



UNSW
SYDNEY

Australia's
Global
University

COMP9321

Data Services Engineering

Term 1, 2019

Week 6 Lecture 2

Supervised Learning

COMP9321 2019T1

Decision Tree

Supervised Learning

Background

- Decision trees have a long history in machine learning
- The first popular algorithm dates back to 1979
- Very popular in many real world problems
 - Intuitive to understand
 - Easy to build

Motivation

- How do people make decisions?
 - Consider a variety of factors
 - Follow a logical path of checks
- An Example
 - Should I eat at this restaurant?
 - If there is no wait
 - Yes
 - If there is short wait and I am hungry
 - Yes
 - Else
 - No

Advantages

Handling of categorical variables

Handling of missing values and unknown labels

Detection of nonlinear relationships

Visualization and interpretation in decision trees

Decision Tree

Goal: Categorization

Given an event, predict its category. Examples:

- Who won a given ball game?
- How should we file a given email?
- What word sense was intended for a given occurrence of a word?

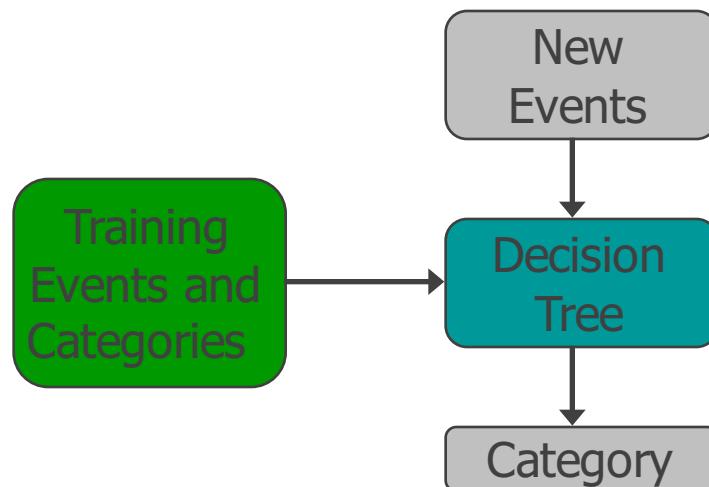
Event = list of features. Examples:

- Ball game: Which players were on offense?
- Email: Who sent the email?
- Disambiguation: What was the preceding word?

Introduction

Use a decision tree to predict categories for new events.

Use training data to build the decision tree.



Introduction

A decision tree has 2 kinds of nodes

1. Each leaf node has a class label, determined by majority vote of training examples reaching that leaf.
2. Each internal node is a question on features. It branches out according to the answers.

Decision Tree for Play Tennis

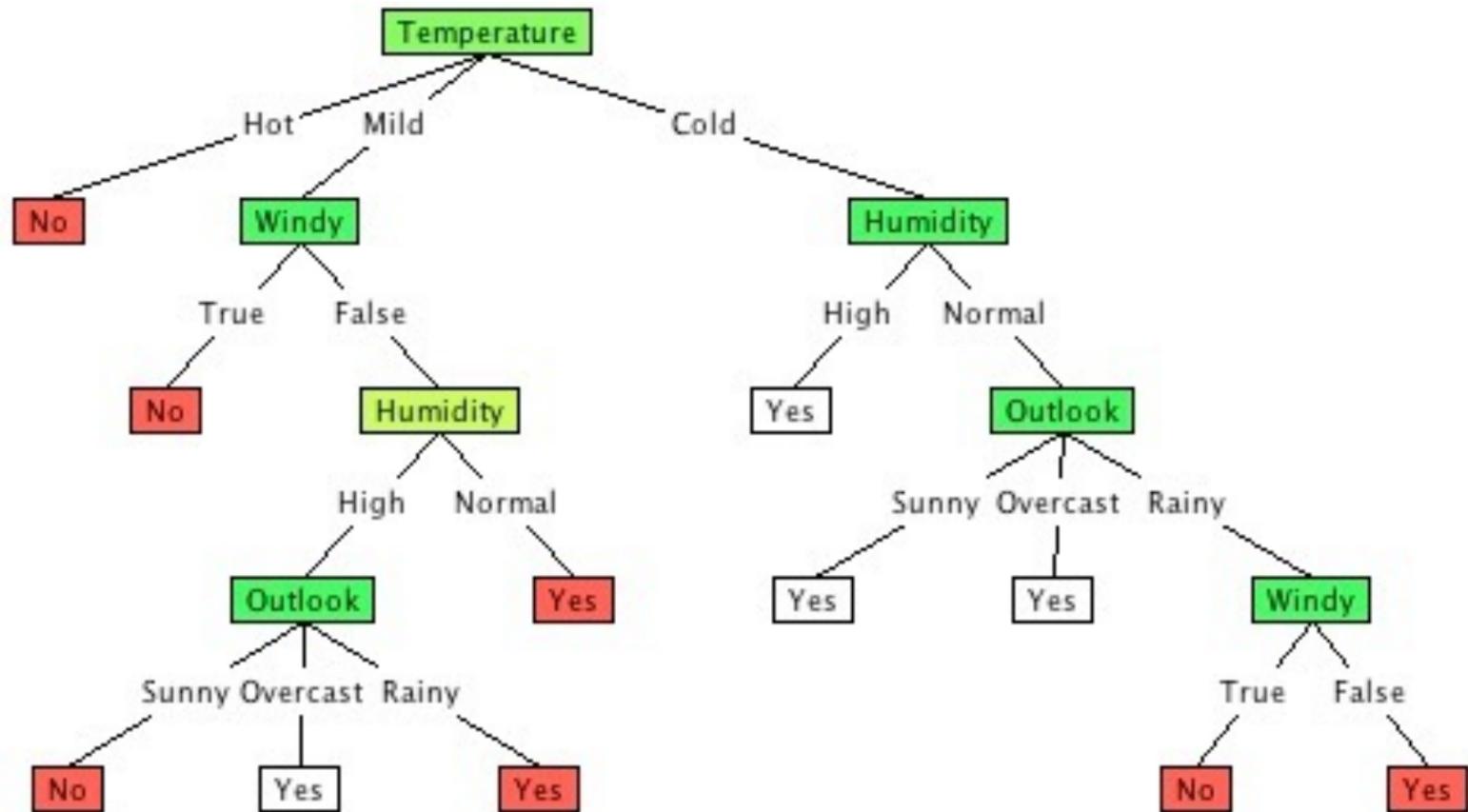
Play Tennis	Outlook	Temperature	Humidity	Windy
No	Sunny	Hot	High	No
No	Sunny	Hot	High	Yes
Yes	Overcast	Hot	High	No
Yes	Rainy	Mild	High	No
Yes	Rainy	Cold	Normal	No

