

# CLASSIFICATION AND VALUE ESTIMATION

## SETTING THE STAGE

“Data science does not replace statistical modeling and data analysis; it augments them.”

(P. Boily)

“Data is not information, information is not knowledge,  
knowledge is not understanding, understanding is not wisdom.”

(attributed to Cliff Stoll in Keeler's *Nothing to Hide: Privacy in the 21st Century*, 2006)

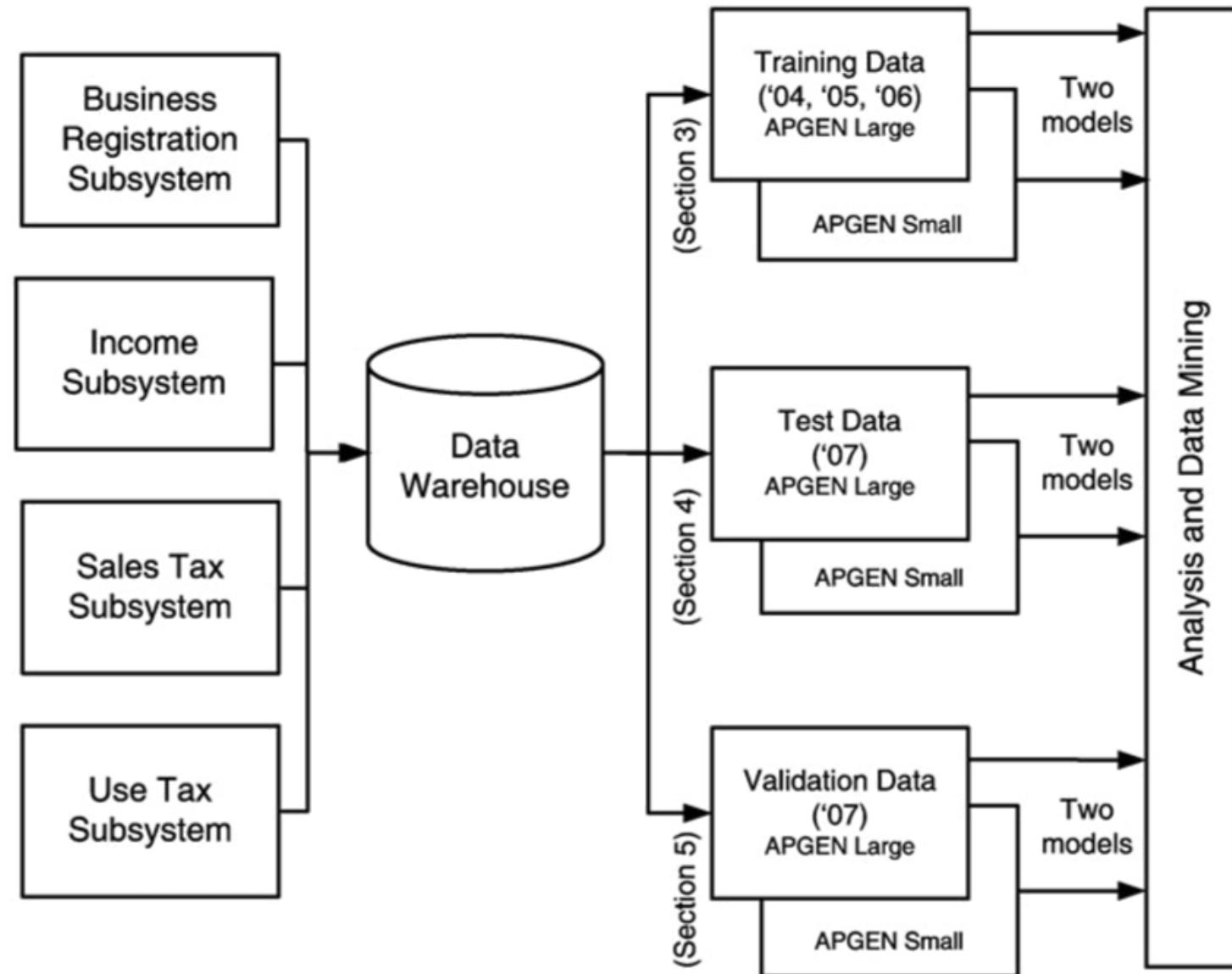
## CASE STUDY: MINNESOTA TAX AUDIT

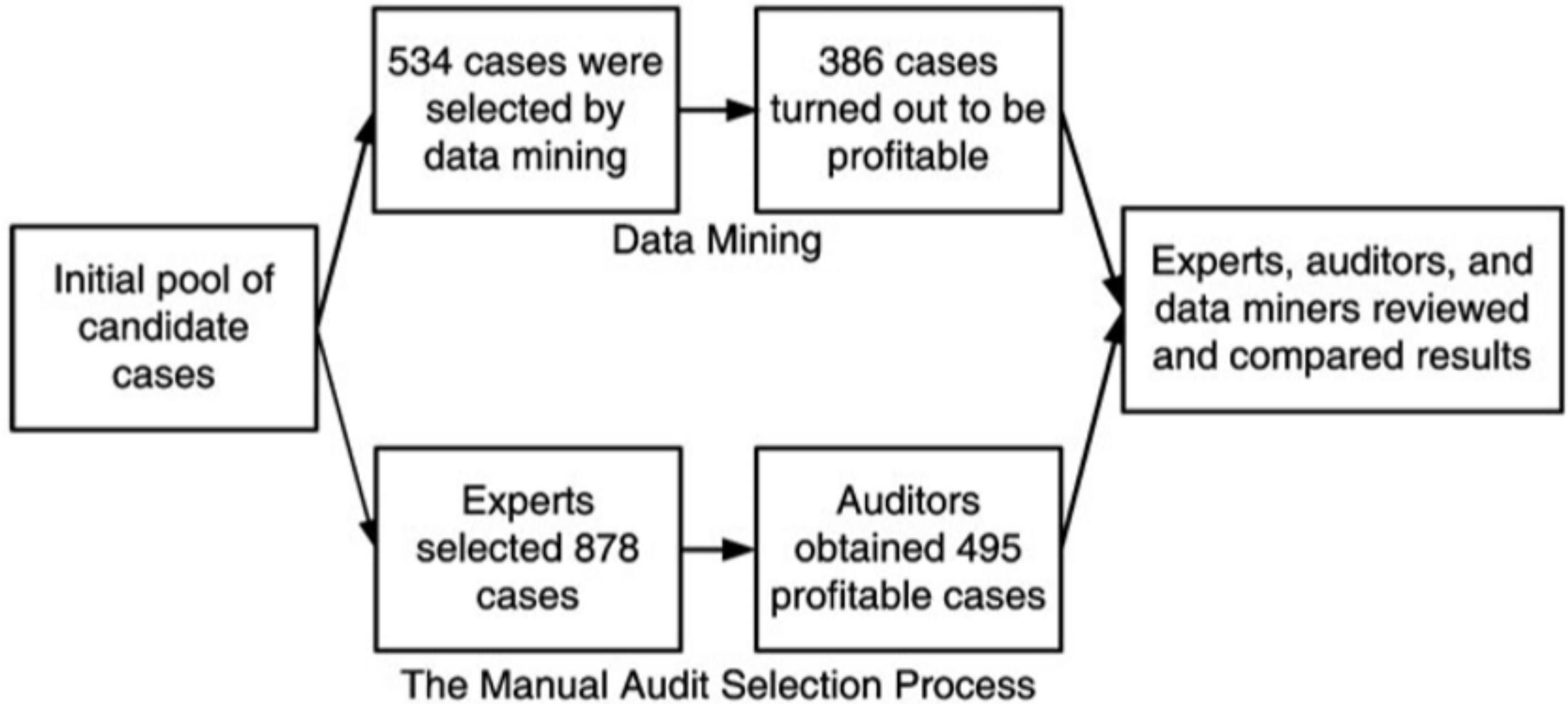
Large gaps between revenue owed (in theory) and revenue collected (in practice) are problematic for governments.

