

15-381: AI
Decision Tree Learning – Part I

October 18, 2009

Manuela Veloso
Chapter 18, Russell and Norvig
Thanks to all past instructors

Carnegie Mellon

Outline

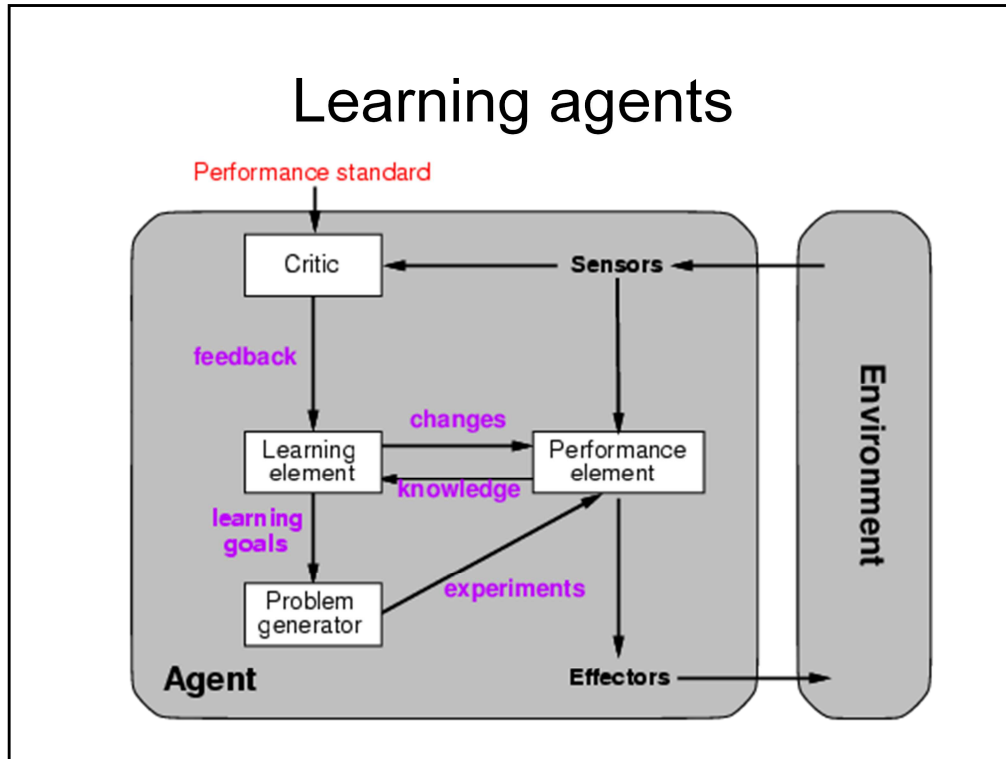
- Learning agents
- Inductive learning
- Decision tree learning

Learning

- Definition?...
- Gathering more knowledge
 - “Knowing more than was known before learning”
- Learning “substitutes” the need to model apriori
- Experience, feedback, refinement
- Learning modifies the agent's decision mechanisms to improve performance

In general, it is hard to define learning, as we don't really know what learning is.

Learning agents



In this example, the learning agent is something that learns by observing and interacting with the environment.

Learning “Element”

- A bit of “magic”:
 - Which **components** of the performance element are to be learned
 - What **feedback** is available to learn these components
 - What **representation** is used for the components
- Type of feedback:
 - **Supervised learning**: correct labels/answers
 - **Unsupervised learning**: correct labels/answers missing
 - **Reinforcement learning**: occasional rewards

Magic == lots of hacks

Need a hypothesis representation that is general enough to express what you want to learn, but specific enough that the search for the correct hypothesis isn't too large.

Inductive Learning

- Simplest form: learn a function from **examples**

f is the **target function**

An **example** is a pair $(x, f(x))$ – *supervised learning*

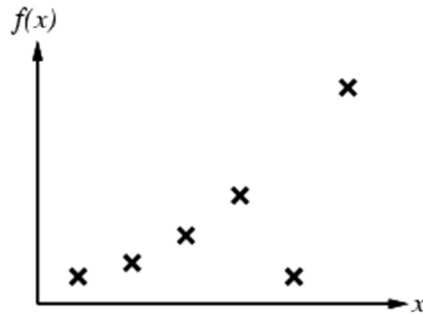
Problem: find a **hypothesis** h

such that $h \approx f$

given a **training set** of examples

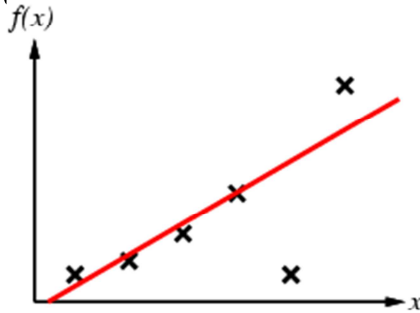
Inductive Learning Method

- 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:



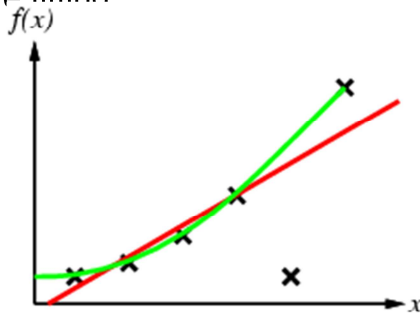
Inductive Learning Method

- 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
-



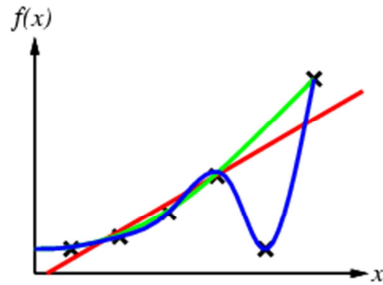
Inductive Learning Method

- 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
-



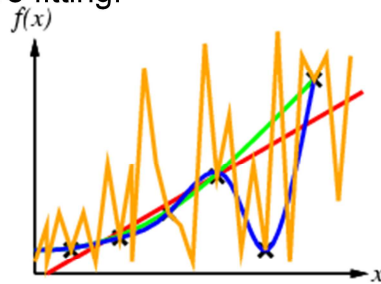
Inductive Learning Method

- 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

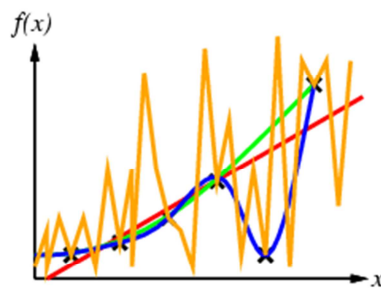
- Construct/adjust h to agree with f on training set
- h is **consistent** if it agrees with f on all examples
- E.g., curve fitting:



Given any number of points, it is possible to find a hypothesis that is consistent for every training point. However we want the hypothesis to also predict points that were not training points. This slide is an example of overfitting. Yes, it is consistent with ever training example, but it is a poor predictor of points not in the training data set.

Inductive Learning Method - Bias

- Hypothesis “form” – **bias** on the learning outcome
- Ockham’s razor: prefer the simplest hypothesis consistent with data



There is an implicit bias in the hypothesis's form. For example, if your hypothesis is a 4th degree polynomial, where you are learning the coefficient of each term, you can only represent curves of degree 4 and less.

Ockham's razor: when two hypothesis perform similarly, prefer the simpler one (ie if a line and a curve both fit points, then prefer a line).

Attribute-Based Data Sets

- Examples described by **attribute values** (Boolean, discrete, continuous)
- Function is class, wait/not wait for table at restaurant

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>Wait</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)

How should data be represented? Here, each row is one example, each column is an attribute.

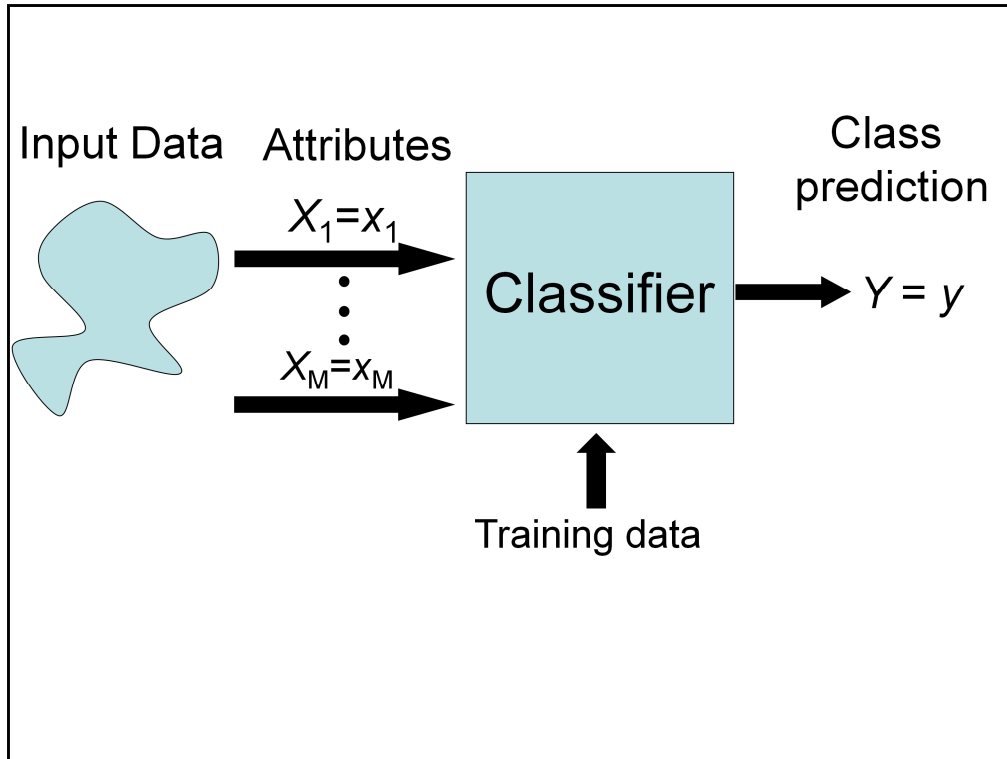
Each attribute takes a value. Boolean attributes take true/false, etc. Some attributes are discrete, some are continuous. As you'll see later on, continuous valued attributes can complicate things.