If temperature is not hot → Play

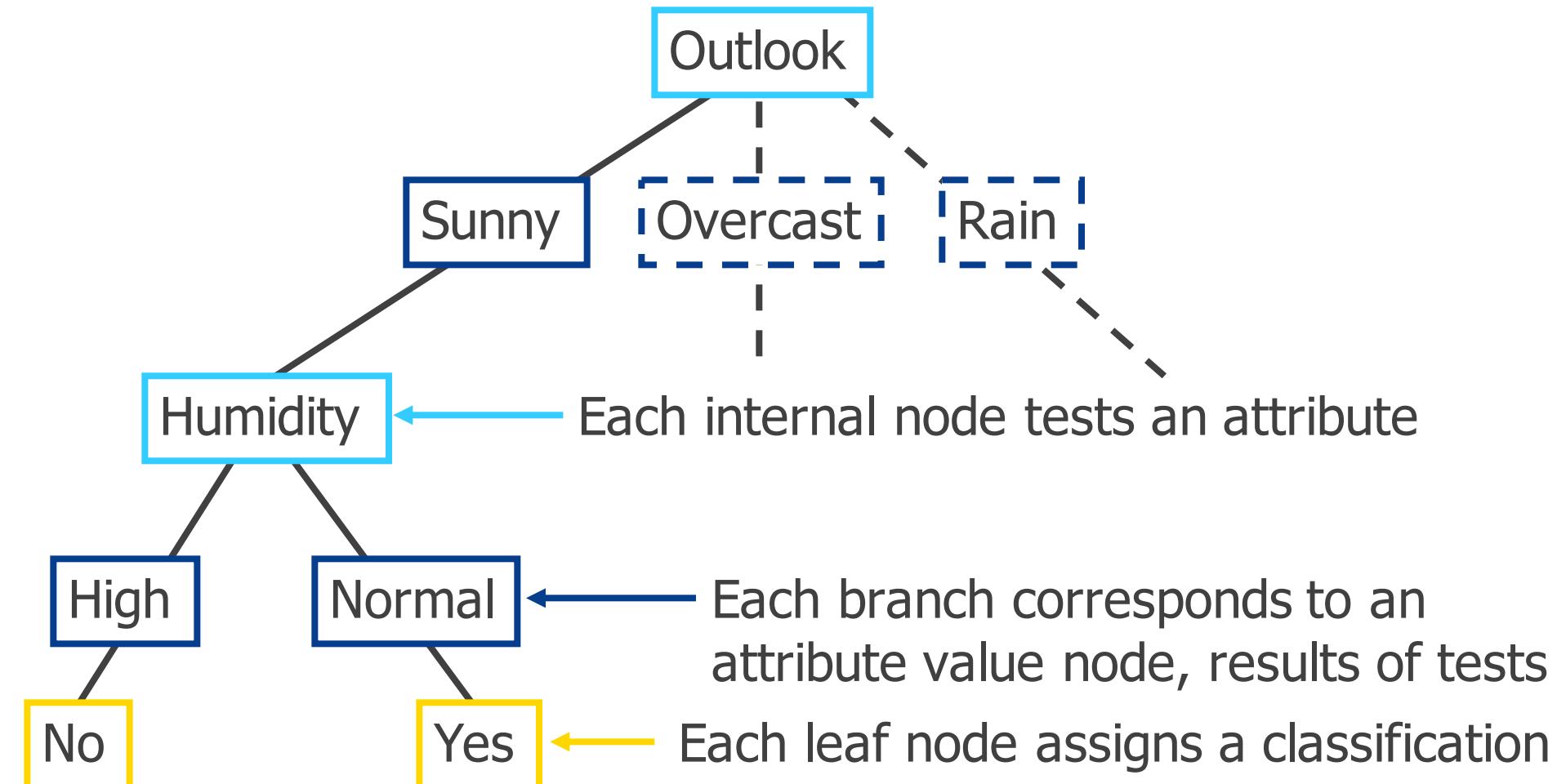
If outlook is overcast → Play

Otherwise → Don't play tennis

Decision Tree for Play Tennis



Decision Tree for PlayTennis - Structure



Function Class

What functions can decision trees model?

Non-linear: very powerful function class

A decision tree can encode any Boolean function

Proof

- Given a truth table for a function
- Construct a path in the tree for each row of the table
- Given a row as input, follow that path to the desired leaf (output)

Problem: exponentially large trees!

Trimming the Tree

Can we produce smaller decision trees for functions?

- Yes (Possible)

Key decision factors

Counter examples

- Parity function: Return 1 on even inputs, 0 on odd inputs
- Majority function: Return 1 if more than half of inputs are 1

Bias-Variance Trade-off

- Bias: Representation Power of Decision Trees
- Variance: require a sample size exponential in depth

What Makes a Good Tree

Small tree:

- Occam's razor: a guideline to help us choose. Simpler is better
- Avoids over-fitting



What Makes a Good Tree

- Occam's razor: a guideline to help us choose. Simpler is better



What Makes a Good Tree

- Occam's razor: a guideline to help us choose. Simpler is better

Prank explanation requires:

1. Human (not observed)
2. Ability to enter house
(unknown)
3. Motive to play prank
(unknown)
4. Leaving no other trace
(observed)

Licking explanation requires:

1. Cat (observed)
2. Licking (observed)

What Makes a Good Tree

Small tree:

- Occam's razor: a guideline to help us choose. Simpler is better
- Avoids over-fitting

A decision tree may be human readable, but not use human logic!

How do we build small trees that accurately capture data?

Learning an optimal decision tree is computationally intractable

Greedy Algorithm

We can get good trees by simple greedy algorithms

Adjustments are usually to fix greedy selection problems

Recursive:

1. Select the “best” variable, and generate child nodes: One for each possible value;
2. Partition samples using the possible values, and assign these subsets of samples to the child nodes;
3. Repeat for each child node until all samples associated with a node that are either all positive or all negative.

Variable Selection

The best variable for partition

- The most informative variable
- Select the variable that is most informative about the labels
- Classification error?

Example: $\mathbb{P}(Y = 1|X_1 = 1) = 0.75$ v.s. $\mathbb{P}(Y = 1|X_2 = 1) = 0.55$

Which to Choose?

Information Theory

The quantification of information

Founded by Claude Shannon

Basic concepts:

Entropy: $H(X) = - \sum_x \mathbb{P}(X = x) \log \mathbb{P}(X = x)$

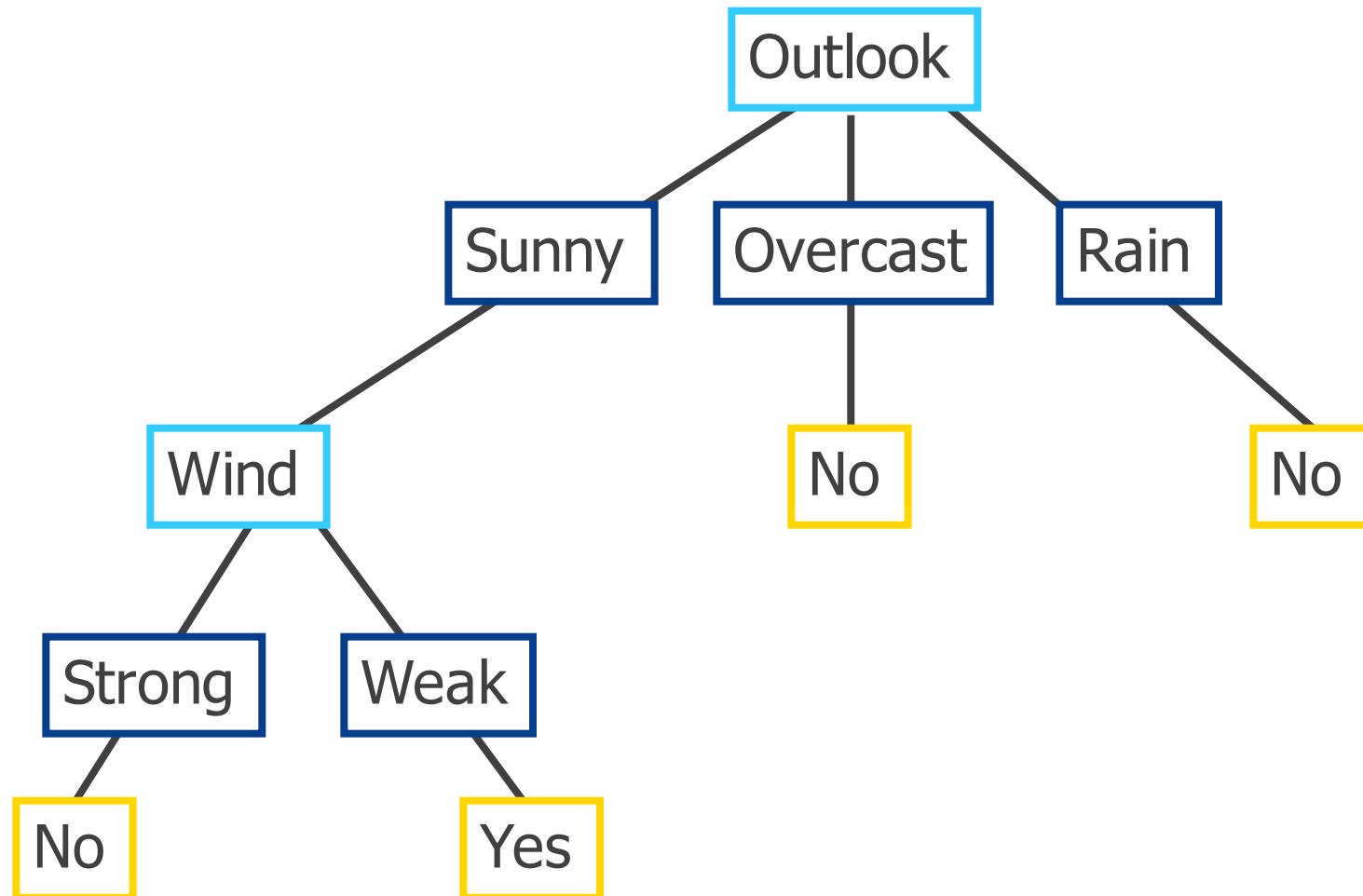
Conditional Entropy: $H(Y|X) = \sum_x \mathbb{P}(X = x) H(Y|X = x)$

