# LEARNING FROM OBSERVATIONS


## CHAPTER 18, SECTIONS 1–3

# Outline

◇ Inductive learning

◇ Decision tree learning

# Learning a model from data

Can involve estimating parameters and/or learning structure of model

Example of parameter estimation: estimating conditional probabilities in Bayesian networks

Practical application: naive Bayes document classifier

# Inductive learning (a.k.a. Science)

Simplest form: learn a function from examples (supervised learning)

$f$ is the target function

An example is a pair $x$, $f(x)$, e.g.,
$$
\begin{array}{c|c|c}
O & O & X \\
\hline
 & X & \\
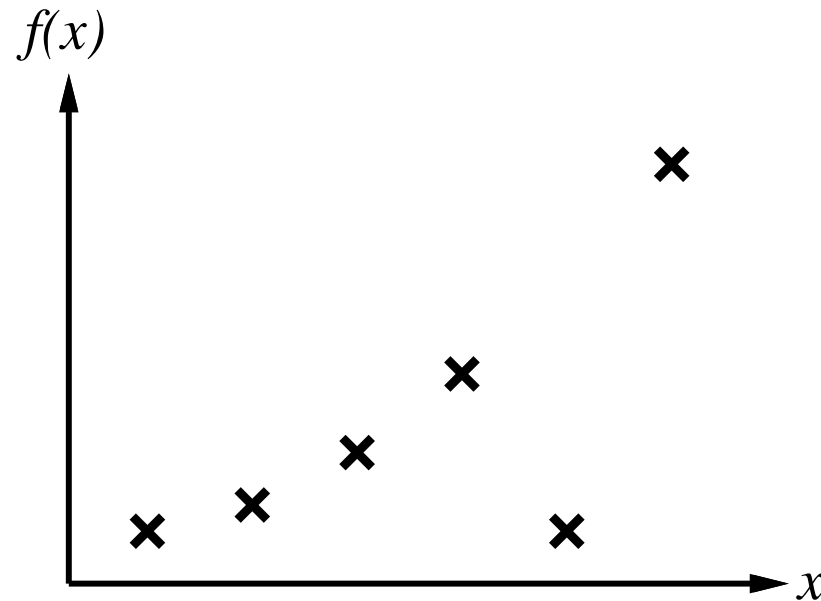\hline
X & & \\
\end{array}
\quad , \quad +1
$$

Problem: find a(n) hypothesis $h$
      such that $h \approx f$
      given a training set of examples

# Inductive learning method

Construct/adjust $h$ to agree with $f$ on training set
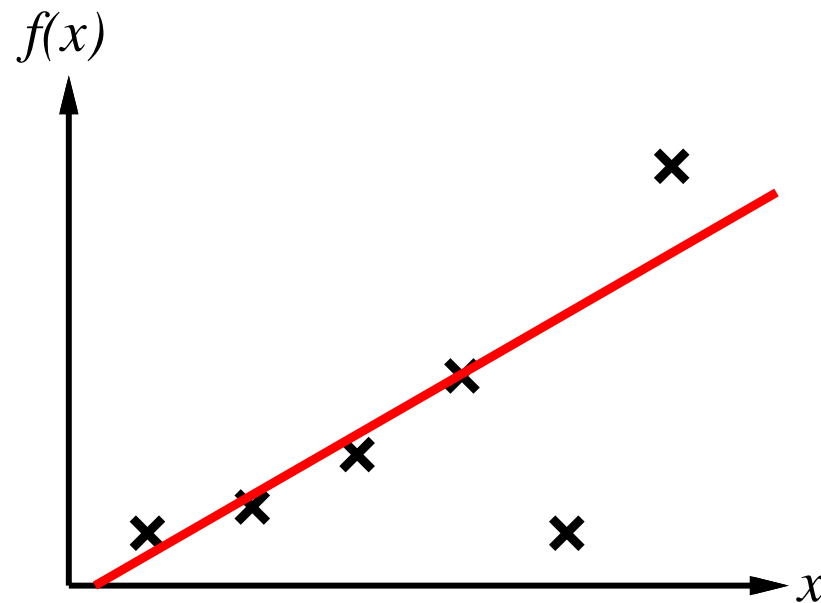($h$ is consistent if it agrees with $f$ on all examples)

