

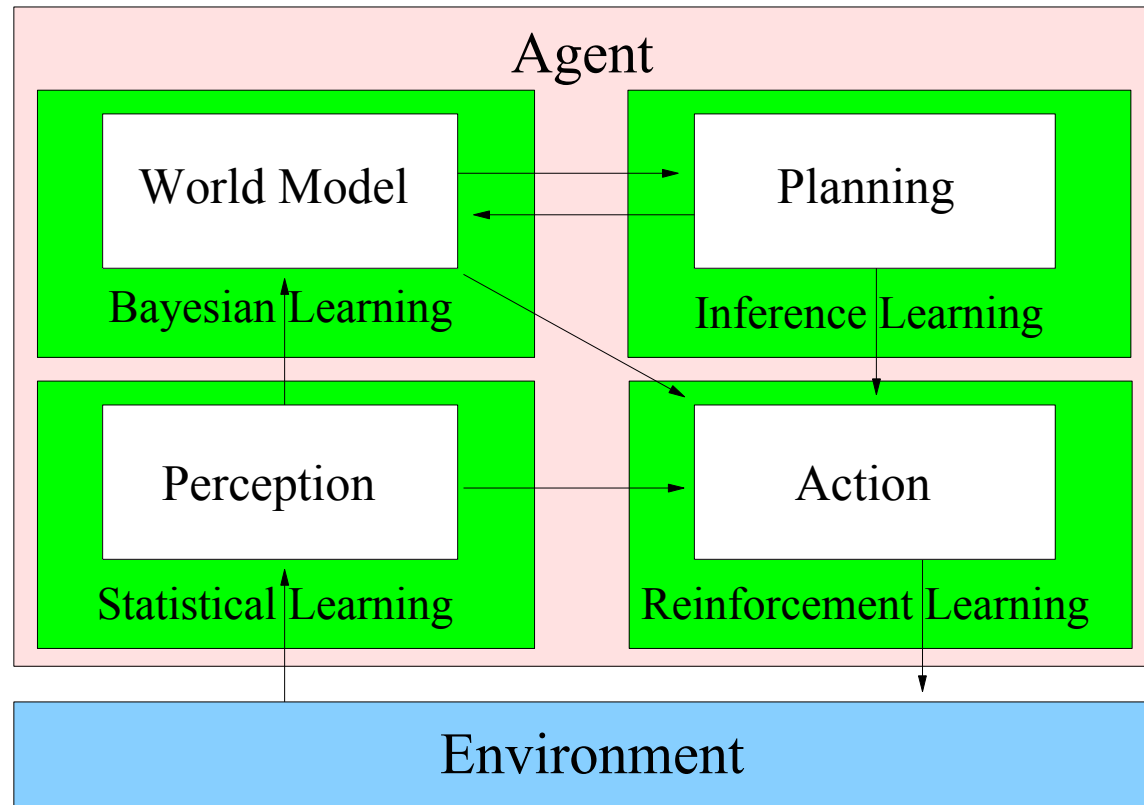
Learning and Decision Trees

COMP3411/9814: Artificial Intelligence

Learning Agents

- Improve performance on future tasks after making observations about the world
- Design of a learning element is affected by
 - Which components of the performance element are to be learned
 - What feedback is available to learn these components
 - What representation is used for the components

Learning Agent



Machine Learning Applications

Data:

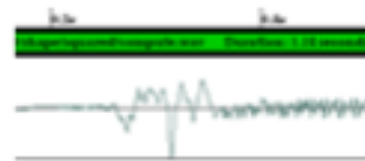
Antenatal aspect	Antenatal aspect	Antenatal aspect
Age 35	Age 35	Age 35
First Pregnancy: no	First Pregnancy: no	First Pregnancy: no
Gravida: 45	Gravida: 45	Gravida: 45
Parity: 10	Parity: 10	Parity: 10
Previous Cesarean Section: no	Previous Cesarean Section: no	Previous Cesarean Section: no
Obstetrical History: 1	Obstetrical History: 1	Obstetrical History: 1
Current Gestation: 3	Current Gestation: 3	Current Gestation: 3
Emergency C-Section: 1	Emergency C-Section: 1	Emergency C-Section: 1

One of 18 learned rules:

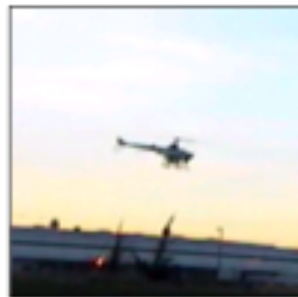
If No previous vaginal delivery, and Abnormal 2nd Trimester Ultrasound, and Malpresentation at admission
Then Probability of Emergency C-Section is 0.6

Over training data: 36/41 = .63,
Over test data: 12/20 = .60

Mining Databases



Speech Recognition



Control learning



Object recognition

Text analysis

Peter H. van Oppen - **Chairman of the Board & Chief Executive Officer**
Mr. van Oppen has served as **Chairman of the Board and Chief Executive Officer of ADG** since its acquisition by Interpoint in 1994 and a **director of ADG** since 1986. Until its acquisition by Crane Co. in October 1996, **Mr. van Oppen** served as **Chairman of the Board, President, Managing and Chief Executive Officer of Interpoint**. Prior to 1985, **Mr. van Oppen** worked as a **consulting manager** at **Price Waterhouse LLP** and at Bain & Company in Boston and London. He has additional experience in medical electronics and venture capital. **Mr. van Oppen** also serves as a **director of Acute Care Medical, Inc.** and **Spacelabs Medical, Inc.**. He holds a B.A. from Whitman College and an M.B.A. from Harvard Business School, where he was a **Baker Scholar**.

Supervised Learning

- Given a **training set** and a **test set**
 - each has set of examples
 - example has **attributes** and a **class, or target value**
 $\langle x_1, x_2, \dots, x_n, y \rangle$ y may be discrete or continuous
- Agent given attributes and class for each example in training set
➡ Try predict class of each example in test set
- Many supervised learning methods, e.g.:
 - Decision Tree
 - Neural Network
 - Support Vector Machine, etc.

Supervised Learning

- Training set:
 - Examples may be presented all at once (batch) or in sequence (online)
 - Examples may be presented at random or in time order (stream)
 - Learner **cannot** use the test set **at all** in defining the model
- Model is evaluated on predicting output for each example in **test set**

Supervised Learning Methodology

1. “Feature engineering” – select relevant features
2. Choose representation of input features and outputs
3. Pre-process to extract features from raw data
4. Choose learning method(s) to evaluate
5. Choose training regime (including parameters)
6. Evaluation
 1. Choose baseline for comparison
 2. Choose type of internal validation, e.g. cross-validation
 3. Sanity check results with human expertise, other benchmark

Supervised Learning Issues

- Framework (decision tree, neural network, SVM, etc.)
 - representation (of inputs and outputs)
 - pre-processing / post-processing
 - training method (perceptron, backpropagation, etc.)
 - generalisation (avoid over-fitting)
 - evaluation (separate training and testing sets)

Inductive Learning

Simplest form: learn a function from examples

f is the **target function**

An **example** is a pair $(x, f(x))$

Problem: find a **hypothesis** h

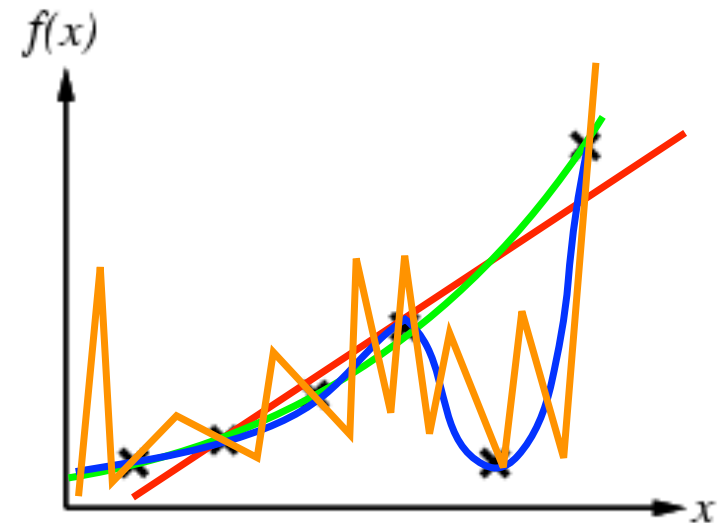
such that $h \approx f$

given a **training set** of examples

- Ignores prior knowledge
 - Assumes examples are given

Inductive Learning Method

- Construct h to agree with f on training set
- (h is **consistent** if it agrees with f in all examples
- E.g., curve fitting
 - Which curve best fits data?



Ockham's razor

“The most likely hypothesis is the simplest one consistent with the data.”

- Occam's razor (sometimes spelled Ockham's razor) is a principle attributed to the 14th-century English logician and Franciscan friar William of Ockham.