Attribute-Based Data Set

- *Type*: drama, comedy, thriller
- *Company*: MGM, Columbia
- *Director*: Bergman, Spielberg, Hitchcock
- *Mood*: stressed, relaxed, normal

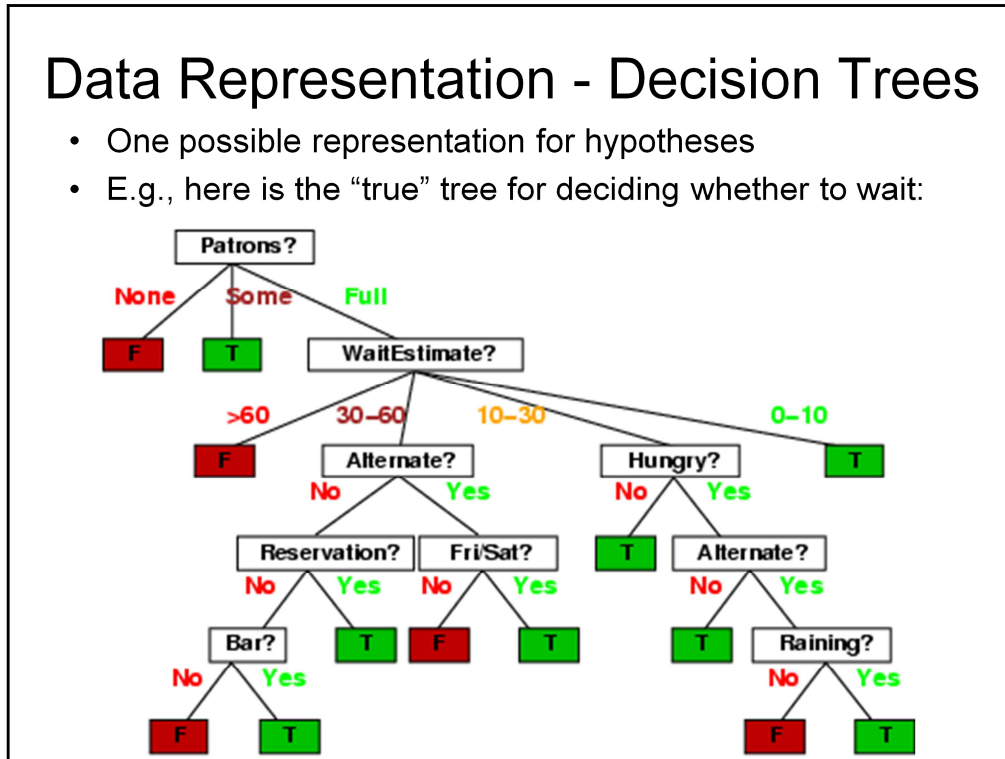
Movie	Type	Company	Director	Mood	Likes-movie?
m_1	thriller	MGM	Bergman	normal	No
m_2	comedy	Columbia	Spielberg	stressed	Yes
m_3	comedy	MGM	Spielberg	relaxed	No
m_4	thriller	MGM	Bergman	relaxed	No
m_5	comedy	MGM	Hitchcock	normal	Yes
m_6	drama	Columbia	Bergman	relaxed	Yes
m_7	drama	Columbia	Bergman	normal	No
m_8	drama	MGM	Spielberg	stressed	No
m_9	drama	MGM	Hitchcock	normal	Yes
m_{10}	comedy	Columbia	Spielberg	relaxed	No
m_{11}	thriller	MGM	Spielberg	normal	No
m_{12}	thriller	Columbia	Hitchcock	relaxed	No

A learner should learn that the type, director and your mood are important to whether or not you liked the movie, while the company is not so important



Data Representation - Decision Trees

- One possible representation for hypotheses
- E.g., here is the “true” tree for deciding whether to wait:



This is a decision tree to decide whether or not you are going to wait for a table at a restaurant.

Here, hypotheses are in the form of decision trees. There are many possible decision trees (hypotheses). You are searching for the simplest decision tree (hypothesis) that is most consistent with the training data

Is this example the simplest?

Expressiveness

- Decision trees can express any function of the input attributes
 - E.g., for Boolean functions, truth table row \rightarrow path to leaf
- Trivially, there is a consistent decision tree for any training set with one path to leaf for each example (unless f nondeterministic in x)
 - But... it probably won't generalize to new examples – goal of learning...
- Goal: Find more “compact” decision trees

Building a Decision Tree

- Three variables:
 - Hair = {blond, dark}
 - Height = {tall, short}
 - Country = {Gromland, Polvia}

Training data:

(B,T,P)

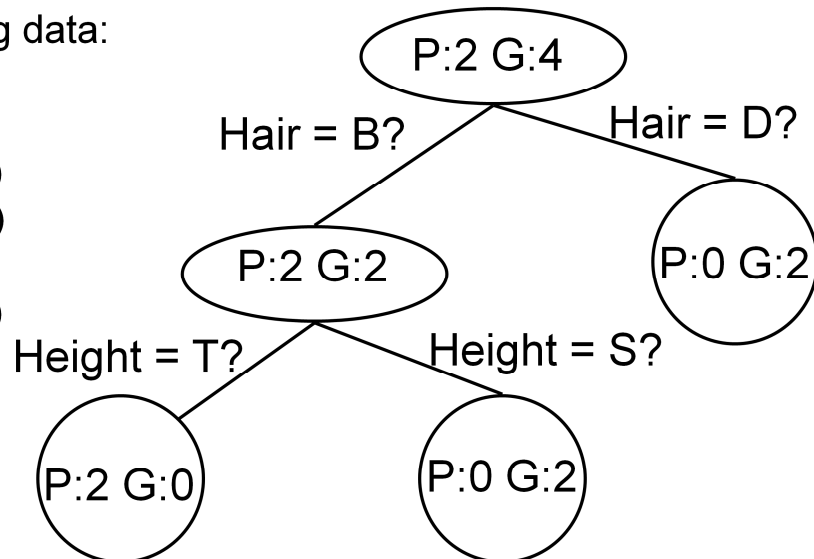
(B,T,P)

(B,S,G)

(D,S,G)

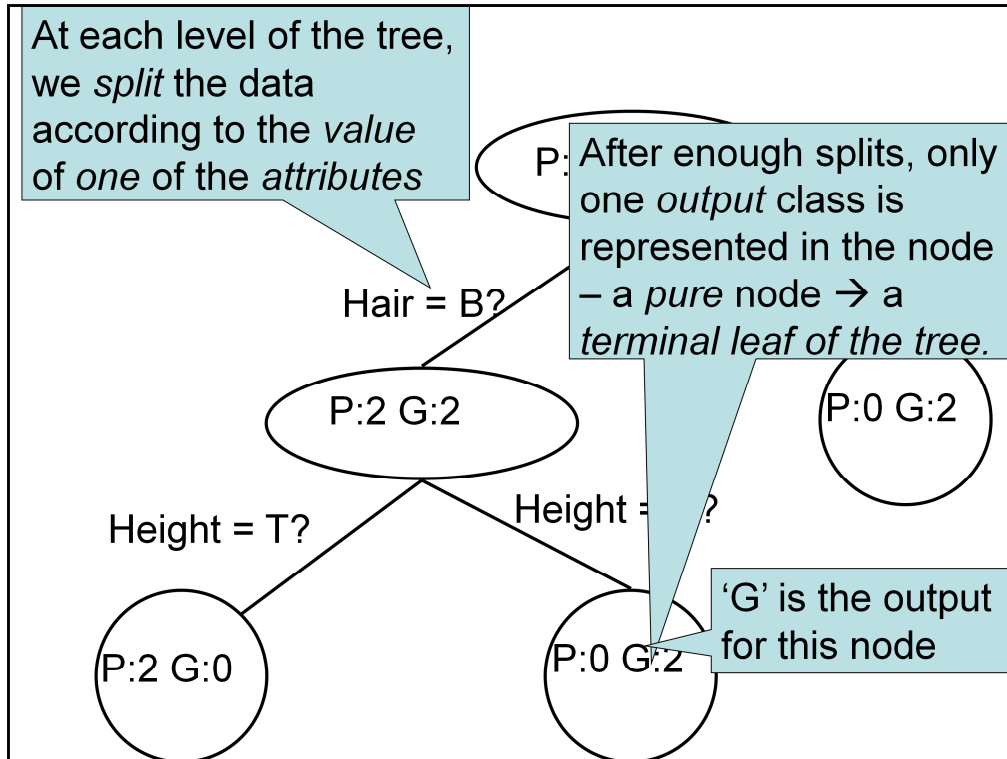
(D,T,G)

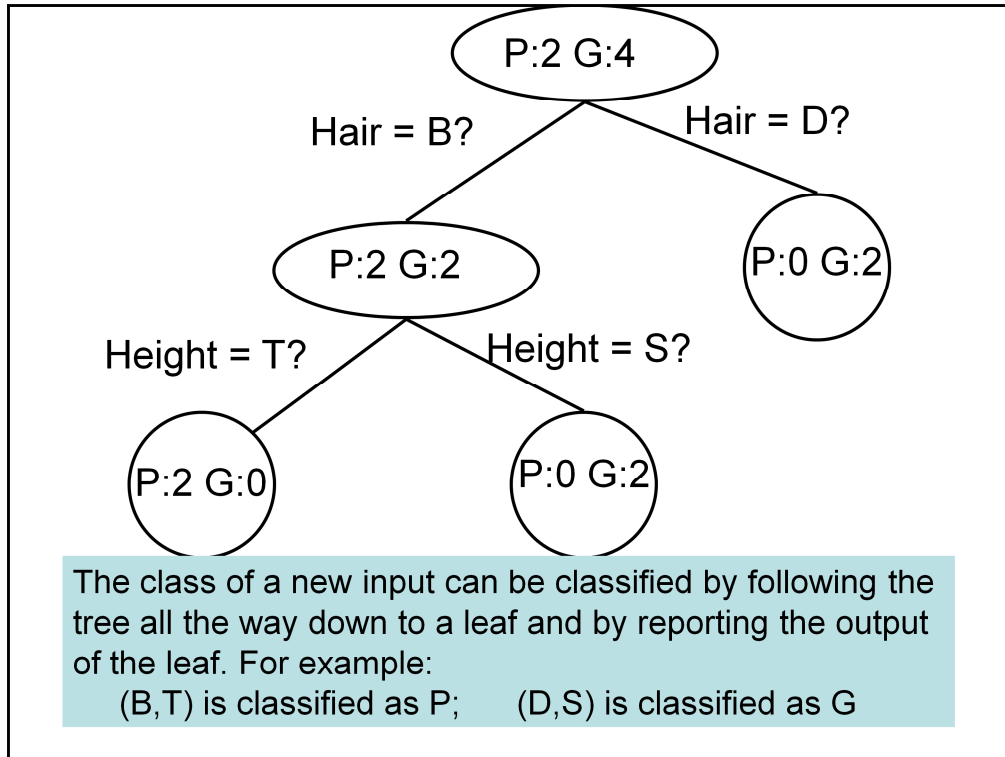
(B,S,G)



