



# Machine Learning, Decision Trees, Overfitting

**Recommended reading: Mitchell, Chapter 3**

Machine Learning 10-701

Tom M. Mitchell  
Center for Automated Learning and Discovery  
Carnegie Mellon University

September 13, 2005

# Machine Learning:

Study of algorithms that

- improve their performance
- at some task
- with experience

# Learning to Predict Emergency C-Sections

[Sims et al., 2000]

Data:

<i>Patient103</i> time=1	<i>Patient103</i> time=2	<i>Patient103</i> time=n
Age: 23	Age: 23	Age: 23
FirstPregnancy: no	FirstPregnancy: no	FirstPregnancy: no
Anemia: no	Anemia: no	Anemia: no
Diabetes: no	Diabetes: YES	Diabetes: no
PreviousPrematureBirth: no	PreviousPrematureBirth: no	PreviousPrematureBirth: no
Ultrasound: ?	Ultrasound: abnormal	Ultrasound: ?
Elective C-Section: ?	Elective C-Section: no	Elective C-Section: no
Emergency C-Section: ?	Emergency C-Section: ?	<b>Emergency C-Section: Yes</b>
...	...	...

9714 patient records,  
each with 215 features

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:  $26/41 = .63$ ,

Over test data:  $12/20 = .60$

# Object Detection

(Prof. H. Schneiderman)