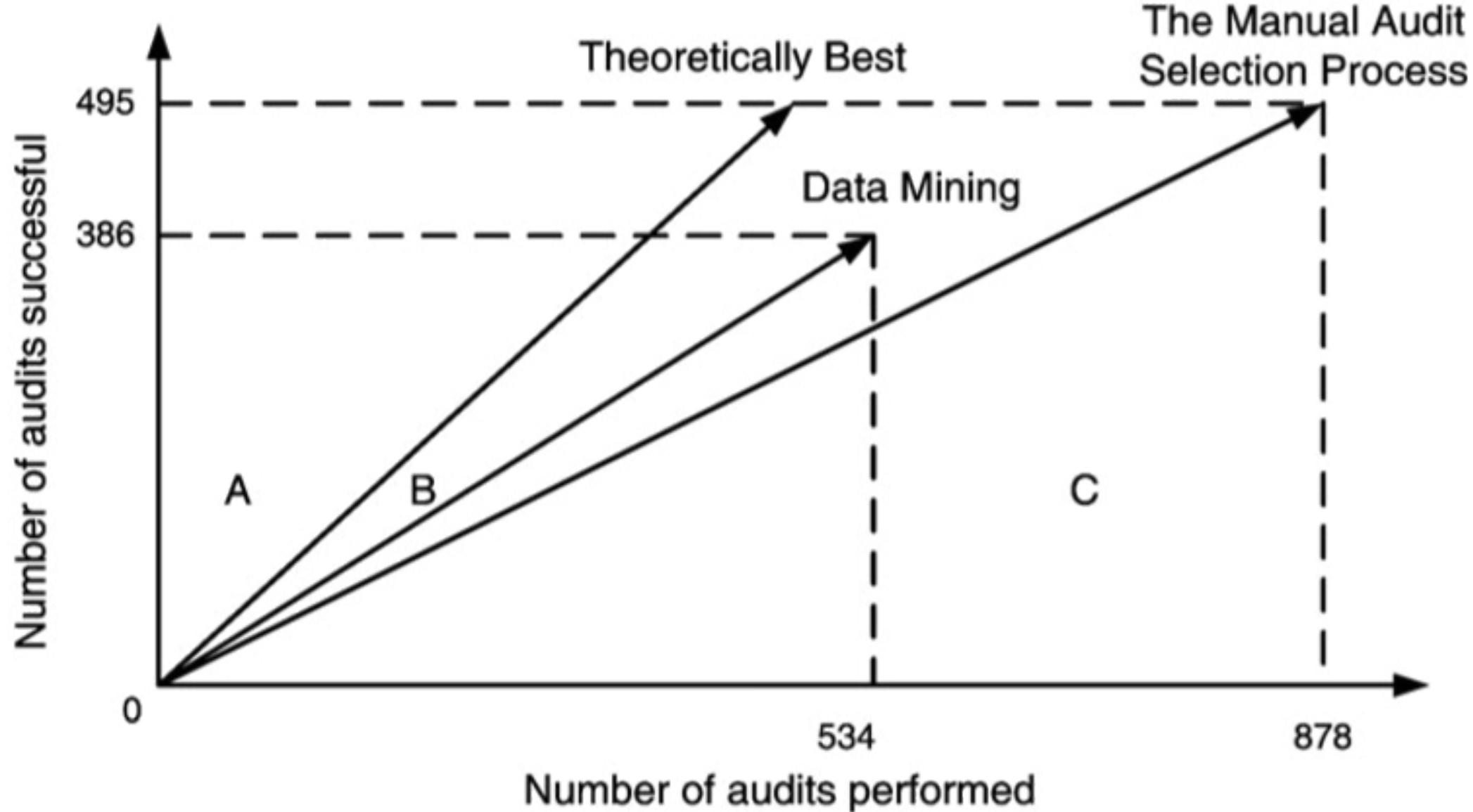
Revenue agencies implement various fraud detection strategies (such as audit reviews) to bridge that gap.

Business audits are costly – are there **algorithms that can predict whether an audit is likely to be successful or a waste of resources?**

Should a tax collection agency seek to maximize its revenues/profits or to ensure compliance?







## CLASSIFICATION OVERVIEW

In **classification**, a sample set of data (the **training** set) is used to determine rules and patterns that divide the data into pre-determined groups, or classes (supervised learning; predictive analytics).

The training data usually consists of a **randomly** selected subset of the **labeled** (target) data.

**Value estimation** (regression) is akin to classification when the target variable is numerical.

# CLASSIFICATION OVERVIEW

In the **testing** phase, the model is used to assign a class to observations for which the label is hidden, but ultimately known (the **testing** set).

The performance of a classification model is evaluated on the testing set, **never** on the training set.

Technical issues include:

- selecting the features to include in the model
- selecting the algorithm
- etc.

# APPLICATIONS

## Medicine and Health Science

- predicting which patient is at risk of suffering a second, fatal heart attack within 30 days based on health factors (blood pressure, age, sinus problems, etc.)

## Social Policies

- predicting the likelihood of requiring assisting housing in old age based on demographic information/survey answers

## Marketing and Business

- predicting which customers are likely to switch to another cell phone company based on demographics and usage

## EXAMPLE

### Scenario:

A motor insurance company has a fraud investigation dept. that studies up to 30% of all claims made, yet money is still getting lost on fraudulent claims.

### Questions: can we predict

- whether a claim is likely to be fraudulent?
- whether a customer is likely to commit fraud in the near future?
- whether an application for a policy is likely to result in a fraudulent claim?
- the amount by which a claim will be reduced if it is fraudulent?

## Testing Set (with labels)

	$Y_1$	$Y_2$	...	$Y_p$	■
02	$x_{02,1}$	$x_{02,2}$	...	$x_{02,p}$	■
03	$x_{03,1}$	$x_{03,2}$	...	$x_{03,p}$	■
05	$x_{05,1}$	$x_{05,2}$	...	$x_{05,p}$	■
06	$x_{06,1}$	$x_{06,2}$	...	$x_{06,p}$	■
07	$x_{07,1}$	$x_{07,2}$	...	$x_{07,p}$	■
08	$x_{08,1}$	$x_{08,2}$	...	$x_{08,p}$	■
09	$x_{09,1}$	$x_{09,2}$	...	$x_{09,p}$	■
11	$x_{11,1}$	$x_{11,2}$	...	$x_{11,p}$	■
...			...		
@@	$x_{@@,1}$	$x_{@@,2}$	...	$x_{@@,p}$	■

## Training Set (with labels)

	$Y_1$	$Y_2$	...	$Y_p$	■
01	$x_{01,1}$	$x_{01,2}$	...	$x_{01,p}$	■
04	$x_{04,1}$	$x_{04,2}$	...	$x_{04,p}$	■
10	$x_{10,1}$	$x_{10,2}$	...	$x_{10,p}$	■
21	$x_{21,1}$	$x_{21,2}$	...	$x_{21,p}$	■
22	$x_{22,1}$	$x_{22,2}$	...	$x_{22,p}$	■
23	$x_{23,1}$	$x_{23,2}$	...	$x_{23,p}$	■
25	$x_{25,1}$	$x_{25,2}$	...	$x_{25,p}$	■
29	$x_{29,1}$	$x_{29,2}$	...	$x_{29,p}$	■
...			...		
**	$x_{**,1}$	$x_{**,2}$	...	$x_{**,p}$	■