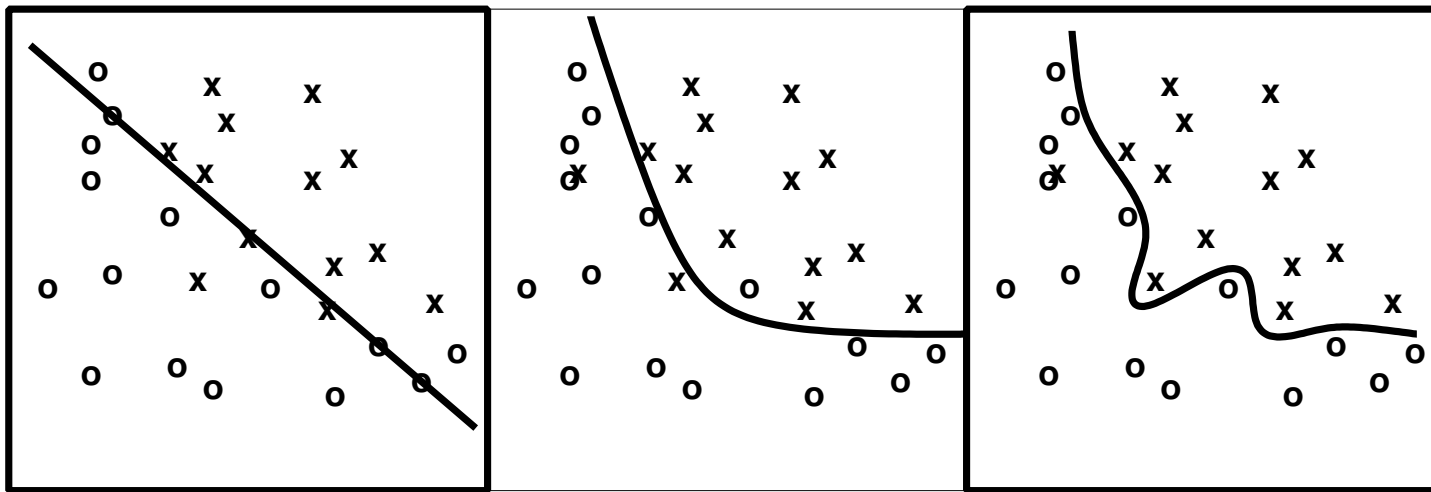
The principle states that the explanation of any phenomenon should make as few assumptions as possible, eliminating those that make no difference in the observable predictions of the explanatory hypothesis or theory.

This is often paraphrased as "All other things being equal, the simplest solution is the best."

Prefer the simplest hypothesis that fits the data

Ockham's razor

“The most likely hypothesis is the simplest one consistent with the data.”



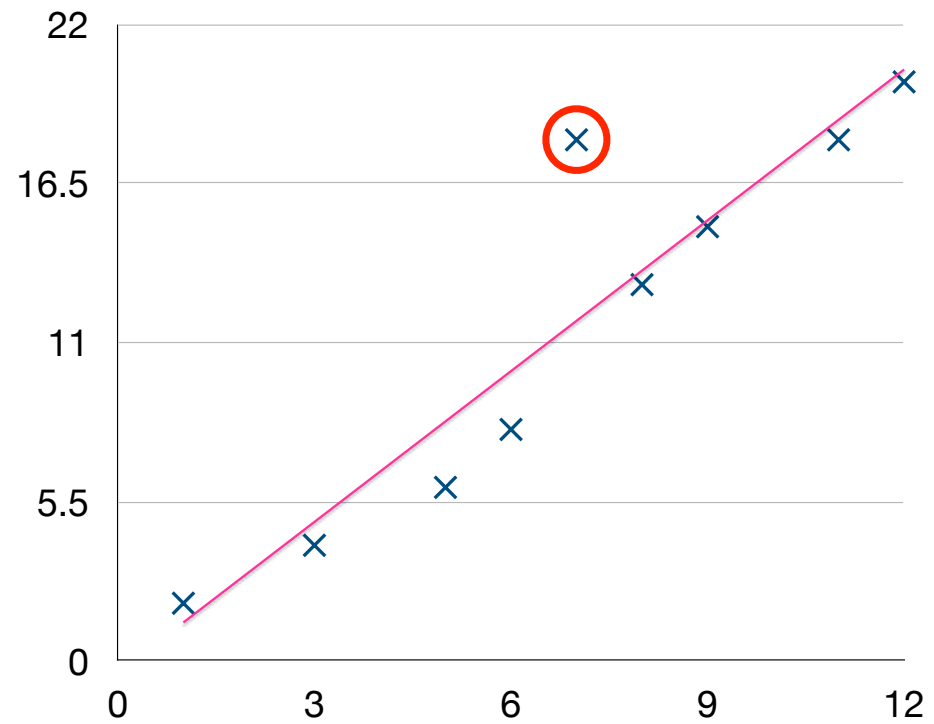
inadequate

good compromise

over-fitting

- Can have **noise** in data
- Need to tradeoff simplicity and accuracy

Outliers



Supervised Learning (again)

Given:

- a set of **training examples**, each with a set of **features** and **target value (or class)**:

$$\langle x_1, x_2, \dots, x_n, y \rangle$$

- a new example, where only the values for the input features are given

predict the target value of the new example.

Learning, Generalisation and Compression

- Could learn by memorising everything
 - Inefficient
 - Can't extract patterns and insights from what is memorised
- Learning usually involved *generalisation*
 - Being able to describe many examples in a simple, general form
 - Many principles shared with compression
 - I.e. learned concept is a compression of original data

Decision Tree Learning

Learning decision trees

Problem: decide whether to wait for a table at a restaurant, based on the following attributes:

1. **Alternate**: is there an alternative restaurant nearby?
2. **Bar**: is there a comfortable bar area to wait in?
3. **Fri/Sat**: is today Friday or Saturday?
4. **Hungry**: are we hungry?
5. **Patrons**: number of people in the restaurant (None, Some, Full)
6. **Price**: price range (\$, \$\$, \$\$\$)
7. **Raining**: is it raining outside?
8. **Reservation**: have we made a reservation?
9. **Type**: kind of restaurant (French, Italian, Thai, Burger)
10. **WaitEstimate**: estimated waiting time (0-10, 10-30, 30-60, >60)

Restaurant Training Data

- Examples described by **attribute values** (Boolean, discrete, continuous)
 - E.g., situations where **I will/won't wait for a table**:

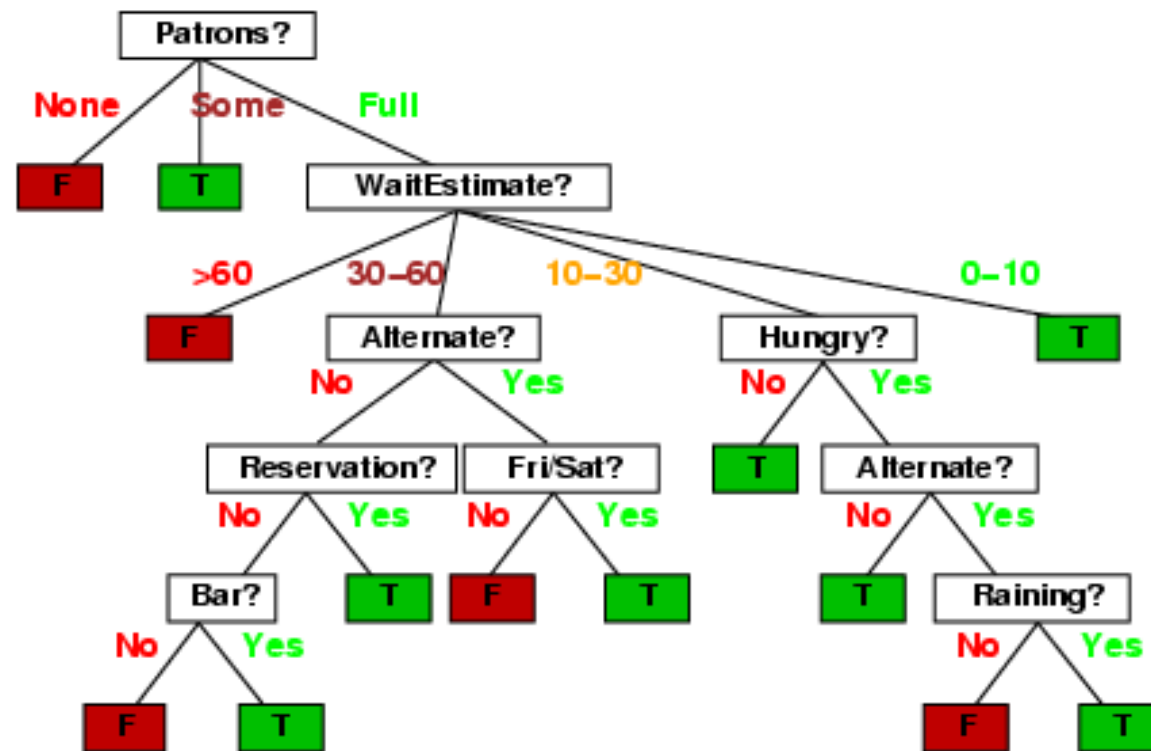
	Alt	Bar	F/S	Hun	Pat	Price	Rain	Res	Type	Est	Wait?
X ₁	T	F	F	T	Some	\$\$\$	F	T	French	0-10	T
X ₂	T	F	F	T	Full	\$	F	F	Thai	30-60	F
X ₃	F	T	F	F	Some	\$	F	F	Burger	0-10	T
X ₄	T	F	T	T	Full	\$	F	F	Thai	10-30	T
X ₅	T	F	T	F	Full	\$\$\$	F	T	French	> 60	F
X ₆	F	T	F	T	Some	\$\$	T	T	Italian	0-10	T
X ₇	F	T	F	F	None	\$	T	F	Burger	0-10	F
X ₈	F	F	F	T	Some	\$\$	T	T	Thai	0-10	T
X ₉	F	T	T	F	Full	\$	T	F	Burger	> 60	F
X ₁₀	T	T	T	T	Full	\$\$\$	F	T	Italian	10-30	F
X ₁₁	F	F	F	F	None	\$	F	F	Thai	0-10	F
X ₁₂	T	T	T	T	Full	\$	F	F	Burger	30-60	T

