat'
- Younger drivers to the left (1-8) and older drivers (9-11) to right
- N9 splits on basis of vehicle value
- N10 \leq \$28.9k giving 15k records and 5.4% of claims
- N11 $>$ \$28.9k+ giving 1.9k records and 8.5% of claims
- Left Split from Root, N2 splits on vehicle body type, on age (N4), then on vehicle value (N6)
- The n value = num of overall population and the y value = probability of claim from a driver in that group

What are they good for?

Cancer Research Screening - Part 1



- Hill et al (2007), models how well cells within an image are segmented, 61 vars with 2019 obs (Training = 1009 & Test = 1010).
 - "Impact of image segmentation on high-content screening data quality for SK-BR-3 cells, Andrew A Hill, Peter LaPan, Yizheng Li and Steve Haney, BMC Bioinformatics 2007, 8:340".
 - b, Well-Segmented (WS)
 - c, WS (e.g. complete nucleus and cytoplasmic region)
 - d, Poorly-Segmented (PS)
 - e, PS (e.g. partial match/es)

Cancer Research Screening Dataset - Part 2

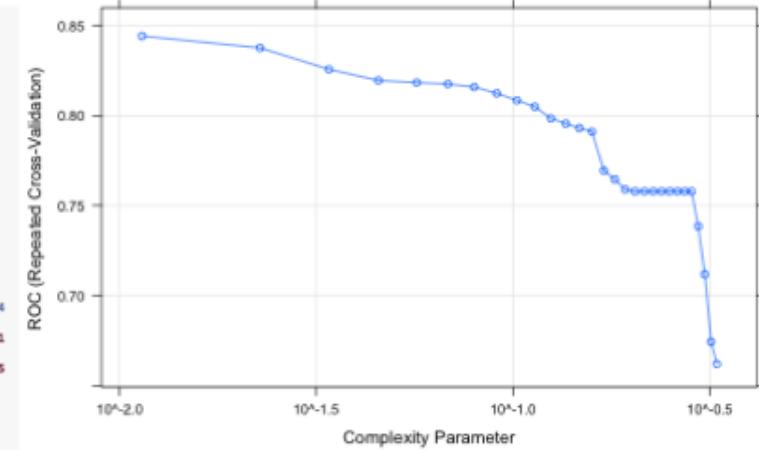
Image Segmentation

```
> table(testing[,1])
  PS   WS 
  664 346
```

```
'data.frame': 1009 obs. of 6 variables:
$ Class      : Factor w/ 2 levels "PS","WS": 1 2 1 2 1 1 1 2 2 2 ...
$ AngleCh1   : num 133.8 106.6 69.2 109.4 104.3 ...
$ AreaCh1    : int 819 431 298 256 258 358 158 315 246 223 ...
$ AvgIntenCh1: num 31.9 28 19.5 18.8 17.6 ...
$ AvgIntenCh2: num 207 116 102 127 125 ...
$ AvgIntenCh3: num 69.9 63.9 28.2 13.6 22.5 ...
```

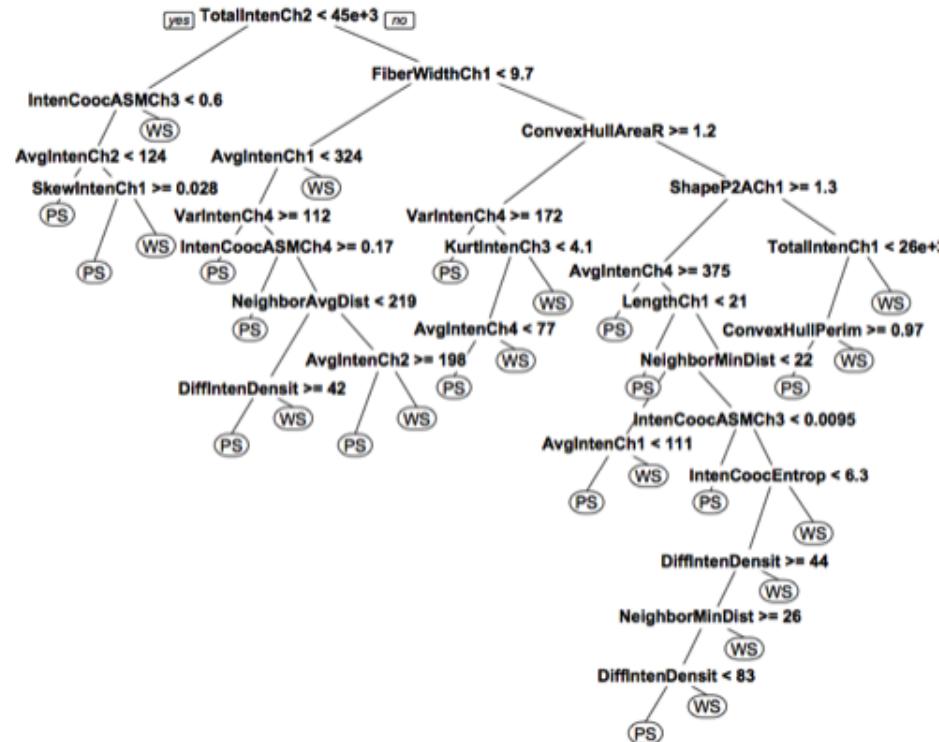
```
> head(training[1:2,])
  Class AngleCh1 AvgIntenCh1 AvgIntenCh2 AvgIntenCh3 AvgIntenCh4 ConvexHullAreaRatioCh1 ConvexHullPerimRatioCh1 DiffIntenDensityCh1
1  PS 133.7528    819   31.92327   286.8785   69.91688   164.1535     1.263158     0.7970881     31.87508
2  WS 133.7528    819   31.92327   286.8785   69.91688   164.1535     1.263158     0.7970881     31.87508
3  PS 133.7528    43.12228   79.30842   6.887592   43.12228   79.30842     0.8754758     0.9354758     32.48771
4  WS 106.6464    431   28.03883   116.3155   63.84175   106.6966     1.053318     0.9354758     32.48771
5  PS 106.6464    35.98577   51.35785   5.883557   35.98577   51.35785     0.8754758     0.9354758     32.48771
6  WS 106.6464    35.98577   51.35785   5.883557   35.98577   51.35785     0.8754758     0.9354758     32.48771
EntropyIntenCh3 EntropyIntenCh4 EqCircDiamCh1 EqEllipseLWRCh1 EqEllipseOblateVolCh1 EqEllipseProlateVolCh1 EqSphereAreaCh1 EqSphereVolCh1
1  6.642761    7.888155   32.38558   1.558394   2232.9855     1432.8246     3278.726     17653.525
2  1.487935    1.352374   64.28238   13.16788   1.487935    1.352374   64.28238   13.16788
3  6.683988    7.144681   23.44892   1.375386   682.1945     583.2584     1727.418     6750.985
FiberAlign2Ch3 FiberAlign2Ch4 FiberLengthCh1 FiberWidthCh1
1  1.398522    1.522316   21.14115   21.14115
2  1.398522    1.522316   21.14115   21.14115
3  1.398522    1.522316   21.14115   21.14115
IntenCoocASMC3 IntenCoocASMC4 IntenCoocContrastCh1 IntenCoocContrastCh4 IntenCoocEntropyCh3 IntenCoocEntropyCh4 IntenCoocMaxCh3
1  0.028051061   0.012594975   8.227953   6.984846   6.822138     7.098988     0.15321477     0.
2  0.028051061   -0.2487691   -0.3387839   -0.2652638   47.21855
3  0.0066862315   0.006141691   14.446874   16.708843   7.588108     7.671478     0.82835852     0.
IntenCoocMaxCh4 KurtIntenCh1 KurtIntenCh3 KurtIntenCh4 LengthCh1
1  0.028051061   0.012594975   8.227953   6.984846   6.822138     7.098988     0.15321477     0.
2  0.028051061   -0.2487691   -0.3387839   -0.2652638   47.21855
3  0.0066862315   0.006141691   14.446874   16.708843   7.588108     7.671478     0.82835852     0.
KurtIntenCh1 KurtIntenCh3 KurtIntenCh4 LengthCh1
1  0.028051061   0.012594975   8.227953   6.984846   6.822138     7.098988     0.15321477     0.
2  0.028051061   -0.2487691   -0.3387839   -0.2652638   47.21855
3  0.0066862315   0.006141691   14.446874   16.708843   7.588108     7.671478     0.82835852     0.
ROC (Repeated Cross-Validation)