E.g., curve fitting:

# Inductive learning method

Construct/adjust $h$ to agree with $f$ on training set
($h$ is consistent if it agrees with $f$ on all examples)
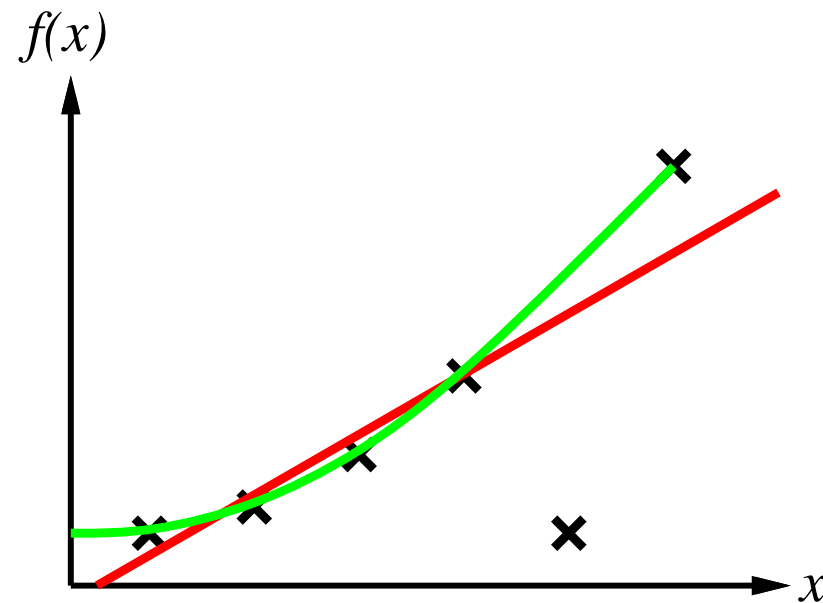
E.g., curve fitting:

# Inductive learning method

Construct/adjust $h$ to agree with $f$ on training set
($h$ is consistent if it agrees with $f$ on all examples)

E.g., curve fitting:

# Inductive learning method

Construct/adjust $h$ to agree with $f$ on training set
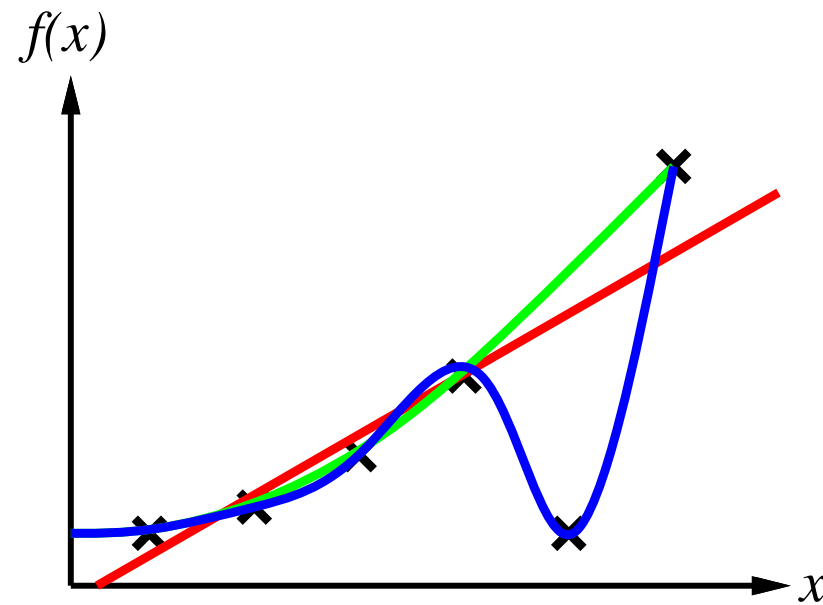($h$ is consistent if it agrees with $f$ on all examples)

