

## LEARNING FROM OBSERVATIONS

CHAPTER 18, SECTIONS 1–3

Chapter 18, Sections 1–3 1

### Outline

- ◇ Learning agents
- ◇ Inductive learning
- ◇ Decision tree learning
- ◇ Measuring learning performance

Chapter 18, Sections 1–3 2

## Learning

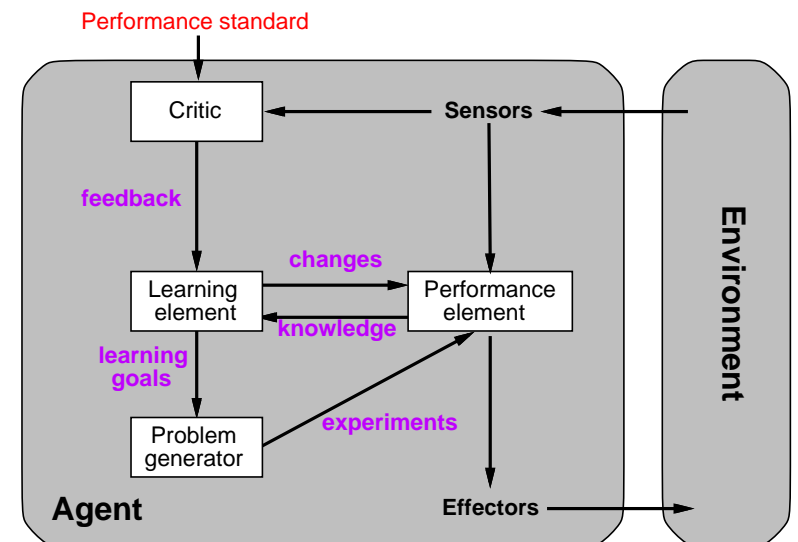
Learning is essential for unknown environments,  
i.e., when designer lacks omniscience

Learning is useful as a system construction method,  
i.e., expose the agent to reality rather than trying to write it down

Learning modifies the agent's decision mechanisms to improve performance

Chapter 18, Sections 1–3 3

### Learning agents



Chapter 18, Sections 1–3 4

## Learning element

Design of learning element is dictated by

- ◇ what type of performance element is used
- ◇ which functional component is to be learned
- ◇ how that functional component is represented
- ◇ what kind of feedback is available

Example scenarios:

Performance element	Component	Representation	Feedback
Alpha-beta search	Eval. fn.	Weighted linear function	Win/loss
Logical agent	Transition model	Successor-state axioms	Outcome
Utility-based agent	Transition model	Dynamic Bayes net	Outcome
Simple reflex agent	Percept-action fn	Neural net	Correct action

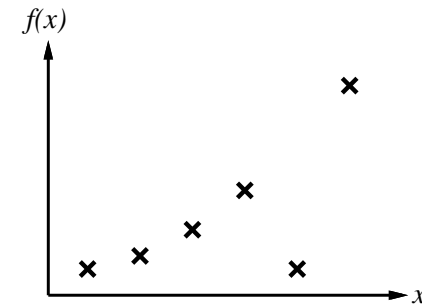
Supervised learning: correct answers for each instance

Reinforcement learning: occasional rewards

## 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 (a.k.a. Science)

Simplest form: learn a function from examples (**tabula rasa**)

$f$  is the **target function**

An **example** is a pair  $x, f(x)$ , e.g.,  $\frac{O|O|X}{X|X|}$ , +1

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

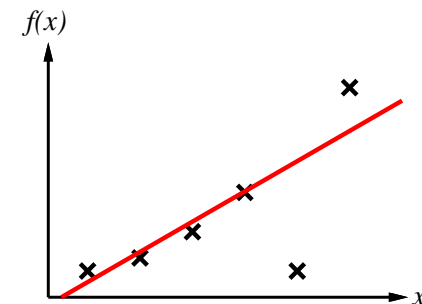
(This is a highly simplified model of real learning:

- Ignores prior knowledge
- Assumes a deterministic, observable “environment”
- Assumes examples are given
- Assumes that the agent wants to learn  $f$ —why?)