Classifier

Model

Classes

## Predictions

	a	p
02	■	■
03	■	■
05	■	■
06	■	■
07	■	■
08	■	■
09	■	■
11	■	■
...	...	...
@@	■	■

Performance Evaluation

Deployment

# CLASSIFICATION SCHEMES

**Logistic Regression**

**Neural Networks**

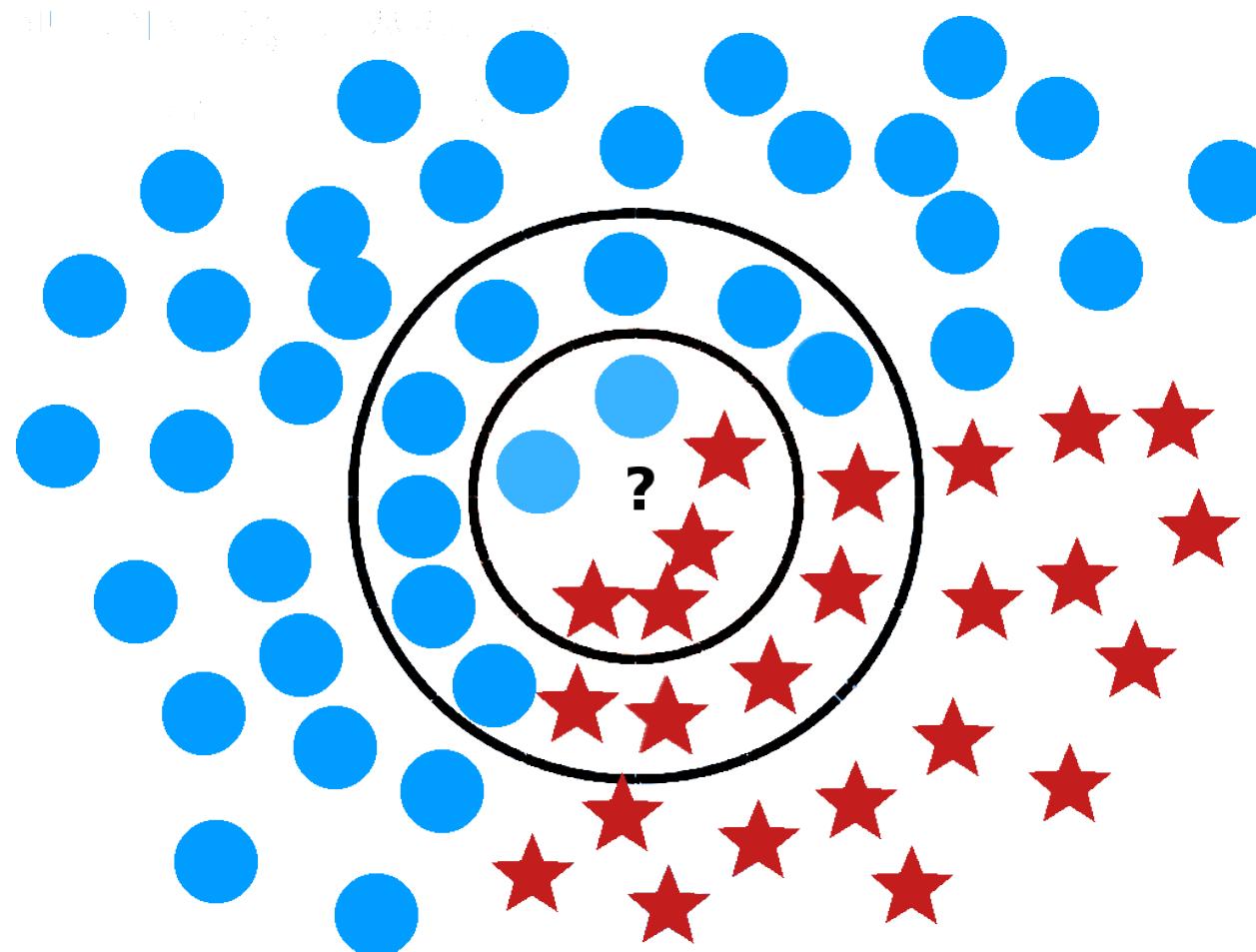
**Decision Trees**

**Naïve Bayes Classifiers**

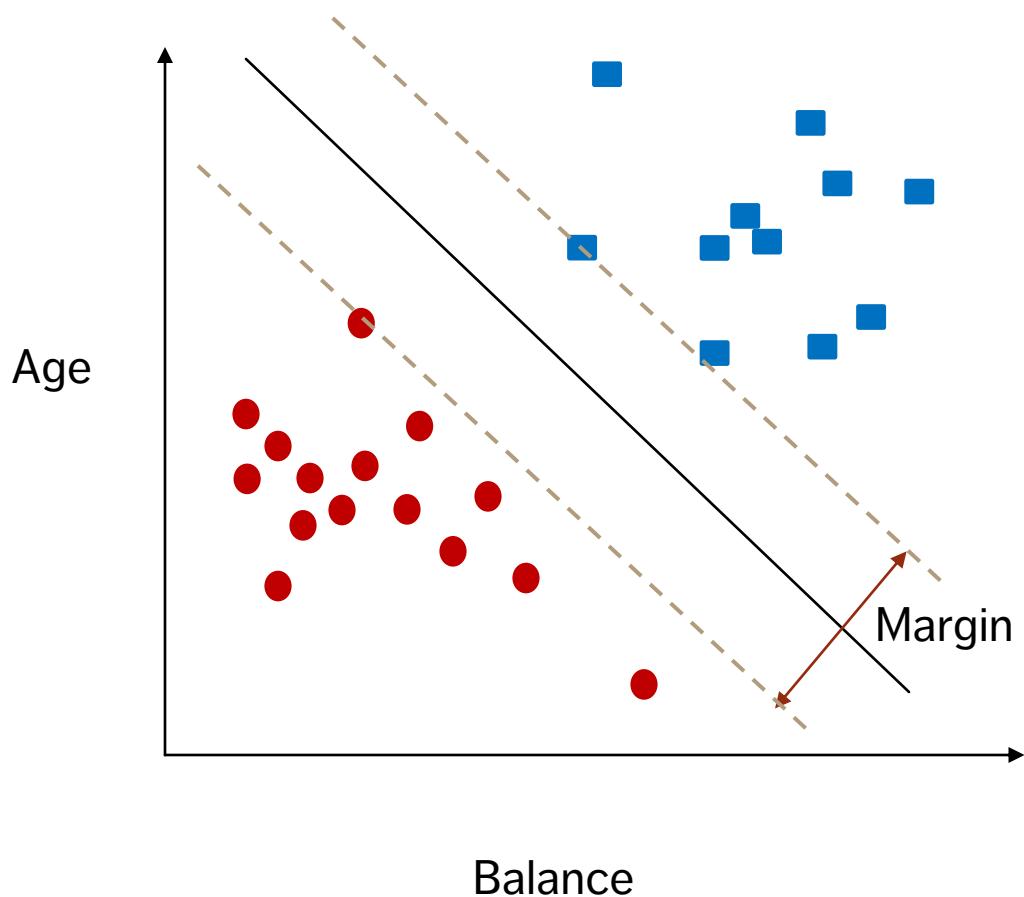