Information Gain: $IG(Y|X) = H(Y) - H(Y|X)$

Select the variable with the highest information gain

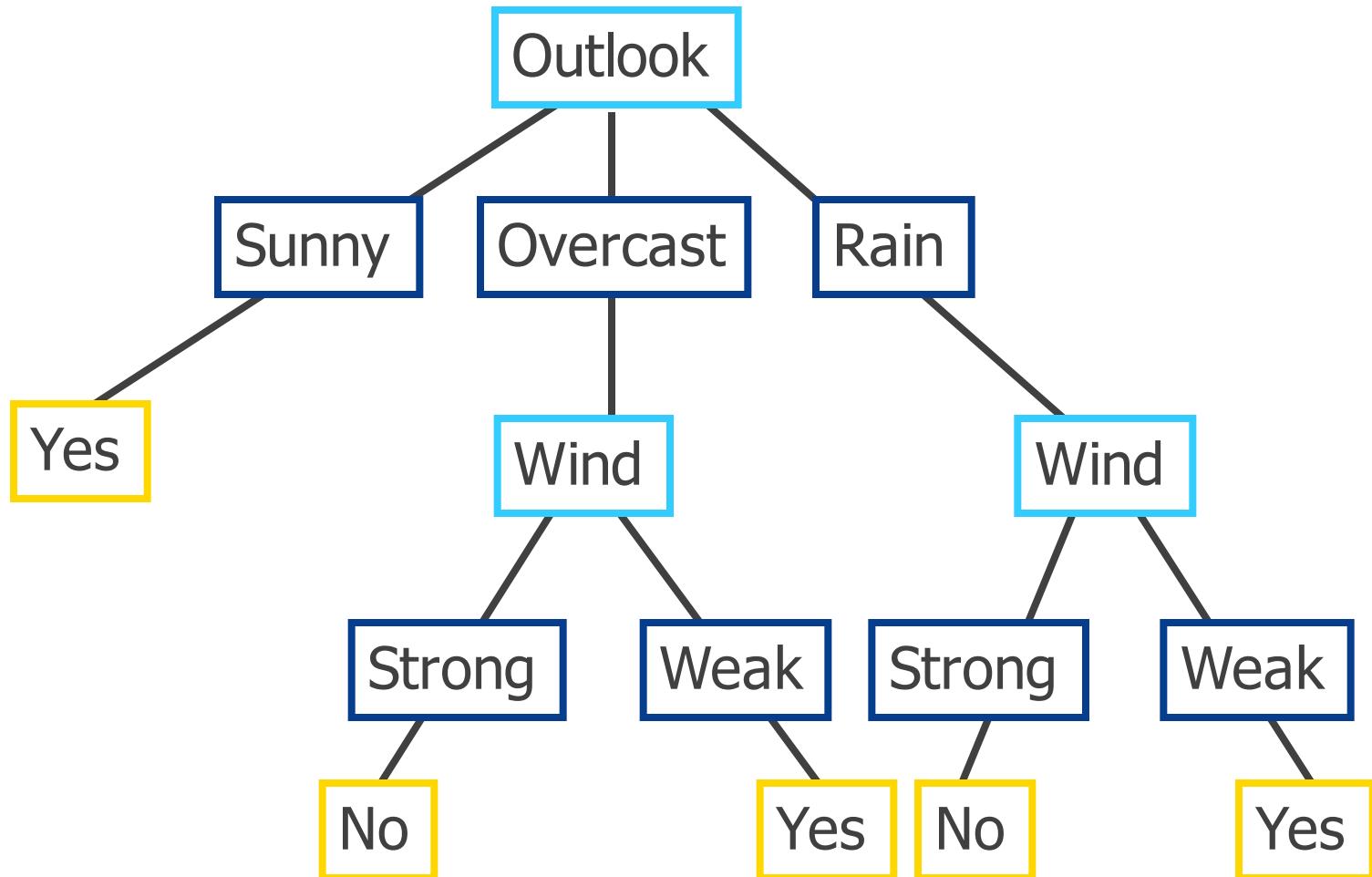
Decision Tree for Conjunction

Outlook=Sunny \wedge Wind=Weak



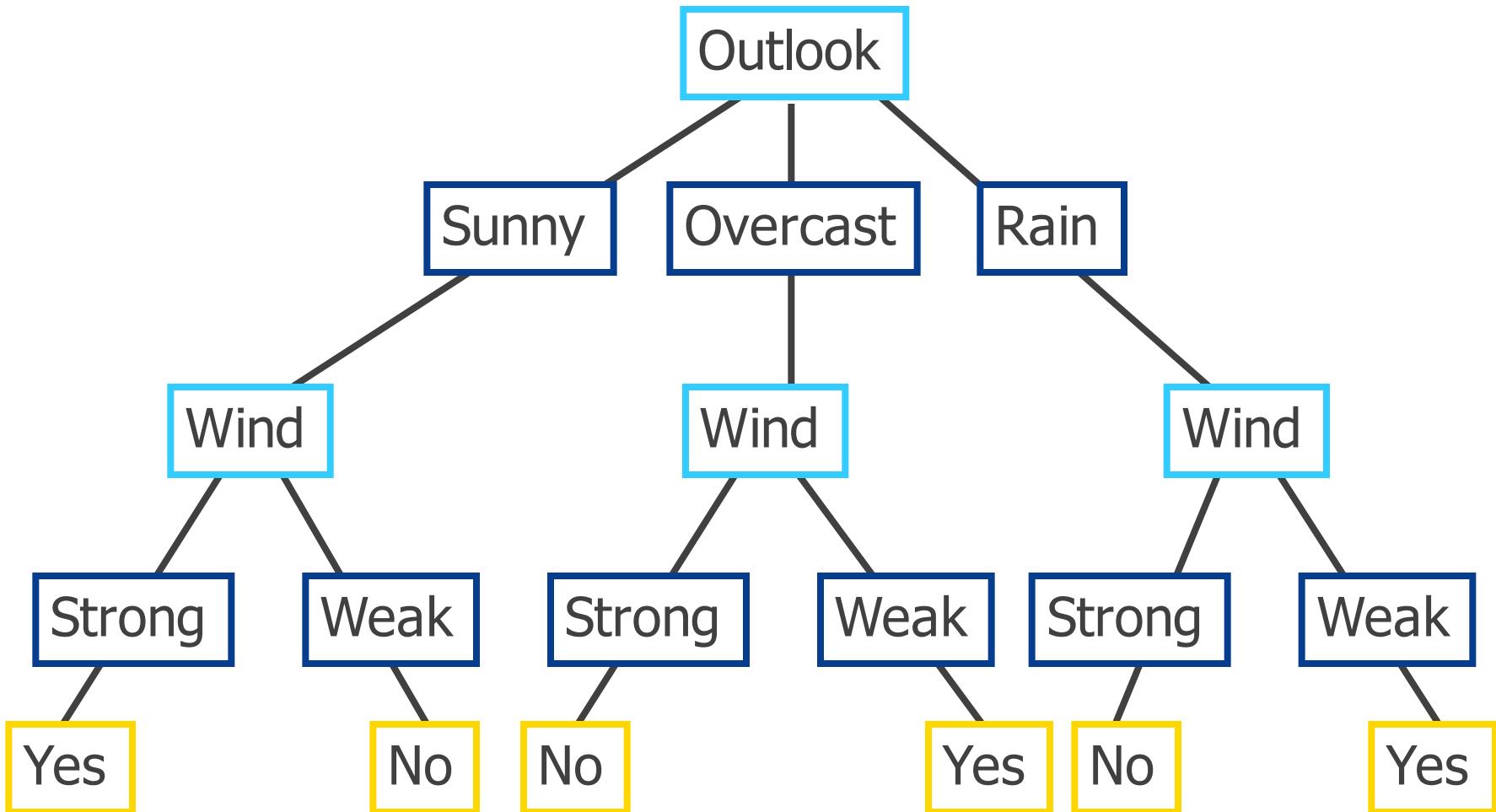
Decision Tree for Disjunction

$\text{Outlook} = \text{Sunny} \vee \text{Wind} = \text{Weak}$



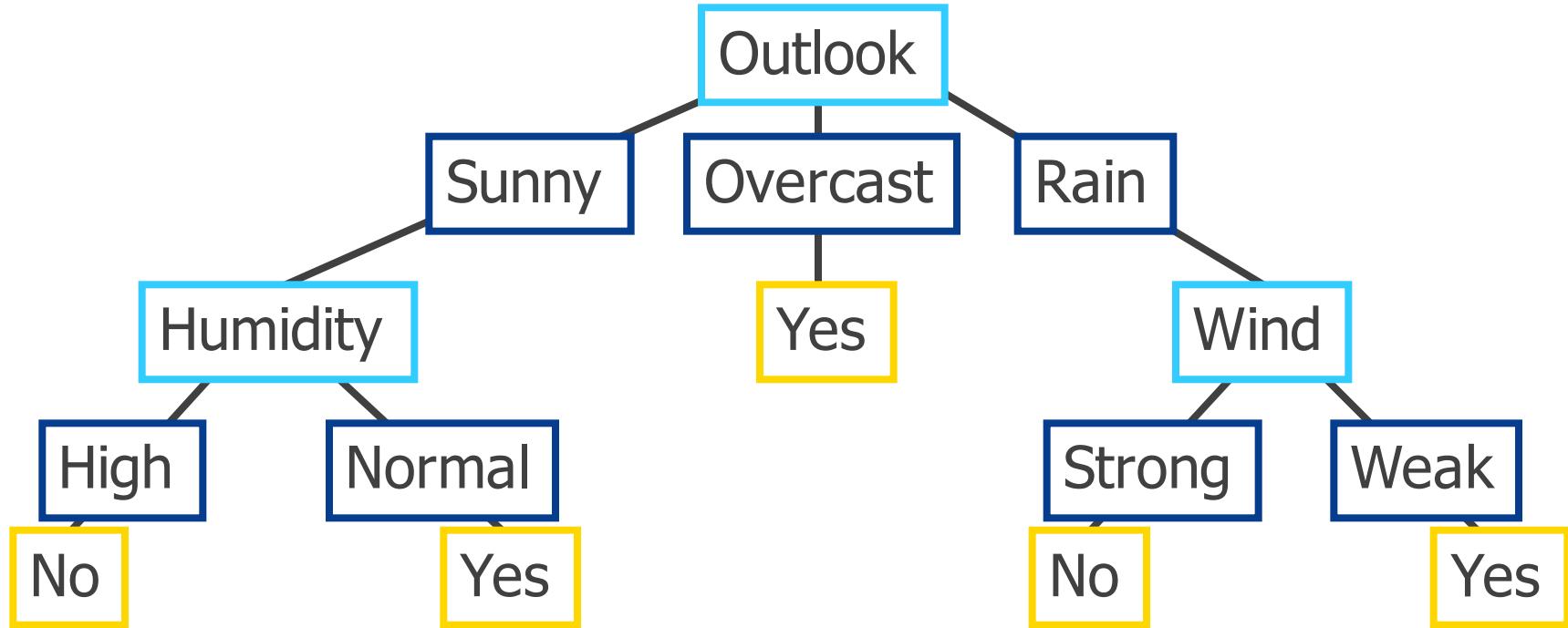
Decision Tree for XOR

Outlook=Sunny XOR Wind=Weak



Decision Tree

- decision trees represent disjunctions of conjunctions


$$\begin{aligned} & (\text{Outlook}=\text{Sunny} \wedge \text{Humidity}=\text{Normal}) \\ \vee & \quad (\text{Outlook}=\text{Overcast}) \\ \vee & \quad (\text{Outlook}=\text{Rain} \wedge \text{Wind}=\text{Weak}) \end{aligned}$$