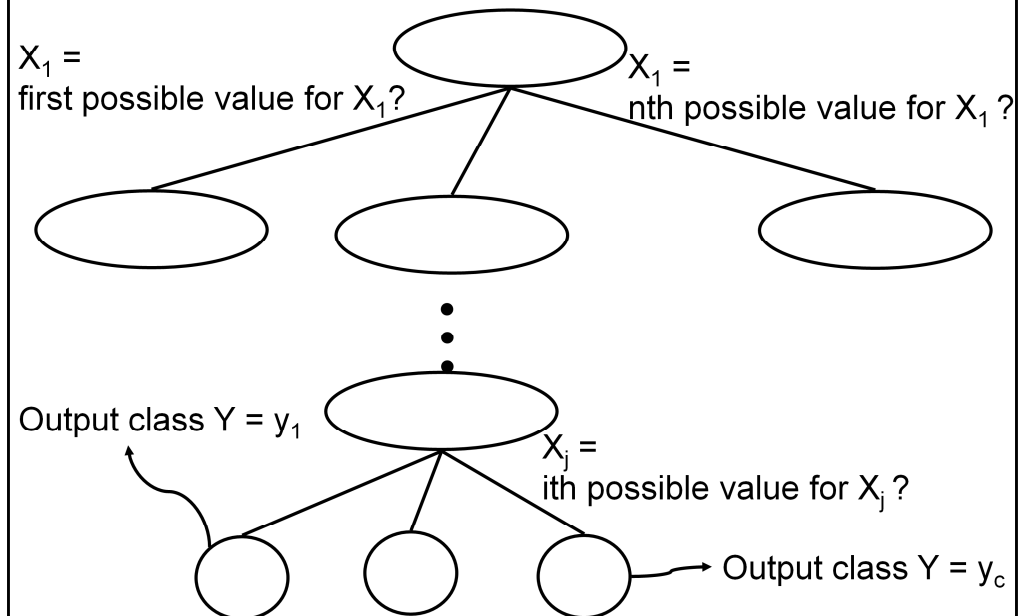
Trying to classify examples as P or G based on their attributes

If tree height is what determines which hypothesis is simpler, there is no simpler, consistent hypothesis





General Decision Tree



Basic Questions

- How **to choose the attribute/value to split** on at each level of the tree?
- When **to stop splitting**? When should a node be declared a leaf?
- If a leaf node is impure, how should the **class label be assigned**?
- If the tree is too large, how can it be **pruned**?

Set of Training Examples

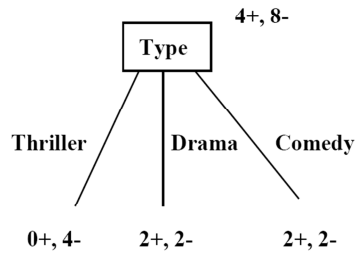
- *Type*: drama, comedy, thriller
- *Company*: MGM, Columbia
- *Director*: Bergman, Spielberg, Hitchcock
- *Mood*: stressed, relaxed, normal

Movie	Type	Company	Director	Mood	Likes-movie?
m_1	thriller	MGM	Bergman	normal	No
m_2	comedy	Columbia	Spielberg	stressed	Yes
m_3	comedy	MGM	Spielberg	relaxed	No
m_4	thriller	MGM	Bergman	relaxed	No
m_5	comedy	MGM	Hitchcock	normal	Yes
m_6	drama	Columbia	Bergman	relaxed	Yes
m_7	drama	Columbia	Bergman	normal	No
m_8	drama	MGM	Spielberg	stressed	No
m_9	drama	MGM	Hitchcock	normal	Yes
m_{10}	comedy	Columbia	Spielberg	relaxed	No
m_{11}	thriller	MGM	Spielberg	normal	No
m_{12}	thriller	Columbia	Hitchcock	relaxed	No

Choosing the “Best” Attribute

- Ideal attribute – partitions examples into all positive or all negative (or from the same class in each partition).
- Attribute that results in higher “discrimination.”

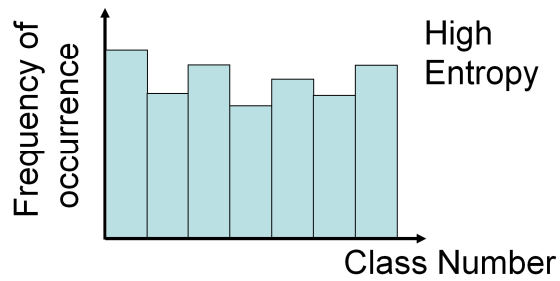
How *good* is the attribute “Type” of movie?



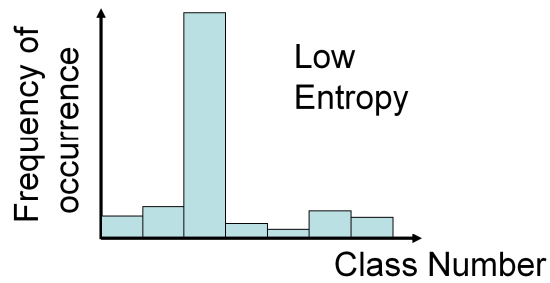
Any attribute splits the data (though some groups might be empty). Here, the type attribute splits the data into three groups.

The optimal attribute would split the data into pure nodes. Here, it would be great if just knowing the type of movie would tell you whether or not the movie was liked.

Entropy (*Shannon and Weaver 1949*)



The entropy captures the degree of “purity” of the data distribution



High entropy because all the information is mixed up, it is chaotic.

Low entropy because the information is organized and orderly.

Entropy *(Shannon and Weaver 1949)*

- In general, entropy is the average number of bits necessary to encode n values:

$$H = -\sum_{i=1}^n P_i \log_2 P_i$$

- P_i = probability of occurrence of value i
 - High entropy → All the classes are (nearly) equally likely
 - Low entropy → A few classes are likely; most of the classes are rarely observed

Know the formula for entropy, it is very, very useful!

Log base conversion from base a to base b: $\log_a(x) = \log_b(x)/\log_b(a)$

Entropy for Attribute Choice

- Measure of *information* provided by the attribute
- **Entropy** of a set of examples S as the information content of S .

$$\text{Entropy}(S) = \sum_{i=1}^c -p_i \log_2 p_i$$

where $p_i = \frac{|S_i|}{|S|}$

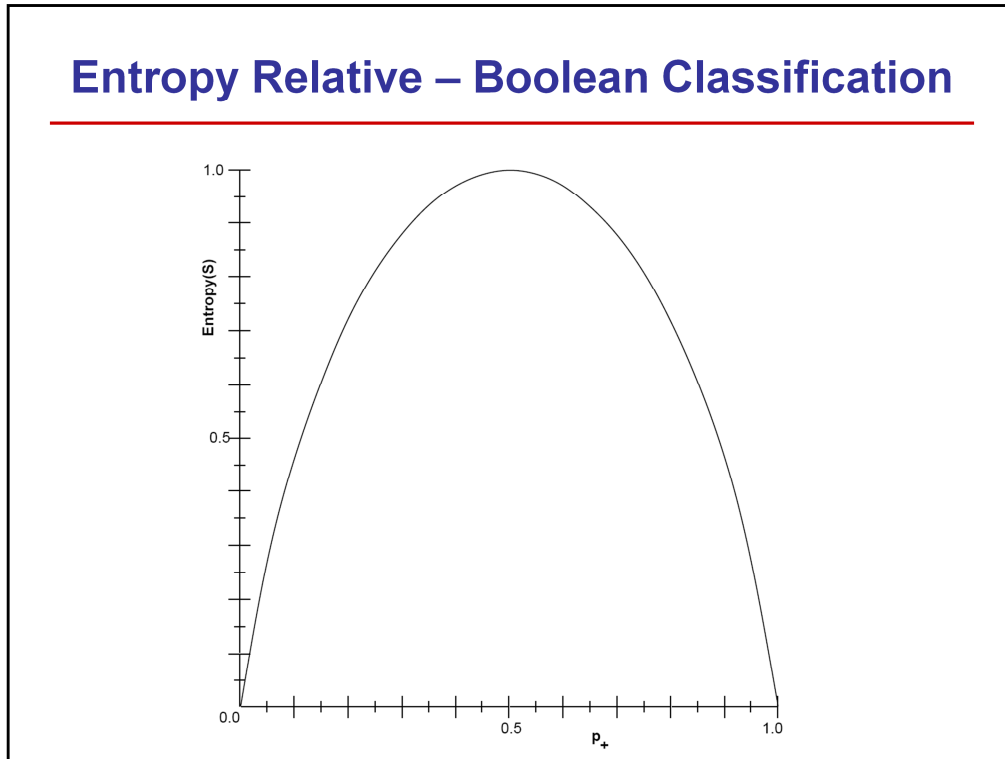
- c classes, S_i size of the data set for class i

Entropy

- Unit – 1 bit of information =
 - the information content of the actual answer when there are two possible answers equally probable.

$$E(S) = -\frac{1}{2} \log_2 \frac{1}{2} - \frac{1}{2} \log_2 \frac{1}{2} = 1$$

Entropy Relative – Boolean Classification



Imagine having a basket full of red and black balls. If p is probability of being black, then when everything is red, p is 0 and entropy is 0. When everything is black, and p is 1, then entropy is 0. If half are red and half are black, p is 0.5 and entropy is 1.

