

SUPERVISED LEARNING

Progress Report

- We've finished Part I: Problem Solving
- We've finished Part II: Reasoning with uncertainty!
- Part III: (Machine) Learning
 - ▣ Supervised Learning
 - ▣ Unsupervised Learning
 - ▣ (Reinforcement Learning)
- Overlaps quite a bit with Part II

Today

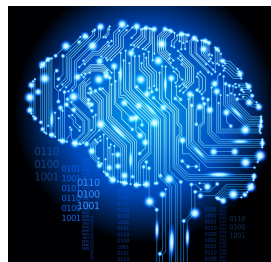
- Reading
 - ▣ We're skipping to AIMA Chapter 18!
 - ▣ AIMA 18.1-18.4

- Goals
 - ▣ Supervised Learning
 - ▣ Naïve Bayes
 - ▣ Decision Trees

Machine Learning

- The term “machine learning” is a bit misleading
 - ▣ Pattern recognition

- We can use machine learning to
 - learn the probabilities for a BN
 - learn the topology of a BN
 - learn heuristic function for games



Subfields of Machine Learning

- Supervised learning
 - ▣ learning with labels
 - ▣ classification, regression, structured prediction
- Unsupervised learning
 - ▣ learning without labels
 - ▣ clustering, projection methods
- Reinforcement learning
 - ▣ learning with rewards
 - ▣ planning

Supervised Learning Terminology

- data set
- instance, input
- features
- label, output
- hypothesis
- hypothesis class
- realizable, consistent

Types of Supervised Learning Tasks

- **Regression**
 - y is a (vector of) real-valued number(s)
 - e.g. price of a commodity, pollution levels, brain activity
- **Classification**
 - y is a discrete (categorical) value
 - e.g. spam or not spam, 5-star ratings
- **Structured prediction**
 - y is a structured object
 - e.g. given sentence predict parse tree, given words in a sentence predict POS tags

Types of Supervised Learning Tasks

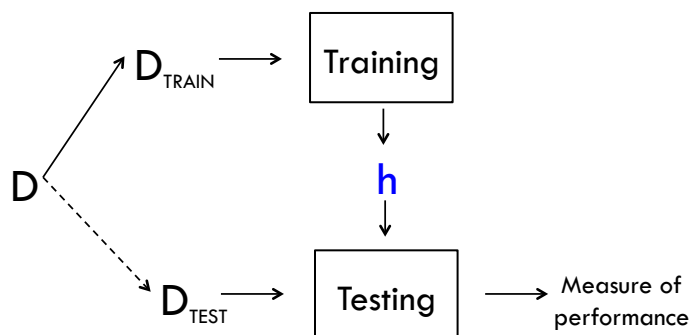
- Supervised learning
 - Spam
 - Digit recognition
 - Rainfall levels in India
 - Pollution index
 - Stock returns
 - User's ratings of movies
 - Genre classification
 - Sentiment analysis
 - Document classification
 - Image recognition
 - Part-of-speech
 - Storm trajectories

0 0 0 0 0 0 0 0 0 0 0 0
 1 1 1 1 1 1 1 1 1 1 1 1
 2 2 2 2 2 2 2 2 2 2 2 2
 3 3 3 3 3 3 3 3 3 3 3 3
 4 4 4 4 4 4 4 4 4 4 4 4
 5 5 5 5 5 5 5 5 5 5 5 5
 6 6 6 6 6 6 6 6 6 6 6 6
 7 7 7 7 7 7 7 7 7 7 7 7
 8 8 8 8 8 8 8 8 8 8 8 8
 9 9 9 9 9 9 9 9 9 9 9 9



So what is learning?

- **Learning** is the process of finding (constructing, searching for) a hypothesis that performs well on the **training data** and **generalizes** well to unseen data (the **test data**)

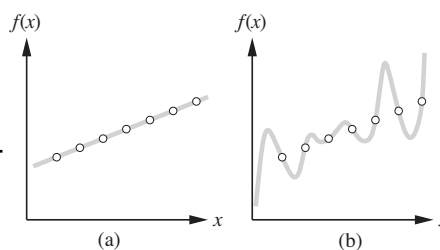


Ockham's Razor

- Ockham's Razor
 - ▣ Prefer the simplest consistent hypothesis

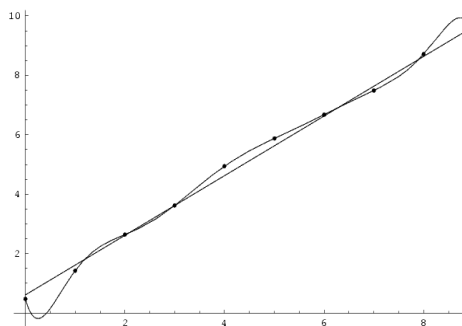
- Example: Curve fitting

- ▣ x is the x-coordinate
- ▣ y is the y-coordinate
- ▣ Both hypotheses are consistent
- ▣ Which is better?



