

# Decision Tree Learning

## Chapter 3

# Recap from Chapter 2

## Decision Tree Learning

### Overview

#### Entropy

#### Information Gain

#### Inductive Bias in ID3

#### Occam's Razor

- Concept learning
- consistent hypothesis ( $h(x) = c(x)$ ) and satisfying hypothesis ( $h(x) = 1$ )
- Find-S algorithm
- Candidate Elimination algorithm
- Inductive Bias

# Topics

## Decision Tree Learning

### Overview

#### Entropy

#### Information Gain

#### Inductive Bias in ID3

#### Occam's Razor

- Decision tree representation
- ID3 learning algorithm
- Entropy, Information gain
- Overfitting

# Decision Tree Definition

## Decision Tree Learning

### Overview

### Entropy

### Information Gain

### Inductive Bias in ID3

### Occam's Razor

Decision tree learning is a method for approximating **discrete-valued** target functions, in which the learned function is represented by a decision tree.

Decision trees classify instances by sorting them down the tree from the **root to some leaf node**, which provides the classification of the instance.

Example:

$$F : \langle Outlook, Humidity, Wind, Temp \rangle \rightarrow PlayTennis?$$

# Decision Tree for *PlayTennis*

Decision Tree  
Learning

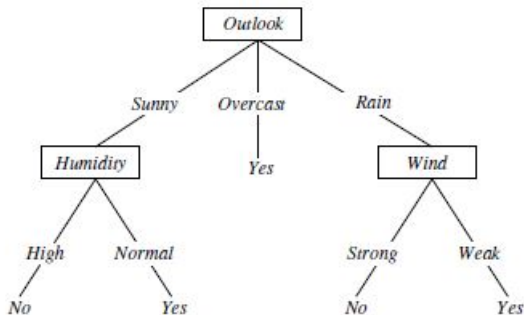
Overview

Entropy

Information  
Gain

Inductive Bias  
in ID3

Occam's  
Razor



This corresponds to

$$(Outlook = Sunny \wedge Humidity = Normal) \vee (Outlook = Overcast) \vee (Outlook = Rain \wedge Wind = Weak)$$

# Decision Trees

## Decision Tree Learning

### Overview

### Entropy

### Information Gain

### Inductive Bias in ID3

### Occam's Razor

Decision tree representation:

- Each internal node tests an attribute
- Each branch corresponds to attribute value
- Each leaf node assigns a classification

How would we represent:

- $\wedge, \vee, \text{XOR}$
- $(A \wedge B) \vee (C \wedge \neg D \wedge E)$
- $M$  of  $N$

The algorithm to construct decision tree (IDB3) is greedy search top down approach.

# Decision Tree Learning

## Decision Tree Learning

### Overview

### Entropy

### Information Gain

### Inductive Bias in ID3

### Occam's Razor

## Problem Setting:

- Set of possible instances  $X$ 
  - each instance  $x$  in  $X$  is a feature vector
  - $x = \langle x_1, x_2 \dots x_n \rangle$ .
- Unknown target function  $f : X \rightarrow Y$ 
  - $Y$  is discrete valued
- Set of function hypotheses  $H = \{h \mid h : X \rightarrow Y\}$ 
  - each hypothesis  $h$  is a decision tree

**Input:** Training examples  $\langle x(i), y(i) \rangle$  of unknown target function  $f$ .

**Output:** Hypothesis  $h \in H$  that best approximates target function  $f$

# Decision Tree

## Decision Tree Learning

### Overview

### Entropy

### Information Gain

### Inductive Bias in ID3

### Occam's Razor

Suppose

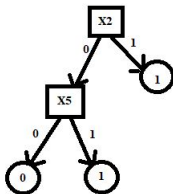
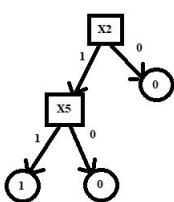
$$X = \langle X_1, \dots, X_n \rangle$$

where  $X_i$  are boolean variables.