**Support Vector Machines**

**Nearest Neighbours Classifiers**

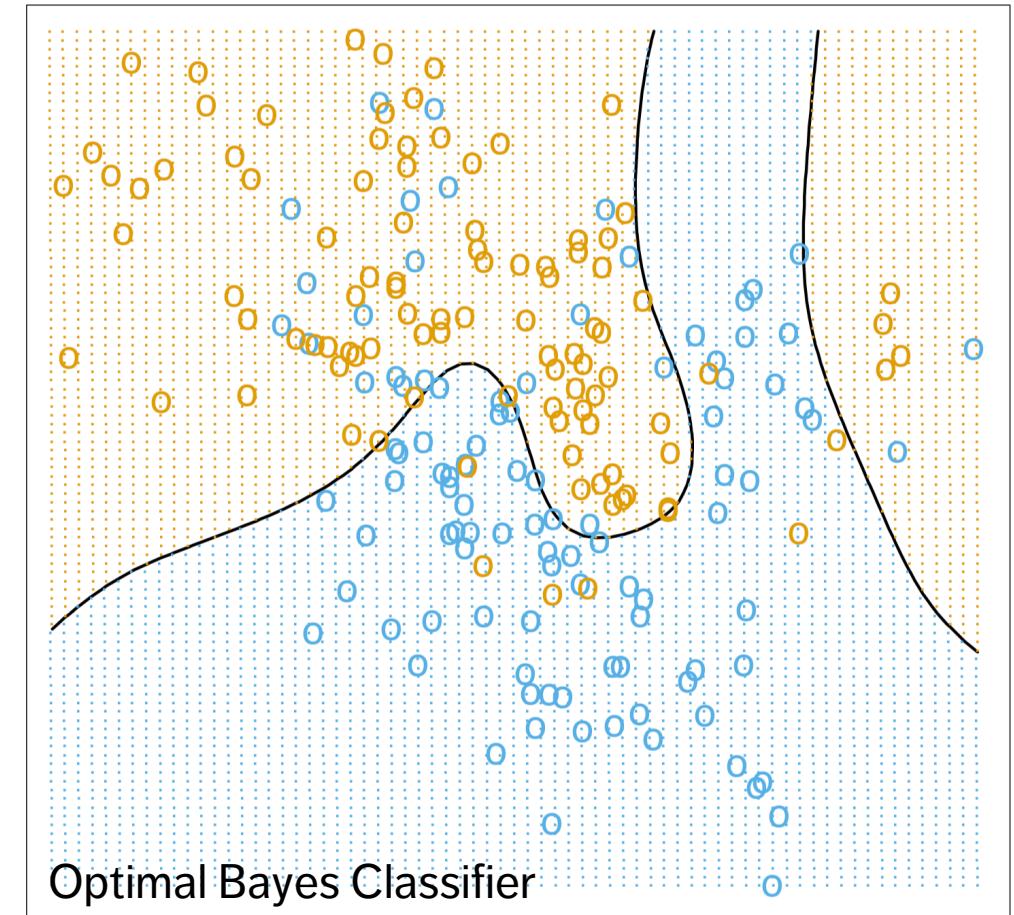
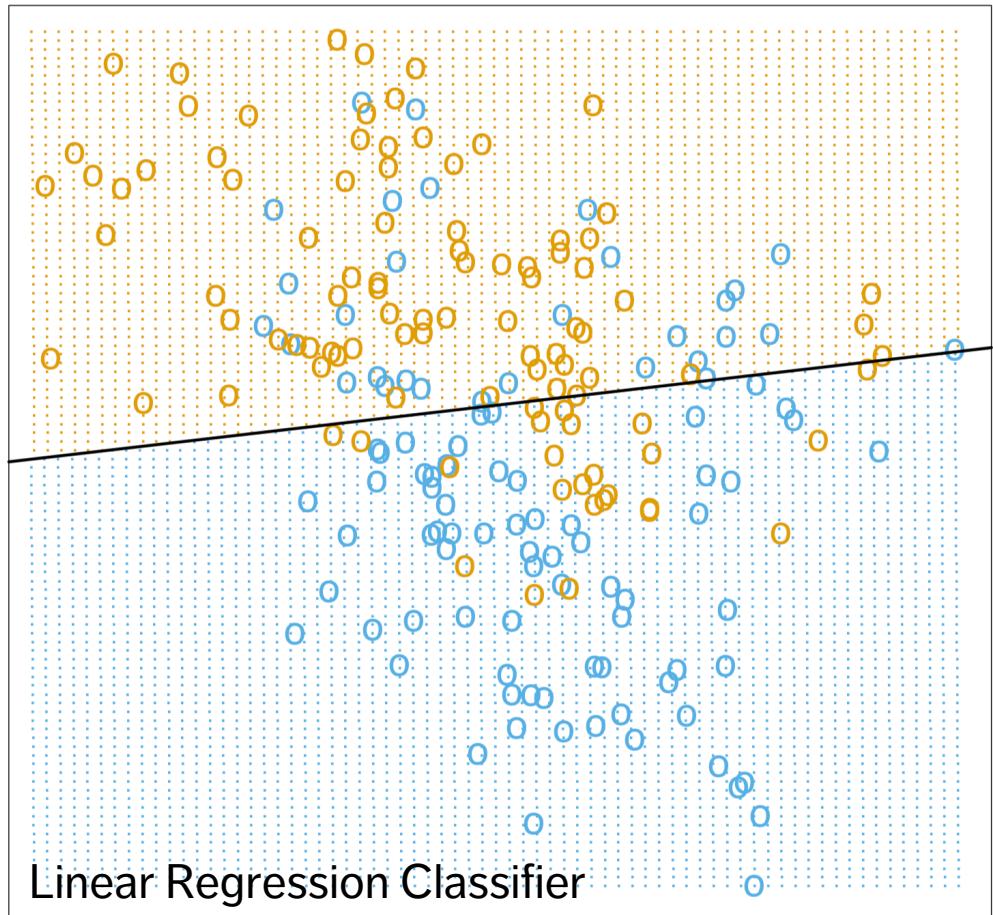
# $k$ NEAREST NEIGHBOURS



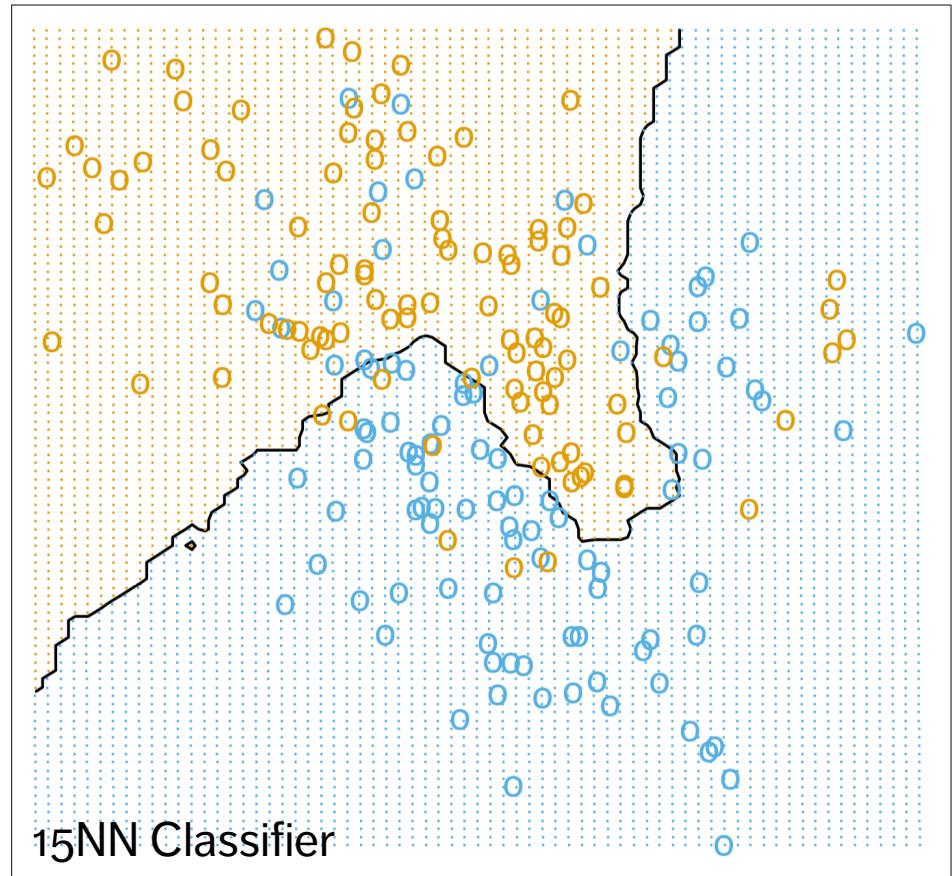
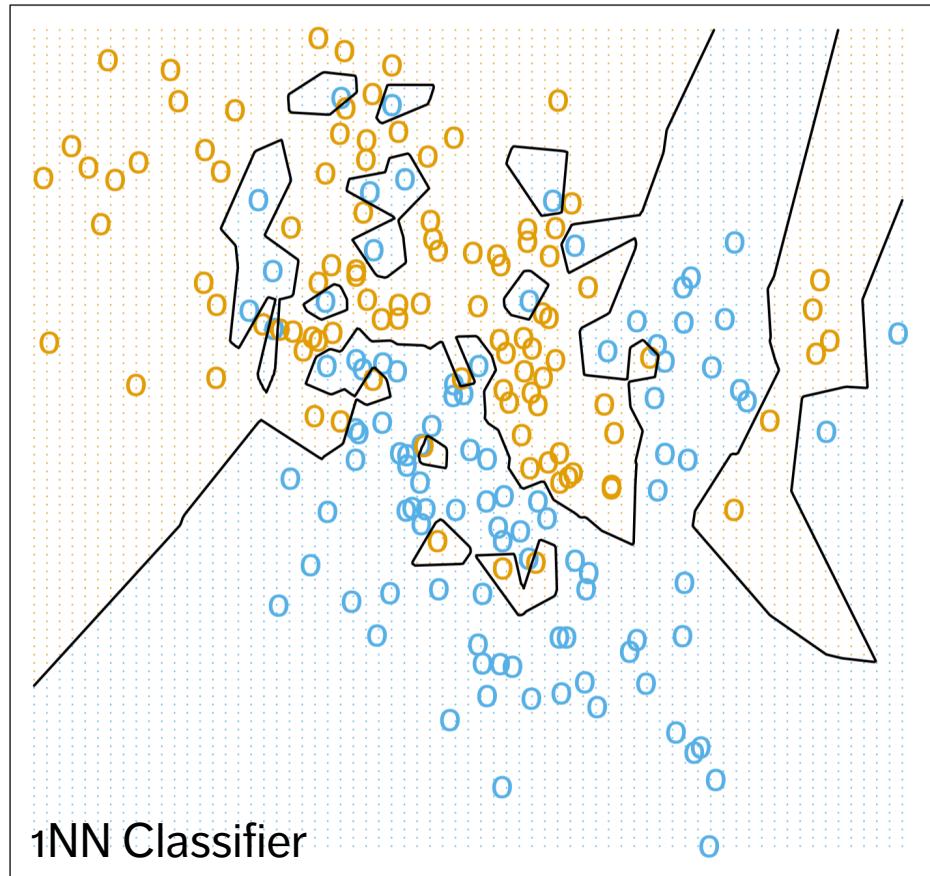
# SUPPORT VECTOR MACHINES