E.g., curve fitting:

# Inductive learning method

Construct/adjust $h$ to agree with $f$ on training set
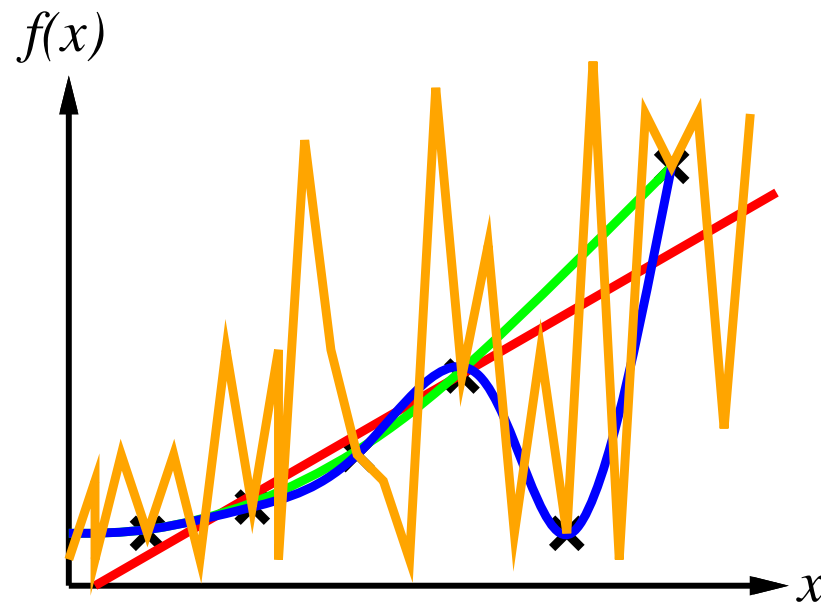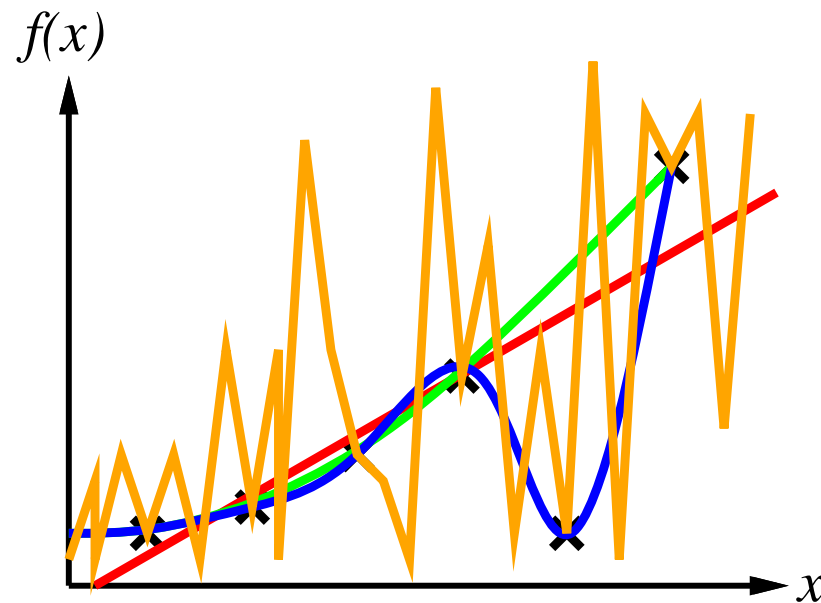($h$ is consistent if it agrees with $f$ on all examples)

E.g., curve fitting:

# Inductive learning method

Construct/adjust $h$ to agree with $f$ on training set
($h$ is consistent if it agrees with $f$ on all examples)