How would you represent

$$Y = X_2 X_5$$

$$Y = X_2 \vee X_5$$





# When to Consider Decision Trees

## Decision Tree Learning

### Overview

### Entropy

### Information Gain

### Inductive Bias in ID3

### Occam's Razor

- Instances describable by attribute–value pairs
- Target function is discrete valued
- Disjunctive hypothesis may be required
- Possibly noisy training data

### Examples:

- Equipment or medical diagnosis
- Credit risk analysis
- Modeling calendar scheduling preferences

# Top-Down Induction of Decision Trees

Decision Tree  
Learning

Overview

Entropy

Information  
Gain

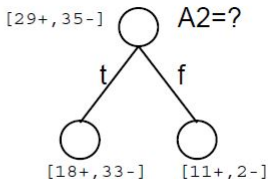
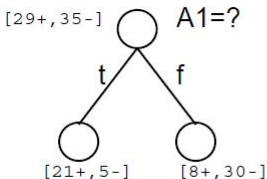
Inductive Bias  
in ID3

Occam's  
Razor

Main loop:

- $A \leftarrow$  the “best” decision attribute for next *node*
- Assign  $A$  as decision attribute for *node*
- For each value of  $A$ , create new descendant of *node*
- Sort training examples to leaf nodes
- If training examples perfectly classified, Then STOP, Else iterate over new leaf nodes

Which attribute is best?



# Entropy

Decision Tree  
Learning

Overview

Entropy

Information  
Gain

Inductive Bias  
in ID3

Occam's  
Razor

$Entropy(S)$  = expected number of bits needed to encode class ( $\oplus$  or  $\ominus$ ) of randomly drawn member of  $S$  (under the optimal, shortest-length code)

Why?

Information theory: optimal length code assigns  $-\log_2 p$  bits to message having probability  $p$ .

So, expected number of bits to encode  $\oplus$  or  $\ominus$  of random member of  $S$ :

$$p_{\oplus}(-\log_2 p_{\oplus}) + p_{\ominus}(-\log_2 p_{\ominus})$$

$$Entropy(S) \equiv -p_{\oplus} \log_2 p_{\oplus} - p_{\ominus} \log_2 p_{\ominus}$$

# Entropy

## Decision Tree Learning

### Overview

### Entropy

### Information Gain

### Inductive Bias in ID3

### Occam's Razor

Suppose  $S$  is a collection of 14 examples of some boolean concept, including 9 positive and 5 negative examples (we adopt the notation  $[9+, 5-]$  to summarize such a sample of data). Then the entropy of  $S$  relative to this boolean classification is

$$\begin{aligned} \text{Entropy}([9+, 5-]) &\equiv -(9/14) \log_2(9/14) - (5/14) \log_2(5/14) \\ &\equiv 0.940 \end{aligned}$$

**Note:** Entropy is 0 if all members of  $S$  belong to the same class. entropy is 1 when the collection contains an equal number of positive and negative examples.

# Entropy

## Decision Tree Learning

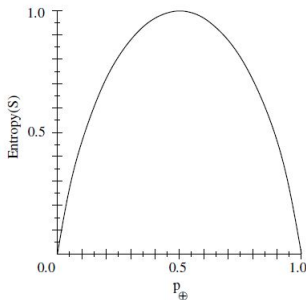
### Overview

### Entropy

### Information Gain

### Inductive Bias in ID3

### Occam's Razor



- $S$  is a sample of training examples
- $p_+$  is the proportion of positive examples in  $S$
- $p_-$  is the proportion of negative examples in  $S$
- Entropy measures the impurity of  $S$

$$Entropy(S) \equiv -p_+ \log_2 p_+ - p_- \log_2 p_-$$

# Information Gain

## Decision Tree Learning

### Overview

### Entropy

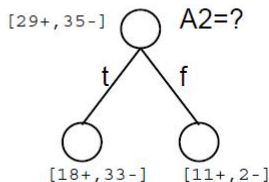
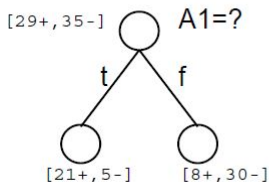
### Information Gain