class
value



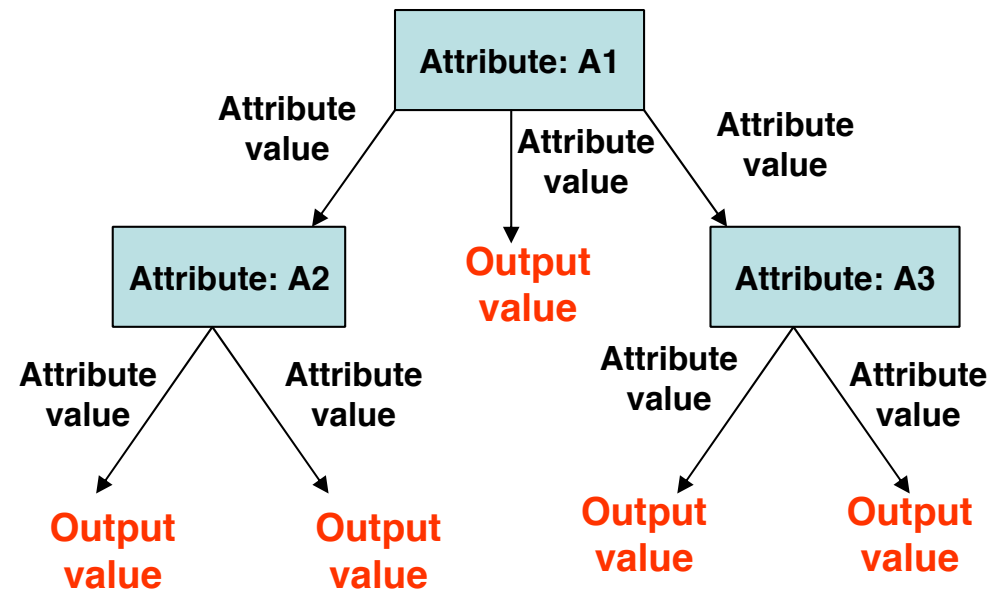
Decision Trees

- One possible representation for hypotheses



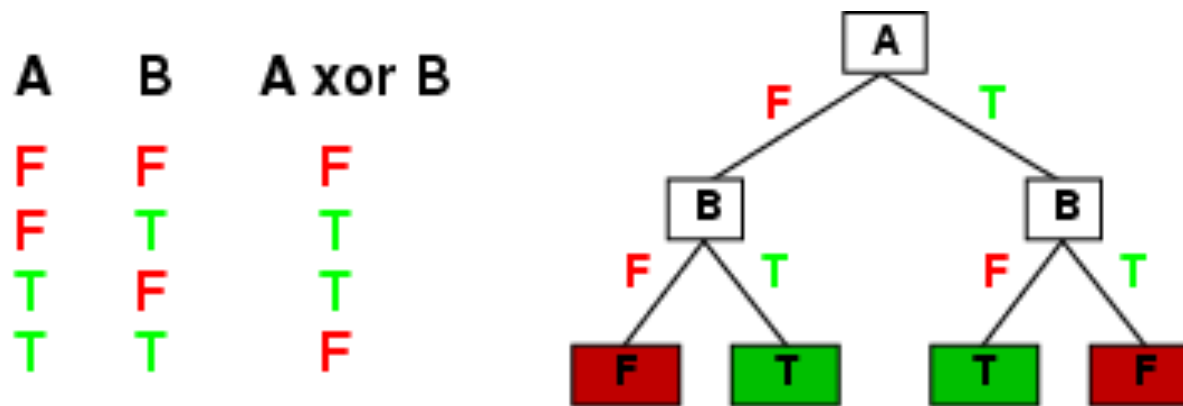
Decision Trees

- Output can be multi-valued, not just binary



Expressiveness

- Decision trees can express any function of the input attributes.
 - E.g., for Boolean functions, truth table row \rightarrow path to leaf:



- There is a consistent decision tree for any training set with one path to leaf for each example (unless f nondeterministic in x) but it probably won't generalise to new examples
- Prefer to find more **compact** decision trees

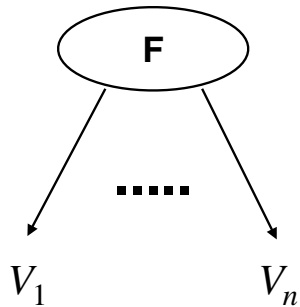
Decision Tree Learning Algorithms

- ID3 Algorithm (Quinlan 1986) and its successors C4.5 and C5.0
 - <http://www.rulequest.com/>
- *Employs a top-down induction*



- Greedy search the space of possible decision trees.
- The algorithm never backtracks to reconsider earlier choices.

ID3 (Quinlan)



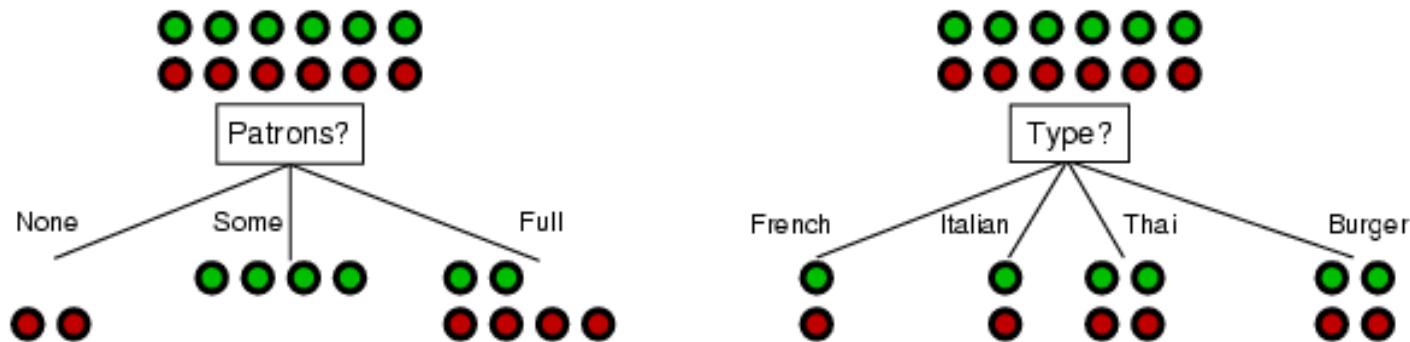
- The algorithm operates over a set of training instances, C .
- If all instances in C are in class P , create a node P and stop, otherwise select a feature F and create a decision node.
- Partition the training instances in C into subsets according to the values of V .
- Apply the algorithm recursively to each of the subsets C_i

Generalisation

- Training data must not be **inconsistent**
 - see later how to handle inconsistent data
- Can split attributes in any order and still produce a tree that correctly classifies all examples in training set
- However, want a tree that is likely to **generalise** to correctly classify (unseen) examples in **test set**.
- Prefer **simpler** hypothesis, i.e. a smaller tree.
- How can we choose attributes to produce a small tree?