# ILLUSTRATION



# ILLUSTRATION



# DECISION TREES

Decision trees are perhaps the most **intuitive** of these methods: classification is achieved by following a path up the tree, from its **root**, through its **branches**, and ending at its **leaves**.



# DECISION TREES

To make a **prediction** for a new instance, follow the path down the tree, reading the prediction directly once a leaf is reached.

Creating the tree and traversing it might be **time-consuming** if there are too many variables.

Prediction accuracy can be a concern in trees whose growth is **unchecked**. In practice, the criterion of **purity** at the leaf-level is linked to bad prediction rates for new instances.

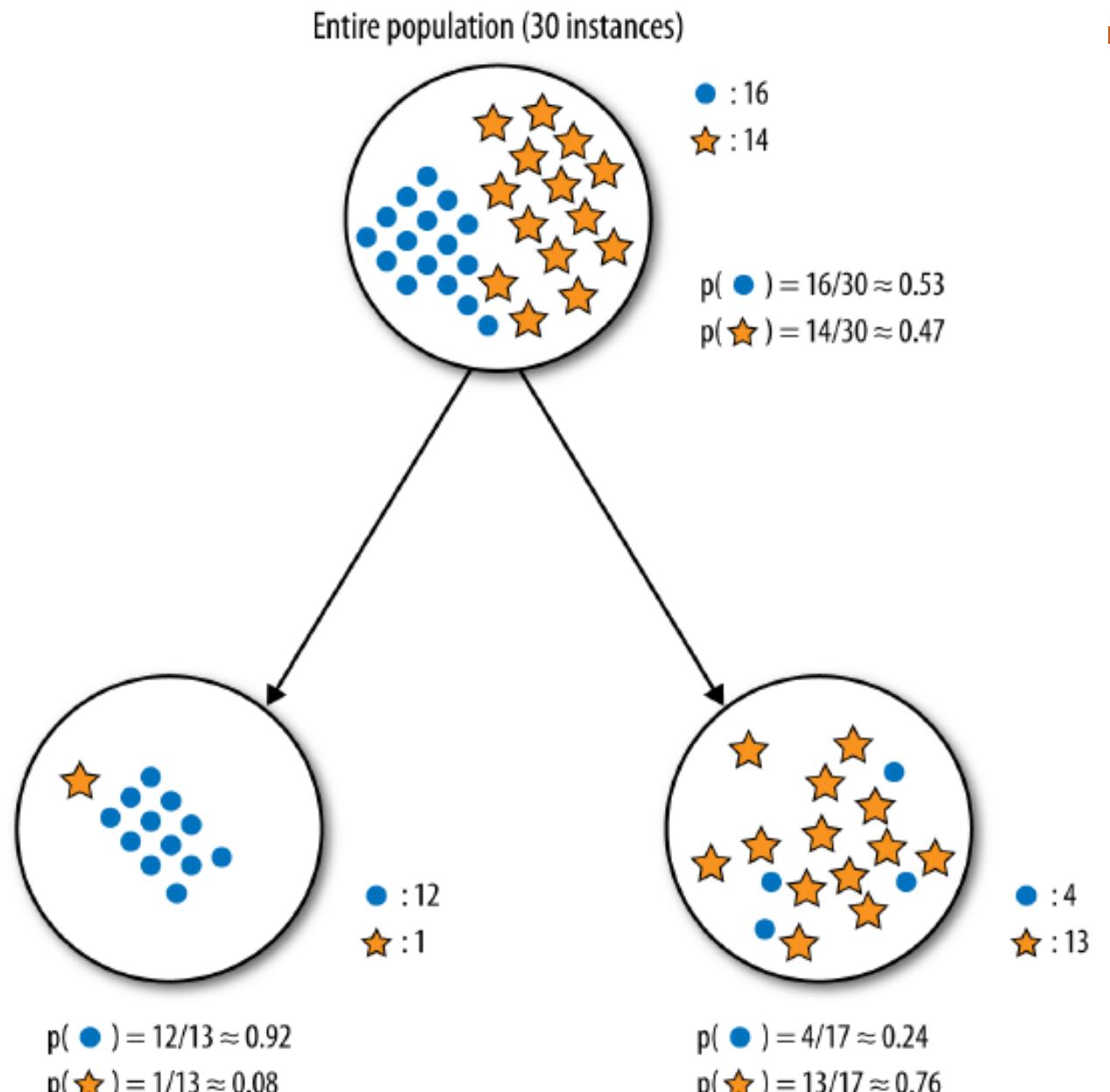
- other criteria are often used to prune trees, which may lead to **impure** leaves (i.e. with non-trivial entropy).

# DECISION TREE ALGORITHM (ID3)

**Task:** grow a decision tree using a training set (a subset of the data for which the correct classification of the target is known).

## Overview:

1. Split the training data (**parent**) set into (**children**) subsets, using the different levels of a particular attribute
2. Compute the **information gain** for each subset
3. Select the **most advantageous** split
4. Repeat for each node until some **leaf** criterion is met (each item in the leaf has the same classification?)