```

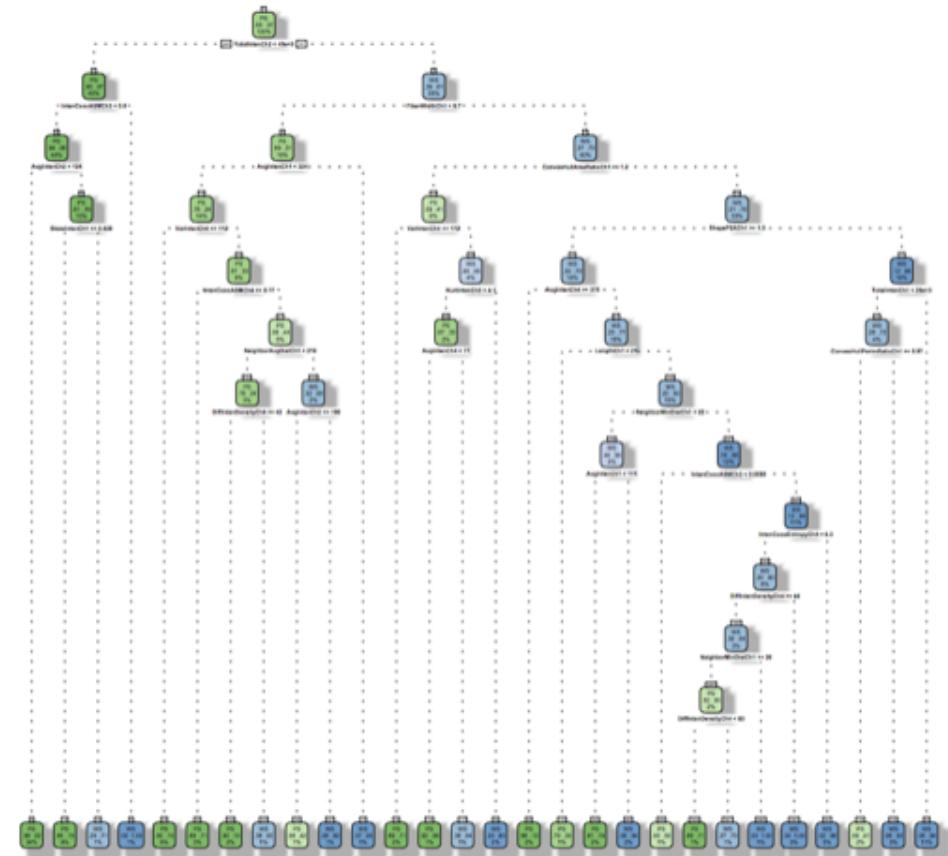


Cancer Research Screening - Part 3

"prp(rpartTune\$finalModel)"

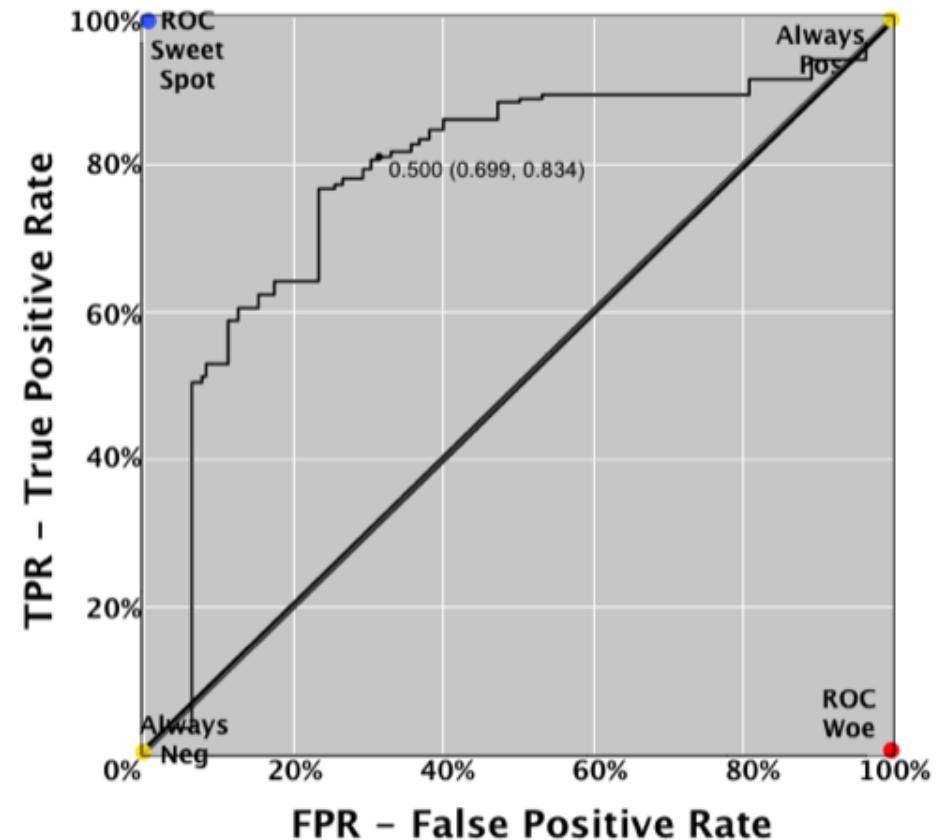


"fancyRpartPlot(rpartTune\$finalModel)"

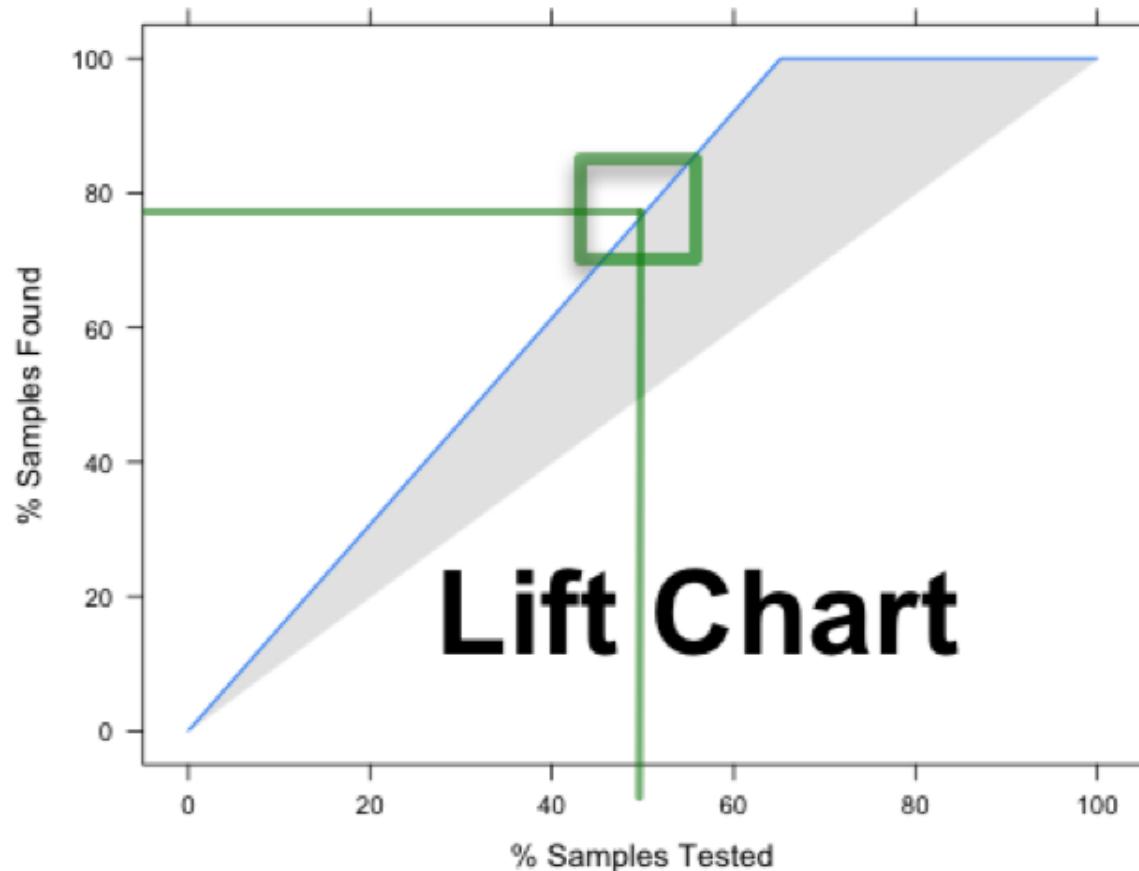


Cancer Research Screening - Part 4

Confusion Matrix		REFERENCE (ACTUAL)	
PREDICTED		Predicted Positive PS	Predicted Negative WS
Positive Examples PS		TP 554 / 55%	TN 104 / 10%
		$\text{Accuracy} = 79\% \\ (55+24)$	
Negative Examples WS		FP 110 / 11%	FN 242 / 24%
		PPos	PNeg
			N



Cancer Research Screening Dataset - Part 5



What are they good for?