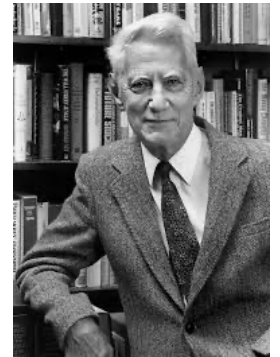
Choosing an attribute

	Alt	Bar	F/S	Hun	Rat	Price	Rain	Res	Type	Est	Wait?
X ₁	T	F	F	T	Some	\$\$\$	F	T	French	0-10	T
X ₂	T	F	F	T	Full	\$	F	F	Thai	30-60	F
X ₃	F	T	F	F	Some	\$	F	F	Burger	0-10	T
X ₄	T	F	T	T	Full	\$	F	F	Thai	10-30	T
X ₅	T	F	T	F	Full	\$\$\$	F	T	French	>60	F
X ₆	F	T	F	T	Some	\$\$	T	T	Italian	0-10	T
X ₇	F	T	F	F	None	\$	T	F	Burger	0-10	F
X ₈	F	F	F	T	Some	\$\$	T	T	Thai	0-10	T
X ₉	F	T	T	F	Full	\$	T	F	Burger	>60	F
X ₁₀	T	T	T	T	Full	\$\$\$	F	T	Italian	10-30	F
X ₁₁	F	F	F	F	None	\$	F	F	Thai	0-10	F
X ₁₂	T	T	T	T	Full	\$	F	F	Burger	30-60	T



- Patrons is a more “informative” attribute than Type, because it splits the examples more nearly into sets that are “all positive” or “all negative”.
- “Informativeness” can be quantified using the mathematical concept of **“entropy”**.
- Tree can be built by minimising entropy at each step

Information Gain



Claude Shannon
(1916 – 2001)

Information gain is based on information theory concept called *Entropy*

In information theory, the **Shannon entropy or information entropy** is a measure of the **uncertainty** associated with a random variable.

- It quantifies the information contained in a message, usually in bits or bits/symbol.
- It is the minimum message length necessary to communicate information.

Entropy

- Entropy is a measure of how much information we **gain** when the class value is revealed to us.
- In decision tree learning, when we partition a dataset by a particular attribute the resulting entropy is lower.
 - An entropy of 0 means all the examples have the same class value
 - An entropy of 1 means that they are randomly distributed
- If the prior probabilities of the n class values are p_1, \dots, p_n then the entropy is

$$H(\langle p_1, \dots, p_n \rangle) = \sum_{i=1} -p_i \log_2 p_i$$

Entropy and Huffman Coding

- Entropy is the number of bits per symbol achieved by a (block) Huffman Coding scheme.
- Suppose we want to encode, into a bit string a long message composed of the two letters A and B , which occur with equal frequency.
- Can be done by assigning $A = 0, B = 1$.
 - I.e., **one bit** (binary digit) is needed to encode each letter.
- Example 1: $H(\langle 0.5, 0.5 \rangle) = 1$ bit

Entropy and Huffman Coding

- Suppose we need to encode a message consisting of the letters A, B and C, and that B and C occur equally often but A occurs twice as often as the other two letters.
- In this case, the most efficient code would be

$$A=0, B=10, C=11.$$

- The average number of bits needed to encode each letter is 1.5 .
- Example 2: $H(\langle 0.5, 0.25, 0.25 \rangle) = 1.5$ bit

Entropy and Huffman Coding

Example 3:

Want a code to distinguish elements of {a, b, c, d} with

$$P(a) = \frac{1}{2}, P(b) = \frac{1}{4}, P(c) = \frac{1}{8}, P(d) = \frac{1}{8}$$

Consider the code:

a	0	b	10	c	110	d	111
---	---	---	----	---	-----	---	-----

Code uses 1 to 3 bits. On average, it uses

$$P(a) \times 1 + P(b) \times 2 + P(c) \times 3 + P(d) \times 3 = \frac{1}{2} + \frac{2}{4} + \frac{3}{8} + \frac{3}{8} = 1\frac{3}{4} \text{ bits}$$

Entropy and Huffman Coding

Information theory

- A **bit** is a binary digit.
- 1 bit can distinguish 2 items
- k bits can distinguish 2^k items
- n items can be distinguished using $\log_2 n$ bits

Entropy and Huffman Coding

If the letters occur in some other proportion, we would need to “block” them together in order to encode them efficiently. But, **the average number of bits required** by the most efficient coding scheme is given by

$$H(\langle p_1, \dots, p_n \rangle) = \sum_{n=1} -p_i \log_2 p_i$$

Entropy

- Suppose we have p positive and n negative examples in a node.
- $H\left(\left\langle \frac{p}{p+n}, \frac{n}{p+n} \right\rangle\right)$ bits needed to encode a new example.
 - e.g. for 12 restaurant examples, $p = n = 6$ so we need 1 bit.
- An attribute splits the examples E into subsets E_i , each of which (we hope) needs less information to complete the classification.

- Let E_i have p_i positive and n_i negative examples
 $H\left(\left\langle \frac{p_i}{p_i+n_i}, \frac{n_i}{p_i+n_i} \right\rangle\right)$ bits needed to encode an example

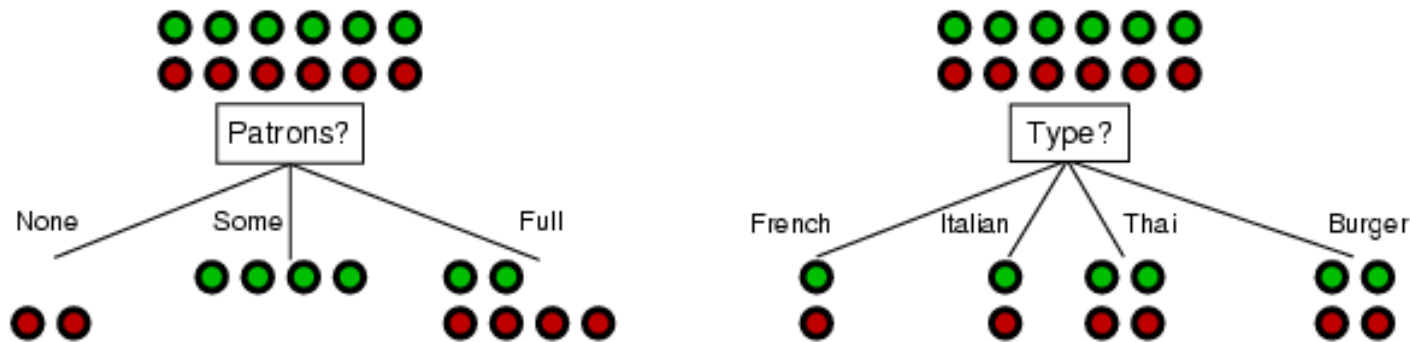
expected number of bits per example over all branches is

$$\sum_i \frac{p_i + n_i}{p + n} H\left(\left\langle \frac{p_i}{p_i + n}, \frac{n_i}{p_i + n_i} \right\rangle\right)$$

Probability of going
down branch i

- For **Patrons**, this is 0.459 bits, for **Type** this is (still) 1 bit → split on **Patrons**

Choosing and Attribute

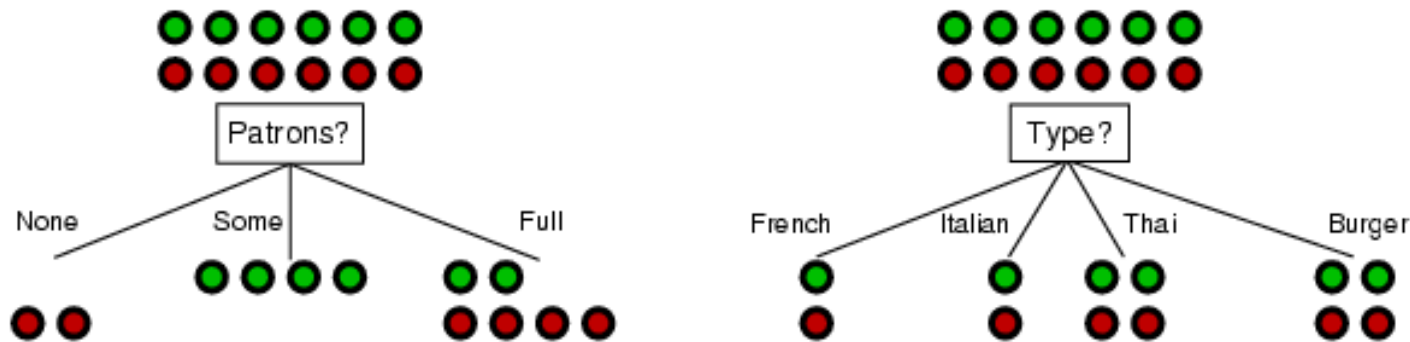


$$\begin{aligned}
 \text{For Patrons, Entropy} &= \frac{1}{6}(0) + \frac{1}{3}(0) + \frac{1}{2} \left[-\frac{1}{3} \log\left(\frac{1}{3}\right) - \frac{2}{3} \log\left(\frac{2}{3}\right) \right] \\
 &= 0 + 0 + \frac{1}{2} \left[\frac{1}{3}(1.585) + \frac{2}{3}(0.585) \right] = 0.459
 \end{aligned}$$

$$\text{For Type, Entropy} = \frac{1}{6}(1) + \frac{1}{6}(1) + \frac{1}{3}(1) + \frac{1}{3}(1) = 1$$

$$\begin{aligned}
 \text{Prior}(\text{none}) &= \frac{2}{12} = \frac{1}{6} \\
 P_{\text{wait}}(\text{false}) &= 1 \\
 -P_{\text{wait}}(\text{false}) \log_2 P_{\text{wait}}(\text{false}) &= 1 \times 0 = 0
 \end{aligned}$$