Entropy – Example

- $|S| = 12$, $c = 2(+,-)$, $|S_+| = 4$, $|S_-| = 8$

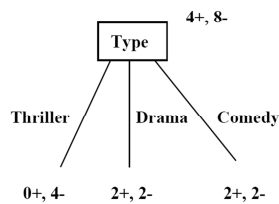
$$\begin{aligned} E(S) &= -\frac{4}{12} \log_2 \frac{4}{12} - \frac{8}{12} \log_2 \frac{8}{12} \\ &= 0.918 \end{aligned}$$

4 positive examples and 8 negative examples. Entropy is still very high.

Choosing the “Best” Attribute

- Attribute with **highest information gain** – expected reduction in entropy of the set S if partitioned according to attribute A .

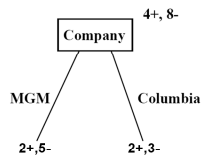
$$\text{Gain}(S,A) = \text{Entropy}(S) - \sum_{i=1}^{\text{values}(A)} \frac{|S_{v_i}|}{|S|} \text{Entropy}(S_{v_i})$$



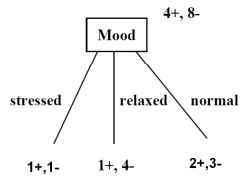
$$\begin{aligned} \text{Gain}(S, \text{Type}) &= E(4^+, 8^-) - \frac{4}{12} E(0^+, 4^-) \\ &\quad - \frac{4}{12} E(2^+, 2^-) \\ &\quad - \frac{4}{12} E(2^+, 2^-) \\ &= 0.252 \end{aligned}$$

Split on the attribute that gives the most information gain (reduces the entropy the most). This is not necessarily optimal, but it is a good heuristic.

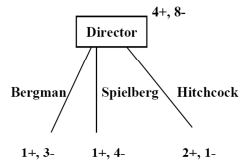
Other Attributes



$$E(4^+, 8^-) - \frac{7}{12} E(2^+, 5^-) - \frac{5}{12} E(2^+, 3^-) = 0.0102$$

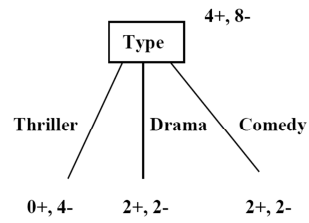


$$E(4^+, 8^-) - \frac{2}{12} E(1^+, 1^-) - \frac{5}{12} E(1^+, 4^-) - \frac{5}{12} E(2^+, 3^-) = 0.046$$



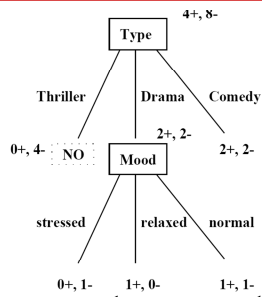
$$E(4^+, 8^-) - \frac{4}{12} E(1^+, 3^-) - \frac{5}{12} E(1^+, 4^-) - \frac{3}{12} E(2^+, 1^-) = 0.118$$

Commitment to One Hypothesis



What is the next *best* attribute? Recursive.

Split Data and Continue

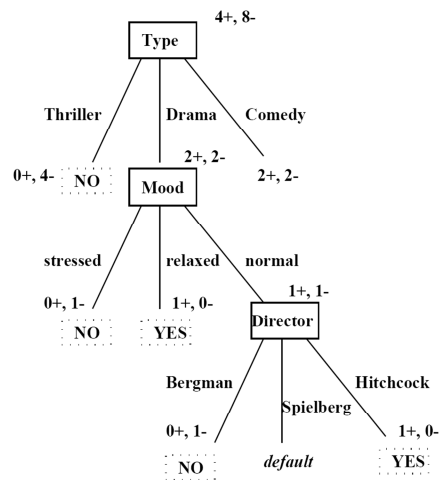


$$\begin{aligned} \text{Gain}(S_d, \text{Mood}) &= E(2^+, 2^-) - \frac{1}{4} E(0^+, 1^-) - \frac{1}{4} E(1^+, 0^-) - \frac{2}{4} E(1^+, 1^-) \\ &= 0.5 \end{aligned}$$

$$\begin{aligned} \text{Gain}(S_d, \text{Company}) &= E(2^+, 2^-) - \frac{2}{4} E(1^+, 1^-) - \frac{2}{4} E(1^+, 1^-) \\ &= 0 \end{aligned}$$

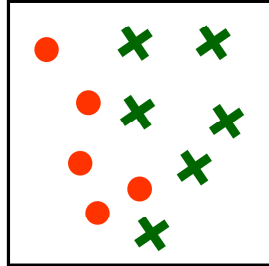
$$\begin{aligned} \text{Gain}(S_d, \text{Director}) &= E(2^+, 2^-) - \frac{2}{4} E(1^+, 1^-) - \frac{1}{4} E(0^+, 1^-) - \frac{2}{4} E(1^+, 0^-) \\ &= 0.5 \end{aligned}$$

Learned Tree



You need some default answer for missing data.

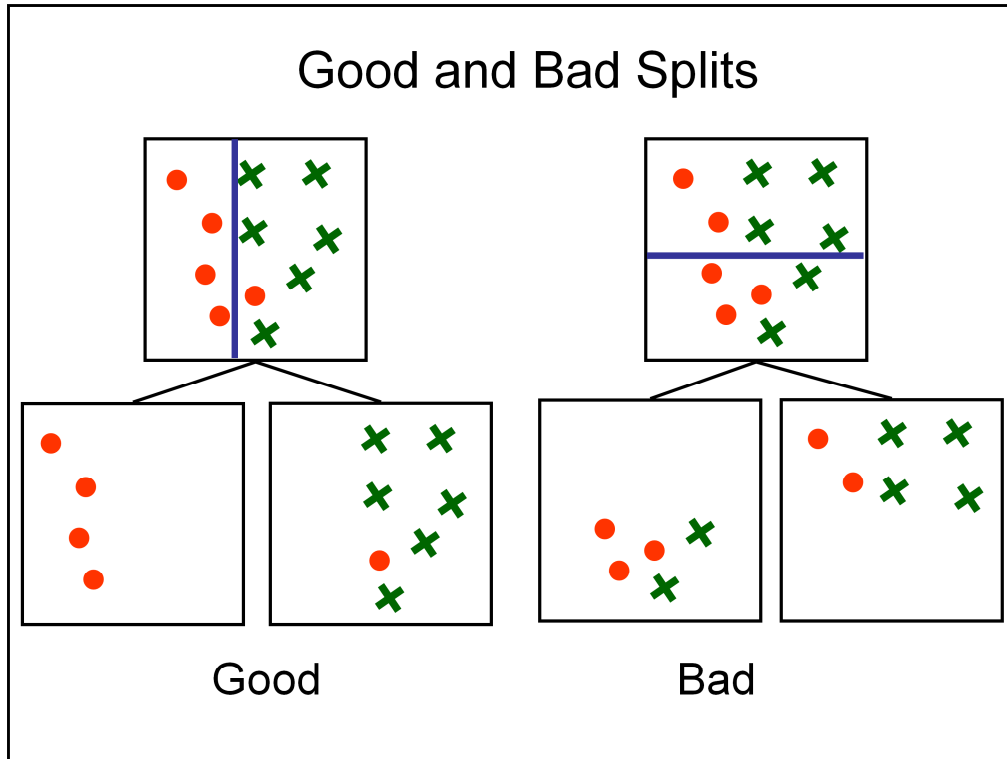
Continuous Case- How to Split?