Predicting the Quality of Wine - Part 1

- Cortez et al (2009), models the quality of wines (Vinho Verde), 14 vars with 4898 obs (Training = 5199 & Test = 1298).
- "Modeling wine preferences by data mining from physicochemical properties, P. Cortez, A. Cerdeira, F. Almeida, T. Matos and J. Reis, Decision Support Systems 2009, 47(4):547-553".
 - Good (quality score is ≥ 6)
 - Bad (quality score is < 6)

```
##  
##  Bad  Good  
##  476   822
```



Predicting the Quality of Wine - Part 2

Quality of Vinho Verde wines

```
> str(wine.df)
'data.frame': 6497 obs. of 15 variables:
 $ fixed.acidity      : num 7.4 7.8 7.8 11.2 7.4 7.4 7.9 7.3 7.8 7.5 ...
 $ volatile.acidity   : num 0.7 0.88 0.76 0.28 0.7 0.66 0.6 0.65 0.58 0.5 ...
 $ citric.acid        : num 0 0 0.04 0.56 0 0 0.06 0 0.02 0.36 ...
 $ residual.sugar     : num 1.9 2.6 2.3 1.9 1.9 1.8 1.6 1.2 2 6.1 ...
 $ chlorides          : num 0.076 0.098 0.092 0.075 0.076 0.075 0.069 0.065 0.073 0.071 ...
 $ free.sulfur.dioxide: num 11 25 15 17 11 13 15 15 9 17 ...
 $ total.sulfur.dioxide: num 34 67 54 60 34 40 59 21 18 102 ...
 $ density            : num 0.998 0.997 0.997 0.998 0.998 ...
 $ pH                 : num 3.51 3.2 3.26 3.16 3.51 3.51 3.3 3.39 3.36 3.35 ...
 $ sulphates          : num 0.56 0.68 0.65 0.58 0.56 0.56 0.46 0.47 0.57 0.8 ...
 $ alcohol             : num 9.4 9.8 9.8 9.8 9.4 9.4 9.4 10 9.5 10.5 ...
 $ quality             : int 5 5 5 6 5 5 5 7 7 5 ...
 $ color               : Factor w/ 2 levels "red","white": 1 1 1 1 1 1 1 1 1 ...
 $ white               : int 0 0 0 0 0 0 0 0 0 ...
 $ good                : Factor w/ 2 levels "Bad","Good": 1 1 1 2 1 1 1 2 2 1 ...

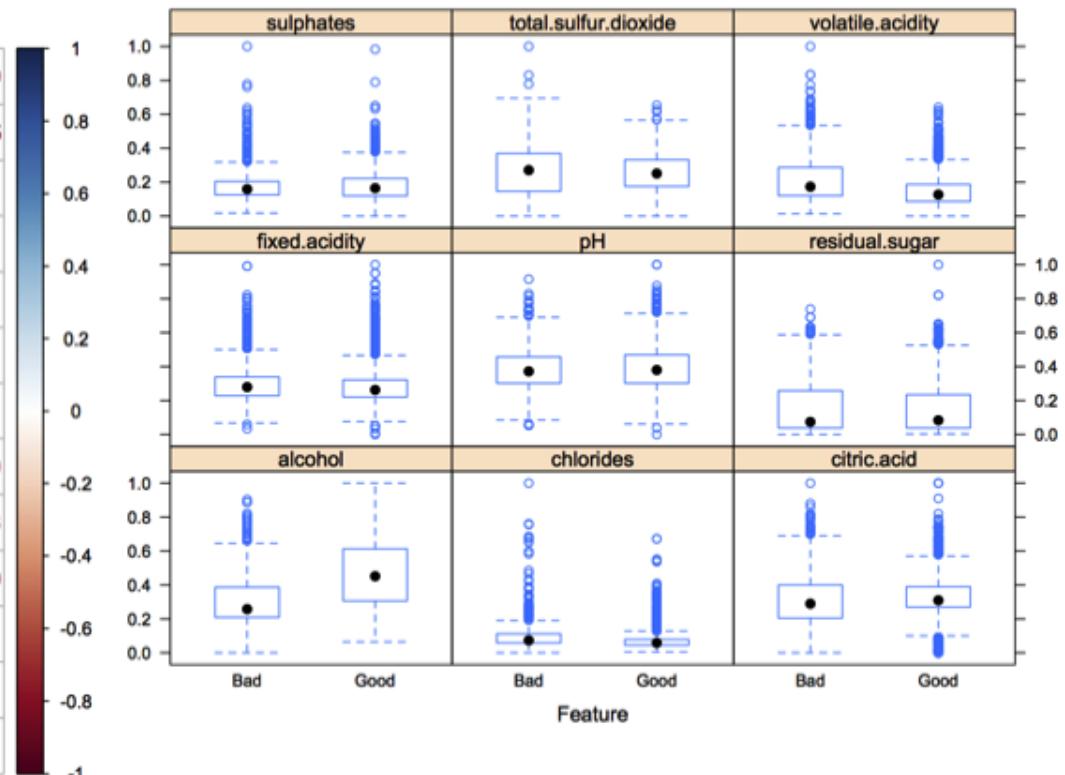
> table(wine.df.test$good)
  Bad Good
    476  822

> head(wine.df)
fixed.acidity volatile.acidity citric.acid residual.sugar chlorides free.sulfur.dioxide total.sulfur.dioxide density pH sulphates alcohol quality color white good
1       7.4         0.70       0.00       1.9      0.076           11              34     0.9978 3.51      0.56      9.4      5 red     0 Bad
2       7.8         0.88       0.00       2.6      0.098           25              67     0.9968 3.20      0.68      9.8      5 red     0 Bad
3       7.8         0.76       0.04       2.3      0.092           15              54     0.9970 3.26      0.65      9.8      5 red     0 Bad
4      11.2         0.28       0.56       1.9      0.075           17              60     0.9988 3.16      0.58      9.8      6 red     0 Good
5       7.4         0.78       0.00       1.9      0.076           11              34     0.9978 3.51      0.56      9.4      5 red     0 Bad
6       7.4         0.66       0.00       1.8      0.075           13              40     0.9978 3.51      0.56      9.4      5 red     0 Bad

> wine.corr
fixed.acidity volatile.acidity citric.acid residual.sugar chlorides free.sulfur.dioxide total.sulfur.dioxide density pH sulphates alcohol quality color white
fixed.acidity 1.0000000 0.21900826 0.32443573 -0.11198128 0.29819477 -0.28273543 -0.32905390 0.45890998 -0.25270047 0.299567744 -0.095451523 -0.07674321 -0.48673983
volatile.acidity 0.21900826 1.00000000 -0.37798132 -0.19680117 0.37712428 -0.35255731 -0.41447619 0.27129565 0.26145444 0.225983680 -0.037640386 -0.26569948 0.65303559
citric.acid 0.32443573 -0.37798132 1.00000000 0.14245123 0.83899881 0.13312581 0.19524198 0.09615393 -0.32980019 0.056197300 -0.010493492 0.08553172 0.18739650
residual.sugar -0.11198128 -0.19680117 0.14245123 1.00000000 -0.12894850 0.48287064 0.49548159 0.55251605 -0.26731984 -0.185927405 -0.359414771 -0.03698848 0.34882181
chlorides 0.29819477 0.37712428 0.03899881 -0.12894850 1.00000000 -0.19504479 -0.27963045 0.36261466 0.84470798 0.395593307 -0.256915580 -0.200665580 -0.51267025
free.sulfur.dioxide -0.28273543 -0.35255731 0.13312581 0.48287064 -0.19504479 1.00000000 0.72093400 0.02571604 -0.14585390 -0.188457249 -0.179838435 0.05546306 0.47164366
total.sulfur.dioxide -0.32905390 -0.41447619 0.19524198 0.49548159 -0.27963045 0.72093400 1.00000000 0.03239451 -0.23841310 -0.275726820 -0.265739639 -0.04138545 0.70835716
density 0.45890998 0.27129565 0.89615393 0.55251605 0.36261466 0.02571604 0.83239451 1.00000000 0.01168608 0.259478495 -0.686745422 -0.30585791 -0.39864532
pH -0.25270047 0.26145444 -0.32980019 -0.26731984 0.04470798 -0.14585390 -0.23841310 0.01168608 1.00000000 0.192123407 0.121248467 0.01958570 -0.32912865
sulphates 0.29956774 0.22598368 0.05619730 -0.18592741 0.39559331 -0.18845725 -0.27572682 0.25947850 0.19212341 1.00000000 0.003029195 0.03848545 -0.48721797
alcohol -0.09545152 -0.03764039 -0.01049349 -0.35941477 -0.25691558 -0.17983843 -0.26573964 -0.68674542 0.12124847 -0.003029195 1.00000000 0.44431852 0.03296955
quality -0.07674321 -0.26569948 0.08553172 -0.03690048 -0.20066558 0.05546306 -0.04138545 -0.30585791 0.01958570 0.030485446 0.444318520 1.00000000 0.11932328
white -0.48673983 -0.65303559 0.18739650 0.34882181 -0.51267025 0.47164366 0.70835716 -0.39864532 -0.32912865 -0.487217970 0.032969551 0.11932328 1.00000000
```