Choosing and Attribute

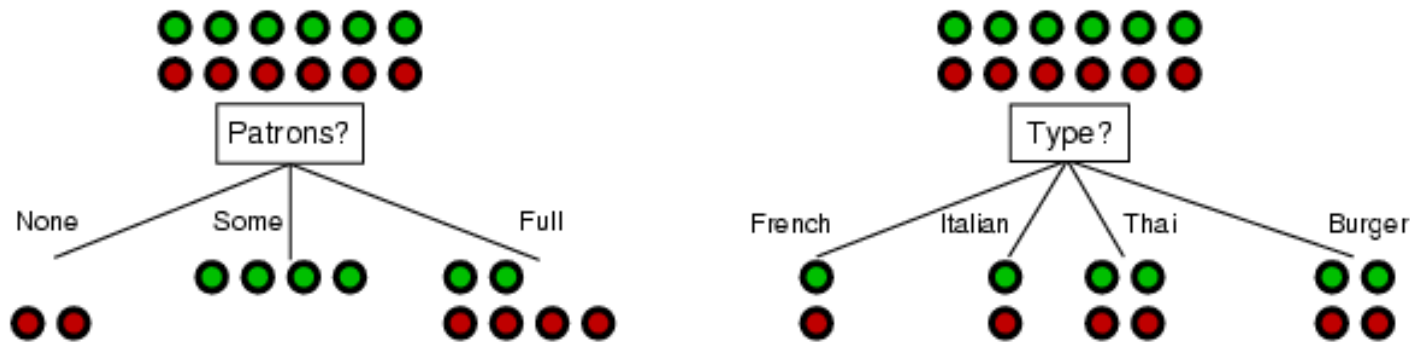


$$\begin{aligned}
 \text{For Patrons, Entropy} &= \frac{1}{6}(0) + \frac{1}{3}(0) + \frac{1}{2} \left[-\frac{1}{3} \log\left(\frac{1}{3}\right) - \frac{2}{3} \log\left(\frac{2}{3}\right) \right] \\
 &= 0 + 0 + \frac{1}{2} \left[\frac{1}{3}(1.585) + \frac{2}{3}(0.585) \right] = 0.459
 \end{aligned}$$

$$\text{For Type, Entropy} = \frac{1}{6}(1) + \frac{1}{6}(1) + \frac{1}{3}(1) + \frac{1}{3}(1) = 1$$

$$\begin{aligned}
 \text{Prior}(\text{some}) &= \frac{4}{12} = \frac{1}{3} \\
 P_{\text{wait}}(\text{true}) &= 1 \\
 -P_{\text{wait}}(\text{true}) \log_2 P_{\text{wait}}(\text{true}) &= 1 \times 0 = 0
 \end{aligned}$$

Choosing and Attribute



$$\begin{aligned}
 \text{For Patrons, Entropy} &= \frac{1}{6}(0) + \frac{1}{3}(0) + \frac{1}{2} \left[-\frac{1}{3} \log\left(\frac{1}{3}\right) - \frac{2}{3} \log\left(\frac{2}{3}\right) \right] \\
 &= 0 + 0 + \frac{1}{2} \left[\frac{1}{3}(1.585) + \frac{2}{3}(0.585) \right] = 0.459
 \end{aligned}$$

$$\text{For Type, Entropy} = \frac{1}{6}(1) + \frac{1}{6}(1) + \frac{1}{3}(1) + \frac{1}{3}(1) = 1$$

$$\text{Prior}(\text{full}) = \frac{6}{12} = \frac{1}{2}$$

$$P_{\text{wait}}(\text{false}) = \frac{2}{6} = \frac{1}{3}$$

$$P_{\text{wait}}(\text{true}) = \frac{4}{6} = \frac{2}{3}$$

$$-P_{\text{wait}}(\text{false}) \log_2 P_{\text{wait}}(\text{false}) = \frac{1}{3} \log\left(\frac{1}{3}\right)$$

$$-P_{\text{wait}}(\text{true}) \log_2 P_{\text{wait}}(\text{true}) = \frac{2}{3} \log\left(\frac{2}{3}\right)$$

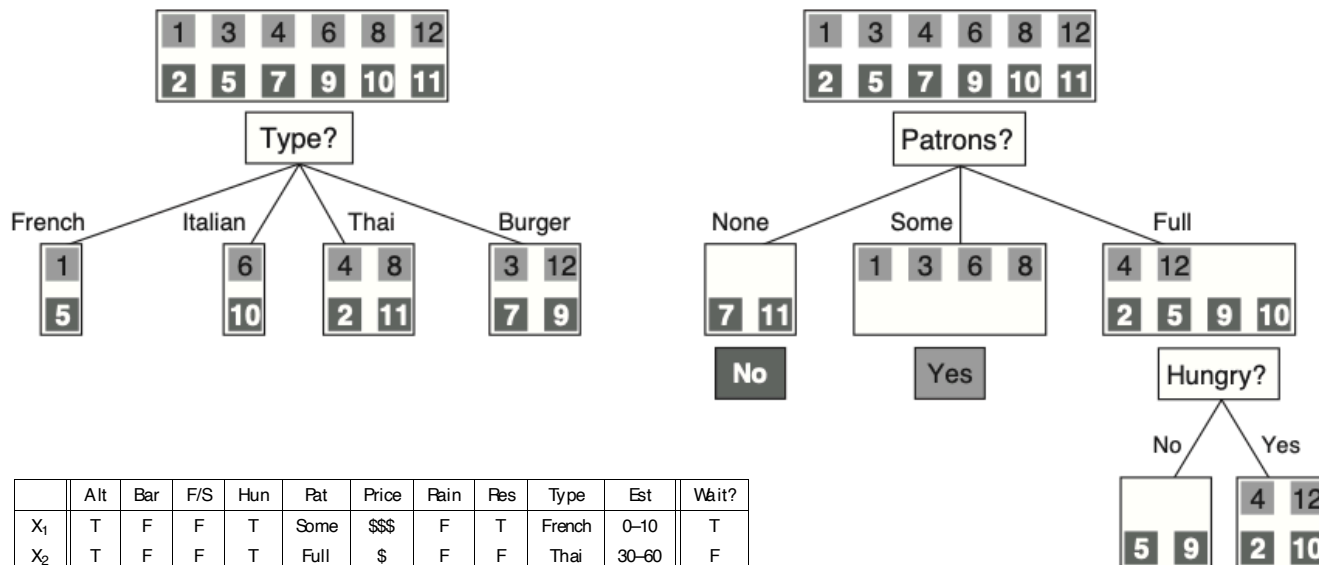
Entropy

- If the prior probabilities of n attribute values are p_1, \dots, p_n , then the entropy of the distribution is

$$H(\langle p_1, \dots, p_n \rangle) = \sum_{i=1}^n -p_i \log_2 p_i$$

- Entropy is a measure of “randomness” (lack of uniformity)
 - Related to prior distribution of some random variable
 - Higher entropy means more randomness
 - “Information” (about distribution) reduces entropy
- Split based on information gain
 - Loss of entropy based on “communicating” value of attribute
 - Related to Shannon’s information theory
 - Measure information gain in bits