Overfitting

- Overfitting
 - ▣ Learner fits itself to noise in the training data failing to generalize well
 - ▣ Causes: noisy data (too little data), overly complex models
- Example: Curve fitting
 - ▣ Which is better?



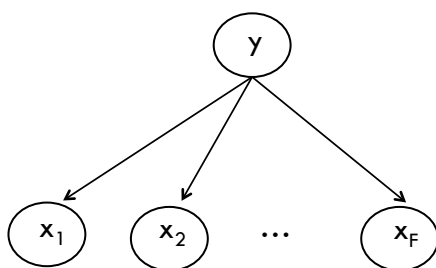
Common Supervised Learning Algorithms

- Graphical models
 - ▣ Naive Bayes classifiers
 - ▣ Bayesian networks
- Decision trees
 - ▣ Random forests (many decision trees)
- Neural Networks
 - ▣ Perceptrons
 - ▣ Artificial neural networks
 - ▣ Deep belief nets
- Max margin classifiers
 - ▣ Support vector machines
- Regression analysis
 - ▣ Logistic regression
 - ▣ Linear regression

Each of these algorithms makes assumptions – these assumptions are known as the **inductive bias**

Naïve Bayes Classifier

- Used for classification
 - ▣ x_i represent symptoms and $y = \{\text{Flu, Appendicitis, ...}\}$
 - ▣ x_i represent words and $y = \{\text{Politics, Sports, Finance, ...}\}$
- **Inductive bias:** features are conditionally independent given label



Naïve Bayes Classifier

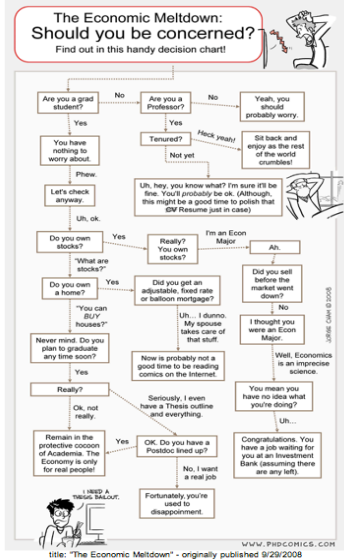
- We're interested in computing the quantity:

$$\begin{aligned}
 & p(Y = y | x_1, x_2, \dots, x_n) \\
 & \propto p(Y = y, x_1, x_2, \dots, x_n) \\
 & = p(x_1 \dots, x_n | Y = y) p(Y = y) \\
 & = p(Y = y) \prod_{i=1}^F p(x_i | Y = y)
 \end{aligned}$$

- How can we compute $p(Y=y)$ and $p(x_i | Y=y)$ from data set D ?

Decision Tree Classifier

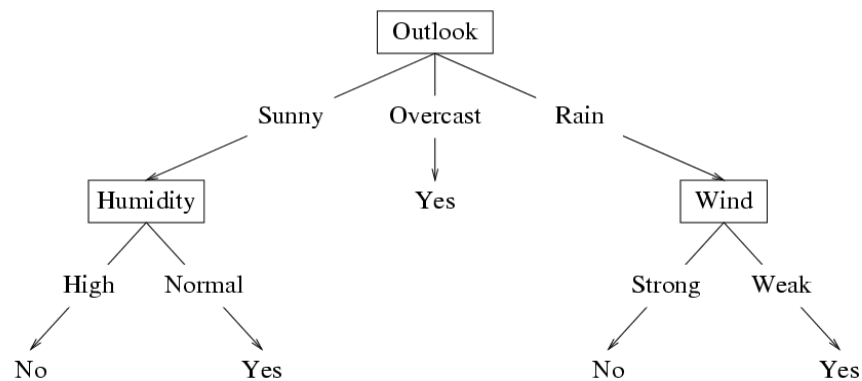
Piled Higher and Deeper by Jorge Cham www.phdcomics.com