Example training images  
for each orientation



# Text Classification



Company home page

VS

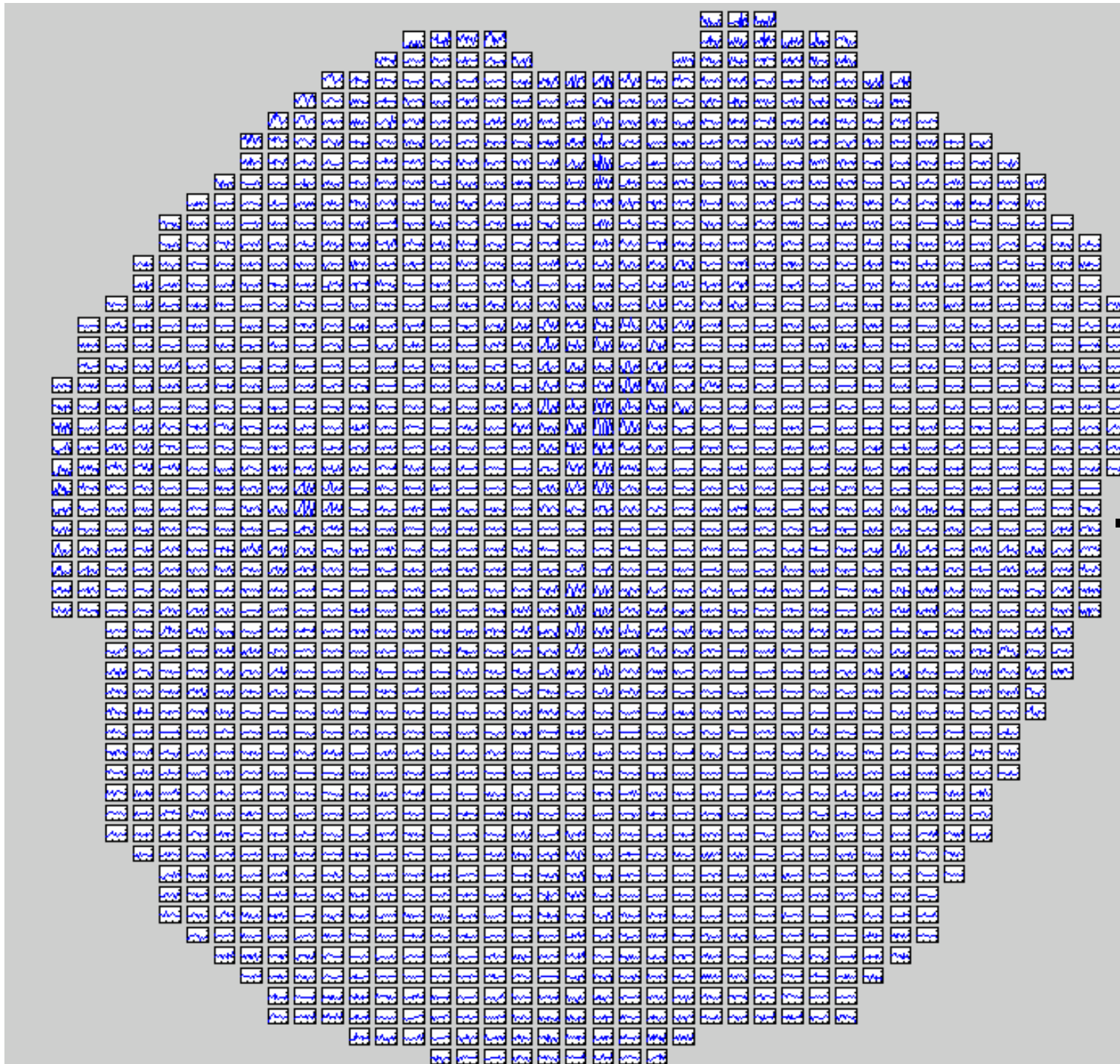
Personal home page

VS

Univeristy home page

VS

...



Reading  
a noun  
(vs verb)

[Rustandi et al.,  
2005]

# Growth of Machine Learning

---

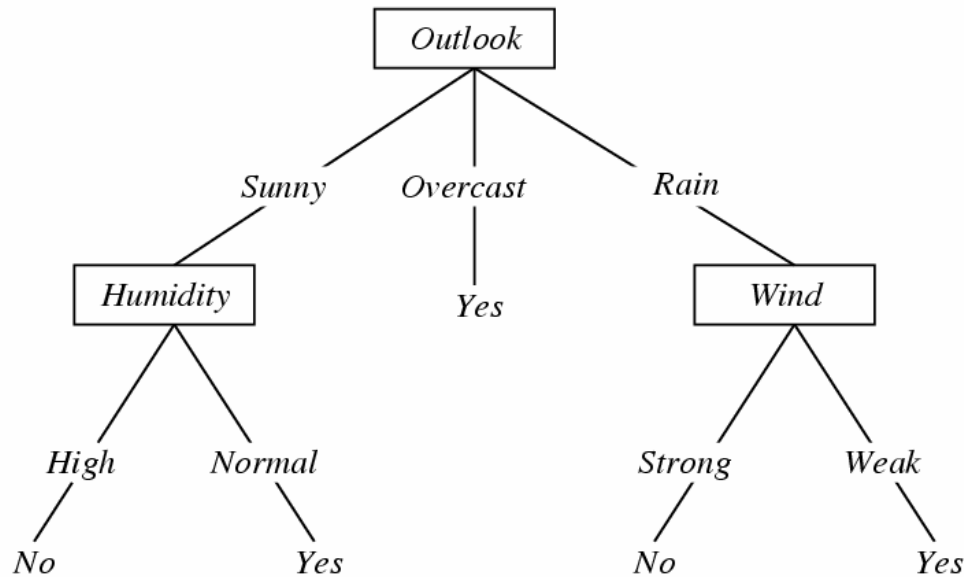
- Machine learning is preferred approach to
  - Speech recognition, Natural language processing
  - Computer vision
  - Medical outcomes analysis
  - Robot control
  - ...
- This trend is accelerating
  - Improved machine learning algorithms
  - Improved data capture, networking, faster computers
  - Software too complex to write by hand
  - New sensors / IO devices
  - Demand for self-customization to user, environment

# Decision tree learning



## Decision Tree for *PlayTennis*

---



How would you  
represent  
 $AB \vee CD(\neg E)$ ?

Each internal node: test one attribute  $X_i$

Each branch from a node: selects one value for  $X_i$

Each leaf node: predict  $Y$  (or  $P(Y|X \in \text{leaf})$ )

# A Tree to Predict C-Section Risk

---

Learned from medical records of 1000 women

Negative examples are C-sections

```
[833+,167-] .83+ .17-
Fetal_Presentation = 1: [822+,116-] .88+ .12-
| Previous_Csection = 0: [767+,81-] .90+ .10-
| | Primiparous = 0: [399+,13-] .97+ .03-
| | Primiparous = 1: [368+,68-] .84+ .16-
| | | Fetal_Distress = 0: [334+,47-] .88+ .12-
| | | | Birth_Weight < 3349: [201+,10.6-] .95+ .05-
| | | | Birth_Weight >= 3349: [133+,36.4-] .78+ .22-
| | | Fetal_Distress = 1: [34+,21-] .62+ .38-
| Previous_Csection = 1: [55+,35-] .61+ .39-
Fetal_Presentation = 2: [3+,29-] .11+ .89-
Fetal_Presentation = 3: [8+,22-] .27+ .73-
```