Choosing Next Attribute

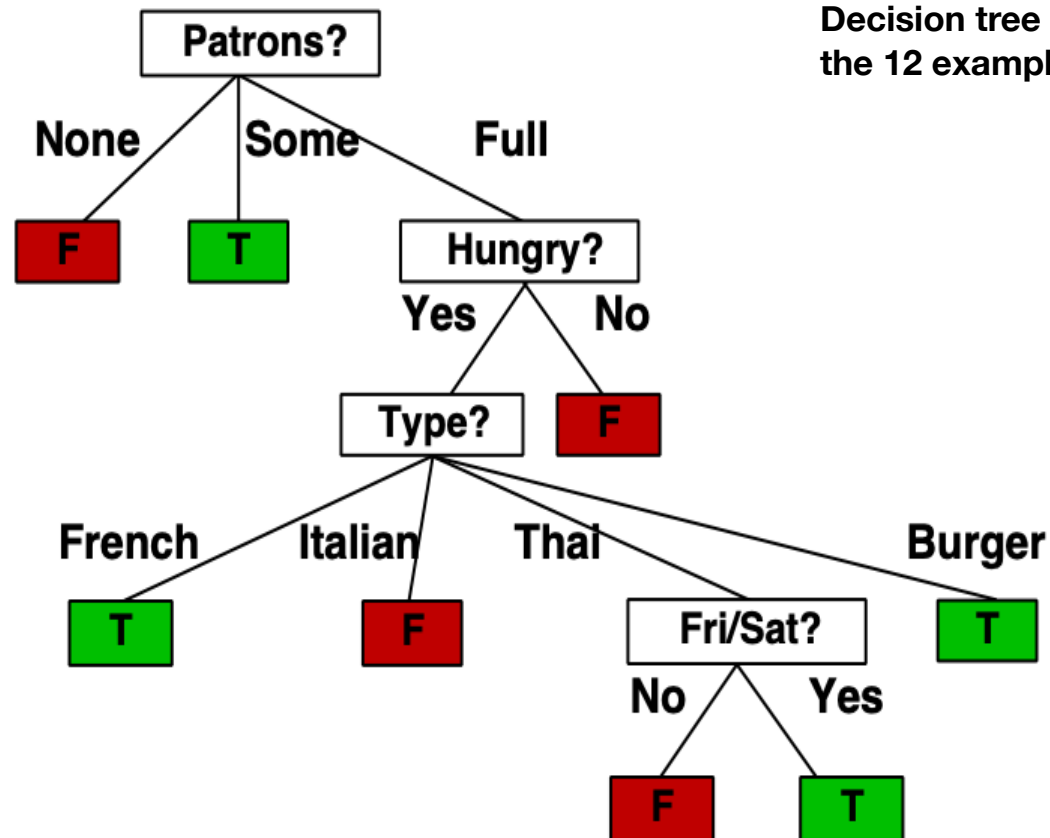


	Alt	Bar	F/S	Hun	Pat	Price	Rain	Res	Type	Est	Wait?
X ₁	T	F	F	T	Some	\$\$\$	F	T	French	0-10	T
X ₂	T	F	F	T	Full	\$	F	F	Thai	30-60	F
X ₃	F	T	F	F	Some	\$	F	F	Burger	0-10	T
X ₄	T	F	T	T	Full	\$	F	F	Thai	10-30	T
X ₅	T	F	T	F	Full	\$\$\$	F	T	French	>60	F
X ₆	F	T	F	T	Some	\$\$	T	T	Italian	0-10	T
X ₇	F	T	F	F	None	\$	T	F	Burger	0-10	F
X ₈	F	F	F	T	Some	\$\$	T	T	Thai	0-10	T
X ₉	F	T	T	F	Full	\$	T	F	Burger	>60	F
X ₁₀	T	T	T	T	Full	\$\$\$	F	T	Italian	10-30	F
X ₁₁	F	F	F	F	None	\$	F	F	Thai	0-10	F
X ₁₂	T	T	T	T	Full	\$	F	F	Burger	30-60	T

$$\sum_i \frac{p_i + n_i}{p + n} H\left(\left\langle \frac{p_i}{p_i + n}, \frac{n_i}{p_i + n_i} \right\rangle\right)$$

After splitting on Patrons,
split on Hungry

Induced Tree

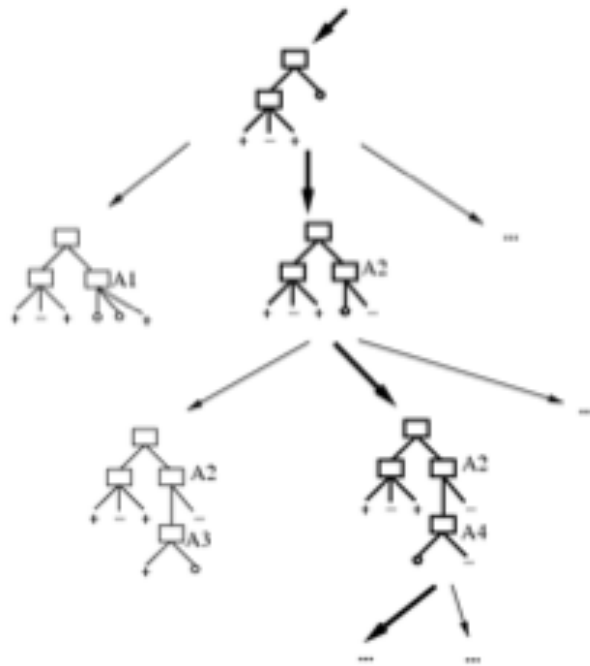


Decision tree learned from the 12 examples.

Decision Trees

- Decision tree learning is a method for [approximating](#) discrete value target functions, in which the learned function is represented by a decision tree.
- Decision trees can also be represented by if-then-else rule
- Decision tree learning is one of the most widely used approach for inductive inference

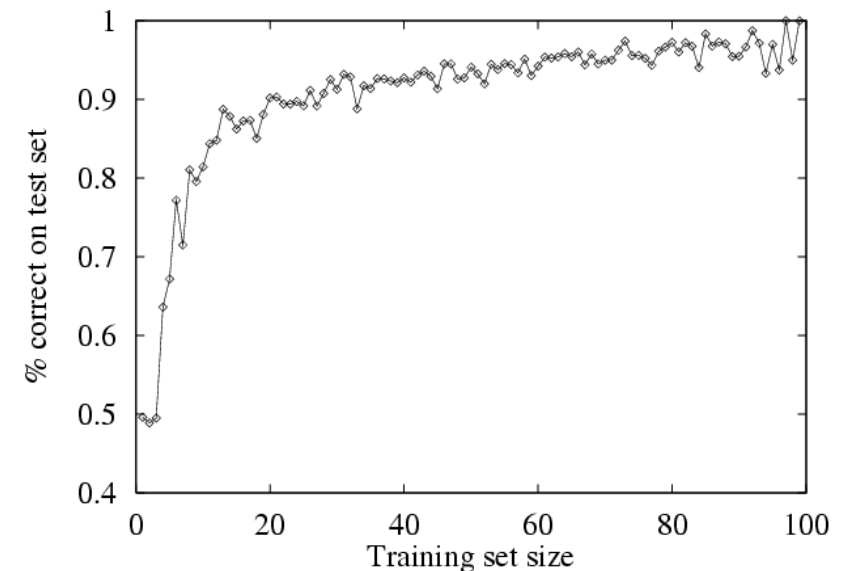
Which Tree Should We Output?



- Decision tree learning performs heuristic search through space of decision trees
- Stops at smallest acceptable tree
- Occam's razor: prefer simplest hypothesis that fits data

Performance measurement

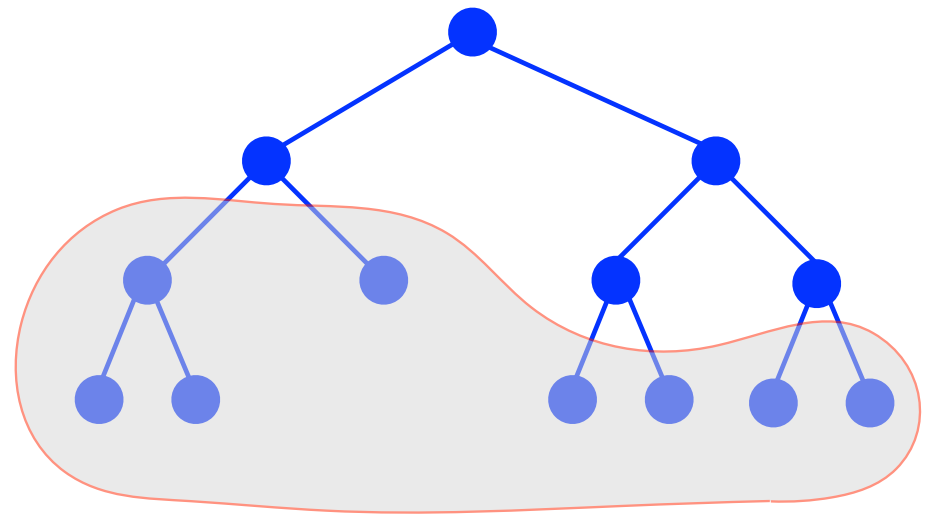
- How do we know that $h \approx f$?
 1. Use theorems of computational/statistical learning theory
 2. Try h on a new **test set** of examples
 - (use **same distribution** over example space as training set)
- **Learning curve** = % correct on test set as a function of training set size



h = hypothesis
 f = target function

Overfitting

- What if training data are noisy?
 - Misclassified examples
 - Incorrect measurements
- Making decision tree classify every example (including noise) could make it less accurate
 - Tries to classify bad examples
 - ➡ Misclassifies correct examples in test set
- This is called *overfitting*
- Can improve accuracy by pruning branches created by try to classify noise.



Tree Pruning To Minimise Error

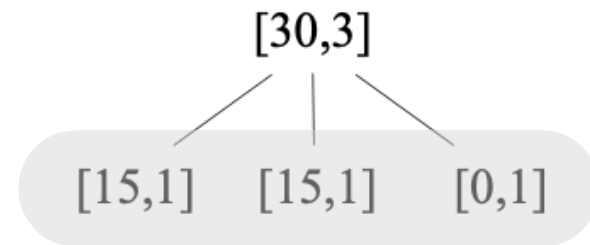
- Maximise expected accuracy, minimise expected probability of error
- For given tree, estimate its expected error
- Also estimate errors for variously pruned tree
- Choose tree with minimal estimated error

How to estimate error?