### Inductive Bias in ID3

### Occam's Razor

$Gain(S, A) =$  expected reduction in entropy due to sorting on  $A$

$$Gain(S, A) \equiv Entropy(S) - \sum_{v \in Values(A)} \frac{|S_v|}{|S|} Entropy(S_v)$$



# Gain - Example

## Decision Tree Learning

### Overview

### Entropy

### Information Gain

### Inductive Bias in ID3

### Occam's Razor

For example, suppose  $S$  is a collection of training-example days described by attributes including *Wind*, which can have the values *Weak* or *Strong*. As before, assume  $S$  is a collection containing 14 examples,  $[9+, 5-]$ . Of these 14 examples, suppose 6 of the positive and 2 of the negative examples have *Wind* = *Weak*, and the remainder have *Wind* = *Strong*. The information gain due to sorting the original 14 examples by the attribute *Wind* may then be calculated as

$$\text{Values}(\text{Wind}) = \text{Weak}, \text{Strong}$$

$$S = [9+, 5-]$$

$$S_{\text{Weak}} \leftarrow [6+, 2-]$$

$$S_{\text{Strong}} \leftarrow [3+, 3-]$$

$$\begin{aligned} \text{Gain}(S, \text{Wind}) &= \text{Entropy}(S) - \sum_{v \in \{\text{Weak}, \text{Strong}\}} \frac{|S_v|}{|S|} \text{Entropy}(S_v) \\ &= \text{Entropy}(S) - (8/14)\text{Entropy}(S_{\text{Weak}}) \\ &\quad - (6/14)\text{Entropy}(S_{\text{Strong}}) \\ &= 0.940 - (8/14)0.811 - (6/14)1.00 \\ &= 0.048 \end{aligned}$$

# Training Examples

## Decision Tree Learning

Overview

Entropy

Information Gain

Inductive Bias in ID3

Occam's Razor

Day	Outlook	Temperature	Humidity	Wind	PlayTennis
D1	Sunny	Hot	High	Weak	No
D2	Sunny	Hot	High	Strong	No
D3	Overcast	Hot	High	Weak	Yes
D4	Rain	Mild	High	Weak	Yes
D5	Rain	Cool	Normal	Weak	Yes
D6	Rain	Cool	Normal	Strong	No
D7	Overcast	Cool	Normal	Strong	Yes
D8	Sunny	Mild	High	Weak	No
D9	Sunny	Cool	Normal	Weak	Yes
D10	Rain	Mild	Normal	Weak	Yes
D11	Sunny	Mild	Normal	Strong	Yes
D12	Overcast	Mild	High	Strong	Yes
D13	Overcast	Hot	Normal	Weak	Yes
D14	Rain	Mild	High	Strong	No



# Selecting the Next Attribute

## Decision Tree Learning

Overview

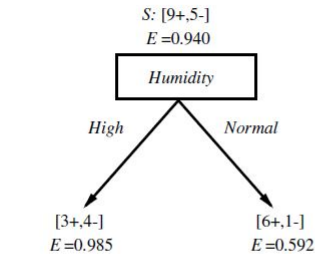
Entropy

Information Gain

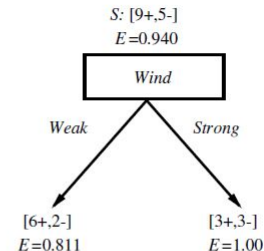
Inductive Bias in ID3

Occam's Razor

Which attribute is the best classifier?



$$\begin{aligned} \text{Gain}(S, \text{Humidity}) &= .940 - (7/14) \cdot .985 - (7/14) \cdot .592 \\ &= .151 \end{aligned}$$



$$\begin{aligned} \text{Gain}(S, \text{Wind}) &= .940 - (8/14) \cdot .811 - (6/14) \cdot 1.0 \\ &= .048 \end{aligned}$$

# Selecting the Next Attribute

## Decision Tree Learning

### Overview

### Entropy

### Information Gain

### Inductive Bias in ID3

### Occam's Razor

Which attribute is the best classifier?