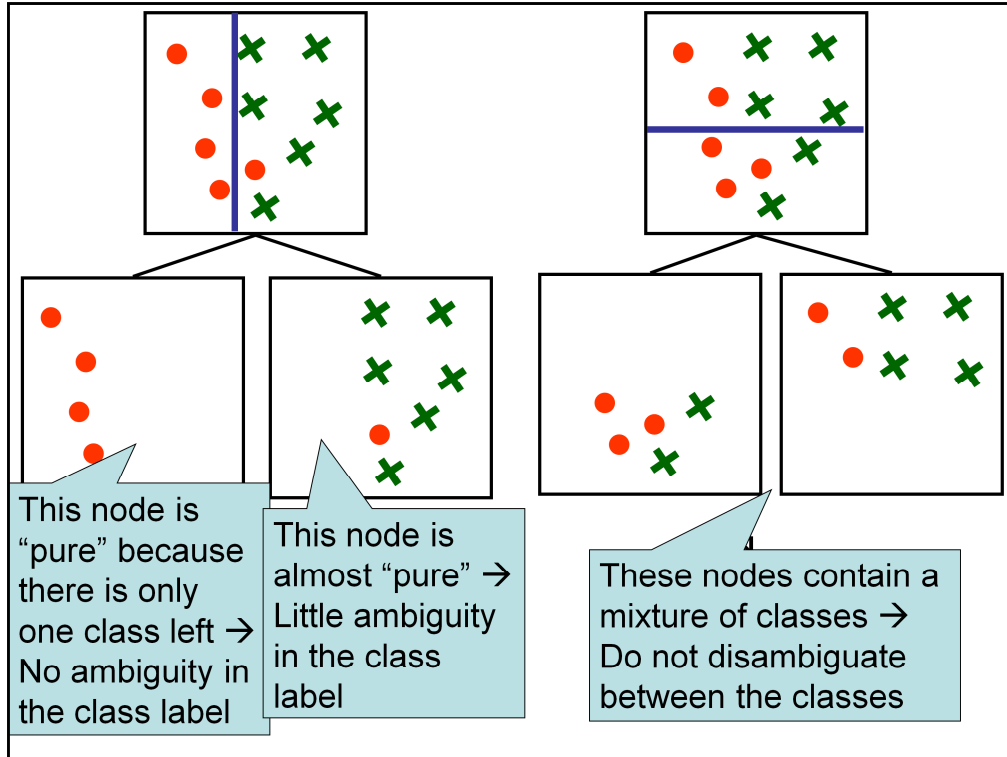


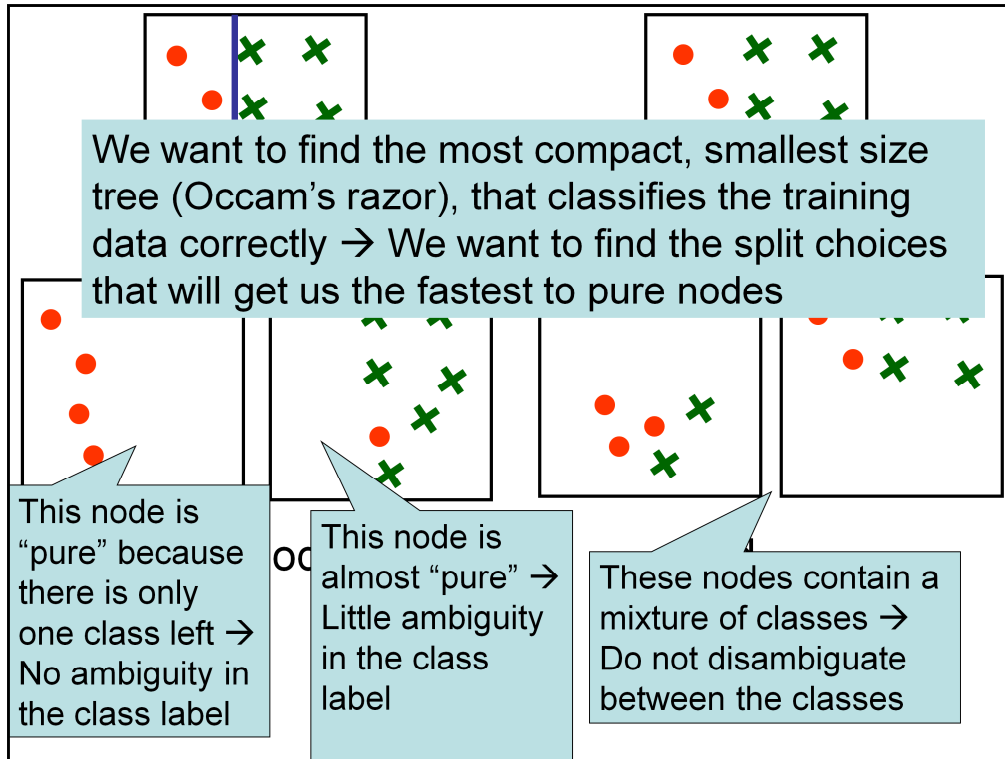
- Two classes (red circles/green crosses)
- Two attributes: X_1 and X_2
- 11 points in training data
- Idea → Construct a decision tree such that the leaf nodes predict correctly the class for all the training examples

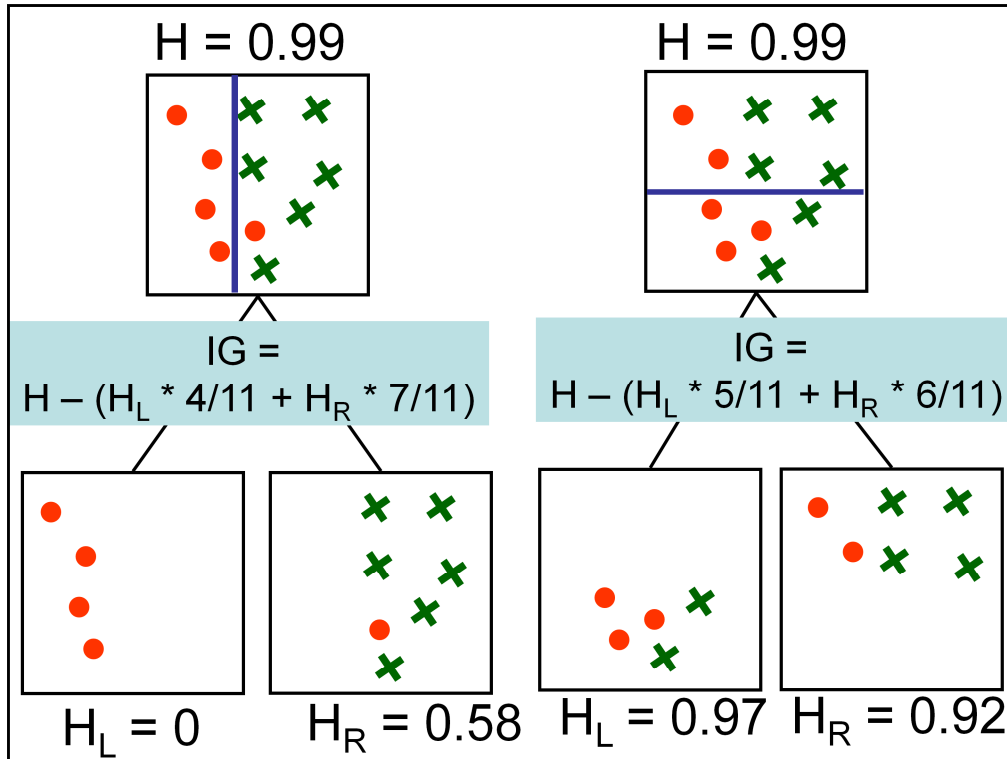
Learn which values of each attribute is the best to split on