- One possibility: use “pruning set”
- Another possibility: estimate error probability from data in tree

Tree Pruning

- The top node of this tree has
 - 30 +ve examples
 - 3 -ve examples



- It is split 3 ways.
- If we decide to prune it, the top node becomes a leaf node and the decision value is the majority class, i.e. +ve.
- Only prune if the expected error is less than the expected error with the children

Laplace Error and Pruning

When a node becomes a leaf, all items will be assigned to the majority class at that node. We can estimate the error rate on the (unseen) test items

using the [Laplace error](#):

$$E = 1 - \frac{n + 1}{N + k}$$

N = total number of training examples

n = number of training examples in the majority class

k = number of classes

If the average Laplace error of the children exceeds that of the parent node, we prune off the children.

How do we get the Laplace Error?

- Suppose a node as 99 +ve and 1 -ve: $\frac{n}{N} = \frac{99}{100}$ is the probability of the majority class.
- If we decide to prune, the expected error will be $1 - \frac{99}{100} = 0.01$
- Good estimate if N is large, but if N is small, better to rely on prior probability of class, i.e. $\frac{1}{k}$
- So the **Laplace error** adds a bias for small N . When N is large, correction is irrelevant

$$E = 1 - \frac{n}{N} + \frac{1}{k}$$

N = total number of training examples

n = number of training examples in the majority class

k = number of classes

Minimal Error Pruning

Should the children of this node be pruned or not?

Left child has class frequencies [7, 3]

$$E = 1 - \frac{n+1}{N+k} = 1 - \frac{7+1}{10+2} = 0.333$$

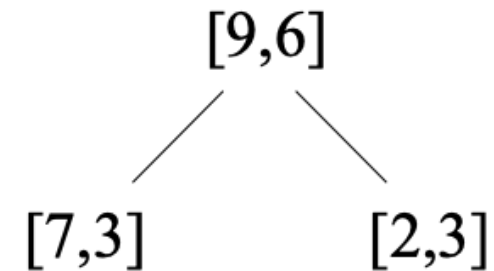
Right child has $E = 0.429$

Parent node has $E = 0.412$

Average for left and right child is:

$$E = \frac{10}{15} \times 0.333 + \frac{5}{15} \times 0.429 = 0.365$$

Since $0.365 < 0.412$, children should NOT be pruned



Minimal Error Pruning

Should the children of this node be pruned or not?

Left child has class frequencies [3, 2]

$$E = 1 - \frac{n+1}{N+k} = 1 - \frac{3+1}{5+2} = 0.429$$

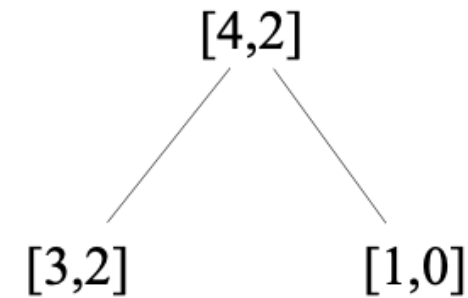
Right child has $E = 0.333$

Parent node has $E = 0.375$

Average for left and right child is:

$$E = \frac{5}{6} \times 0.429 + \frac{1}{6} \times 0.333 = 0.413$$

Since $0.413 > 0.375$, children should be pruned

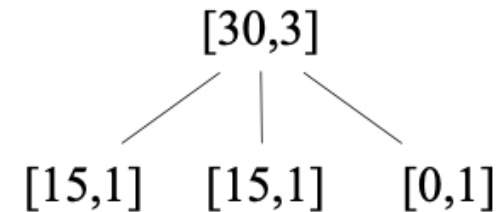


Minimal Error Pruning

Should the children of this node be pruned or not?

Left child has class frequencies [15, 1]

$$E = 1 - \frac{n+1}{N+k} = 1 - \frac{15+1}{16+2} = 0.111$$



Right child has $E = 0.333$

Parent node has $E = \frac{4}{35} = 0.114$

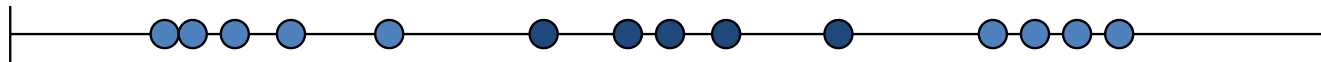
Average for left and right child is:

$$E = \frac{16}{33} \times 0.111 + \frac{16}{33} \times 0.111 + \frac{1}{33} \times 0.333 = 0.118$$

Since $0.118 > 0.114$, children should be pruned

Numerical Attributes

- ID3 algorithm is designed for attributes that have discrete values.
 - How can we handle attributes with continuous numerical values?
- ➡ Must discretise values.



Problems Suitable for Decision Trees

- Instances are represented by attribute-value pairs
- Instances are described by a fixed set of attributes (e.g., Temperature) and their values (e.g., Hot).
 - Easiest domains for decision tree learning is when each attribute takes on a small number of disjoint possible values (e.g., Hot, Mild, Cold).
 - Extensions allow handling real-valued attributes as well (e.g., representing temperature numerically).
- The target function has discrete output values
- The training data
 - The training data may contain errors
 - The training data may contain missing attribute values

Some TDIDT Systems

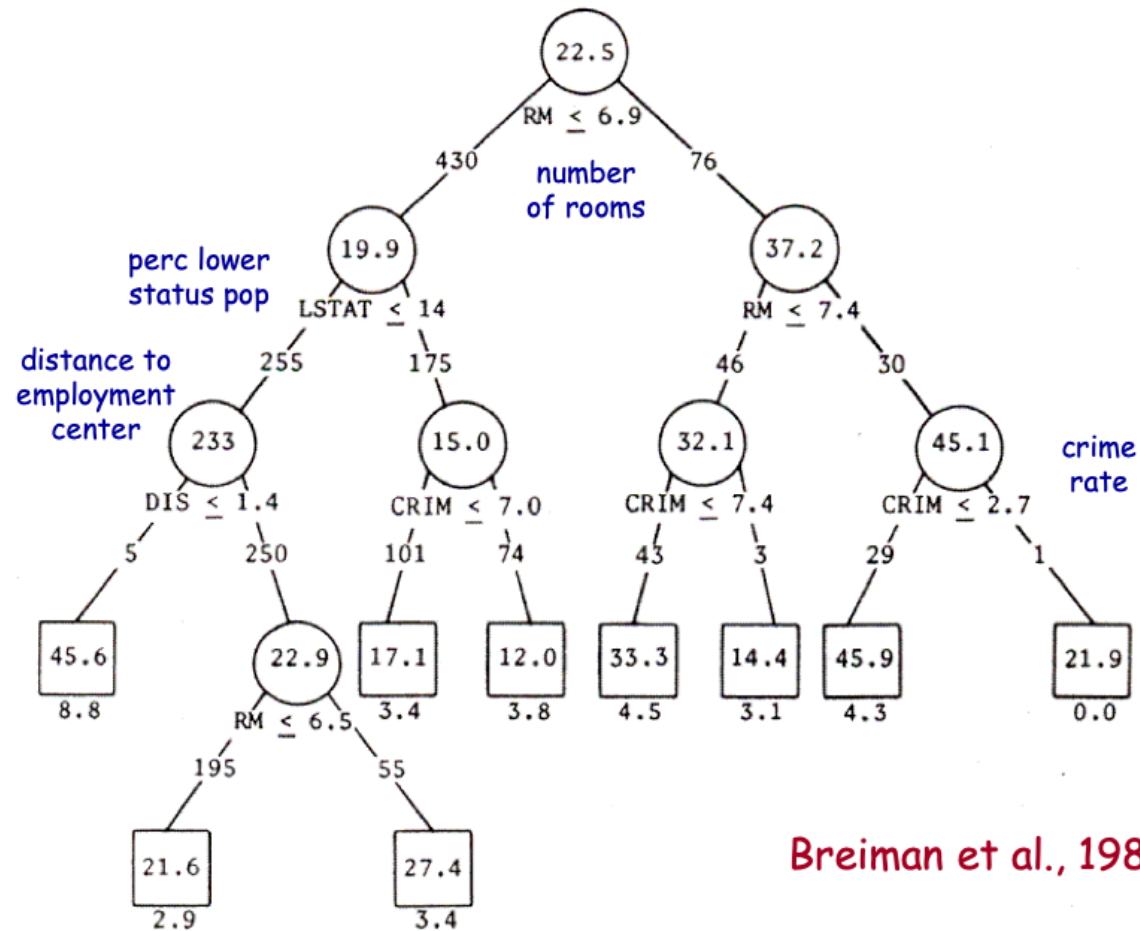
- ID3 (Quinlan 79)
- CART (Breiman et al. 84)
- Assistant (Cestnik et al. 87)
- C4.5 (Quinlan 93)
- C5 (Quinlan 97)
- J48 in WEKA (Witten, Frank 2000 - ...)
- scikit-learn (for continuous values)

TDIDT - Top-Down Induction of Decision Trees

Inducing Regression Trees

- Like induction of decision trees
- Regression trees useful in continuous domains,
 - e.g. predict biomass of algae
- Decision trees: discrete class
- Regression trees: continuous, real-valued class

Example: Boston Housing Values



Breiman et al., 1984

Some Systems that Induce Regression Trees

- CART (Breiman et al. 1984)
- RETIS (Karalič 1992)
- M5 (Quinlan 1993)
- WEKA (Witten and Frank 2000-...)