# Top-Down Induction of Decision Trees

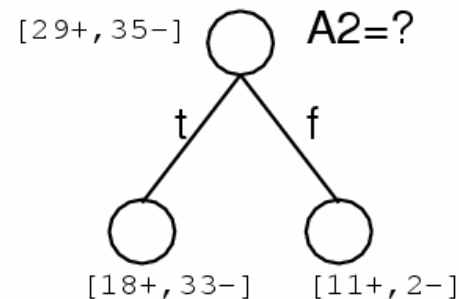
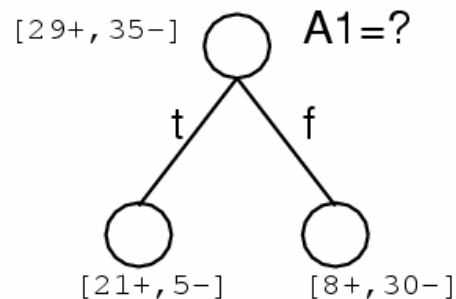
[ID3, C4.5, ...]

***node*** = Root

Main loop:

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

Which attribute is best?



# Entropy

Entropy  $H(X)$  of a random variable  $X$

$$H(X) = - \sum_{i=1}^n P(X = i) \log_2 P(X = i)$$

$H(X)$  is the expected number of bits needed to encode a randomly drawn value of  $X$  (under most efficient code)

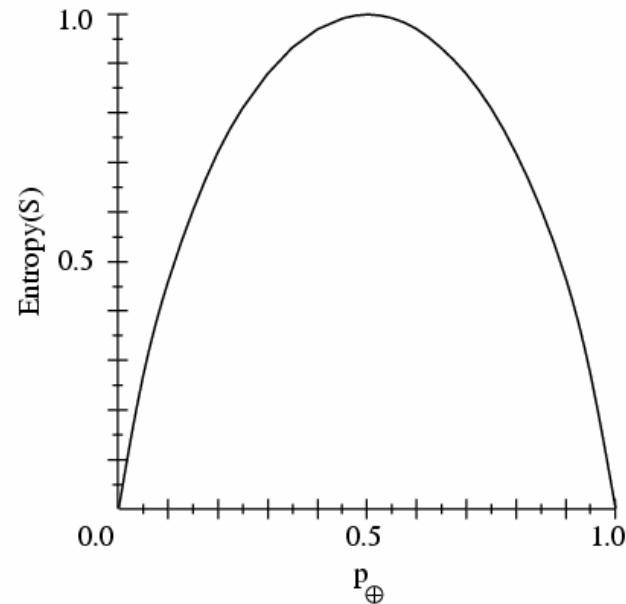
Why? Information theory:

- Most efficient code assigns  $-\log_2 P(X=i)$  bits to encode the message  $X=i$
- So, expected number of bits is:

$$\sum_{i=1}^n P(X = i)(-\log_2 P(X = i))$$

# Sample Entropy

---



- $S$  is a sample of training examples
- $p_{\oplus}$  is the proportion of positive examples in  $S$
- $p_{\ominus}$  is the proportion of negative examples in  $S$
- Entropy measures the impurity of  $S$

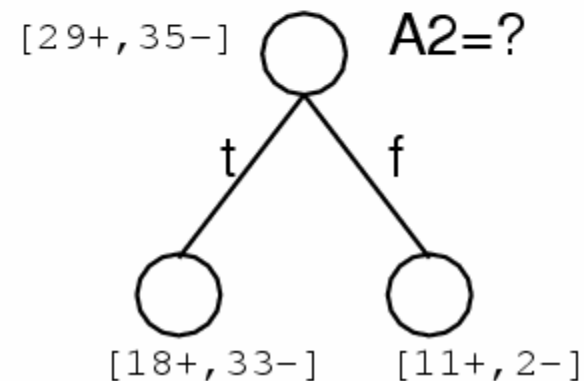
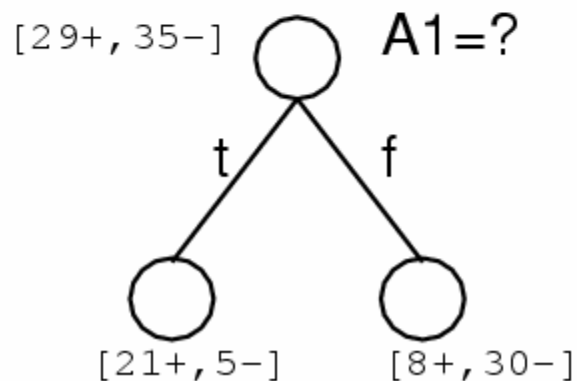
$$Entropy(S) \equiv -p_{\oplus} \log_2 p_{\oplus} - p_{\ominus} \log_2 p_{\ominus}$$