When to consider Decision Trees

Instances describable by attribute-value pairs

Target function is discrete valued

Disjunctive hypothesis may be required

Possibly noisy training data

Missing attribute values

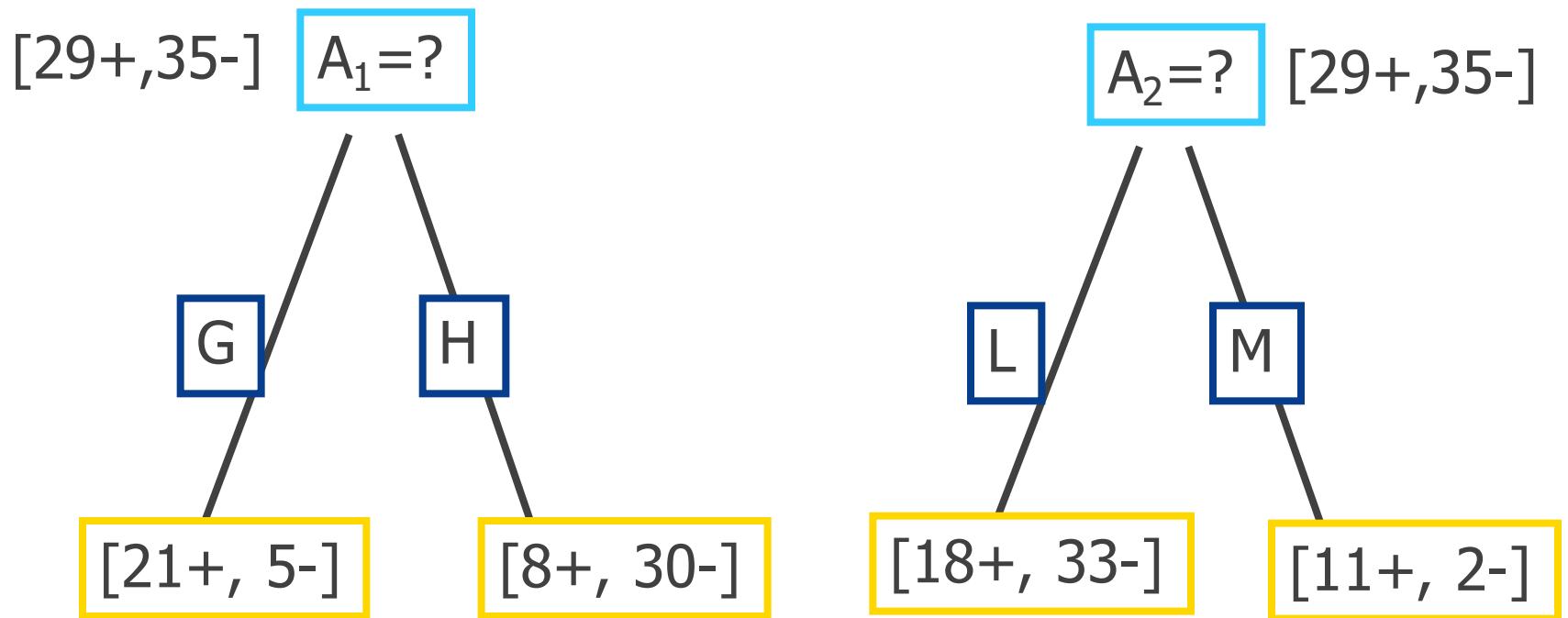
Examples:

- Medical diagnosis
- Credit risk analysis
- Object classification for robot manipulator (Tan 1993)

Top-Down Induction of Decision Trees ID3

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
4. Sort training examples to leaf node according to the attribute value of the branch
5. If all training examples are perfectly classified (same value of target attribute) stop, else iterate over new leaf nodes.

Which attribute is best?