$$\text{Gain}(S, \text{Outlook}) = 0.246$$

$$\text{Gain}(S, \text{Humidity}) = 0.151$$

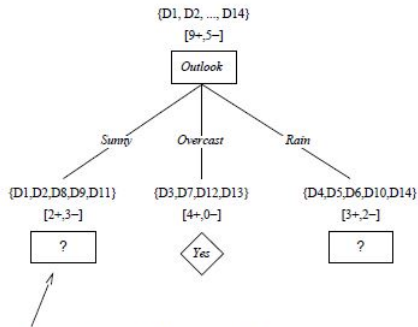
$$\text{Gain}(S, \text{Wind}) = 0.048$$

$$\text{Gain}(S, \text{Temperature}) = 0.029$$

**Note:** Process of selecting new attribute continues for each new leaf node until either of two conditions is met:

- (1) every attribute has already been included along this path through the tree, or
- (2) the training examples associated with this leaf node all have the same target attribute value

## Which attribute to be tested



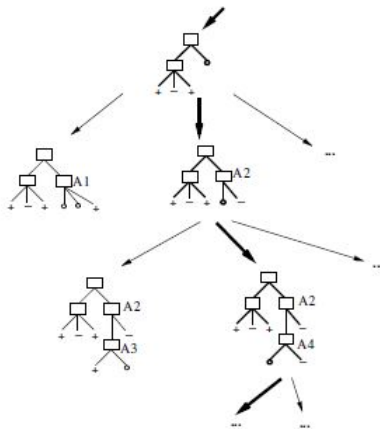
$$S_{\text{sunny}} = \{D1, D2, D8, D9, D11\}$$

$$\text{Gain}(S_{\text{sunny}}, \text{Humidity}) = .970 - (3/5) 0.0 - (2/5) 0.0 = .970$$

$$\text{Gain}(S_{\text{sunny}}, \text{Temperature}) = .970 - (2/5) 0.0 - (2/5) 1.0 - (1/5) 0.0 = .570$$

$$\text{Gain}(S_{\text{sunny}}, \text{Wind}) = .970 - (2/5) 1.0 - (3/5) .918 = .019$$

## Hypothesis Space Search by ID3



# Hypothesis Space Search by ID3

## Decision Tree Learning

### Overview

### Entropy

### Information Gain

### Inductive Bias in ID3

### Occam's Razor

- Hypothesis space is complete!
  - Target function surely in there...
- Outputs a single hypothesis (which one?)
  - Can't play 20 questions...
- No back tracking
  - Local minima...
- Statistically-based search choices
  - Robust to noisy data...
- Inductive bias: approx “prefer shortest tree”

# Inductive Bias in ID3

## Decision Tree Learning

### Overview

### Entropy

### Information Gain

### Inductive Bias in ID3

### Occam's Razor

Note  $H$  is the power set of instances  $X$

→ Unbiased?

Not really...

- Preference for short trees, and for those with high information gain attributes near the root
- Bias is a *preference* for some hypotheses, rather than a *restriction* of hypothesis space  $H$
- Occam's razor: prefer the shortest hypothesis that fits the data

# Occam's Razor

Decision Tree  
Learning

Overview

Entropy

Information  
Gain

Inductive Bias  
in ID3

Occam's  
Razor

## Why prefer short hypotheses?

### Argument in favor:

- Fewer short hyps. than long hyps.
- a short hyp that fits data unlikely to be coincidence
- a long hyp that fits data might be coincidence

### Argument opposed:

- There are many ways to define small sets of hyps
- e.g., all trees with a prime number of nodes that use attributes beginning with "Z"
- What's so special about small sets based on size of hypothesis??

# Overfitting in Decision Tree Learning

## Decision Tree Learning

Overview

Entropy

Information Gain

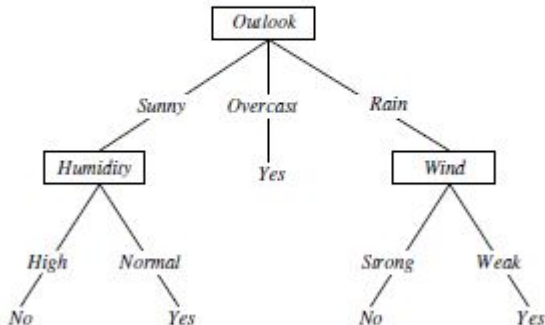
Inductive Bias in ID3

Occam's Razor

Consider adding noisy training example #15:

*Sunny, Hot, Normal, Strong, PlayTennis = No*

What effect on earlier tree?