E.g., curve fitting:



Ockham's razor: maximize a combination of consistency and simplicity

# Attribute-based representations

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

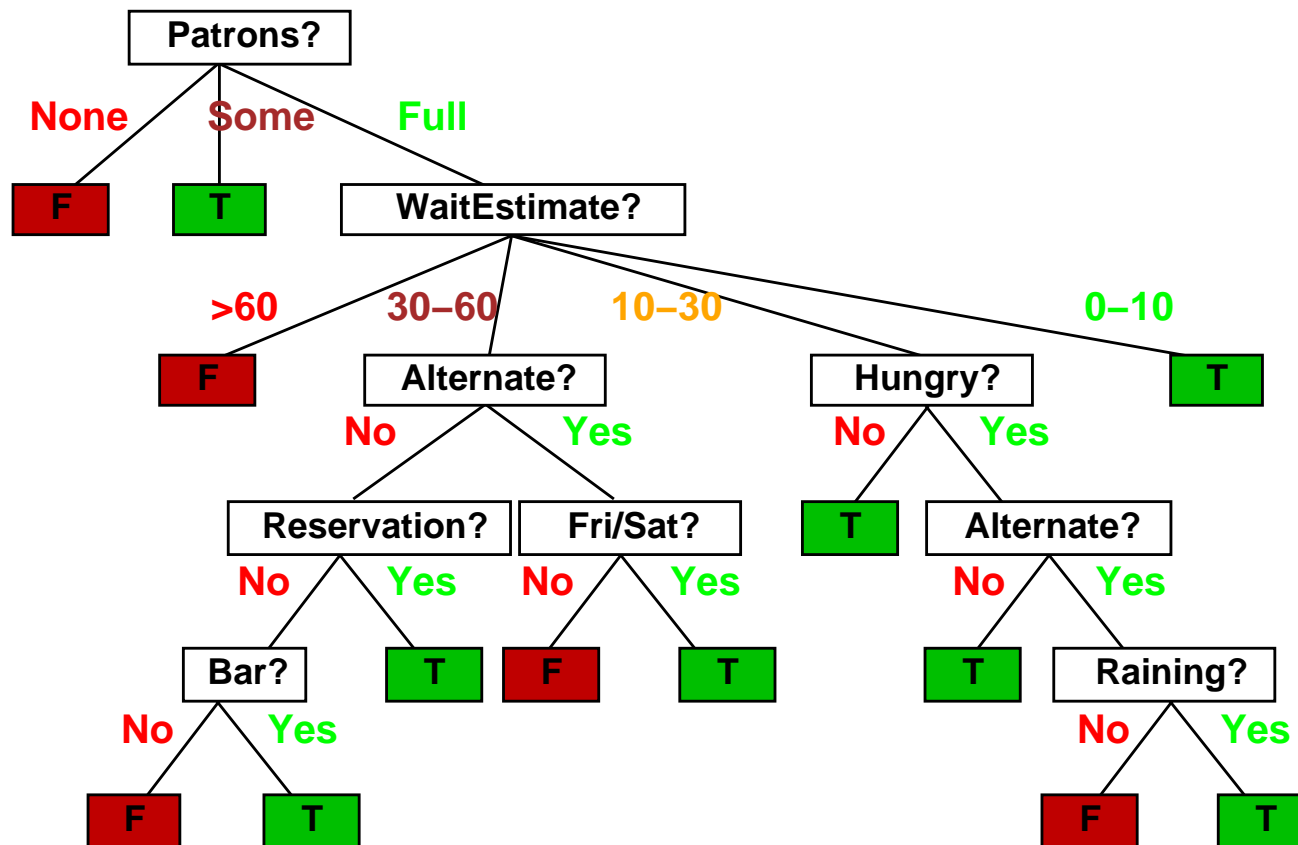| Example | Attributes | | | | | | | | | | Target |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | Alt | Bar | Fri | Hun | Pat | Price | Rain | Res | Type | Est | WillWait |
| $X_1$ | T | F | F | T | Some | $$$ | F | T | French | 0–10 | T |
| $X_2$ | T | F | F | T | Full | $ | F | F | Thai | 30–60 | F |
| $X_3$ | F | T | F | F | Some | $ | F | F | Burger | 0–10 | T |
| $X_4$ | T | F | T | T | Full | $ | F | F | Thai | 10–30 | T |
| $X_5$ | T | F | T | F | Full | $$$ | 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 |