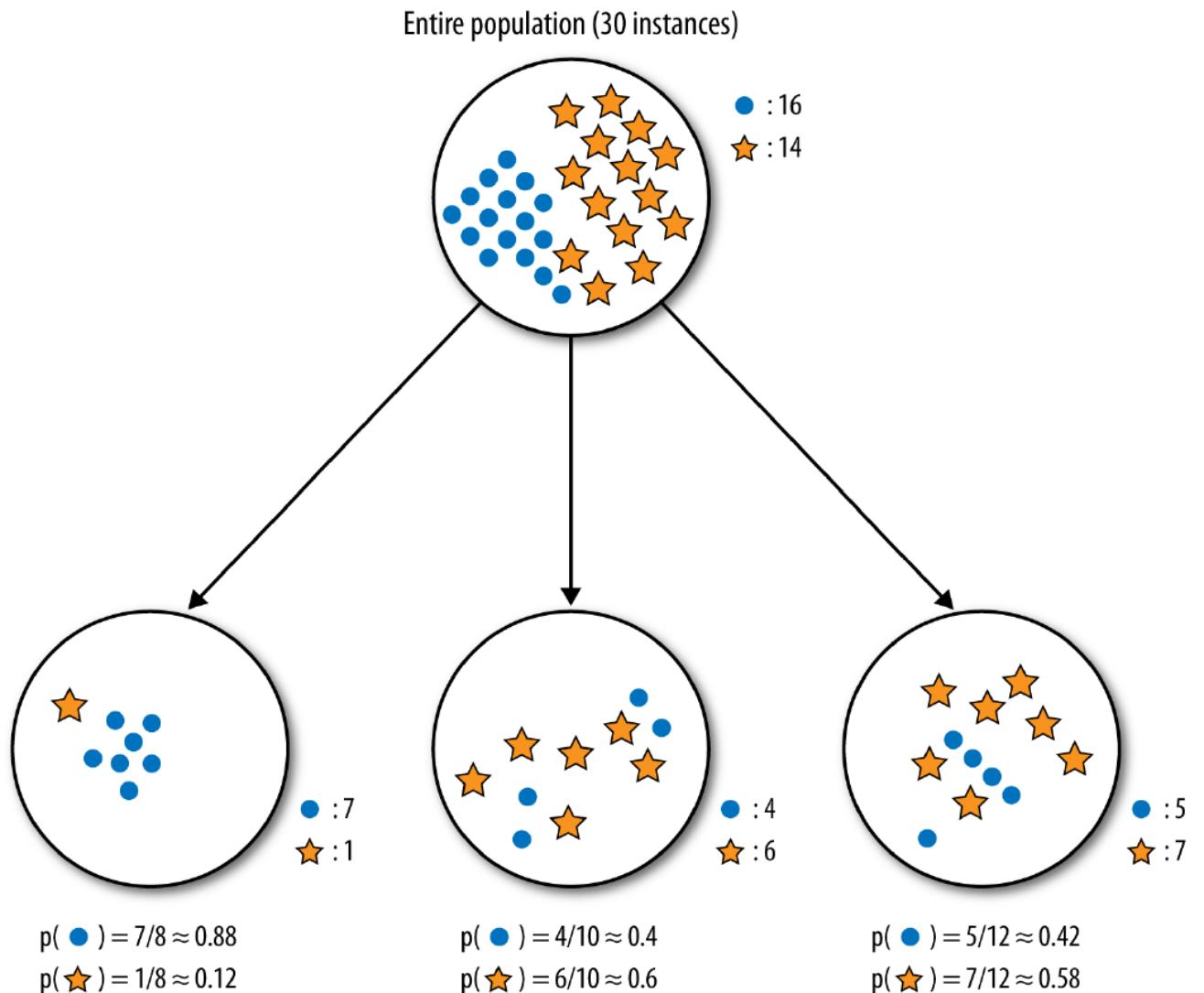


$$\begin{aligned}E(S) &= -p_o \log p_o - p_* \log p_* \\&= -\frac{16}{30} \log \frac{16}{30} - \frac{14}{30} \log \frac{14}{30} \approx 0.99\end{aligned}$$

$$\begin{aligned}E(L) &= -p_o \log p_o - p_* \log p_* \\&= -\frac{12}{13} \log \frac{12}{13} - \frac{1}{13} \log \frac{1}{13} \approx 0.39\end{aligned}$$

$$\begin{aligned}E(R) &= -p_o \log p_o - p_* \log p_* \\&= -\frac{4}{17} \log \frac{4}{17} - \frac{13}{17} \log \frac{13}{17} \approx 0.79\end{aligned}$$

$$\begin{aligned}IG &= E(S) - \frac{1}{30}[q_L E(L) + q_R E(R)] \\&\approx 0.99 - \frac{1}{30}[13(0.39) + 17(0.79)] \\&\approx \mathbf{0.37}\end{aligned}$$



# DECISION TREES STRENGTHS AND LIMITATIONS

**White box** model

Can be used with **incomplete** datasets

**Built-in** feature selection

Makes **no assumption** about

**Not as accurate** as other algorithms (usually)

**Not robust**

Particularly vulnerable to **overfitting**

Optimal decision tree learning is **NP-complete**

Biased towards categorical features with high number of levels

# DECISION TREES NOTES

## Splitting Metrics

- information gain, Gini impurity, variance reduction, etc.

## Common Algorithms

- Iterative Dichotomiser 3, C4.0, C4.5, CHAID, MARS, conditional inference trees, CART

Decision trees can also be combined together using boosting algorithms (**AdaBoost**) or **Random Forests**, providing a type of voting procedure (Ensemble Learning).

## OTHER POINTS TO PONDER

Classification is linked to **probability estimation**

- approaches based on regression models could prove fruitful

**Rare occurrences** (often more interesting or important) continue to plague classification attempts

- historical data at Fukushima's nuclear reactor prior to the meltdown could not have been used to learn about meltdowns

**No Free-Lunch Theorem:** no classifier works best for all data.

With big datasets, algorithms must also consider efficiency.

# PERFORMANCE EVALUATION

Classifiers are evaluated on the testing set.

Ideally, a good classifier would have high rates of both **True Positives** (TP) and **True Negatives** (TN), and low rates of both **False Positives** (FP, Type I error) and **False Negatives** (FN, Type II error).

Evaluation metrics mean very little on their own: context requires comparison with other classifiers, and other evaluation metrics.

		Predicted		Total	79.0%
Actuals	A	54	10		
	B	6	11	17	21.0%
Total	60	21	81		
	74.1%	25.9%			

Classification Rates	
Sensitivity:	0.84
Specificity:	0.65
Precision:	0.90
Negative Predictive Value:	0.52
False Positive Rate:	0.35
False Discovery Rate:	0.10
False Negative Rate:	0.16

Performance Metrics	
Accuracy:	0.80
F1-Score:	0.87
Informedness (ROC):	0.49
Markedness:	0.42
M.C.C.:	0.46
Pearson's chi2:	0.01
Hist. Stat:	0.10

		Predicted		Total	66.7%
Actuals	A	54	0		
	B	16	11	27	33.3%
Total	70	11	81		
	86.4%	13.6%			

Classification Rates	
Sensitivity:	1.00
Specificity:	0.41
Precision:	0.77
Negative Predictive Value:	1.00
False Positive Rate:	0.59
False Discovery Rate:	0.23
False Negative Rate:	0.00

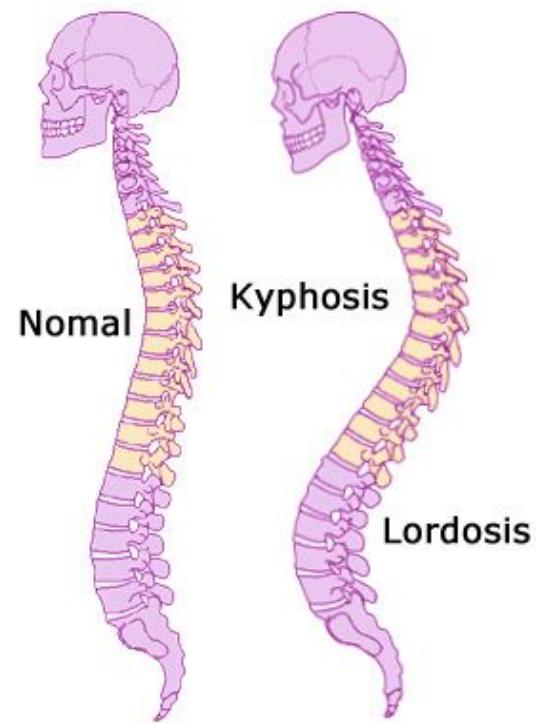
Performance Metrics	
Accuracy:	0.80
F1-Score:	0.87
Informedness (ROC):	0.41
Markedness:	0.77
M.C.C.:	0.56
Pearson's chi2:	0.33
Hist. Stat:	0.40

## EXAMPLE – KYPHOSIS DATASET

Kyphosis is a medical condition related to the excessive convex curvature of the spine. Corrective spinal surgery is at times performed on children.