# Avoiding Overfitting

## Decision Tree Learning

### Overview

### Entropy

### Information Gain

### Inductive Bias in ID3

### Occam's Razor

How can we avoid overfitting?

- stop growing when data split not statistically significant
- grow full tree, then post-prune

How to select “best” tree:

- Measure performance over training data
- Measure performance over separate validation data set
- MDL: minimize
$$size(tree) + size(misclassifications(tree))$$

# Reduced-Error Pruning

## Decision Tree Learning

### Overview

### Entropy

### Information Gain

### Inductive Bias in ID3

### Occam's Razor

Split data into *training* and *validation* set

Do until further pruning is harmful:

- Evaluate impact on *validation* set of pruning each possible node (plus those below it)
- Greedily remove the one that most improves *validation* set accuracy
- produces smallest version of most accurate subtree
- What if data is limited?

# Rule Post-Pruning

## Decision Tree Learning

Overview

Entropy

Information Gain

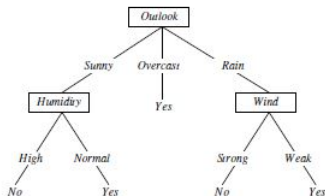
Inductive Bias in ID3

Occam's Razor

- Convert tree to equivalent set of rules
- Prune each rule independently of others
- Sort final rules into desired sequence for use

# Converting A Tree to Rules

## Decision Tree Learning



IF  $(Outlook = Sunny) \wedge (Humidity = High)$   
THEN  $PlayTennis = No$

IF  $(Outlook = Sunny) \wedge (Humidity = Normal)$   
THEN  $PlayTennis = Yes$

# Continuous Valued Attributes

## Decision Tree Learning

### Overview

### Entropy

### Information Gain

### Inductive Bias in ID3

### Occam's Razor

Create a discrete attribute to test continuous

- $Temperature = 82.5$
- $(Temperature > 72.3) = t, f$

<i>Temperature:</i>	40	48	60	72	80	90
<i>PlayTennis:</i>	No	No	Yes	Yes	Yes	No

# Attributes with Many Values

Problem:

- If attribute has many values, *Gain* will select it
- Imagine using *Date = Jun\_3\_1996* as attribute

One approach: use *GainRatio* instead

$$\text{GainRatio}(S, A) \equiv \frac{\text{Gain}(S, A)}{\text{SplitInformation}(S, A)}$$

$$\text{SplitInformation}(S, A) \equiv - \sum_{i=1}^c \frac{|S_i|}{|S|} \log_2 \frac{|S_i|}{|S|}$$

where  $S_i$  is subset of  $S$  for which  $A$  has value  $v_i$

# Attributes with Costs

Consider

- medical diagnosis, *BloodTest* has cost \$150
- robotics, *Width\_from\_1ft* has cost 23 sec.

How to learn a consistent tree with low expected cost?

One approach: replace gain by

- Tan and Schlimmer (1990)

$$\frac{Gain^2(S, A)}{Cost(A)}.$$

- Nunez (1988)

$$\frac{2^{Gain(S,A)} - 1}{(Cost(A) + 1)^w}$$

where  $w \in [0, 1]$  determines importance of cost

# Unknown Attribute Values

## Decision Tree Learning

### Overview

### Entropy

### Information Gain

### Inductive Bias in ID3

### Occam's Razor

What if some examples missing values of  $A$ ?

Use training example anyway, sort through tree

- If node  $n$  tests  $A$ , assign most common value of  $A$  among other examples sorted to node  $n$
- assign most common value of  $A$  among other examples with same target value
- assign probability  $p_i$  to each possible value  $v_i$  of  $A$ 
  - assign fraction  $p_i$  of example to each descendant in tree