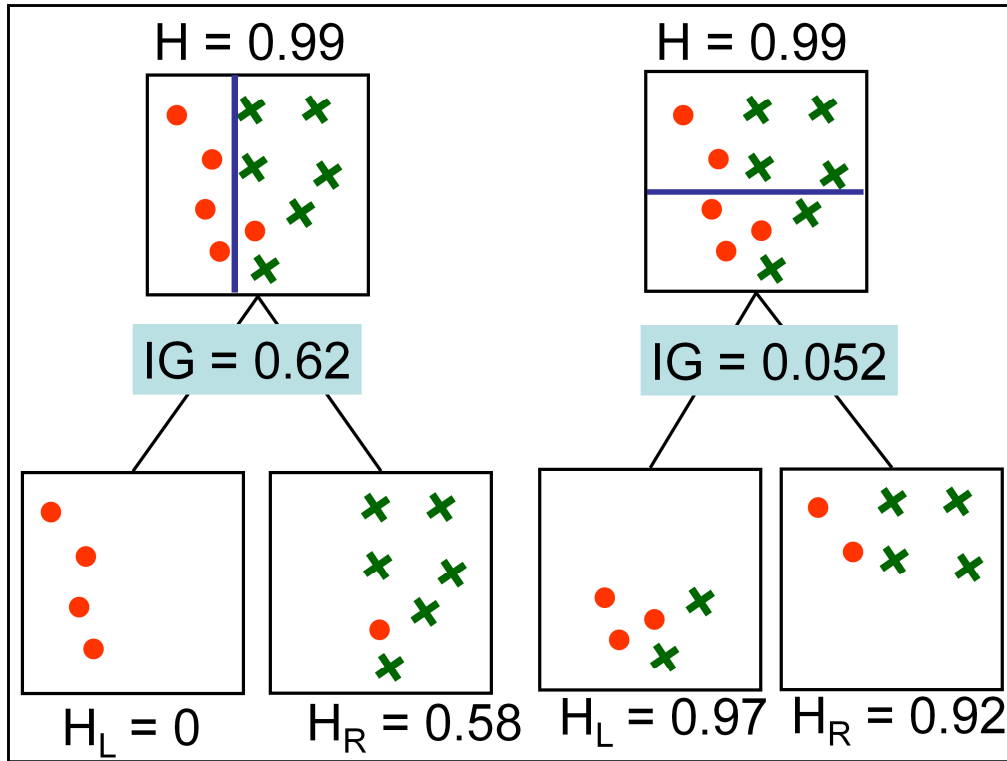
Good and Bad Splits

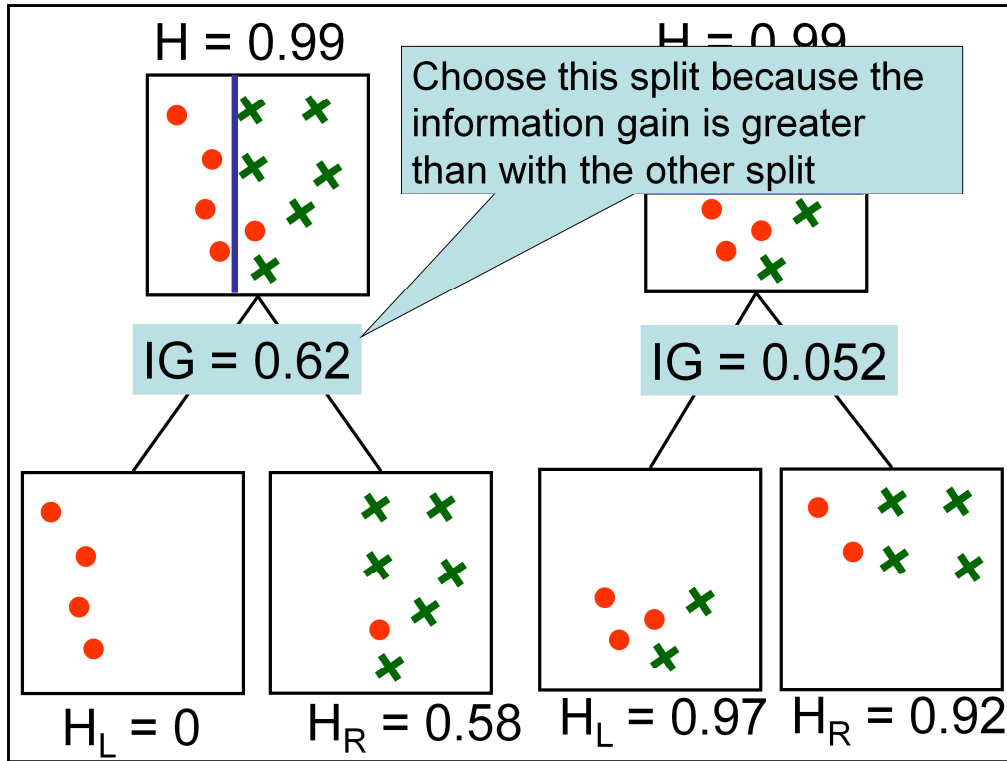




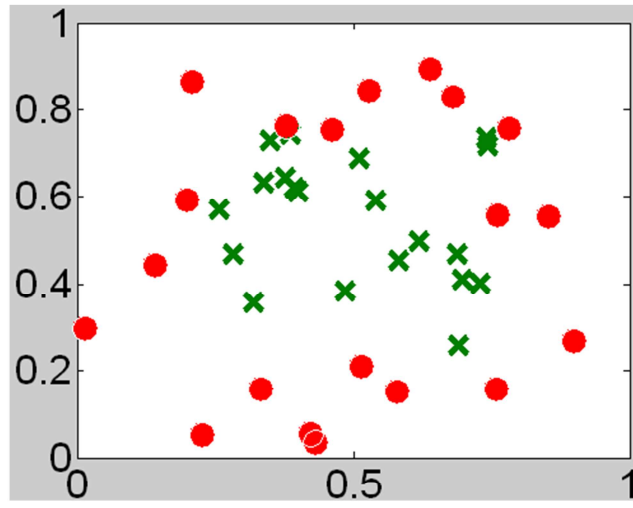




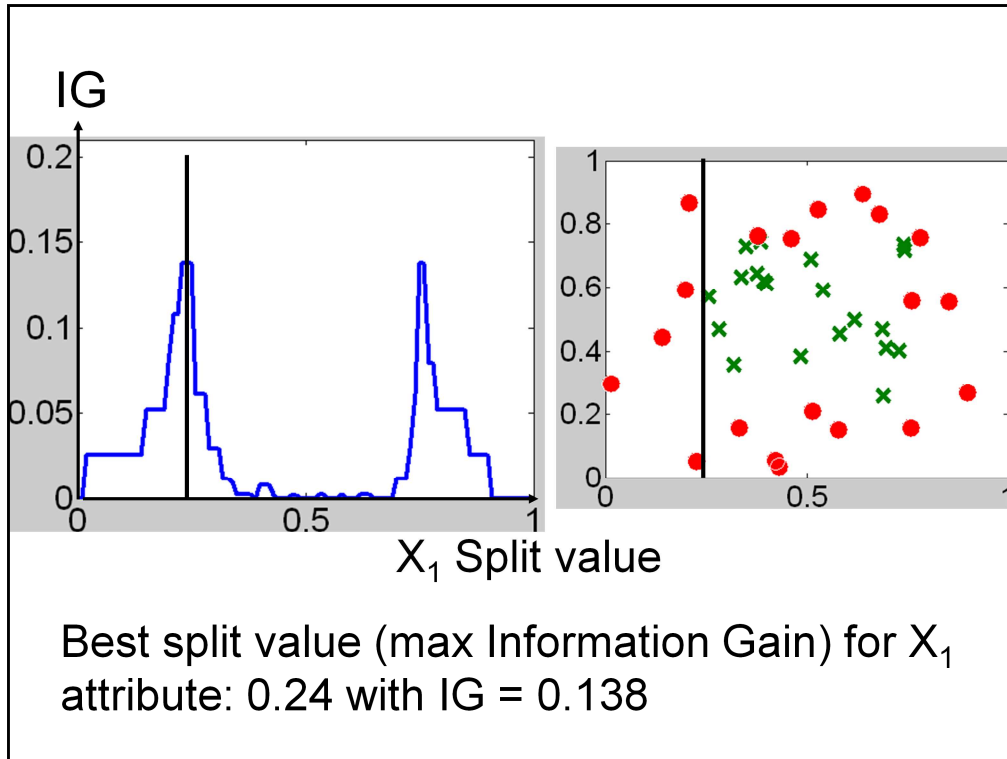


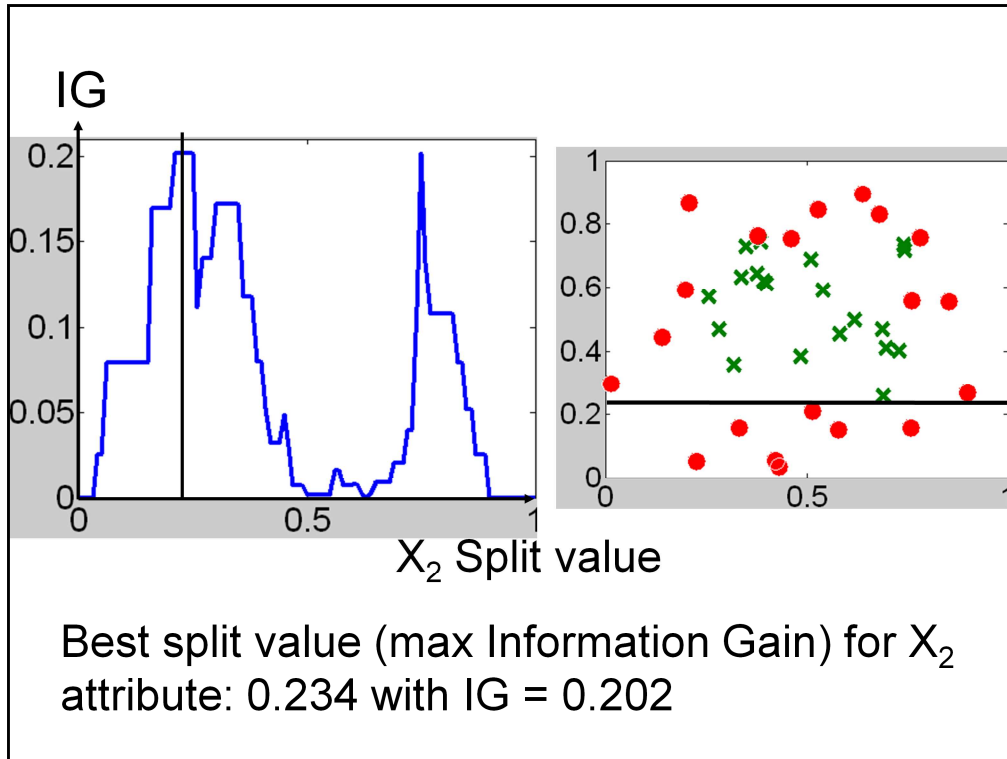