The dataset has 81 observations and 4 attributes:

- **kyphosis** (absent or present after operation)
- **age** (at time of operation, in months)
- **number** (of vertebrae involved)
- **start** (topmost vertebra operated on)

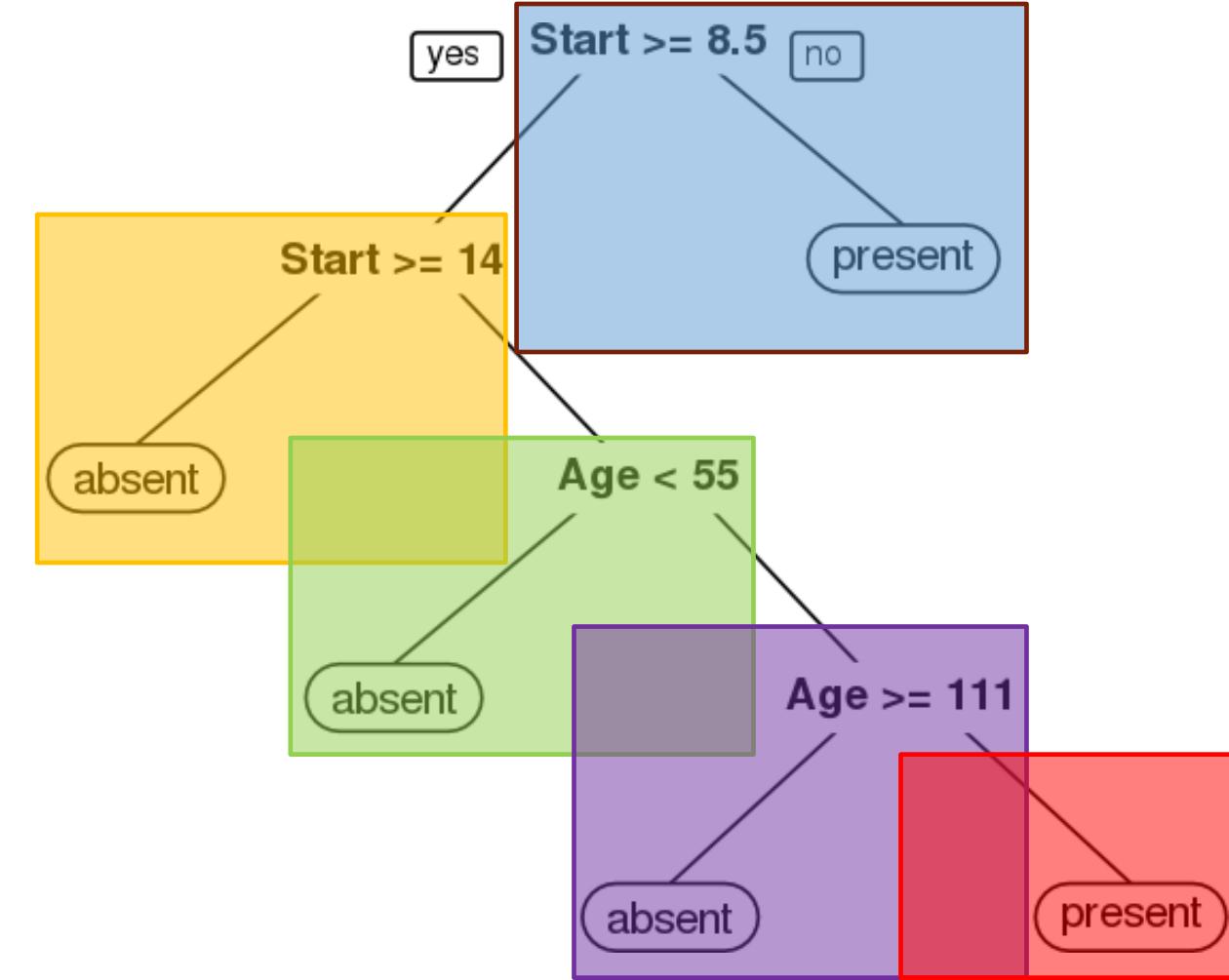
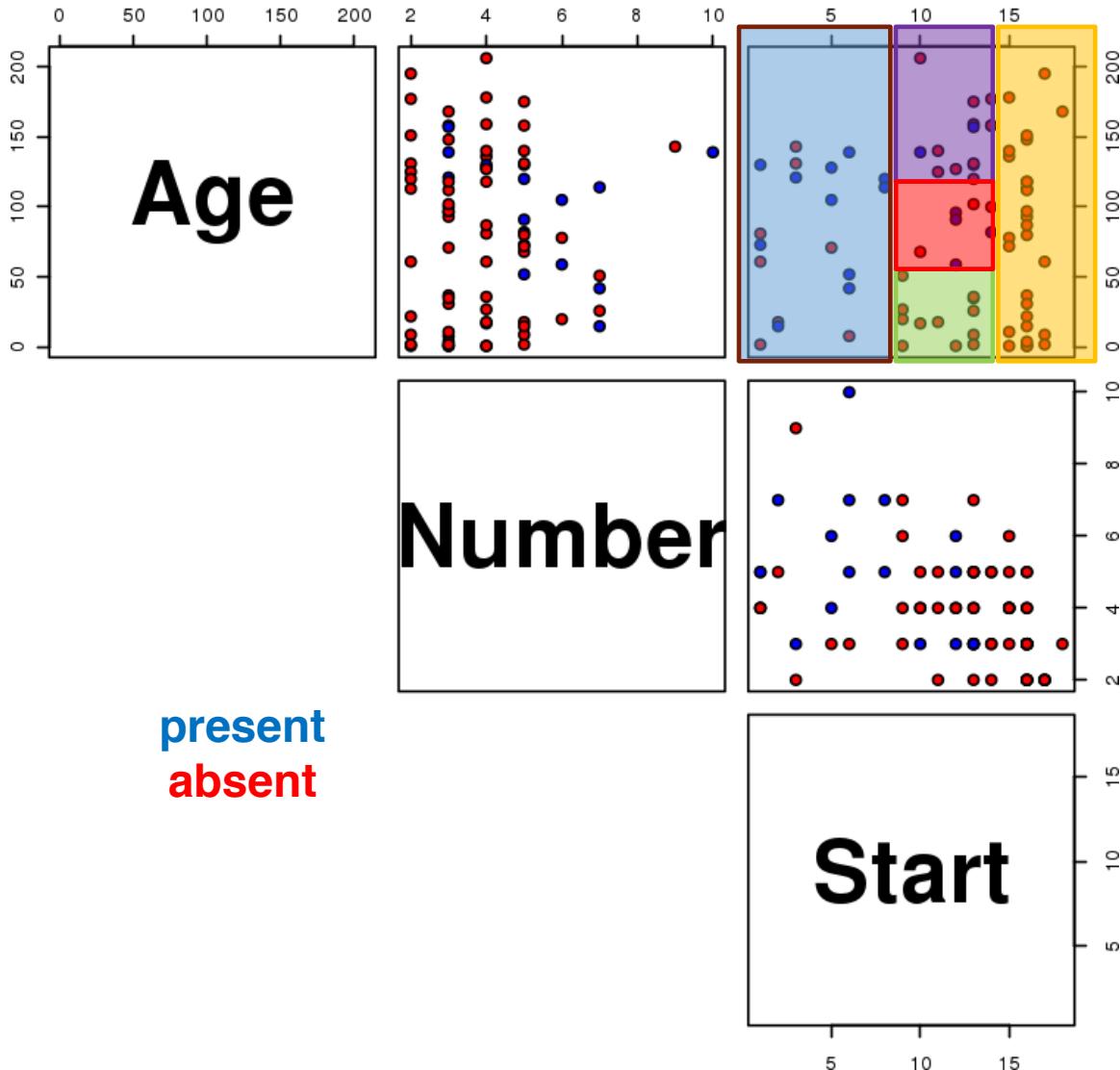