Predicting the Quality of Wine - Part 3

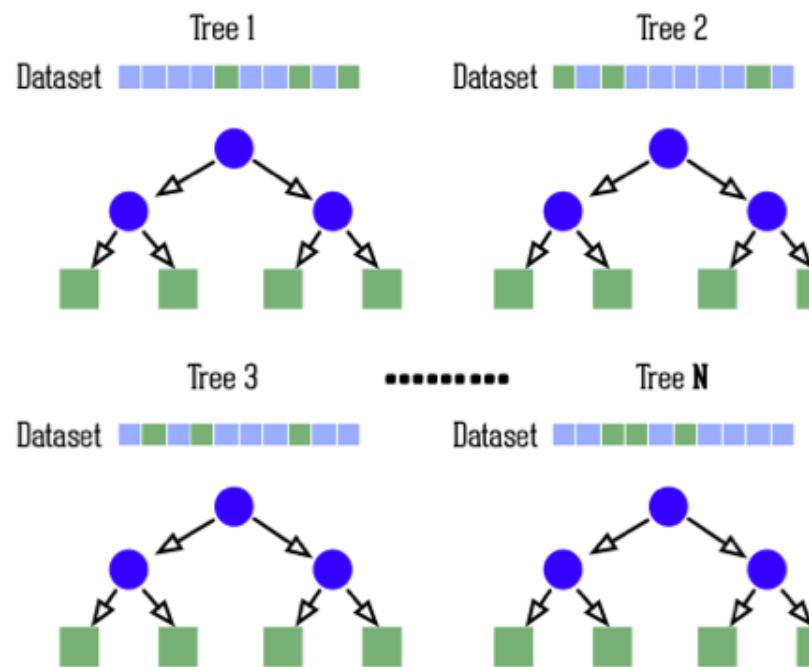
	fixed.acidity	volatile.acidity	citric.acid	residual.sugar	chlorides	free.sulfur.dioxide	total.sulfur.dioxide	density	pH	sulphates	alcohol	quality	white
fixed.acidity	1	0.22	0.32	-0.11	0.3	-0.28	-0.33	0.46	-0.25	0.3	-0.1	-0.08	-0.49
volatile.acidity	0.22	1	-0.38	-0.2	0.38	-0.35	-0.41	0.27	0.26	0.23	0.0	-0.27	-0.65
citric.acid	0.32	-0.38	1	0.14	0.04	0.13	0.2	0.1	-0.33	0.06	0.0	0.09	0.19
residual.sugar	-0.11	-0.2	0.14	1	-0.13	0.4	0.5	0.55	-0.27	-0.19	-0.36	0.0	0.35
chlorides	0.3	0.38	0.04	-0.13	1	-0.2	-0.28	0.36	0.04	0.4	-0.26	-0.2	-0.51
free.sulfur.dioxide	-0.28	-0.35	0.13	0.4	-0.2	1	0.72	0.03	-0.15	-0.19	-0.18	0.00	0.47
total.sulfur.dioxide	-0.33	-0.41	0.2	0.5	-0.28	0.72	1	0.03	-0.24	-0.28	-0.27	-0.04	0.7
density	0.46	0.27	0.1	0.55	0.36	0.03	0.03	1	0.26	-0.69	-0.31	-0.39	
pH	-0.25	0.26	-0.33	-0.27	0.04	-0.15	-0.24	0.03	1	0.19	0.12	0.00	-0.33
sulphates	0.3	0.23	0.06	-0.19	0.4	-0.19	-0.28	0.26	0.19	1	0.00	0.00	-0.49
alcohol	-0.1	-0.04	0.0	-0.36	-0.26	-0.18	-0.27	-0.69	0.12	0.00	1	0.44	0.01
quality	-0.08	-0.27	0.09	-0.04	-0.2	0.06	-0.04	-0.31	0.00	0.00	0.44	1	0.12
white	-0.49	-0.65	0.19	0.35	-0.51	0.47	0.7	-0.39	-0.33	-0.49	0.03	0.12	1



Predicting the Quality of Wine - Part 4 - Problems with Trees

- Deal with irrelevant inputs
 - No data preprocessing required
 - Scalable computation (fast to build)
 - Tolerant with missing values (little loss of accuracy)
 - Only a few tunable parameters (easy to learn)
 - Allows for human understandable graphic representation
- Data fragmentation for high-dimensional sparse data set (over-fitting)
 - Difficult to fit to a trend / piece-wise constant model
 - Highly influenced by changes to the data set and local optima (deep trees might be questionable as the errors propagate down)

Aside - How does a random forest work?

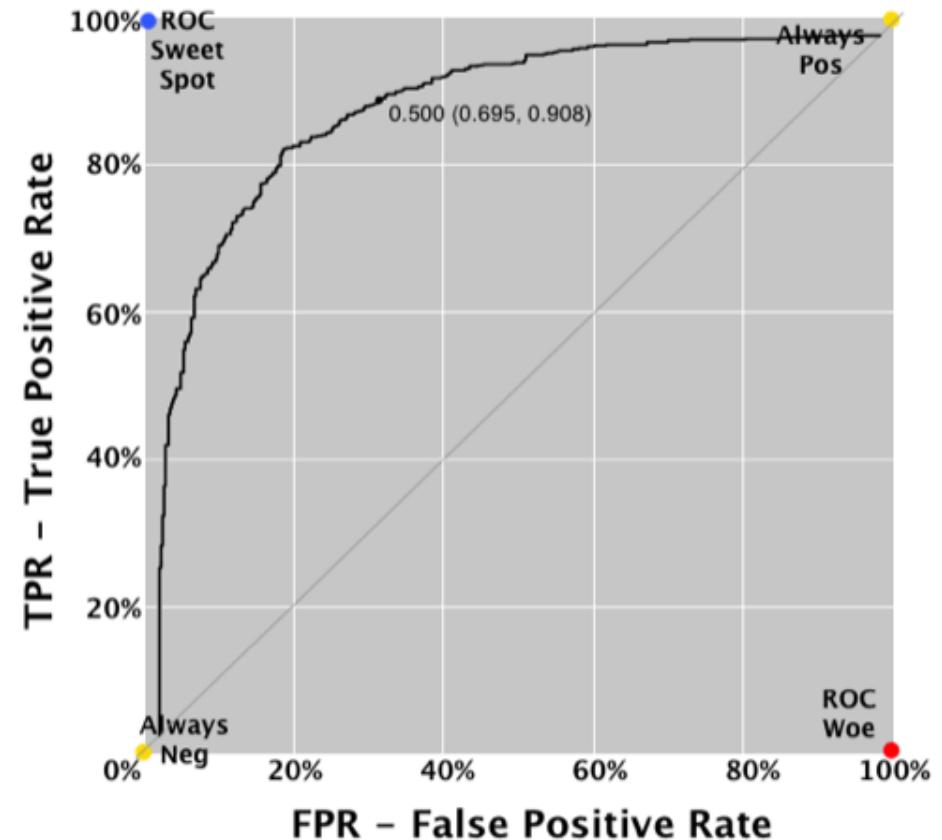


Traverse each tree and at each node in a tree:

- i. Select m random predictor variables from available set
- ii. Use the variable with best split [use objective function]
- iii. Move to next node in tree

Predicting the Quality of Wine - Part 5 - Random Forest

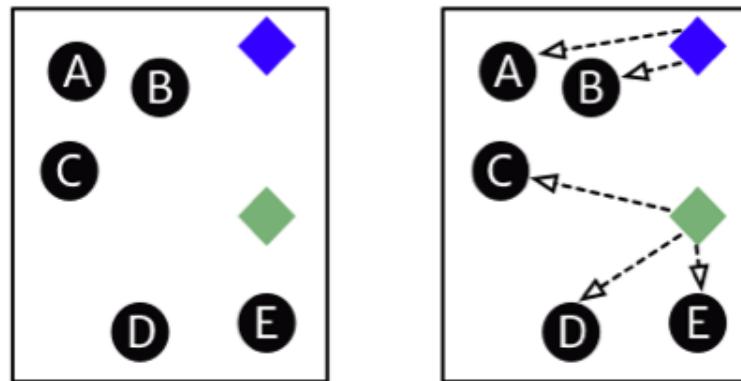
Confusion Matrix		REFERENCE (ACTUAL)		
PREDICTED	Predicted Positive Bad	Predicted Negative Good		
	Positive Examples Bad	TP 331 / 25%	TN 76 / 6%	Pos
Negative Examples Good	FP 145 / 11%	FN 746 / 58%	Neg	
		PPos	PNeg	N



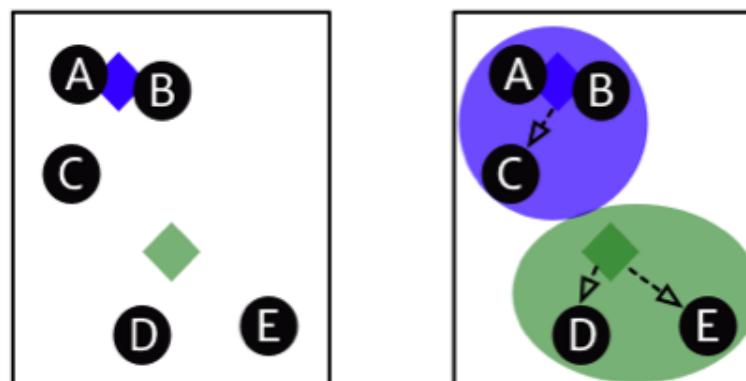
Predicting the Quality of Wine - Part 6 - Other ML methods

- K-nearest neighbors
 - Unsupervised learning / non-target based learning
 - Distance matrix / cluster analysis using Euclidean distances.
- Neural Nets
 - Looking at basic feed forward simple 3-layer network (input, 'processing', output)
 - Each node / neuron is a set of numerical parameters / weights tuned by the learning algorithm used
- Support Vector Machines
 - Supervised learning
 - non-probabilistic binary linear classifier / nonlinear classifiers by applying the kernel trick
 - constructs a hyper-plane/s in a high-dimensional space

Aside - How does k nearest neighbors work?



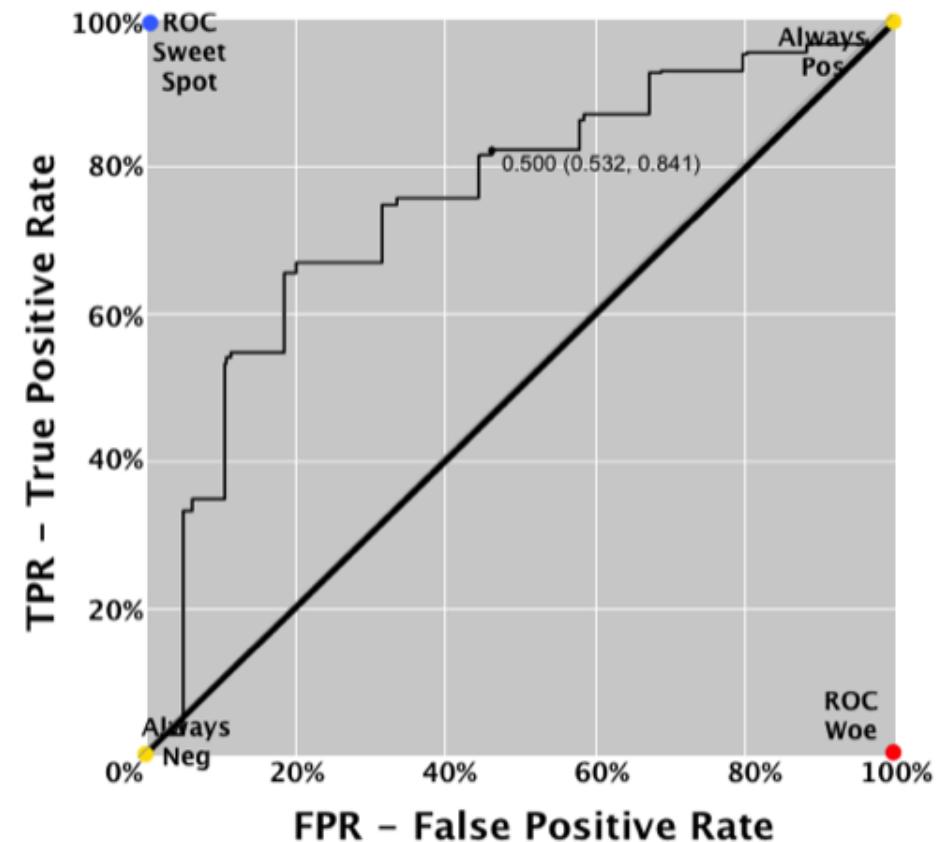
K=2 random centroids, Each item assigned to a centroid, then centroids moved to average location and re-assigned, iterate until assignments stop changing.



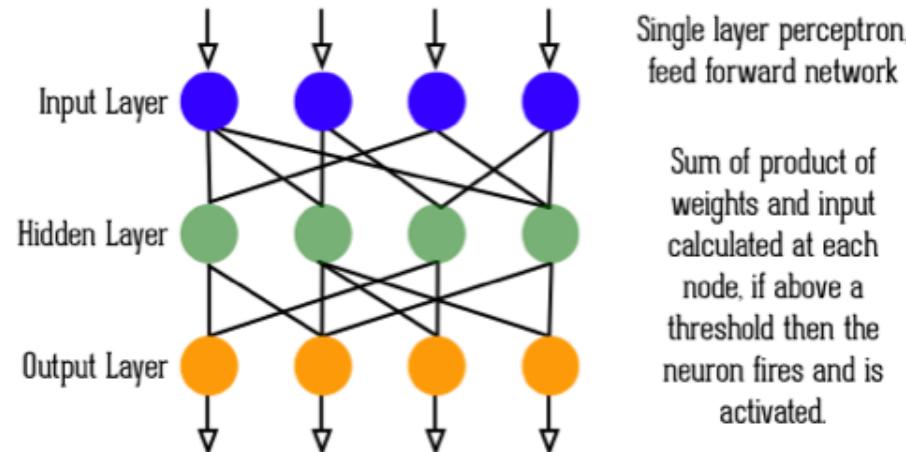
Predicting the Quality of Wine - Part 7 - kNN

		REFERENCE (ACTUAL)		
		Predicted Positive Bad	Predicted Negative Good	
				Pos
PREDICTED	Positive Examples Bad	TP 257 / 20%	TN 132 / 10%	Pos
	Negative Examples Good	FP 219 / 17%	FN 690 / 53%	Neg
		PPos	PNeg	N

*Accuracy = 73%
(20+53)*



Aside - How do neural networks work?