## Inductive learning method

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

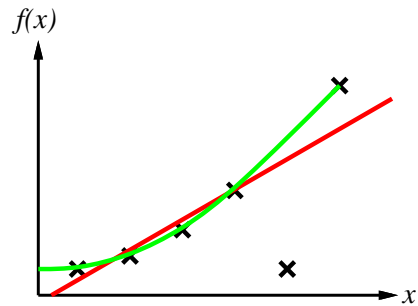
E.g., curve fitting:



## Inductive learning method

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

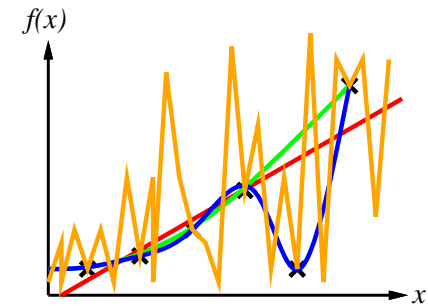
E.g., curve fitting:



## Inductive learning method

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

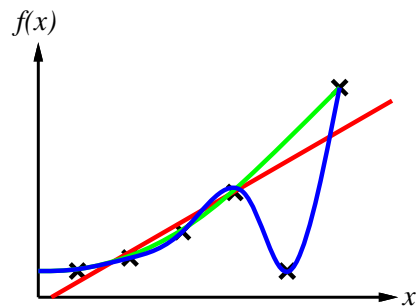
E.g., curve fitting:



## Inductive learning method

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

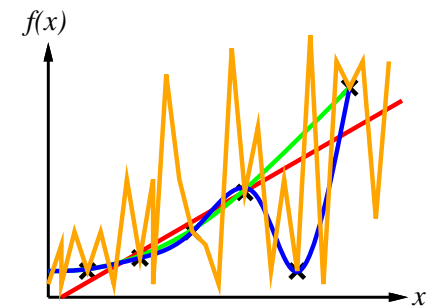
E.g., curve fitting:



## 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

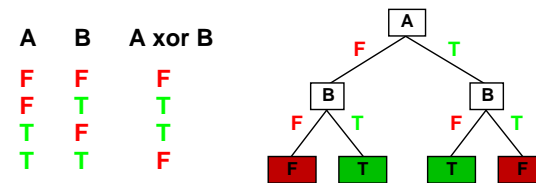
Examples described by **attribute values** (Boolean, discrete, continuous, etc.)  
 E.g., situations where I will/won't wait for a table:

Example	Attributes										Target
	<i>Alt</i>	<i>Bar</i>	<i>Fri</i>	<i>Hun</i>	<i>Pat</i>	<i>Price</i>	<i>Rain</i>	<i>Res</i>	<i>Type</i>	<i>Est</i>	<i>WillWait</i>
$X_1$	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)

## Expressiveness

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



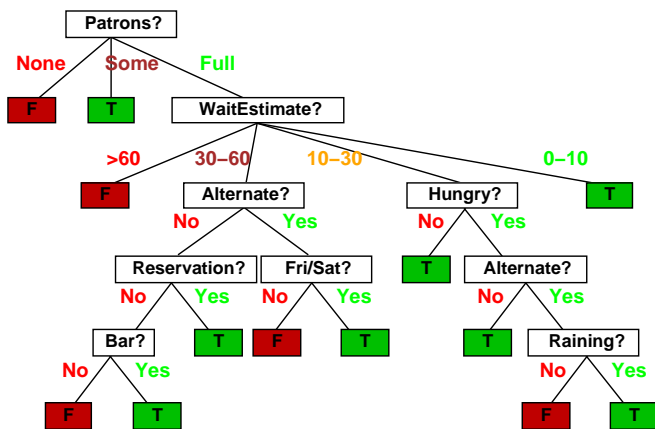
Trivially, there is a consistent decision tree for any training set  
 w/ one path to leaf for each example (unless  $f$  nondeterministic in  $x$ )  
 but it probably won't generalize to new examples

Prefer to find more **compact** decision trees

## Decision trees

One possible representation for hypotheses

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



## 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 =  $2^{2^n}$

## Hypothesis spaces

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

= number of Boolean functions

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

## 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

⇒  $3^n$  distinct conjunctive hypotheses

More expressive hypothesis space

– increases chance that target function can be expressed 😊

– increases number of hypotheses consistent w/ training set

⇒ may get worse predictions 😞

## 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

⇒  $3^n$  distinct conjunctive hypotheses

## Decision tree learning

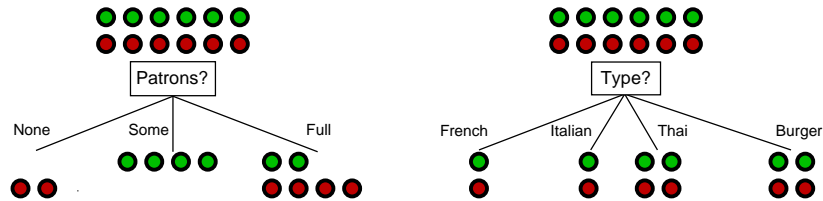
Aim: find a small tree consistent with the training examples