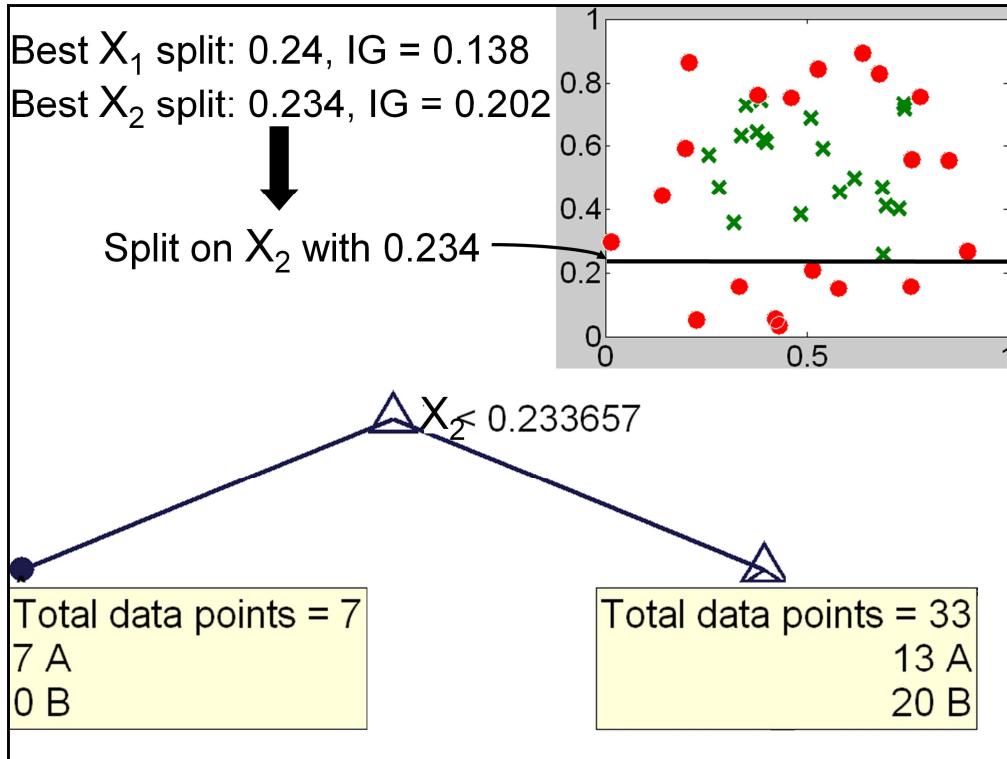
More Complete Example

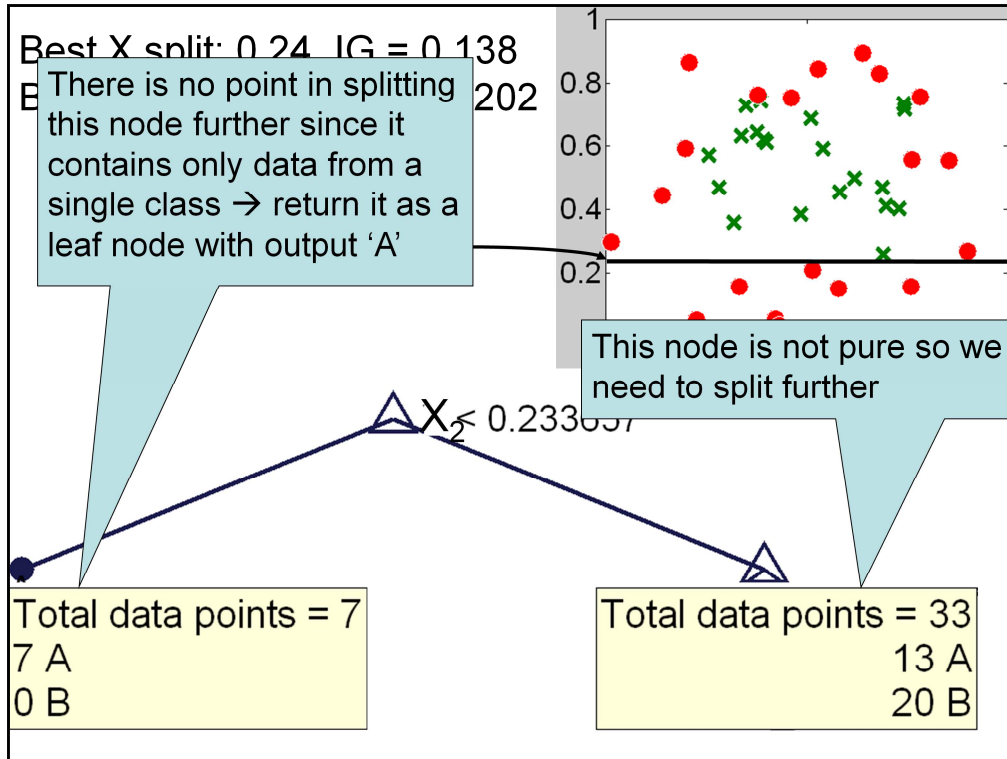


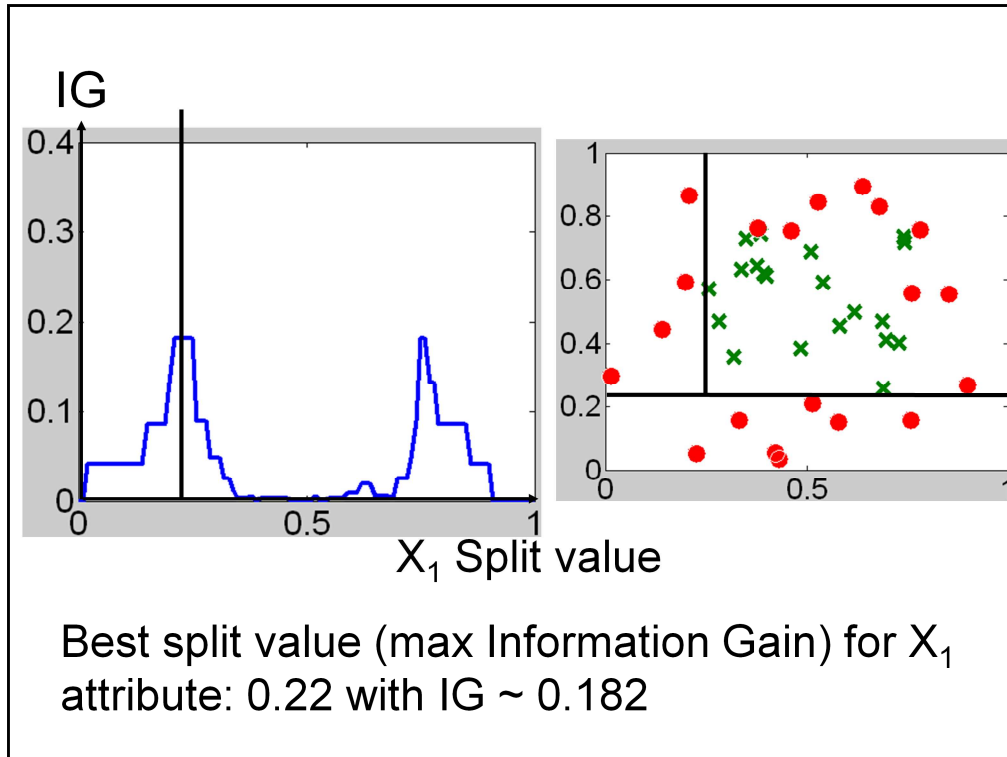
- = 20 training examples from class A
 - ✕ = 20 training examples from class B
- Attributes = X_1 and X_2 coordinates