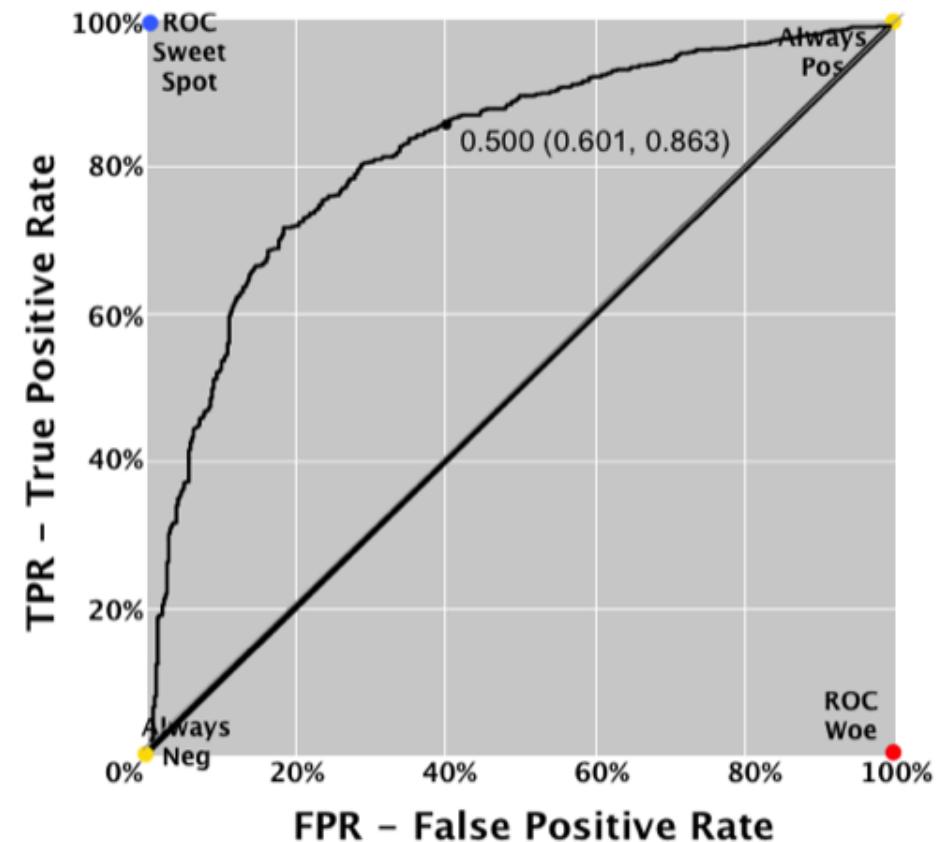
An example of how this can be used:



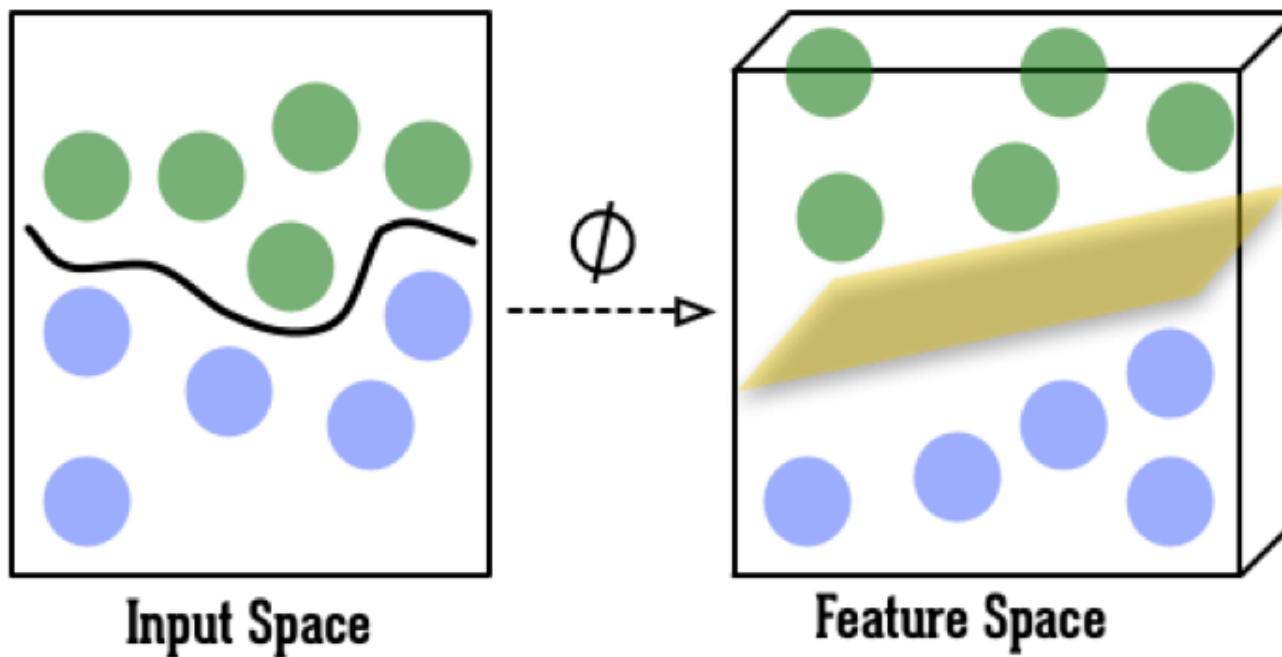
Predicting the Quality of Wine - Part 8 - NNET

Confusion Matrix		REFERENCE (ACTUAL)		
PREDICTED	Positive Examples Bad	Predicted Positive Bad	Predicted Negative Good	
	Positive Examples Bad	TP 286 / 22%	TN 113 / 9%	Pos
	Negative Examples Good	FP 190 / 14%	FN 709 / 55%	Neg
		PPos	PNeg	N

*Accuracy = 77%
(22+55)*



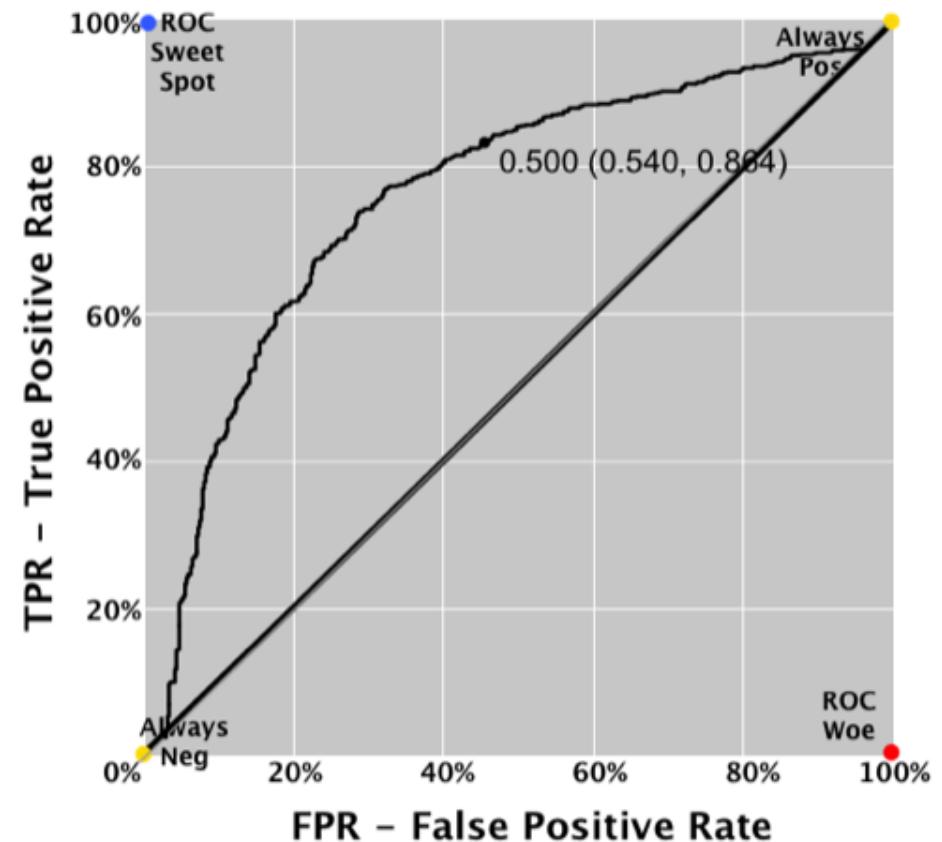
Aside - How do support vector machines work?