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(examplesi, 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 Theory

- ◇ Consider communicating two messages (A and B) between two parties
- ◇ Bits are used to measure message size
- ◇ If  $P(A) = 1$  and  $P(B) = 0$ , how many bits are needed?
- ◇ If  $P(A) = .5$  and  $P(B) = .5$ , how many bits are needed?
- ◇ Information:  $I(P(v_1), \dots, P(v_n)) = \sum_{i=1}^n -P(v_i) \log_2 P(v_i)$
- ◇  $I(1, 0) = 0$  bit
- ◇  $I(0.5, 0.5) = -0.5 \times \log_2 0.5 - 0.5 \times \log_2 0.5 = 1$  bit

## Information Theory

- ◇ Consider communicating two messages (A and B) between two parties
- ◇ Bits are used to measure message size
- ◇ If  $P(A) = 1$  and  $P(B) = 0$ , how many bits are needed?
- ◇ If  $P(A) = .5$  and  $P(B) = .5$ , how many bits are needed?

## Information Theory

- ◇ Consider communicating two messages (A and B) between two parties
- ◇ Bits are used to measure message size
- ◇ If  $P(A) = 1$  and  $P(B) = 0$ , how many bits are needed?
- ◇ If  $P(A) = .5$  and  $P(B) = .5$ , how many bits are needed?
- ◇ Information:  $I(P(v_1), \dots, P(v_n)) = \sum_{i=1}^n -P(v_i) \log_2 P(v_i)$
- ◇  $I(1, 0) = 0$  bit
- ◇  $I(0.5, 0.5) = -0.5 \times \log_2 0.5 - 0.5 \times \log_2 0.5 = 1$  bit
- ◇  $I$  measures the information content for communication (or uncertainty in what is already known)
- ◇ The more one knows, the less to be communicated, the smaller is  $I$
- ◇ The less one knows, the more to be communicated, the larger is  $I$

## Using Information Theory

- ◇  $(P(pos), P(neg))$ : probabilities of positive and negative
- ◇ Attribute *color*: black (1,0), white (0,1)
- ◇ Attribute *size*: large (.5,.5), small (.5,.5)

## Using Information Theory

- ◇  $(P(pos), P(neg))$ : probabilities of positive and negative
- ◇ Attribute *color*: black (1,0), white (0,1)
- ◇ Attribute *size*: large (.5,.5), small (.5,.5)
  
- ◇ Before selecting an attribute
  - $p$  = number of positive examples,  $n$  = number of negative examples
  - Estimating probabilities:  $P(pos) = \frac{p}{p+n}$ ,  $P(neg) = \frac{n}{p+n}$
  - $Before() = I(P(pos), P(neg))$

## Selecting an Attribute

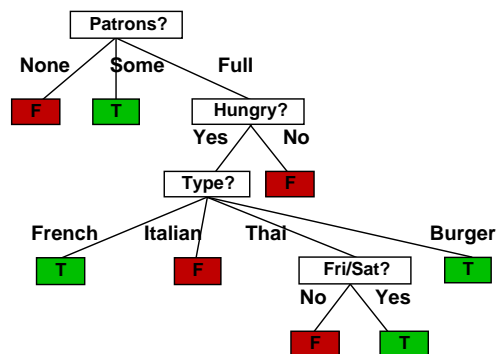
- ◇ Evaluating an attribute (e.g., color)
  - $p_i$  = number of positive examples for value  $i$  (e.g., black),  $n_i$  = number of negative ones
  - Estimating probabilities for value  $i$ :  $P_i(pos) = \frac{p_i}{p_i+n_i}$ ,  $P_i(neg) = \frac{n_i}{p_i+n_i}$
  - $v$  values for attribute  $A$  (e.g., 2 for color)
  - $Remainder(A) = After(A) = \sum_{i=1}^v \frac{p_i+n_i}{p+n} I(P_i(pos), P_i(neg))$  [expected information]

## Selecting an Attribute

- ◇ Evaluating an attribute (e.g., color)
  - $p_i$  = number of positive examples for value  $i$  (e.g., black),  $n_i$  = number of negative ones
  - Estimating probabilities for value  $i$ :  $P_i(pos) = \frac{p_i}{p_i+n_i}$ ,  $P_i(neg) = \frac{n_i}{p_i+n_i}$
  - $v$  values for attribute  $A$  (e.g., 2 for color)
  - $Remainder(A) = After(A) = \sum_{i=1}^v \frac{p_i+n_i}{p+n} I(P_i(pos), P_i(neg))$  [expected information]
  