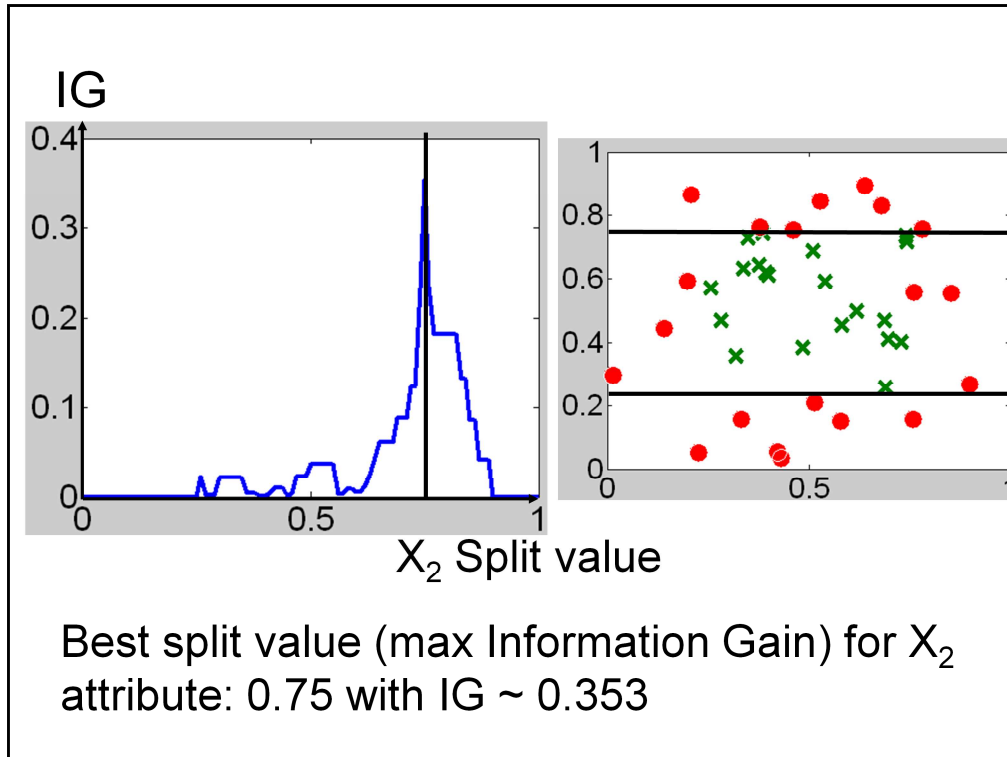


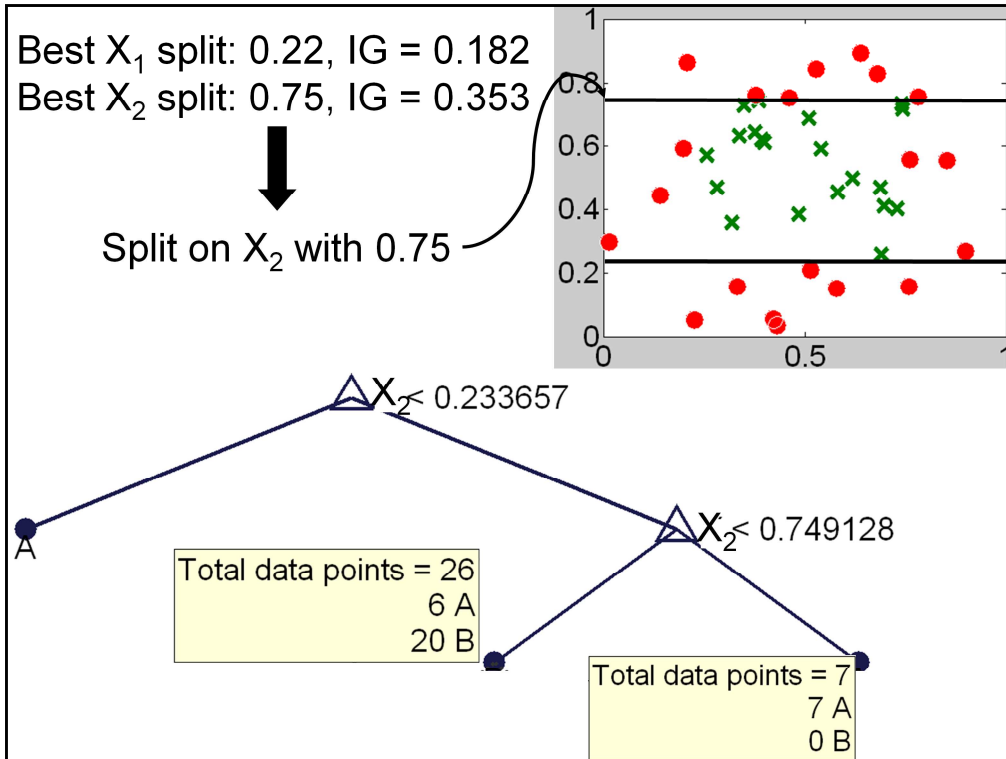


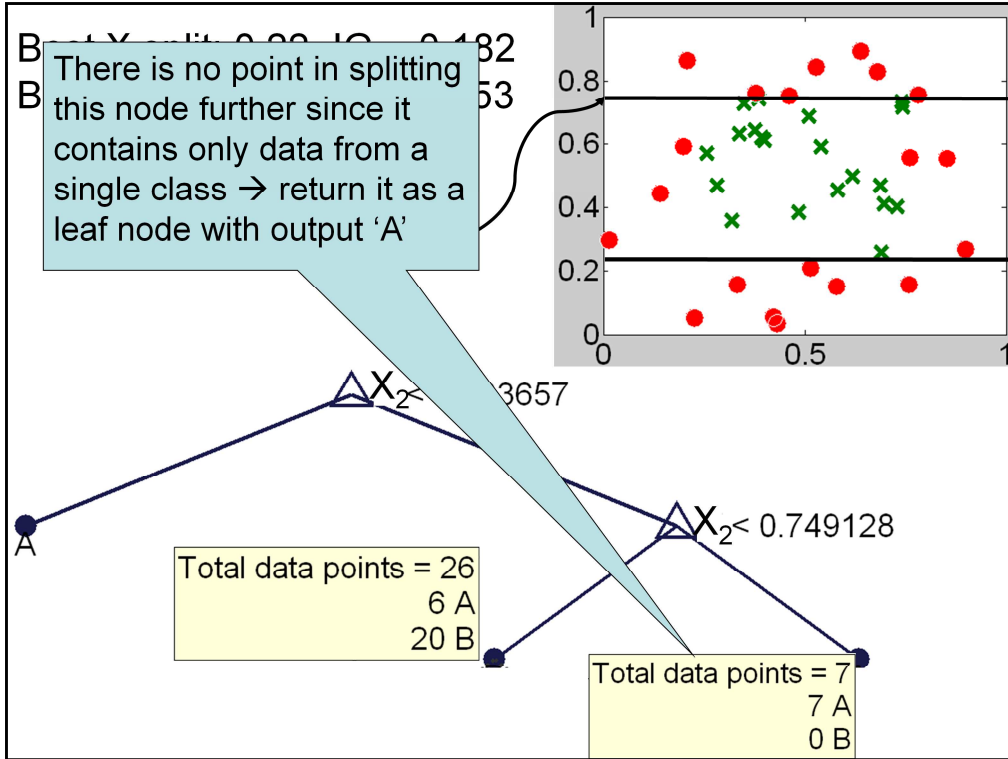


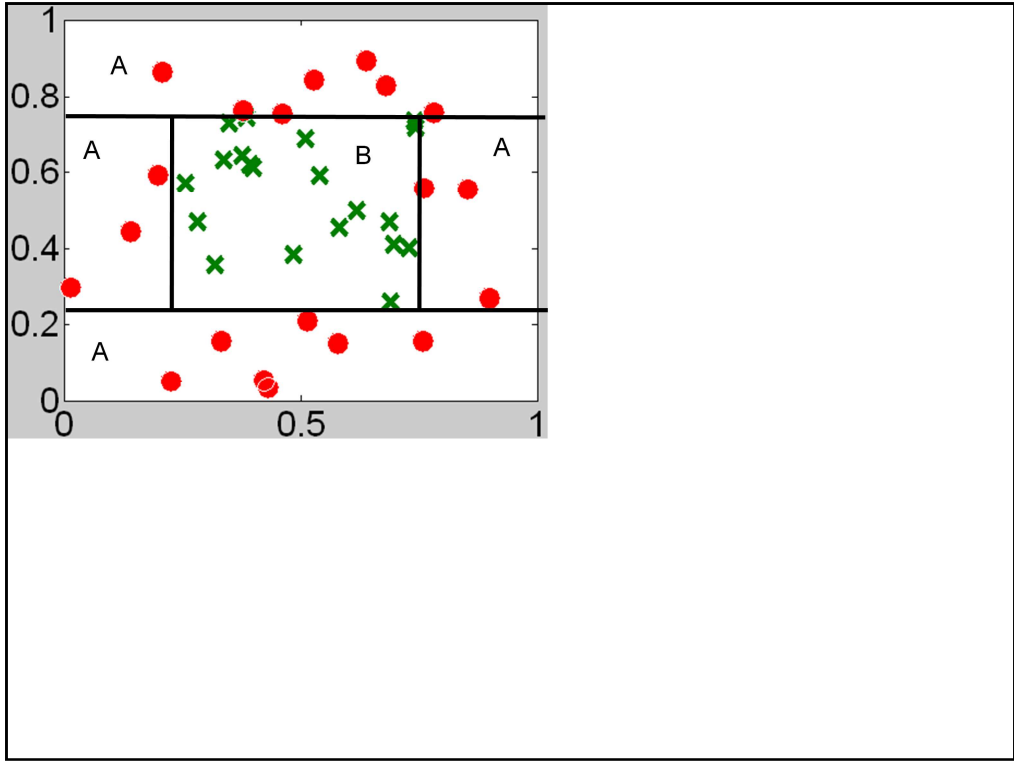


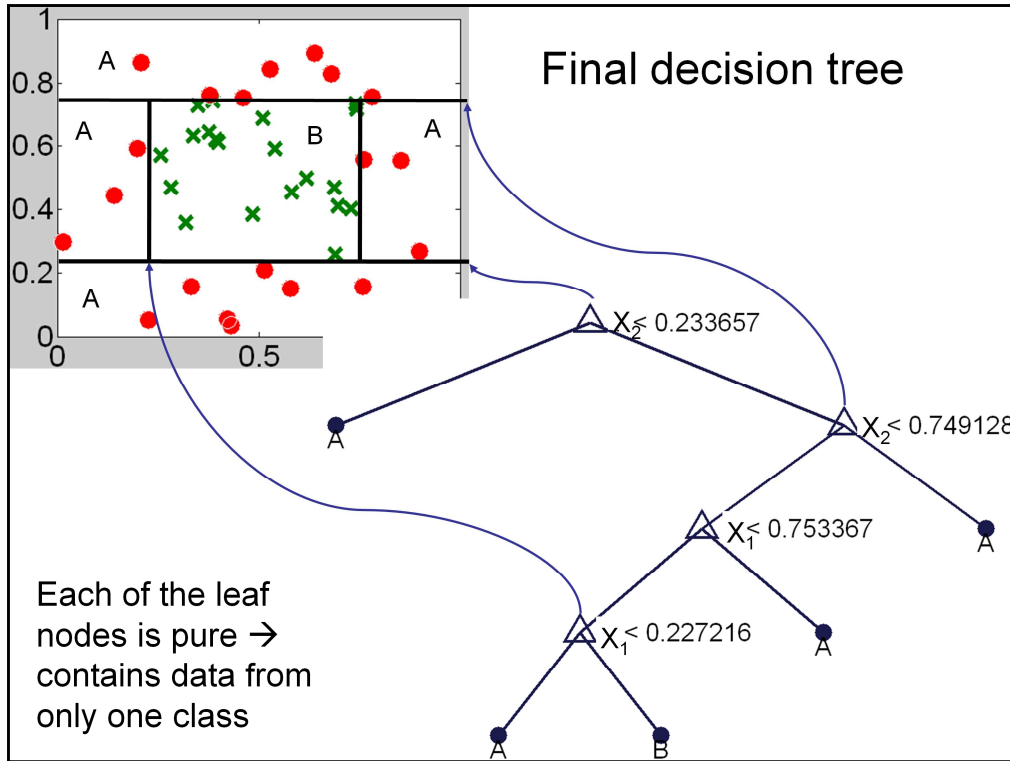


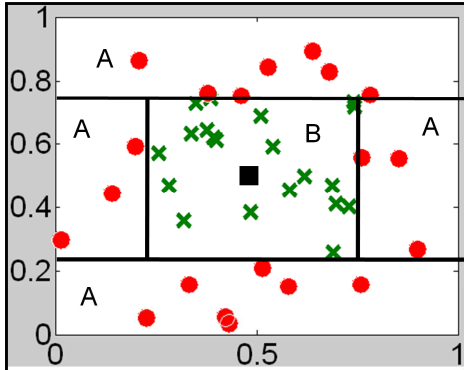






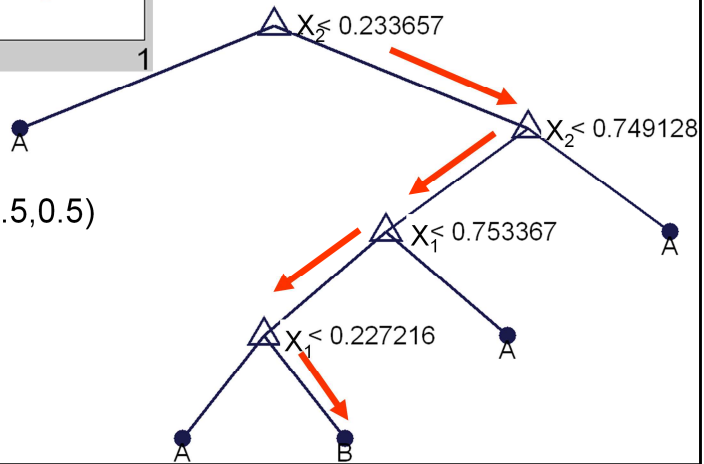






Final decision tree
 Given an input $(X,Y) \rightarrow$
 Follow the tree down to a
 leaf.
 Return corresponding
 output class for this leaf

Example $(X,Y) = (0.5,0.5)$



Induction of Decision Trees

```
ID3 (examples, attributes)
if there are no examples
  then return default
else
  if all examples are members of the same class
    then return the class
  else
    if there are no attributes
      then return Most_Common_Class (examples)
    else
      best_attribute ← Choose_Best (examples, attributes)
      root ← Create_Node_Test (best_attribute)
      for each value  $v_i$  of best_attribute
        examples $v_i$  ← subset of examples with best_attribute =  $v_i$ 
        subtree ← ID3(examples $v_i$ , attributes - best_attribute)
        set subtree as a child of root with label  $v_i$ 
return root
```

Basic Questions

- How to choose the attribute/value to split on at each level of the tree?
- When to stop splitting? When should a node be declared a leaf?
- If a leaf node is impure, how should the class label be assigned?
- If the tree is too large, how can it be pruned?

Comments on Tree Termination

- Zero entropy in data set, perfect classification
 - Type = thriller, 0+ 4-
 - Type = drama, Mood = stressed, 0+ 1-
 - Type = drama, Mood = relaxed, 1+ 0-
 - Type = drama, Mood = normal, Director = Bergman, 0+ 1-
 - Type = drama, Mood = normal, Director = Hitchcock, 1+ 0-
- No examples available
 - Type = drama, Mood = normal, Director = Spielberg
- Indistinguishable data
 - Type = comedy, attribute Company:
 - ❖ m2, comedy, Columbia, Yes, and m10, comedy, Columbia, No
 - ❖ m3, comedy, MGM, No, and m5, comedy, MGM, Yes
 - Type = comedy, attribute Director:
 - ❖ m2, Spielberg, Yes, and m3, Spielberg, No
 - ❖ m5, Hitchcock, Yes, and m10, Spielberg, No
 - (Type = comedy, attribute Mood, could have split?).

Errors

Split labeled data D into *training* and *validation* sets

- Training set error - fraction of training examples for which the decision tree disagrees with the true classification
- Validation set error - fraction of testing examples - from given labeled examples - for which the decision tree disagrees with the true classification

Overfitting

Tree too specialized – “perfectly” fit the training data.

A tree T *overfits* the training data, if there is an alternate T' such that

- for training set, error with $T <$ error with T'
- for complete D , error with $T >$ error with T'

- Node Pruning
- Rule Post-pruning
- Cross validation

Reduced-Error Pruning

- Remove subtree, make node a leaf node, assign the most common classification of the examples of that node.
- Check validation set error
- Continue pruning until error does not increase with respect to the error of the unpruned tree
- Rule post-pruning: represent tree as set of rules; remove rules independently; check validation set error; stop with same criteria

Other Issues

- Attributes with many values
 - Information gain may select it.
 - Consider splitting value and use the ratio between Gain and the splitting value
- Attributes with cost
 - Use other metrics
- Attributes with missing values
 - Infer most common value from examples at node

Summary

- Inductive learning – supervised learning – classification
- Decision trees represent hypotheses
- DT learning driven by information gain heuristic
- Recursive algorithm to build decision tree
- Next class:
 - More on continuous values
 - Missing, noisy attributes
 - Overfitting, pruning
 - Different attribute selection criteria