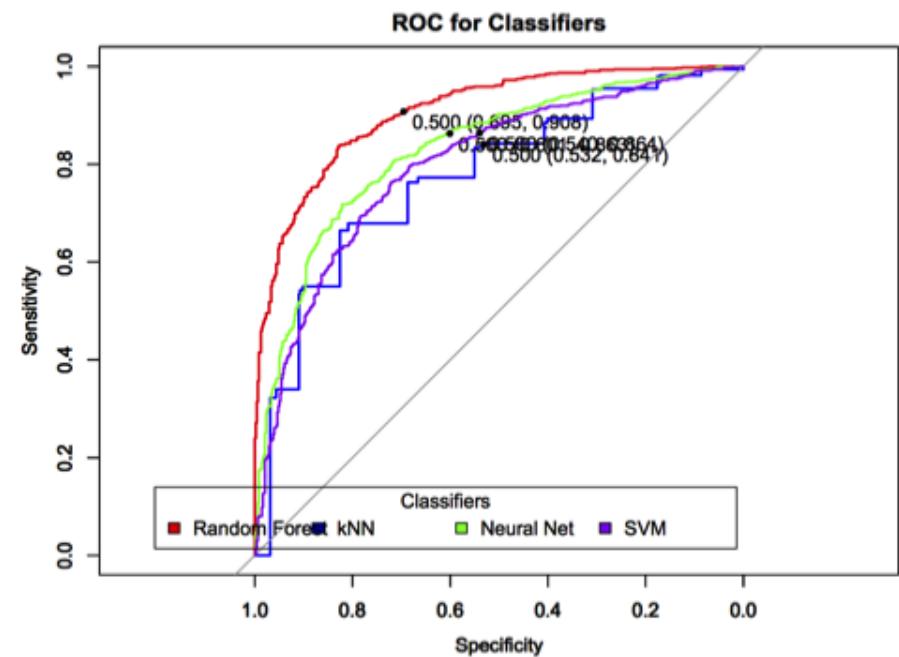
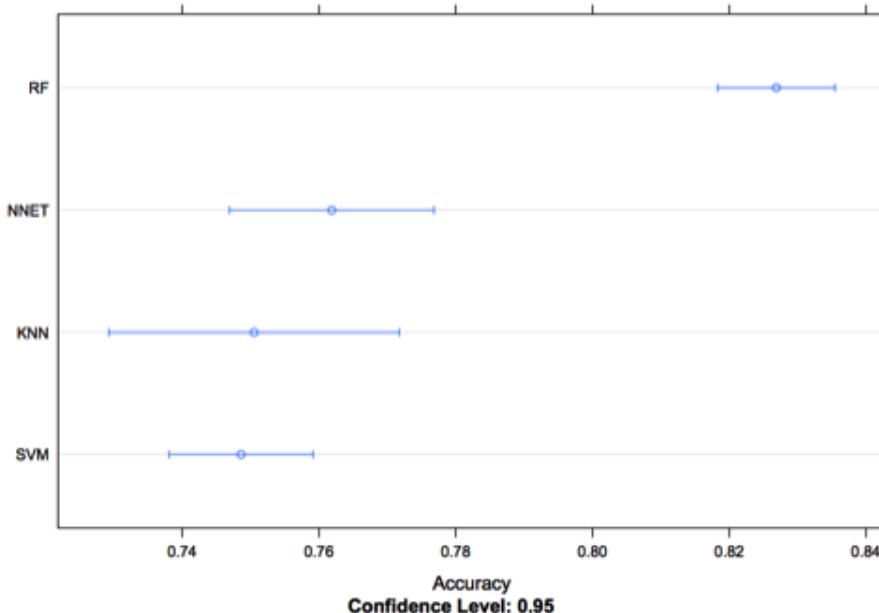
A list of numbers [a n-dimensional vector] and transform the points into higher dimensions so it is easier to separate them using a [n-1] dimensional hyperplane.

Predicting the Quality of Wine - Part 9 - SVN

Confusion Matrix		REFERENCE (ACTUAL)		
PREDICTED	Positive Examples Bad	Predicted Positive Bad	Predicted Negative Good	
	Negative Examples Good	FP 219 / 16%	TN 112 / 9%	Pos
		TP 257 / 20%	FN 710 / 55%	Accuracy = 75% $(20+55)$
		PPos	PNeg	N



Predicting the Quality of Wine - Part 10 - All Results



What are they not good for ?

Predicting the Extramarital Affairs

- Fair, R.C. et al (1978), models the possibility of affairs, 9 vars with 601 obs (Training = 481 & Test = 120).
- "A Theory of Extramarital Affairs, Fair, R.C., Journal of Political Economy 1978, 86:45-61".
 - Yes (affairs is ≥ 1 in last 6 months)
 - No (affairs is < 1 in last 6 months)

```
##  
##  No Yes  
##  90   30
```

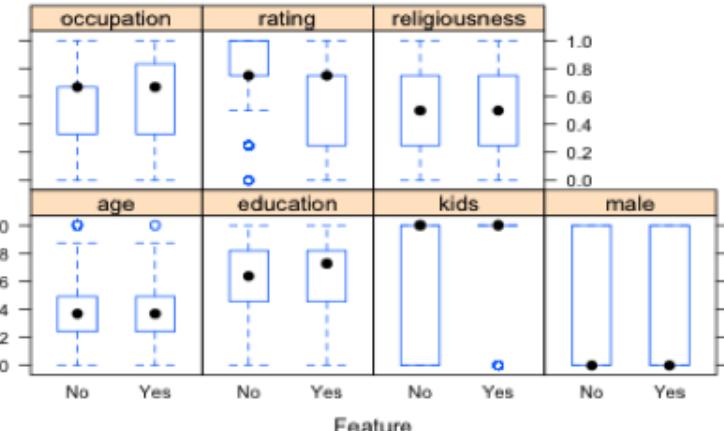


Extramarital Dataset

```
> str(Affairs)
'data.frame': 601 obs. of 12 variables:
$ affairs      : num  0 0 0 0 0 0 0 0 0 ...
$ gender       : Factor w/ 2 levels "female","male": 2 1 1 2 2 1 1 2 1 2 ...
$ age          : num  37 27 32 57 22 32 22 57 32 22 ...
$ yearsmarried : num  10 4 15 15 0.75 1.5 0.75 15 15 1.5 ...
$ children     : Factor w/ 2 levels "no","yes": 1 1 2 2 1 1 1 2 2 1 ...
$ religiousness: int  3 4 1 5 2 2 2 2 4 4 ...
$ education    : num  18 14 12 18 17 17 12 14 16 14 ...
$ occupation   : int  7 6 1 6 6 5 1 4 1 4 ...
$ rating       : int  4 4 4 5 3 5 3 4 2 5 ...
$ male          : num  1 0 0 1 1 0 0 1 0 1 ...
$ kids          : num  0 0 1 1 0 0 0 1 1 0 ...
$ hadaffair    : Factor w/ 2 levels "No","Yes": 1 1 1 1 1 1 1 1 1 1 ...
```