Decision Tree Classifier

	Day	Outlook	Temp.	Humidity	Wind	PlayTennis	
$x_1 \rightarrow$	D1	Sunny	Hot	High	Weak	No	$\leftarrow Y_1$
$x_2 \rightarrow$	D2	Sunny	Hot	High	Strong	No	$\leftarrow Y_2$
$x_3 \rightarrow$	D3	Overcast	Hot	High	Weak	Yes	$\leftarrow Y_3$
	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	

Decision Tree Classifier

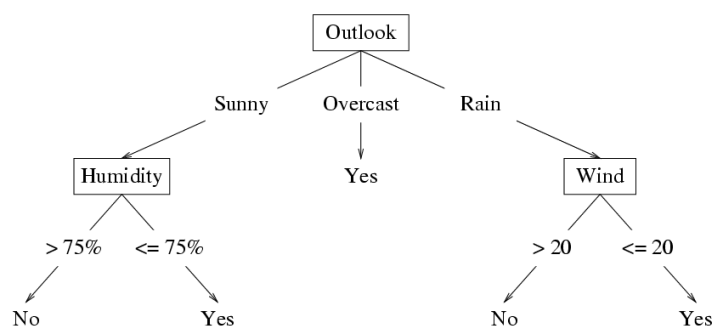


Decision Tree Classifier

- Decision trees are best suited to problems where
 - ▣ Each attribute is discrete
 - ▣ The label y is discrete
 - ▣ The hypothesis can be expressed using conjunctions (AND) and disjunctions (OR)
 - ▣ The training data may contain errors
 - ▣ The training data may contain missing attribute values

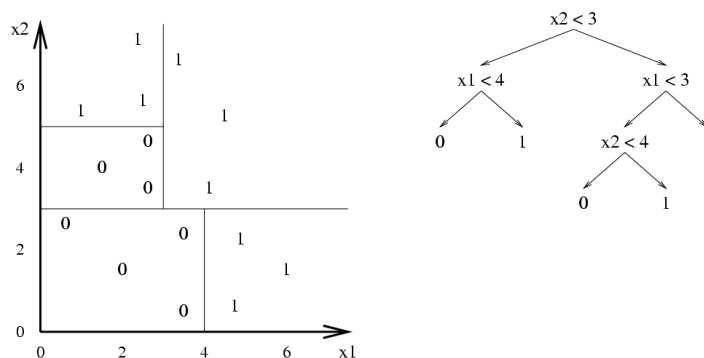
Decision Tree Classifier

- If the features are continuous, internal nodes may test the value of a feature against a threshold

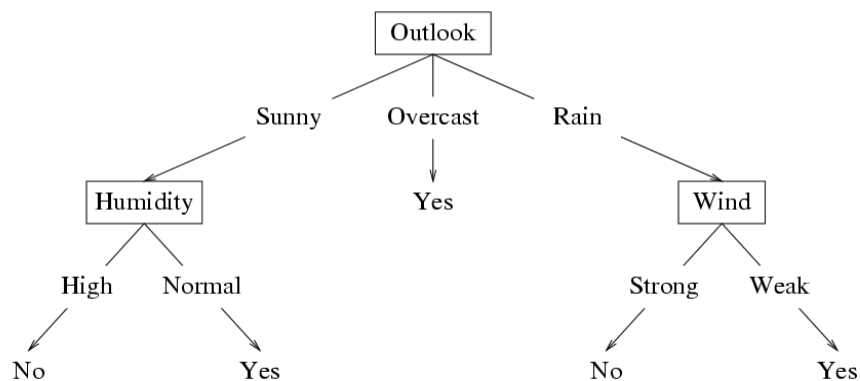


Decision Tree Classifier

- Learns axis-parallel decision boundaries, i.e. divides feature space into hyper-rectangles



Learning a Decision Tree



Learning a Decision Tree

```

function DECISION-TREE-LEARNING (examples, attributes, parents) returns a tree
if examples is empty return MAJORITY_VOTE(parents)
else if all examples have same classification return classification
else if attributes is empty return MAJORITY_VOTE(examples)
else
  A ← CHOOSE-BEST-ATTRIBUTE (examples)
  tree ← a new decision tree with root A
  for each value  $v_k$  of A
     $S_k$  ← examples with value  $v_k$  for attribute A
    subtree ← DECISION-TREE-LEARNING( $S_k$ , attributes-A, examples)
    add branch to tree with label ( $A=v_k$ ) and subtree
  return tree
  