Entropy

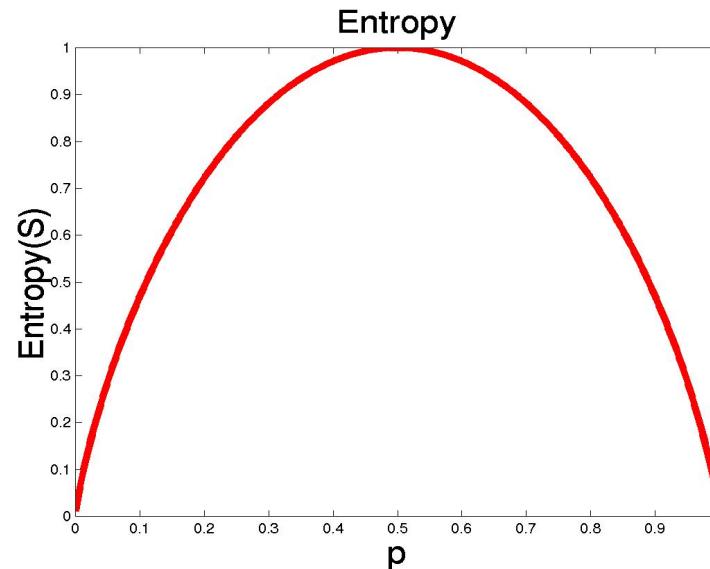
S is a sample of training examples

p_+ is the proportion of positive examples

p_- is the proportion of negative examples

Entropy measures the impurity of S

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



Entropy

Entropy(S) = expected number of bits needed to encode class (+ or -) of randomly drawn members of S (under the optimal, shortest length-code)

Why?

Information theory optimal length code assign
 $-\log_2 p$ bits to messages having probability p .

So the expected number of bits to encode
(+ or -) of random member of S :

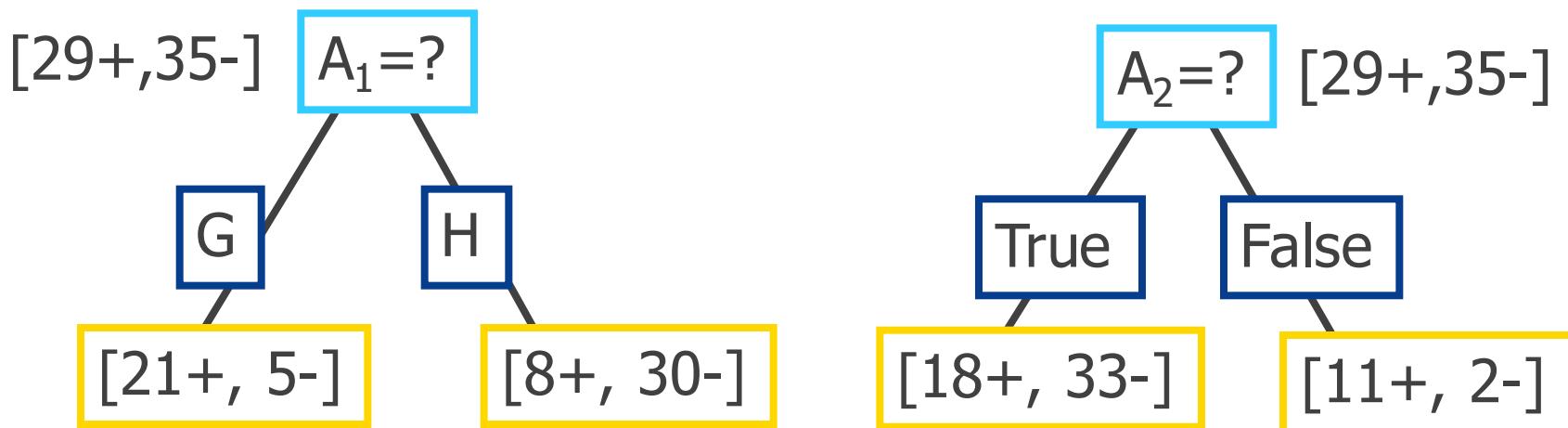
- $-p_+ \log_2 p_+ - p_- \log_2 p_-$

Information Gain (S=E)

$\text{Gain}(S, A)$: expected reduction in entropy due to sorting S on attribute A

$$\begin{aligned}\text{Entropy}([29+, 35-]) &= -29/64 \log_2 29/64 - 35/64 \log_2 35/64 \\ &= 0.99\end{aligned}$$

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



Information Gain

$$\text{Entropy}([21+, 5-]) = 0.71$$

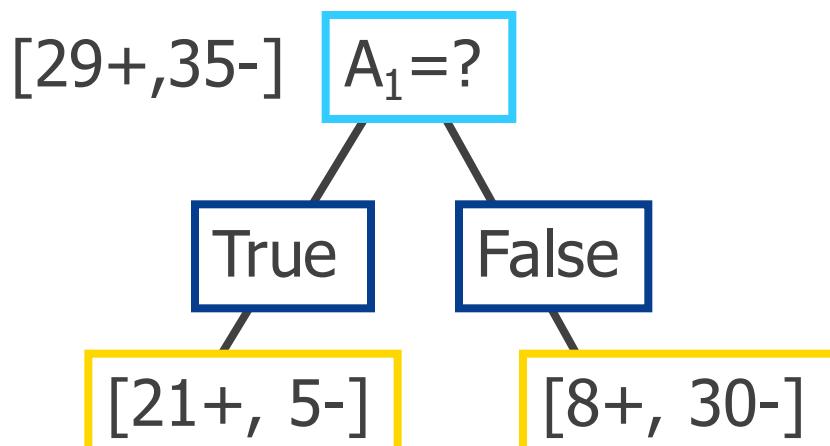
$$\text{Entropy}([8+, 30-]) = 0.74$$

$$\text{Gain}(S, A1) = \text{Entropy}(S)$$

$$-26/64 * \text{Entropy}([21+, 5-])$$

$$-38/64 * \text{Entropy}([8+, 30-])$$

$$= 0.27$$



$$\text{Entropy}([18+, 33-]) = 0.94$$

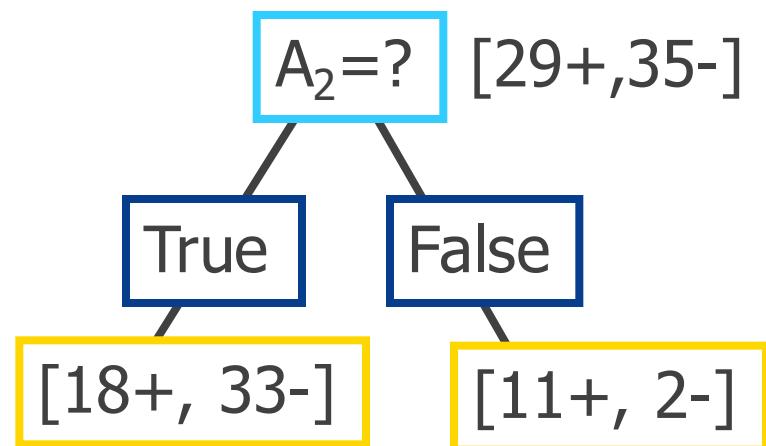
$$\text{Entropy}([11+, 2-]) = 0.62$$

$$\text{Gain}(S, A2) = \text{Entropy}(S)$$

$$-51/64 * \text{Entropy}([18+, 33-])$$

$$-13/64 * \text{Entropy}([11+, 2-])$$

$$= 0.12$$



Training Examples

Day	Outlook	Temp.	Humidity	Wind	Play Tennis
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	Weak	Yes
D8	Sunny	Mild	High	Weak	No
D9	Sunny	Cold	Normal	Weak	Yes
D10	Rain	Mild	Normal	Strong	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

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

$$E=0.940$$

Humidity

High

$$[3+, 4-]$$

$$E=0.985$$

Normal

$$[6+, 1-]$$

$$E=0.592$$

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

$$E=0.940$$

Wind

Weak

$$[6+, 2-]$$

$$E=0.811$$

Strong

$$[3+, 3-]$$

Gain(S, Humidity)

$$=0.940 - (7/14) * 0.985$$

$$- (7/14) * 0.592$$

$$=0.151$$

Humidity provides greater info. gain than Wind, w.r.t target classification.