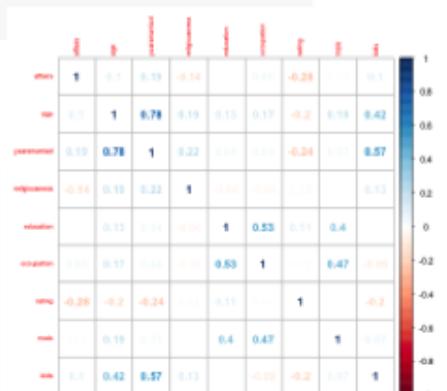
```
> head(Affairs)
affairs gender age yearsmarried children religiousness education occupation rating male kids hadaffair
4      0 male   37      10.00      no      3      18      7      4      1      0      No
5      0 female  27       4.00      no      4      14      6      4      0      0      No
11     0 female  32      15.00      yes     1      12      1      4      0      1      No
16     0 male    57      15.00      yes     5      18      6      5      1      1      No
23     0 male    22       0.75      no      2      17      6      3      1      0      No
29     0 female  32       1.50      no      2      17      5      5      0      0      No
```

```
> affairs.corr
affairs      affairs      age  yearsmarried  religiousness  education  occupation  rating  male  kids  hadaffair
affairs  1.000000000  0.0952372  0.18684169 -0.144501345 -0.002437441  0.04961176 -0.279512403  0.011736251  0.104010057
age        0.095237204  1.0000000  0.77754585  0.193776931  0.134596015  0.16641254 -0.198999899  0.190641080  0.421930815
yearsmarried  0.186841686  0.7775458  1.000000000  0.218260672  0.040002716  0.04459281 -0.243118827  0.030282521  0.572857364
religiousness -0.144501345  0.1937769  0.21826067  1.000000000 -0.042571079 -0.03972232  0.024295777  0.007679445  0.129351259
education    -0.002437441  0.1345960  0.040002727 -0.042571079  1.000000000  0.53360524  0.109303473  0.397504680 -0.006985882
occupation   0.049611758  0.1664125  0.04459201 -0.039722324  0.533605242  1.000000000  0.017422273  0.467923152 -0.092727118
rating       -0.279512403 -0.1989999 -0.24311883  0.024295777  0.109303473  0.01742227  1.000000000 -0.007523748 -0.196275616
male         0.011736251  0.1906411  0.03028252  0.007679445  0.397504680  0.46792315 -0.007523748  1.000000000  0.069222338
kids          0.0180057  0.4219308  0.57285736  0.129351259 -0.006985882 -0.09272712 -0.196275616  0.069222338  1.000000000
```



```
> table(affairs.df.test$hadaffair)
```

	No	Yes
	90	30



Predicting the Extramarital Affairs - RF & NB

Random Forest

```
##             Reference  
## Prediction No Yes  
##          No 90 30  
##          Yes 0 0
```

```
## Accuracy  
## 0.75
```

Naive Bayes

```
##             Reference  
## Prediction No Yes  
##          No 88 29  
##          Yes 2 1
```

```
## Accuracy  
## 0.75
```

Other related tools: OpenRefine (formerly Google Refine) / Rattle

Google refine
A power tool for working with messy data.

Create Project Start Over Configure Parsing Options Project name: goodine.csv Create Project ▾

	fixed acidity	volatile acidity	citric acid	residual sugar	chlorides	free sulfur dioxide	total sulfur dioxide	density	pH	sulphates	alcohol	quality	color	white
1.	7.4	0.7	0	1.9	0.076	11	34	0.9978	3.51	0.56	9.4	5	red	0
2.	7.8	0.66	0	2.6	0.098	25	67	0.9968	3.2	0.68	9.8	5	red	0
3.	7.8	0.76	0.04	2.3	0.092	15	54	0.997	3.26	0.65	9.8	5	red	0
4.	11.2	0.28	0.56	1.9	0.075	17	60	0.998	3.16	0.58	9.8	6	red	0
5.	7.4	0.7	0	1.9	0.076	11	34	0.9978	3.51	0.56	9.4	5	red	0
6.	7.4	0.66	0	1.8	0.075	13	40	0.9978	3.51	0.56	9.4	5	red	0
7.	7.9	0.6	0.06	1.6	0.099	15	59	0.9964	3.3	0.66	9.4	5	red	0
8.	7.3	0.65	0	1.2	0.065	15	21	0.9946	3.39	0.47	10	7	red	0
9.	7.8	0.66	0.62	2	0.073	9	18	0.9968	3.36	0.57	9.5	7	red	0
10.	7.5	0.5	0.36	6.1	0.071	17	102	0.9978	3.35	0.8	10.5	5	red	0
11.	8.7	0.66	0.08	1.6	0.097	15	65	0.9909	3.29	0.54	9.2	5	red	0
12.	7.5	0.5	0.36	6.1	0.071	17	102	0.9978	3.35	0.8	10.5	5	red	0
13.	8.6	0.616	0	1.6	0.089	16	59	0.9943	3.58	0.52	9.9	5	red	0
14.	7.8	0.61	0.29	1.6	0.114	9	29	0.9974	3.26	1.66	9.1	5	red	0
15.	8.9	0.62	0.18	3.8	0.176	52	145	0.9986	3.16	0.88	9.2	5	red	0
16.	8.9	0.62	0.19	3.9	0.17	51	148	0.9986	3.17	0.93	9.2	5	red	0
17.	8.5	0.28	0.56	1.8	0.092	35	103	0.9969	3.3	0.75	10.5	7	red	0

Parse data as Character encoding: Update Preview

CSV / TSV / separator-based files Line-based text files Fixed-width field text files PC-Axis text files JSON files RDF/NDX files XML files Open Document Format spreadsheets (.ods) RDF/XML files

Character encoding: Columns are separated by: Ignore first: 0 line(s) at beginning of file

Ignore first: 0 line(s) at beginning of file Parse next: 1 line(s) as column headers Discard initial: 0 row(s) of data Load at most: 0 row(s) of data

Parse next: 1 line(s) as column headers Discard initial: 0 row(s) of data Load at most: 0 row(s) of data

Parse-cell text into numbers, dates, ... Store blank rows Quotation marks are used to enclose cells containing column separators Store blank cells as nulls Store file source (.file names, URLs) in each row

R Data Miner - [Rattle]

Project Tools Settings Help Rattle Version 2.6.6 togaware.com

Execute New Open Save Report Export Stop Quit

Data Explore Test Transform Cluster Associate Model Evaluate Log

Source: Spreadsheet ARFF ODBC R Dataset RData File Library Corpus Script

Filename: (None) Separator: Decimal: Header:

Partition: 70/15/15 Seed: 42 View Edit

Input Ignore Weight Calculator: Target Data Type: Auto Categorical Numeric Survival

Welcome to Rattle (rattle.togaware.com).

Rattle is a free graphical user interface for Data Mining, developed using R. R is a free software environment for statistical computing and graphics. Together they provide a sophisticated environments for data mining, statistical analyses, and data visualisation.

See the Help menu for extensive support in using Rattle. The Togaware Desktop Data Mining Survival Guide includes Rattle documentation and is available from datamining.togaware.com

Rattle is licensed under the GNU General Public License, Version 2. Rattle comes with ABSOLUTELY NO WARRANTY. See Help -> About for details.

Rattle Version 2.6.6. Copyright 2006-2011 Togaware Pty Ltd
Rattle is a registered trademark of Togaware Pty Ltd

To Begin: Choose the data source, specify the details, then click the Execute button.

Other related tools: Command Line Utilities

- <http://www.gregreda.com/2013/07/15/unix-commands-for-data-science/>
 - sed / awk
 - head / tail
 - wc (word count)
 - grep
 - sort / uniq
- <http://blog.comsysto.com/2013/04/25/data-analysis-with-the-unix-shell/>
 - join
 - Gnuplot
- <http://jeroenjanssens.com/2013/09/19/seven-command-line-tools-for-data-science.html>
 - <http://csvkit.readthedocs.org/en/latest/>
 - <https://github.com/jehiah/json2csv>
 - <http://stedolan.github.io/jq/>
 - <https://github.com/jeroenjanssens/data-science-toolbox/blob/master/sample>
 - https://github.com/bitly/data_hacks
 - <https://github.com/jeroenjanssens/data-science-toolbox/blob/master/Rio>
 - <https://github.com/parmentf/xml2json>

A (incomplete) tour of the packages in R

- caret
- party
- rpart
- rpart.plot
- AppliedPredictiveModeling
- randomForest
- corrplot
- arules
- arulesViz
- C50
- pROC
- corrplot
- kernlab
- rattle
- RColorBrewer
- corrgram
- ElemStatLearn
- car

In Summary

An idea of some of the types of classifiers available in ML.

What a confusion matrix and ROC means for a classifier and how to interpret them

An idea of how to test a set of techniques and parameters to help you find the best model for your data

Slides, Data, Scripts are all on GH:

<https://github.com/braz/DublinR-ML-treesandforests>