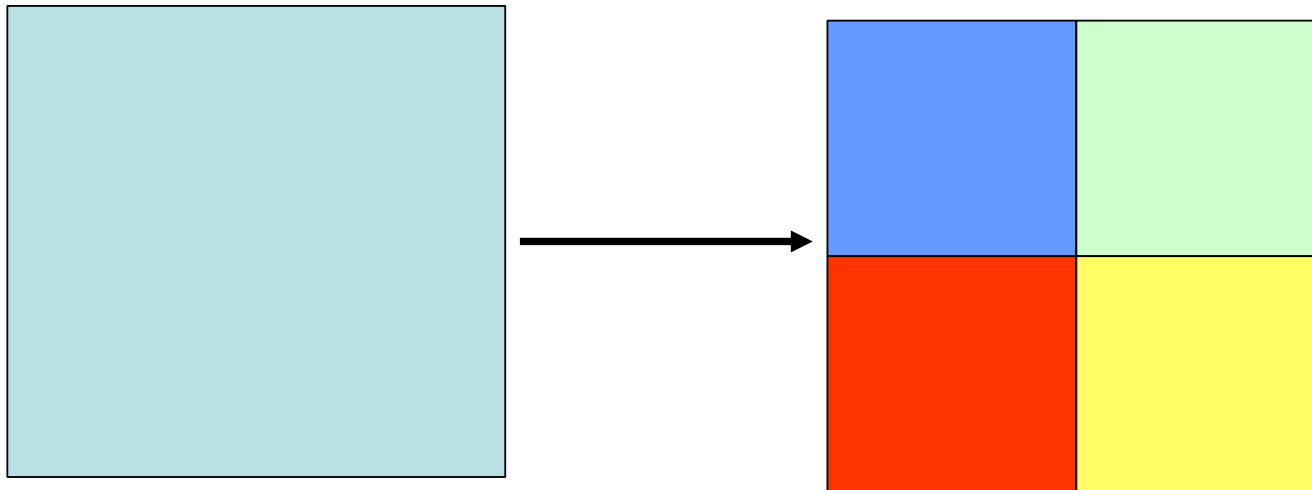
Training and Testing

- For classification problems, a classifier's performance is measured in terms of the *error rate*.
- The classifier predicts the class of each instance: if it is correct, that is counted as a *success*; if not, it is an *error*.
- The error rate is just the proportion of errors made over a whole set of instances, and it measures the overall performance of the classifier.

Training and Testing

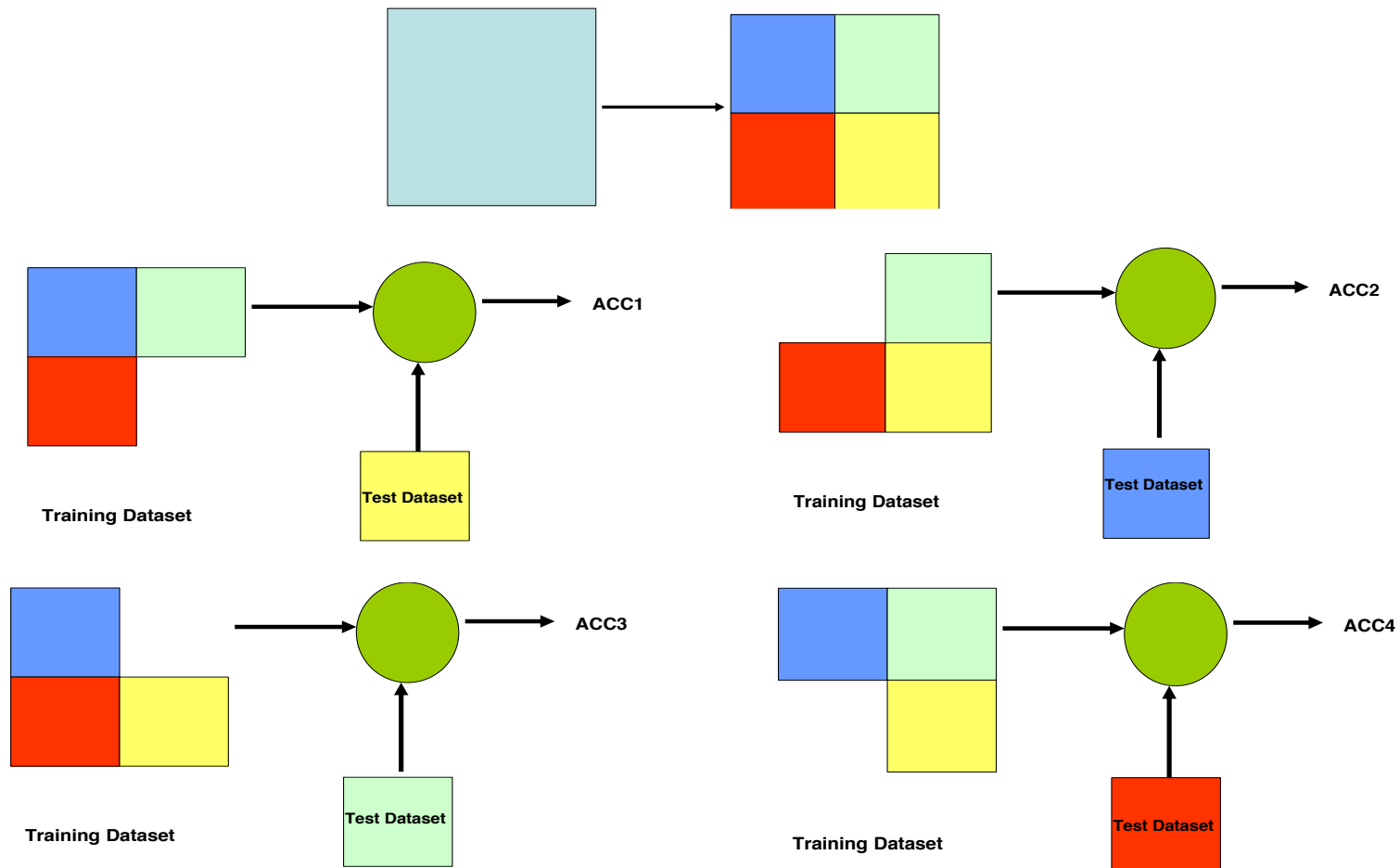
- Self-consistency test: when training and test dataset are the same
- Hold out strategy: reserve some examples for testing and use the rest for training (set part of that aside for validation, if required)
- K-fold Cross validation:
 - If don't have many examples
 - Partition dataset randomly into k equal-sized sets.
 - Train and test classifier k times using one set for testing and other $k - 1$ sets for training

4-Fold Cross-validation



$$\text{ACC} = (\text{ACC1} + \text{ACC2} + \text{ACC3} + \text{ACC4}) / 4$$

4-Fold Cross-validation



K-Fold Cross Validation

- Train multiple times, leaving out a disjoint subset of data each time for test.
- Average the test set accuracies.

Partition data into K disjoint subsets

for $k \in 1..K$:

$testData \leftarrow k_{th}$ subset

$h \leftarrow$ classifier trained on all data except for $testData$

$accuracy(k) =$ accuracy of h on $testData$

end

$FinalAccuracy =$ mean of the K recorded test set accuracies

Leave-One-Out Cross Validation

- This is just k-fold cross validation leaving out one example each iteration. Average the test set accuracies

Partition data into K disjoint subsets *each containing one example*

for $k \in 1..K$:

$testData \leftarrow k_{th}$ subset

$h \leftarrow$ classifier trained on all data except for $testData$

$accuracy(k) =$ accuracy of h on $testData$

end

$FinalAccuracy =$ mean of the K recorded test set accuracies

Summary

- Supervised Learning
 - Training set and test set
 - Try to predict target value based on input attributes
- Ockham's Razor
 - Tradeoff between simplicity and accuracy
- Decision Trees
 - Improve generalisation by building a smaller tree (using entropy)
 - Prune nodes based on Laplace error
 - Other ways to prune Decision Trees

References

- Poole & Mackworth, Artificial Intelligence: Foundations of Computational Agents, Chapter 7.
- Russell & Norvig, *Artificial Intelligence: a Modern Approach*, Chapter 18.1, 18.2, 18.3