Classification of examples is positive (T) or negative (F)

# Decision trees

One possible representation for hypotheses

E.g., here is the "true" tree for deciding whether to wait:

# Expressiveness

Decision trees can express any function of the input attributes.

E.g., for Boolean functions, truth table row $\rightarrow$ path to leaf:

| A | B | A xor B |
|---|---|---------|
| F | F | F |
| F | T | T |
| T | F | T |
| T | T | F |



Trivially, $\exists$ a consistent decision tree for any training set
w/ one path to leaf for each example (unless $f$ nondeterministic in $x$)
but it probably won't generalize to new examples

Prefer to find more **compact** decision trees

# Hypothesis spaces

How many distinct decision trees with $n$ Boolean attributes??

# Hypothesis spaces

How many distinct decision trees with $n$ Boolean attributes??

= number of Boolean functions

# Hypothesis spaces

How many distinct decision trees with $n$ Boolean attributes??

$=$ number of Boolean functions
$=$ number of distinct truth tables with $2^n$ rows

# Hypothesis spaces

How many distinct decision trees with $n$ Boolean attributes??

= number of Boolean functions
= number of distinct truth tables with $2^n$ rows = $2^{2^n}$

# Hypothesis spaces

How many distinct decision trees with $n$ Boolean attributes??

= number of Boolean functions
= number of distinct truth tables with $2^n$ rows = $2^{2^n}$

E.g., with 6 Boolean attributes, there are 18,446,744,073,709,551,616 trees

# Hypothesis spaces

How many distinct decision trees with $n$ Boolean attributes??

= number of Boolean functions
= number of distinct truth tables with $2^n$ rows = $2^{2^n}$

E.g., with 6 Boolean attributes, there are 18,446,744,073,709,551,616 trees

How many purely conjunctive hypotheses (e.g., $Hungry \wedge \neg Rain$)??

# Hypothesis spaces

How many distinct decision trees with $n$ Boolean attributes??

= number of Boolean functions
= number of distinct truth tables with $2^n$ rows = $2^{2^n}$

E.g., with 6 Boolean attributes, there are 18,446,744,073,709,551,616 trees

How many purely conjunctive hypotheses (e.g., $Hungry \wedge \neg Rain$)??

Each attribute can be in (positive), in (negative), or out
$\Rightarrow$   $3^n$ distinct conjunctive hypotheses

More expressive hypothesis space
– increases chance that target function can be expressed
– increases number of hypotheses consistent w/ training set
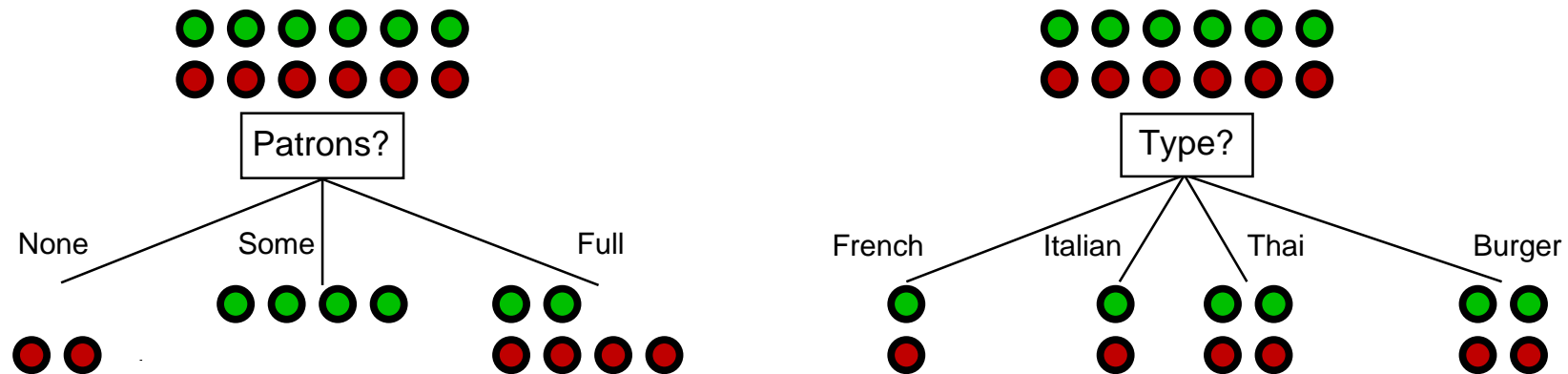$\Rightarrow$   may get worse predictions