- ◇ "Information Gain" (reduction in uncertainty of what is known)
  - $Gain(A) = Before() - After(A)$  [ $Before()$  has more uncertainty]
  - Choose attribute  $A$  with the largest  $Gain(A)$

## 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**)

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
  - divide into two disjoint subsets: training and test sets
  - prediction accuracy: accuracy on the (unseen) test set

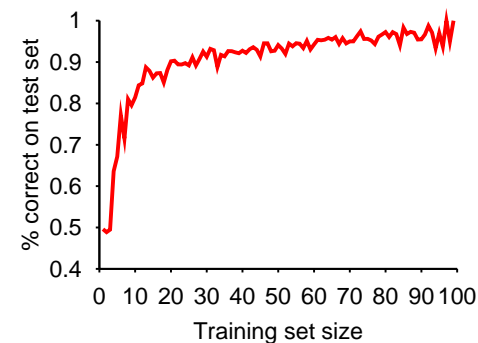
## Performance measurement

How do we know that  $h \approx f$ ?

How about measuring the accuracy of  $h$  on the examples that were used to learn  $h$ ?

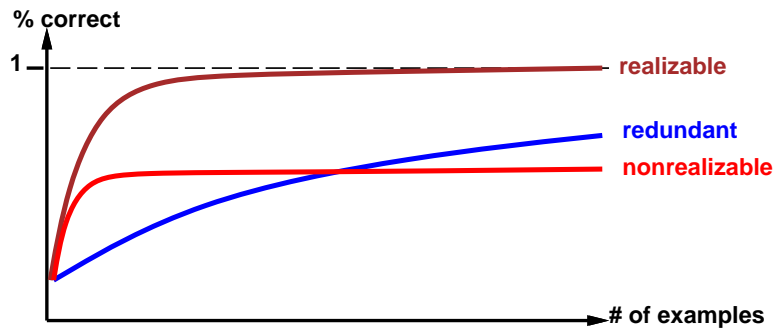
## Performance measurement

**Learning curve** = % correct on test set as a function of training set size



## Learning curve

- **realizable** (can express target function) vs. **non-realizable**  
non-realizability can be due to:
  - missing attributes
  - and/or restricted hypothesis class (e.g., thresholded linear function)
- redundant expressiveness (e.g., loads of irrelevant attributes)



Chapter 18, Sections 1-3 37

## Overfitting

- More attributes  $\Rightarrow$  larger hypothesis space
- Larger hypothesis space can lead to more hypotheses that represent meaningless “regularity/patterns”
- Overfitting: high accuracy on training set, but low accuracy on test set—low prediction accuracy
- Select the attribute with the largest information gain
  - however, is the gain significant?
  - (statistical) significance test
- Pruning
  - do not include the attribute if information gain is not statistically significant
  - potentially, less than 100% accurate on the training set, why?
  - however, improved prediction accuracy on the test set

Chapter 18, Sections 1-3 39

## Irrelevant Attributes

- Consider adding the attribute: Date (month and day)
- How can it affect the learned tree?

Chapter 18, Sections 1-3 38

## Significance Test

- “Null hypothesis” (in statistics): attribute is irrelevant (gain is not significant)
- “Alternative hypothesis”: attribute is relevant
- Calculating the deviation
  - expected  $\hat{p}_i = p \times \frac{p_i + n_i}{p + n}$
  - expected  $\hat{n}_i = n \times \frac{p_i + n_i}{p + n}$
  - Deviation (from expected):
$$D = \sum_{i=1}^v \frac{(p_i - \hat{p}_i)^2}{\hat{p}_i} + \frac{(n_i - \hat{n}_i)^2}{\hat{n}_i}$$
  - $D$  is  $\chi^2$  (chi-squared) distributed with  $v - 1$  degrees of freedom
  - $\chi^2$  Test in statistics
- With a confidence level (e.g. 95%), if  $D > \chi^2$  value, attribute is relevant (Null hypothesis is rejected)

Chapter 18, Sections 1-3 40

## Additional Issues

- ◇ Missing attribute values.
- ◇ Gain() biases to attributes with more values.
- ◇ Continuous-valued (numeric) attributes have infinite number of values.

## Summary

Learning needed for unknown environments, lazy designers

Learning agent = performance element + learning element

Learning method depends on type of performance element, available feedback, type of component to be improved, and its representation

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

Decision tree learning using information gain

Learning performance = prediction accuracy measured on test set

## Learning as search

- ◇ What is the state space in learning decision trees?
- ◇ State-space formulation
  - State: a decision tree
  - Initial state: an empty decision tree
  - Action: select an attribute and add it to the tree
  - Goal test: all examples in each leaf have the same classification
- ◇ What kind of search is it?