# Information Gain

---

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

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



# Training Examples

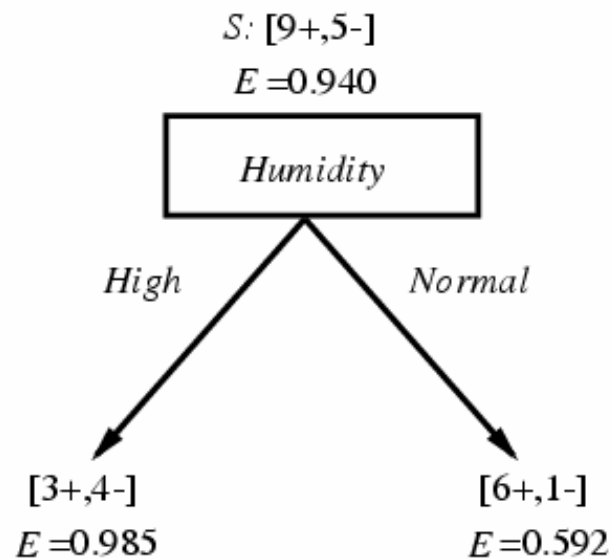
---

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

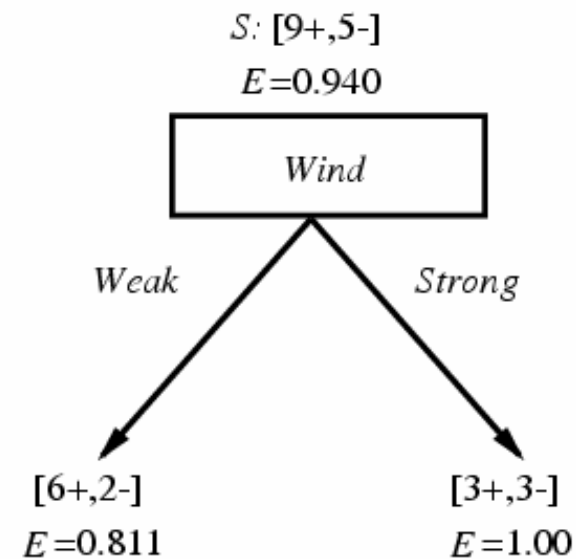
# Selecting the Next Attribute

---

Which attribute is the best classifier?