# Decision tree learning

Aim: find a small tree consistent with the training examples

Idea: (recursively) choose "most significant" attribute as root of (sub)tree

---

**function** DTL(*examples, attributes, default*) **returns** a decision tree

  **if** *examples* is empty **then return** *default*
  **else if** all *examples* have the same classification **then return** the classification
  **else if** *attributes* is empty **then return** MODE(*examples*)
  **else**
      *best* ← CHOOSE-ATTRIBUTE(*attributes, examples*)
      *tree* ← a new decision tree with root test *best*
      **for each** value $v_i$ of *best* **do**
          $examples_i$ ← {elements of *examples* with *best* = $v_i$}
          *subtree* ← DTL($examples_i$, *attributes* − *best*, MODE(*examples*))
          add a branch to *tree* with label $v_i$ and subtree *subtree*
      **return** *tree*

---

# Choosing an attribute

Idea: a good attribute splits the examples into subsets that are (ideally) "all positive" or "all negative"



*Patrons?* is a better choice—gives **information** about the classification

# Information

Information answers questions

The more clueless I am about the answer initially, the more information is contained in the answer

Scale: 1 bit = answer to Boolean question with prior $\langle 0.5, 0.5 \rangle$

Information in an answer when prior is $\langle P_1, \ldots, P_n \rangle$ is

$$H(\langle P_1, \ldots, P_n \rangle) = \sum_{i=1}^{n} - P_i \log_2 P_i$$

(also called entropy of the prior)

# Information contd.

Suppose we have $p$ positive and $n$ negative examples at the root
$$\Rightarrow \quad H(\langle p/(p+n), n/(p+n)\rangle) \text{ bits needed to classify a new example}$$
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
$$\Rightarrow \quad H(\langle p_i/(p_i+n_i), n_i/(p_i+n_i)\rangle) \text{ bits needed to classify a new example}$$
$\Rightarrow$ **expected** number of bits per example over all branches is