Gain(S, Wind)

$$E=1.0 = 0.940 - (8/14) * 0.811$$

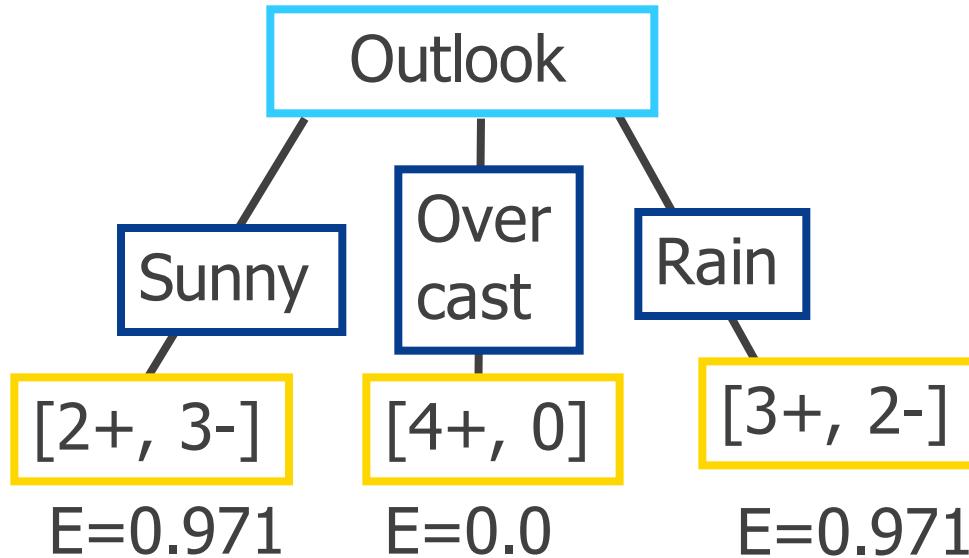
$$- (6/14) * 1.0$$

$$=0.048$$

Selecting the Next Attribute

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

$$E=0.940$$



$$\begin{aligned} \text{Gain}(S, \text{Outlook}) &= 0.940 - (5/14) * 0.971 \\ &\quad - (4/14) * 0.0 - (5/14) * 0.0971 \\ &= 0.247 \end{aligned}$$

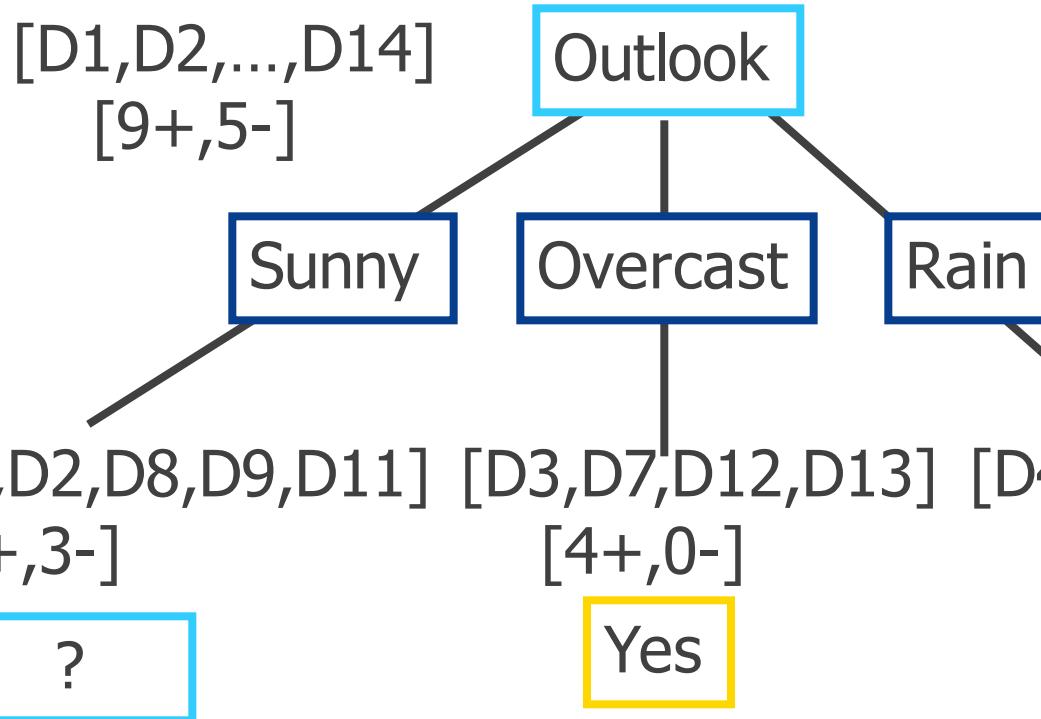
Selecting the Next Attribute

The information gain values for the 4 attributes are:

- $\text{Gain}(S, \text{Outlook}) = 0.247$
- $\text{Gain}(S, \text{Humidity}) = 0.151$
- $\text{Gain}(S, \text{Wind}) = 0.048$
- $\text{Gain}(S, \text{Temperature}) = 0.029$

where S denotes the collection of training examples

ID3 Algorithm



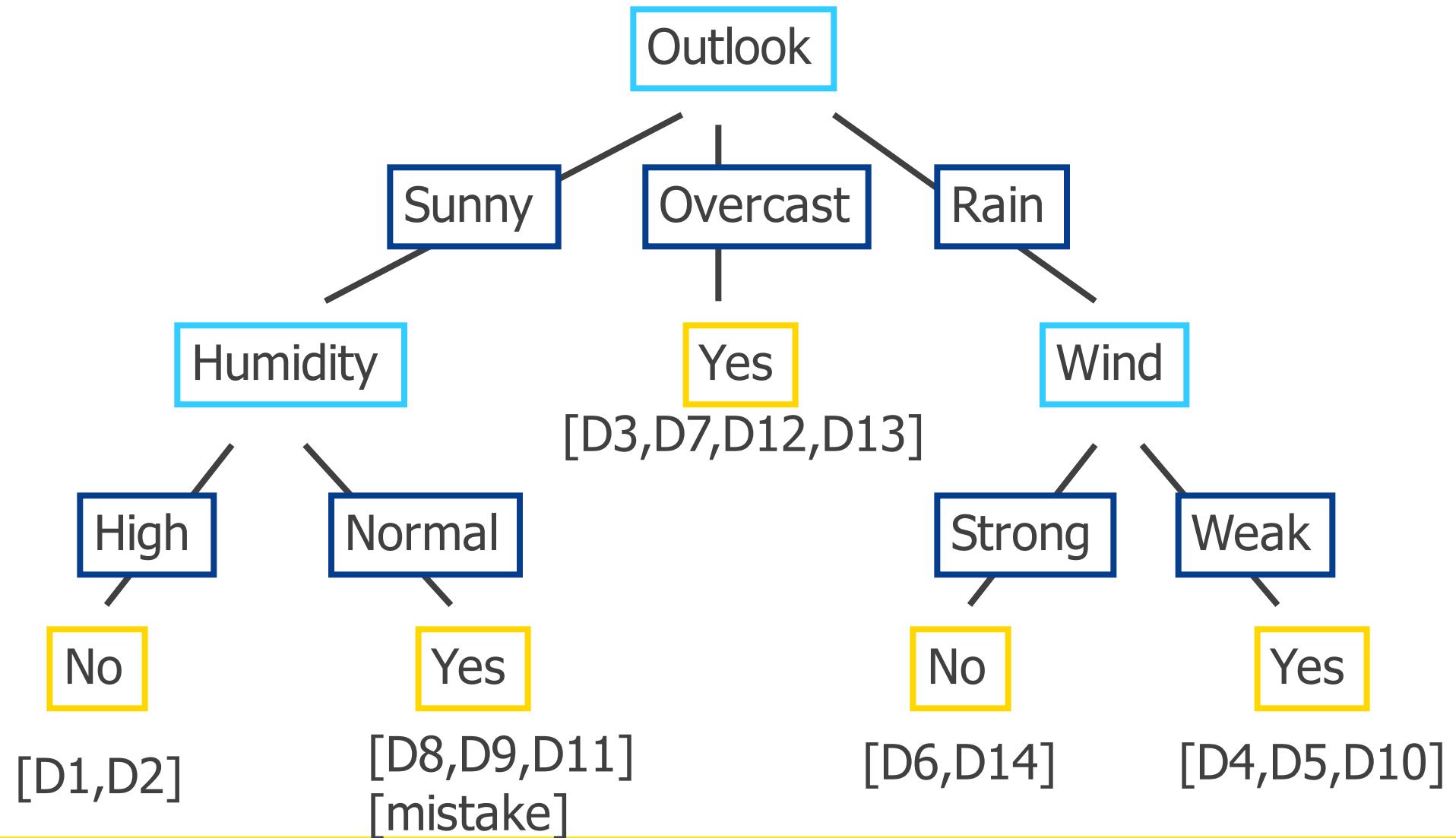
$S_{\text{sunny}} = [D1, D2, D8, D9, D11] [D3, D7, D12, D13] [D4, D5, D6, D10, D14]$
 $[2+, 3-]$ $[4+, 0-]$ $[3+, 2-]$