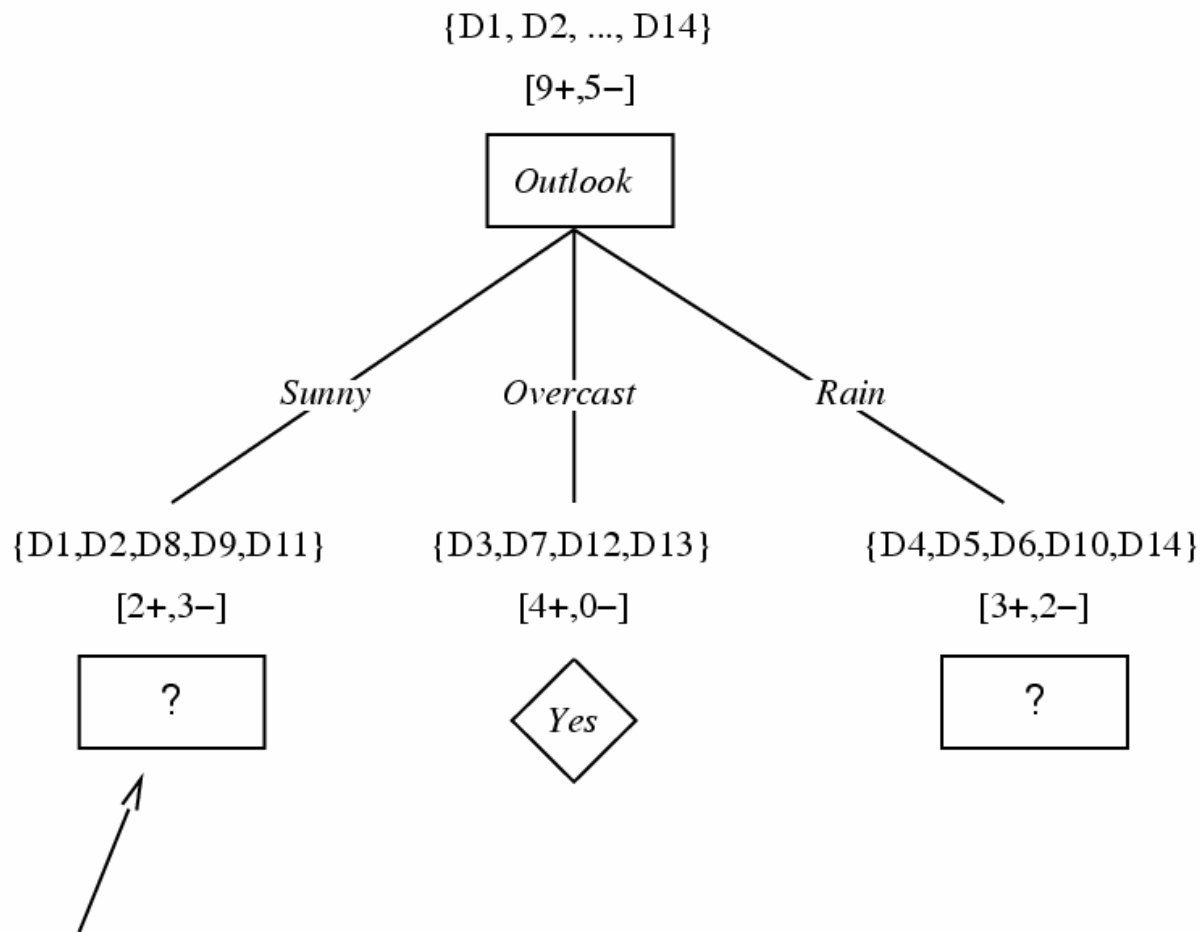


$$\begin{aligned} \text{Gain}(S, \text{Humidity}) &= .940 - (7/14).985 - (7/14).592 \\ &= .151 \end{aligned}$$



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





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

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

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

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

# 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??

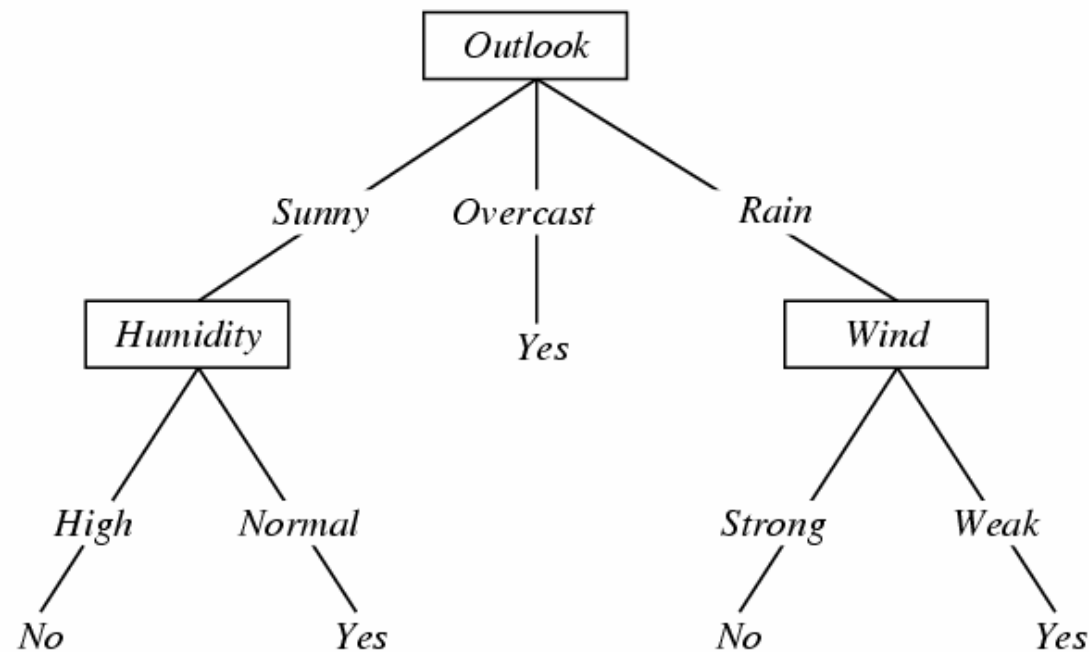
# Overfitting in Decision Trees

---

Consider adding noisy training example #15:

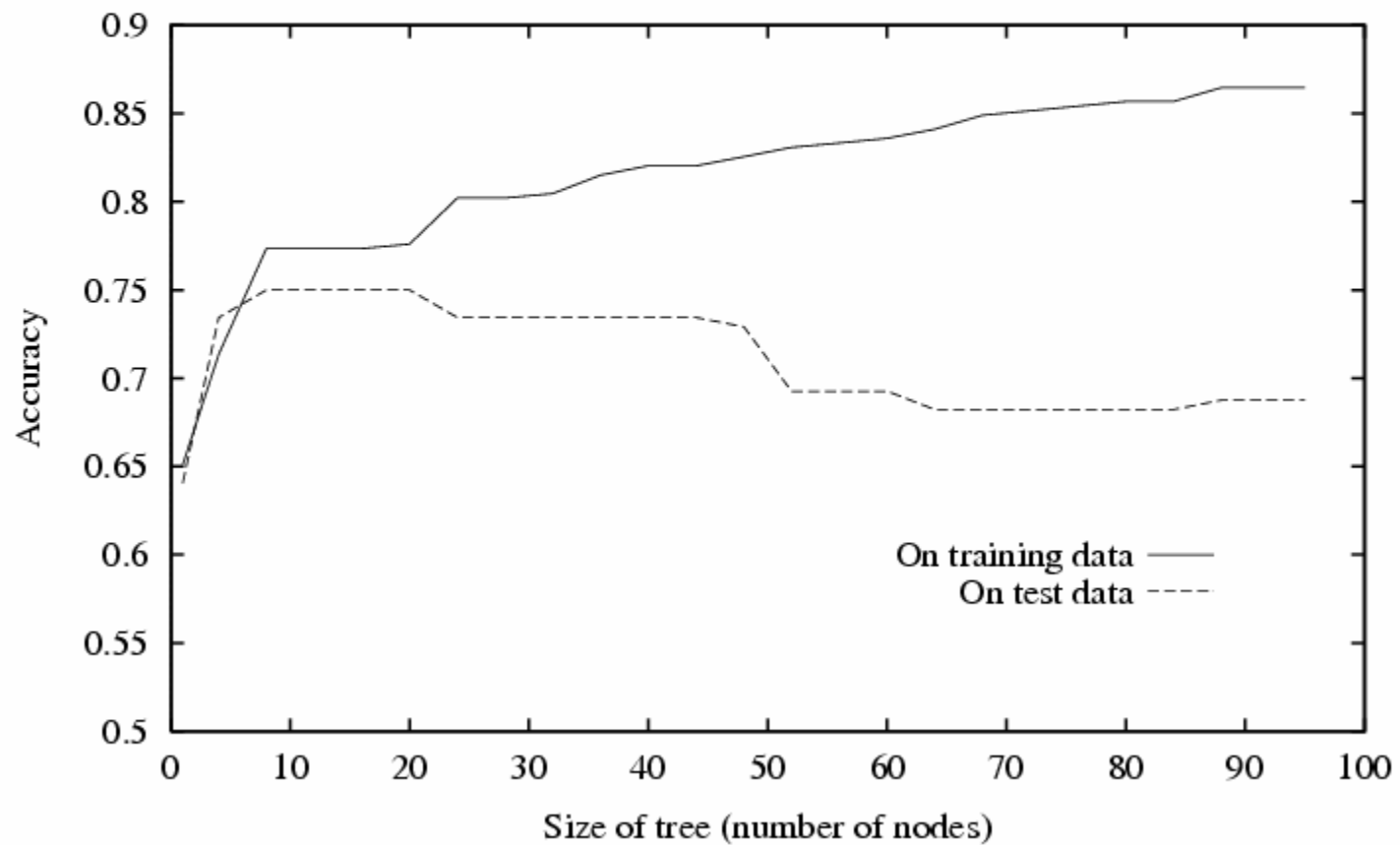
*Sunny, Hot, Normal, Strong, PlayTennis = No*

What effect on earlier tree?