## EXAMPLE – KYPHOSIS DATASET

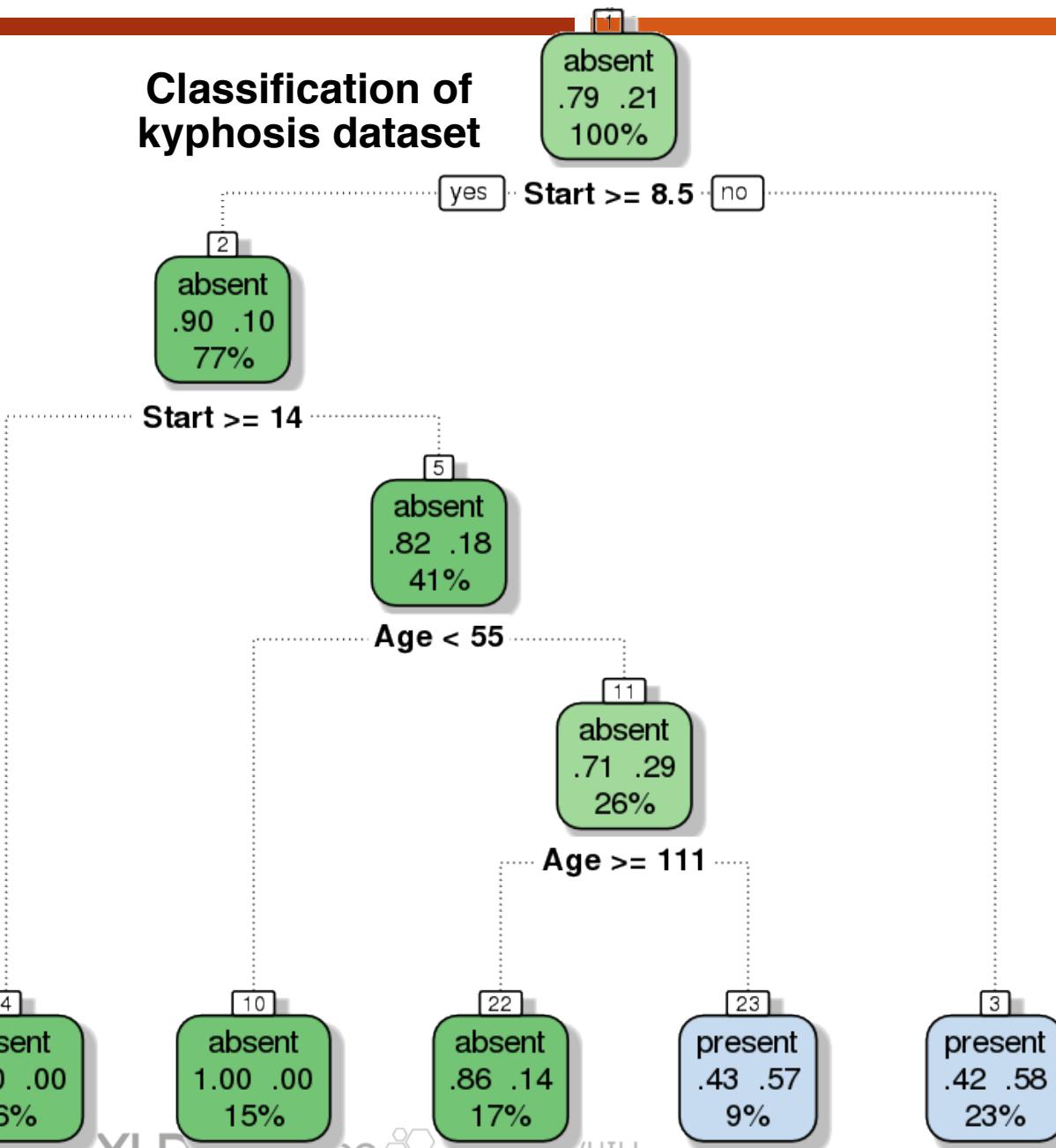
The question of interest for this natural dataset is how the three explanatory attributes might impact the operation's success.

We use the rpart implementation of CART to generate candidate decision trees.

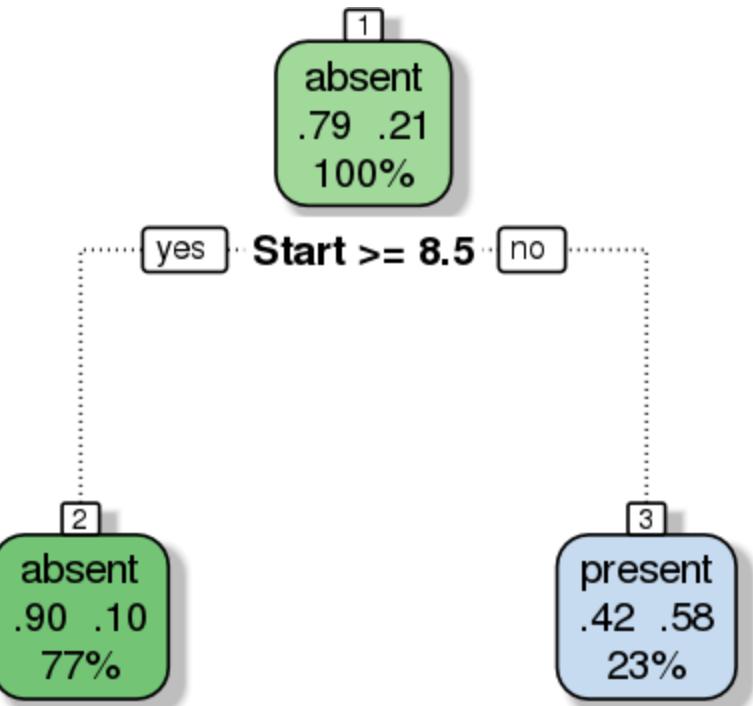
Strictly speaking, this is not a predictive supervised task as we treat the entire dataset as a training set (there are no hold-out testing observations for the time being).



## Classification of kypnosis dataset



## Pruned classification of kypnosis dataset



## EXAMPLE – KYPHOSIS DATASET

We train a model on 50 observations (selected randomly) and evaluate the performance on the remaining 31 observations.

		Predicted		Total	83.9%
		A	B		
Actuals	A	23	3	26	83.9%
	B	3	2	5	16.1%
Total		26	5	31	
		83.9%	16.1%		

Classification Rates		Performance Metrics
	Sensitivity: 0.88	Accuracy: 0.81
	Specificity: 0.40	F1-Score: 0.88
	Precision: 0.88	Informedness (ROC): 0.28
	Negative Predictive Value: 0.40	Markedness: 0.28
	False Positive Rate: 0.60	M.C.C.: 0.28
	False Discovery Rate: 0.12	Pearson's chi2: 0.00
	False Negative Rate: 0.12	Hist. Stat: 0.00

Is this a good model?

## PERFORMANCE EVALUATION

**In both the categorical and numerical estimation problem,** isolated performance metric does not provide enough of a rationale for model validation, unless it has first been normalized.

There is (a lot) more to be said on the topic of model selection.