$$\text{Gain}(S_{\text{sunny}}, \text{Humidity}) = 0.970 - (3/5)0.0 - 2/5(0.0) = 0.970$$

$$\text{Gain}(S_{\text{sunny}}, \text{Temp.}) = 0.970 - (2/5)0.0 - 2/5(1.0) - (1/5)0.0 = 0.570$$

$$\text{Gain}(S_{\text{sunny}}, \text{Wind}) = 0.970 - (2/5)1.0 - 3/5(0.918) = 0.019$$

ID3 Algorithm



Occam's Razor

"If two theories explain the facts equally well, then the simpler theory is to be preferred"

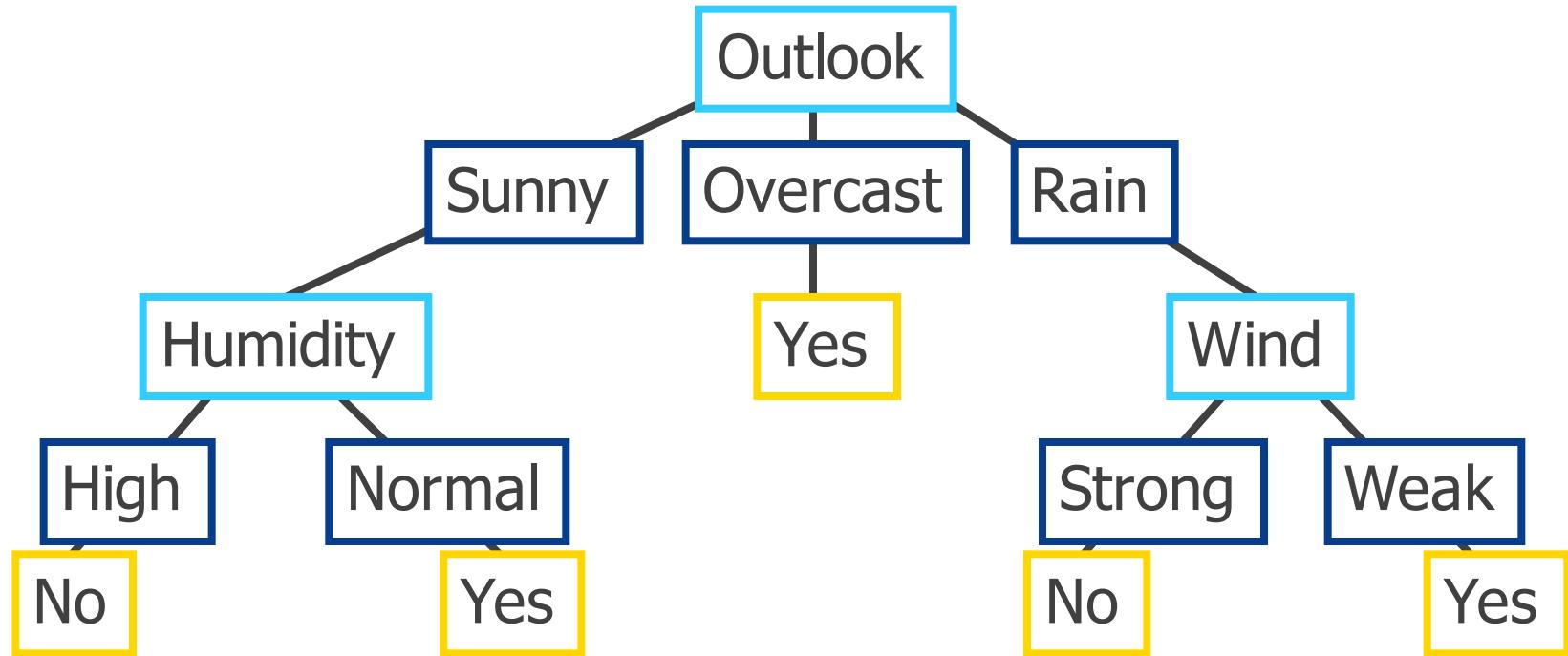
Arguments in favor:

- Fewer short hypotheses than long hypotheses
- A short hypothesis that fits the data is unlikely to be a coincidence
- A long hypothesis that fits the data might be a coincidence

Arguments opposed:

- There are many ways to define small sets of hypotheses

Converting a Tree to Rules



- R₁: If (Outlook=Sunny) \wedge (Humidity=High) Then PlayTennis>No
- R₂: If (Outlook=Sunny) \wedge (Humidity=Normal) Then PlayTennis=Yes
- R₃: If (Outlook=Overcast) Then PlayTennis=Yes
- R₄: If (Outlook=Rain) \wedge (Wind=Strong) Then PlayTennis>No
- R₅: If (Outlook=Rain) \wedge (Wind=Weak) Then PlayTennis=Yes

Continuous Valued Attributes

Create a discrete attribute to test continuous

Temperature = 24.50C

(Temperature > 20.00C) = {true, false}

Where to set the threshold?

Temperature	15 ^o C	18 ^o C	19 ^o C	22 ^o C	24 ^o C	27 ^o C
PlayTennis	No	No	Yes	Yes	Yes	No

Unknown Attribute Values

What if some examples have 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

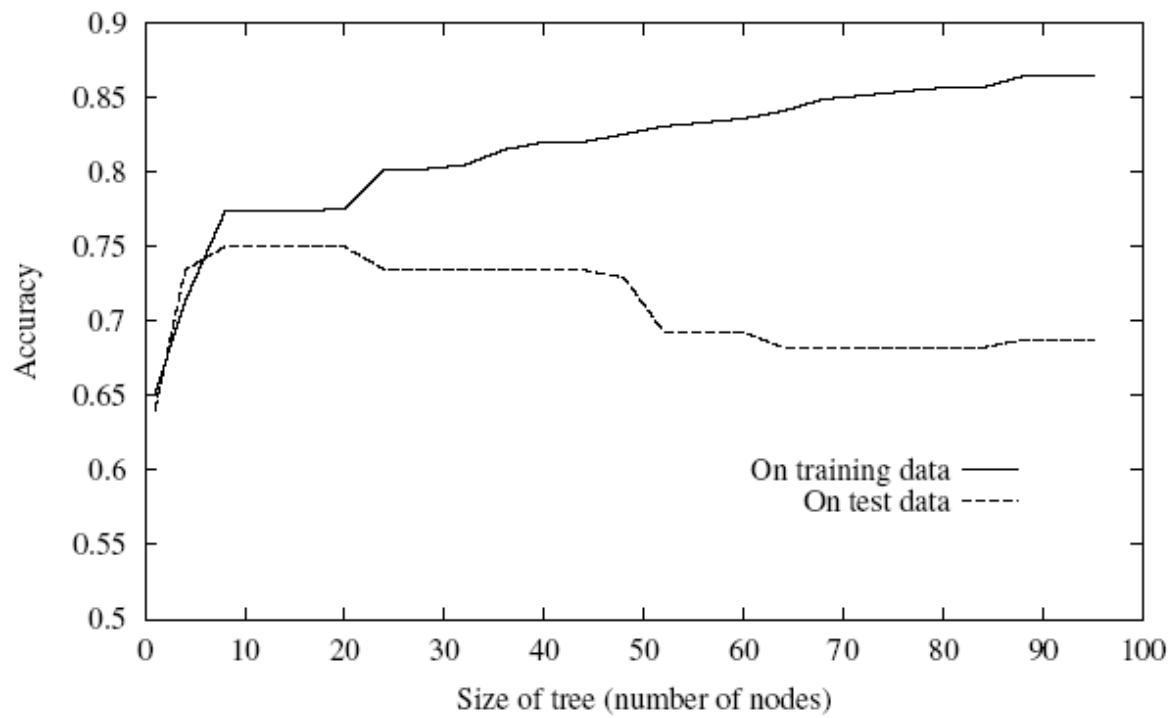
Assign probability π_i to each possible value v_i of A