# Overfitting in Decision Tree Learning

---



# Avoiding Overfitting

---

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  
 $size(tree) + size(misclassifications(tree))$

# Minimum Description Length Principle

---

Occam's razor: prefer the shortest hypothesis

MDL: prefer the hypothesis  $h$  that minimizes

$$h_{MDL} = \operatorname{argmin}_{h \in H} L_{C_1}(h) + L_{C_2}(D|h)$$

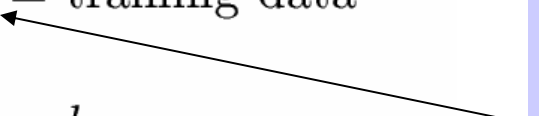
where  $L_C(x)$  is the description length of  $x$  under encoding  $C$

---

Example:  $H$  = decision trees,  $D$  = training data labels

- $L_{C_1}(h)$  is # bits to describe tree  $h$
- $L_{C_2}(D|h)$  is # bits to describe  $D$  given  $h$ 
  - Note  $L_{C_2}(D|h) = 0$  if examples classified perfectly by  $h$ . Need only describe exceptions
- Hence  $h_{MDL}$  trades off tree size for training errors

Assume  $X$   
values known,  
labels  $Y$   
encoded



# Minimum Description Length Principle

---

$$\begin{aligned} h_{MAP} &= \arg \max_{h \in H} P(D|h)P(h) \\ &= \arg \max_{h \in H} \log_2 P(D|h) + \log_2 P(h) \\ &= \arg \min_{h \in H} -\log_2 P(D|h) - \log_2 P(h) \quad (1) \end{aligned}$$

Interesting fact from information theory:

The optimal (shortest expected coding length) code for an event with probability  $p$  is  $-\log_2 p$  bits.

So interpret (1):

- $-\log_2 P(h)$  is length of  $h$  under optimal code
- $-\log_2 P(D|h)$  is length of  $D$  given  $h$  under optimal code

→ prefer the hypothesis that minimizes

$$\text{length}(h) + \text{length}(\text{misclassifications})$$

# Reduced-Error Pruning

---

Split data into *training* and *validation* set

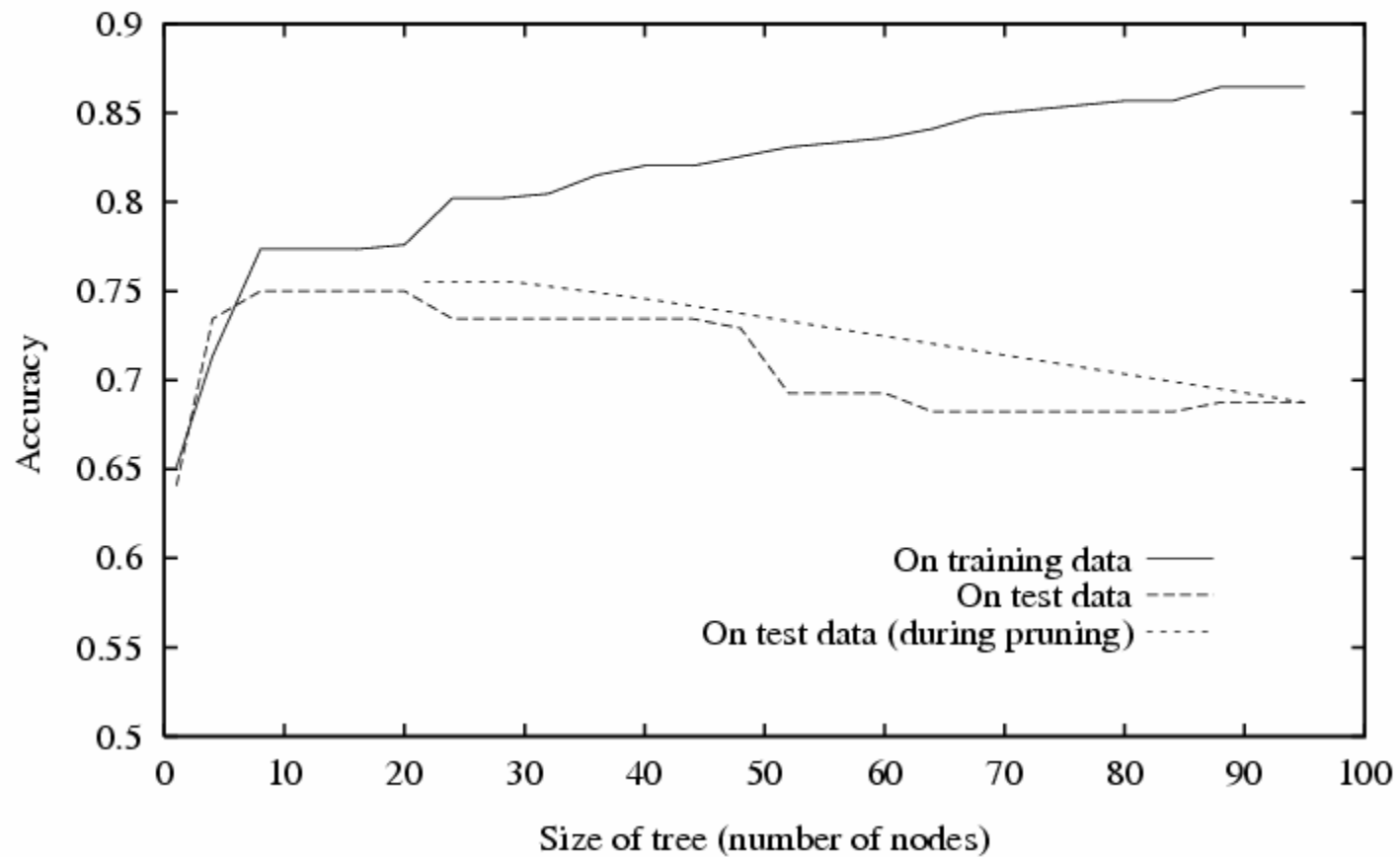
Do until further pruning is harmful:

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



# Effect of Reduced-Error Pruning

---



# Rule Post-Pruning

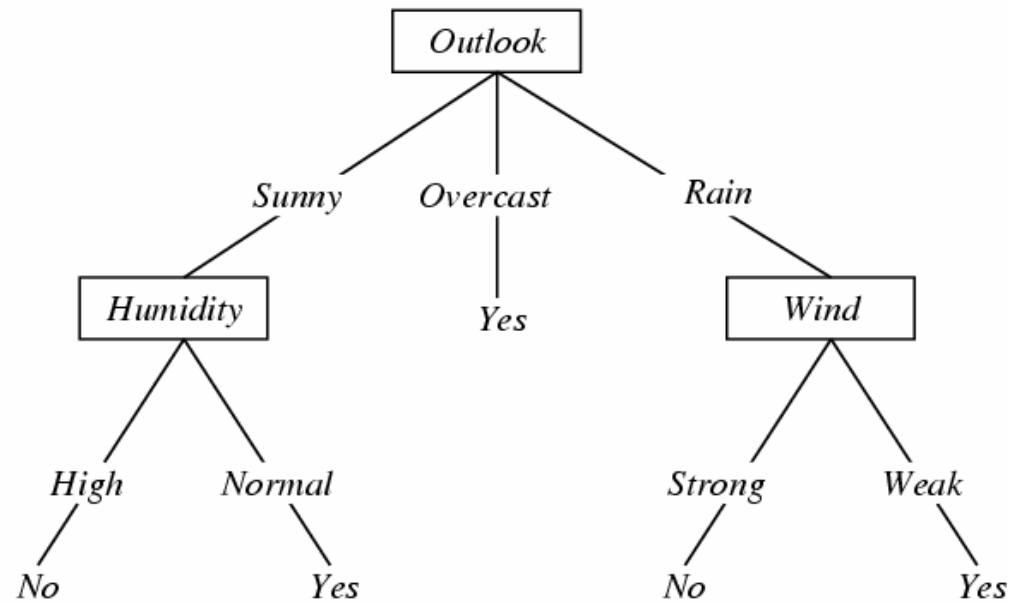
---

1. Convert tree to equivalent set of rules
2. Prune each rule independently of others
3. Sort final rules into desired sequence for use

Perhaps most frequently used method (e.g., C4.5)

# Converting A Tree to Rules

---



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

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

# Continuous Valued Attributes

---

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

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

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

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

# Unknown Attribute Values

---

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

Classify new examples in same fashion

# What you should know:

---

- Well posed function approximation problems:
  - Instance space,  $X$
  - Sample of labeled training data,  $D = \{ \langle x_i, y_i \rangle \}$
  - Hypothesis space,  $H = \{ f: X \rightarrow Y \}$
- Learning is a search/optimization problem over  $H$ 
  - Various objective functions
  - Today: minimize training error (0-1 loss)
- Decision tree learning
  - Greedy top-down learning of decision trees (ID3, C4.5, ...)
  - Overfitting and tree/rule post-pruning
  - Extensions...