```

Choosing the best attribute

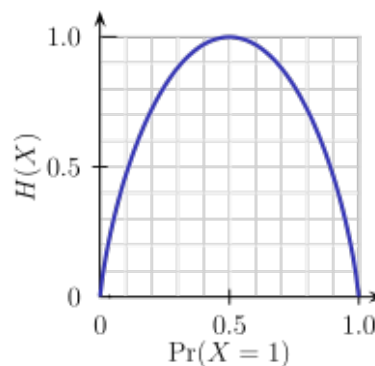
- Splitting on a **good** attribute
 - ▣ After the split, the examples at each branch have the same classification

- Splitting on a **bad** attribute
 - ▣ After the split, the examples at each branch have the same proportion of positive and negative examples

- We will use entropy and information gain to formalize what we mean by *good* and *bad* attributes

Entropy

- **Entropy** measures the uncertainty of a random variable
 - ▣ How many bits are needed to efficiently encode the possible values (outcomes) of a random variable?
- Introduced by Shannon in 1948 paper
- Example: flipping a coin
 - ▣ A completely biased coin requires 0 bits of entropy
 - ▣ A fair coin requires 1 bit of entropy
 - ▣ How many bits are need to encode the outcome of flipping a fair coin twice?



Entropy and Information Gain

- Let A be a random variable with values v_k
- Each value v_k occurs with probability $p(v_k)$
- Then the entropy of A is defined as

$$\begin{aligned}
 H(A) &= \sum_k p(v_k) \log_2 \left(\frac{1}{p(v_k)} \right) \\
 &= - \sum_k p(v_k) \log_2 p(v_k)
 \end{aligned}$$

- (Apply this notion of entropy to choosing the best attribute)

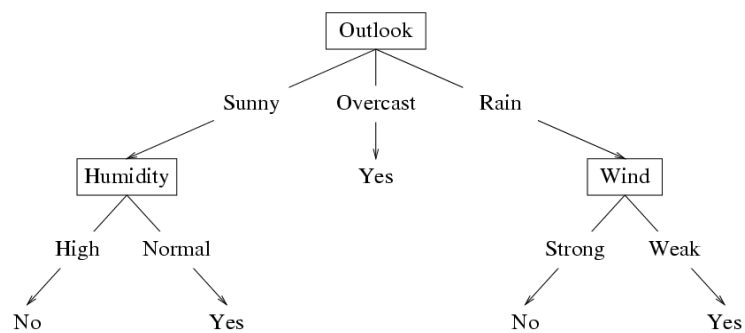
Entropy and Information Gain

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

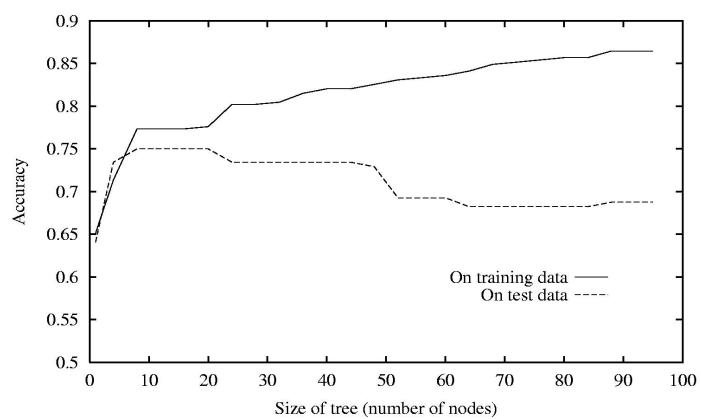
Day	Outlook	Temp.	Humidity	Wind	PlayTennis
D1	Sunny	Hot	High	Weak	No
D2	Sunny	Hot	High	Strong	No
D3	Overcast	Hot	High	Weak	Yes
D4	Rain	Mild	High	Weak	Yes
D5	Rain	Cool	Normal	Weak	Yes
D6	Rain	Cool	Normal	Strong	No
D7	Overcast	Cool	Normal	Strong	Yes
D8	Sunny	Mild	High	Weak	No
D9	Sunny	Cool	Normal	Weak	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

Decision Trees: additional considerations

- Overfitting can cause additional nodes to be added to the decision tree



Decision Trees: additional considerations



Decision Trees: additional considerations

- Overfitting
 - ▣ Can post-process the learned decision tree and prune using significance testing at final nodes
 - ▣ Cross-validation using validity set
- Continuous or integer-valued attributes
 - ▣ Use ranges
- Continuous label y
 - ▣ Combination of splitting and linear regression