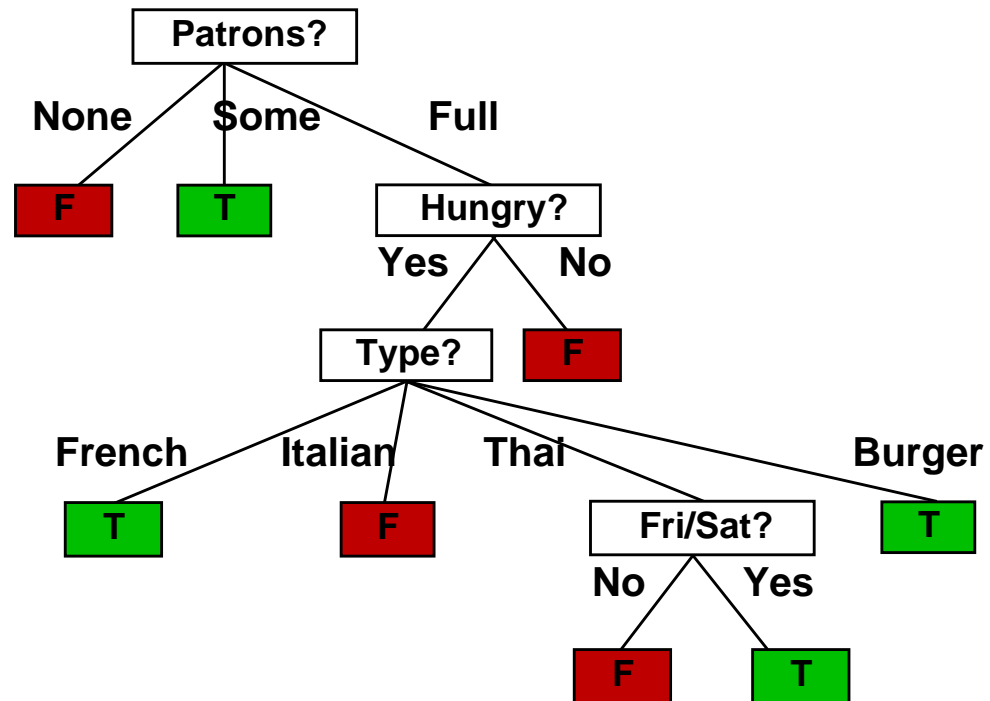$$\Sigma_i \ \frac{p_i + n_i}{p + n} \ H(\langle p_i/(p_i + n_i), n_i/(p_i + n_i)\rangle)$$

For $Patrons?$, this is 0.459 bits, for $Type$ this is (still) 1 bit

$\Rightarrow$ choose the attribute that minimizes the remaining information needed

# Example contd.

Decision tree learned from the 12 examples:



Substantially simpler than "true" tree—a more complex hypothesis isn't justified by small amount of data

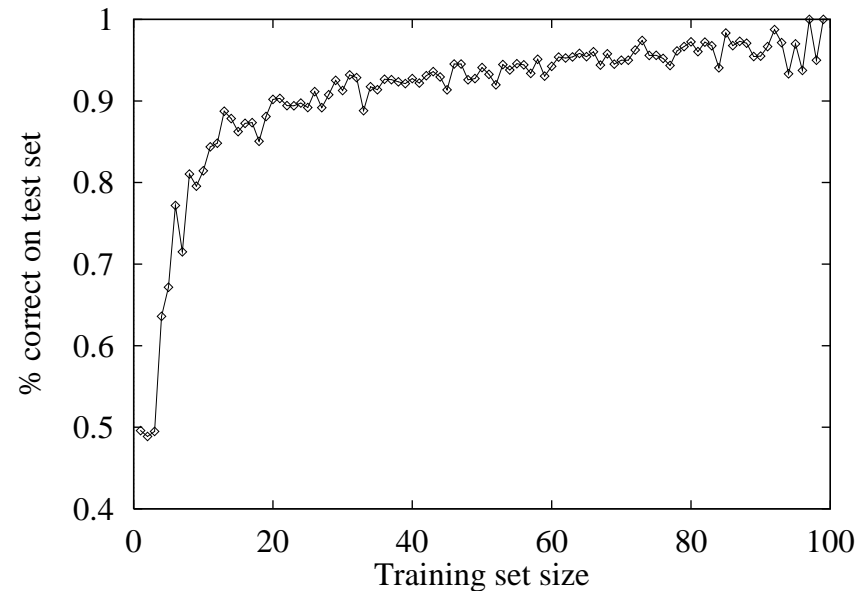# Performance measurement

How do we know that $h \approx f$? (Hume's **Problem of Induction**)

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

# Summary

Learning needed for unknown environments, lazy designers

For supervised learning, the aim is to find a simple hypothesis approximately consistent with training examples

Decision tree learning using information gain

Learning performance = prediction accuracy measured on test set