- Assign fraction π_i of example to each descendant in tree

Classify new examples in the same fashion

Overfitting

One of the biggest problems with decision trees is Overfitting



Avoid Overfitting

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

Select “best” tree:

measure performance over training data

measure performance over separate validation data set

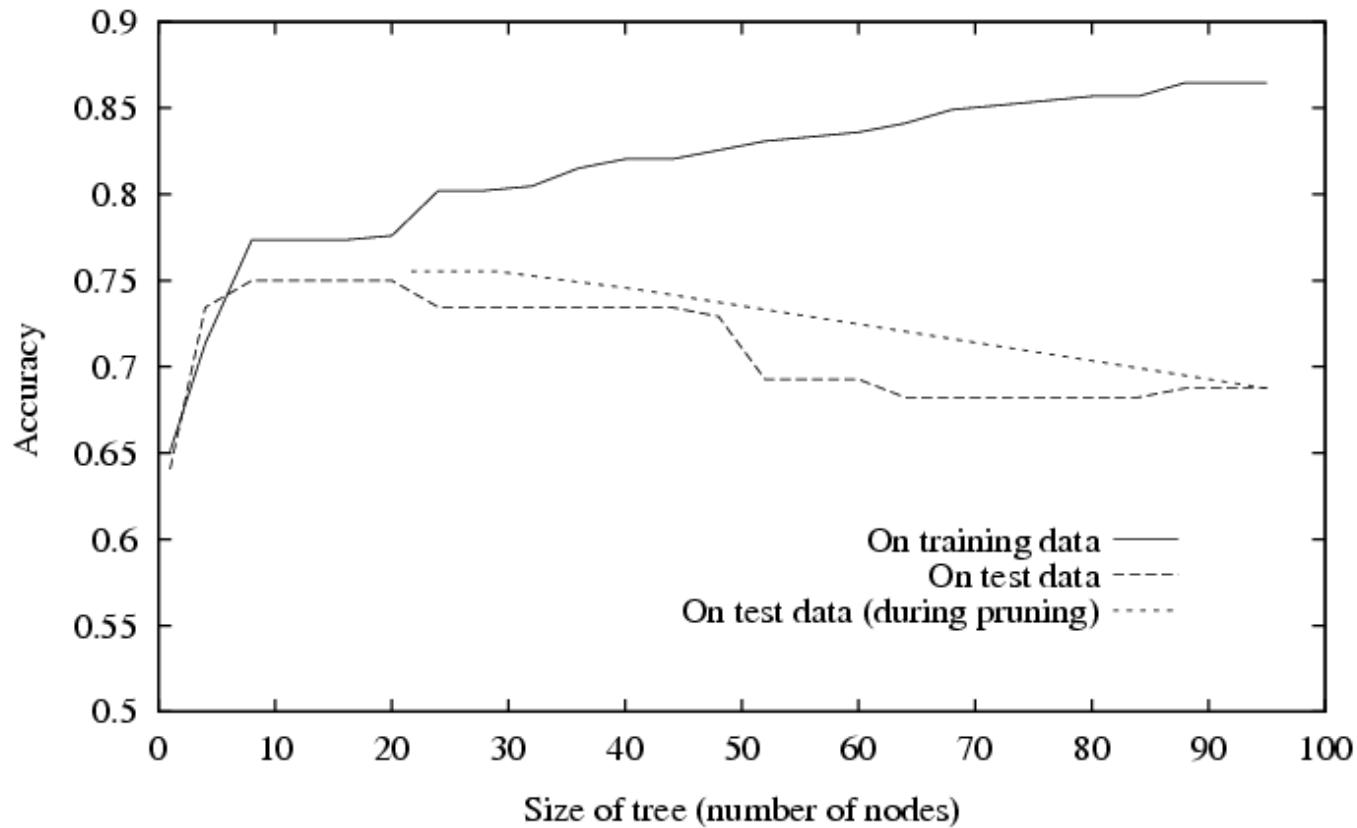
$\min(|\text{tree}| + |\text{misclassifications(tree)}|)$

Avoid Overfitting

Idea 1: Stop growing the tree when the RSS doesn't drop by more than a threshold with any new cut.

Idea 2: Prune a large tree from the leaves to the root.
Weakest link pruning:

Effect of Reduced Error Pruning



Cross-Validation

Idea 3: Find the optimal subtree by cross validation. →
There are too many possibilities – harder than best subsets!

Estimate the accuracy of an hypothesis induced by a supervised learning algorithm

Predict the accuracy of an hypothesis over future unseen instances

Select the optimal hypothesis from a given set of alternative hypotheses

- Pruning decision trees
- Model selection
- Feature selection

Combining multiple classifiers (bagging/boosting)

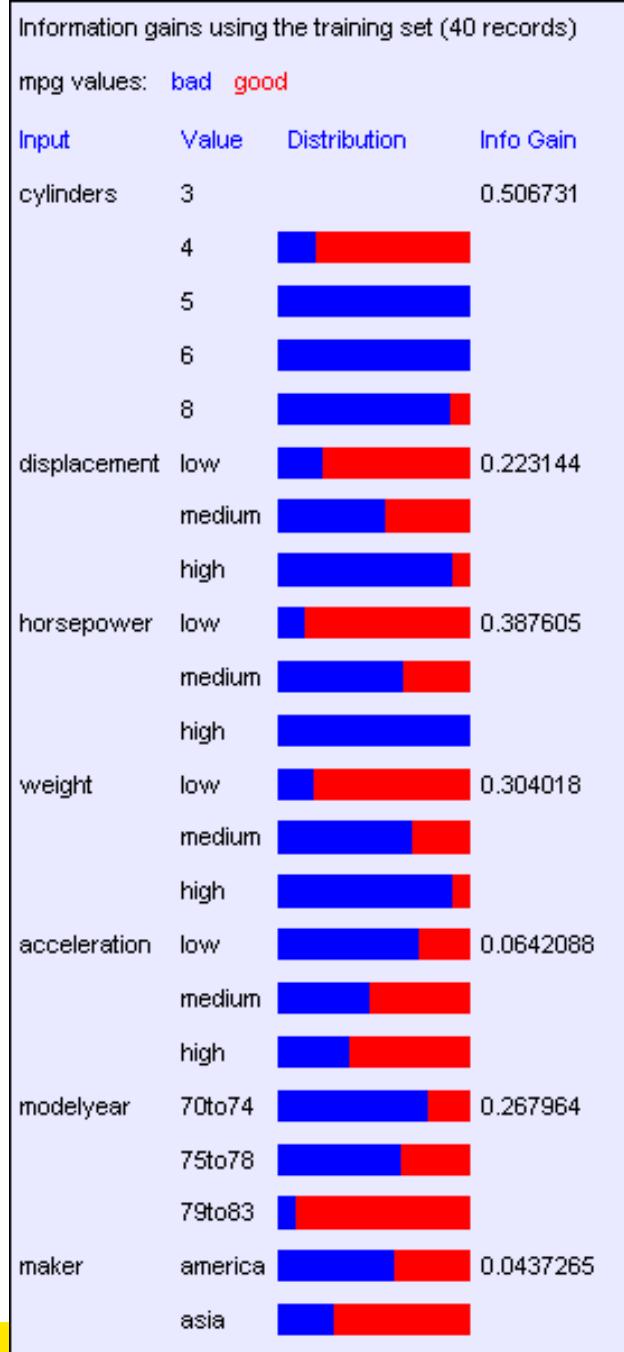
Automobile Miles-per-gallon prediction



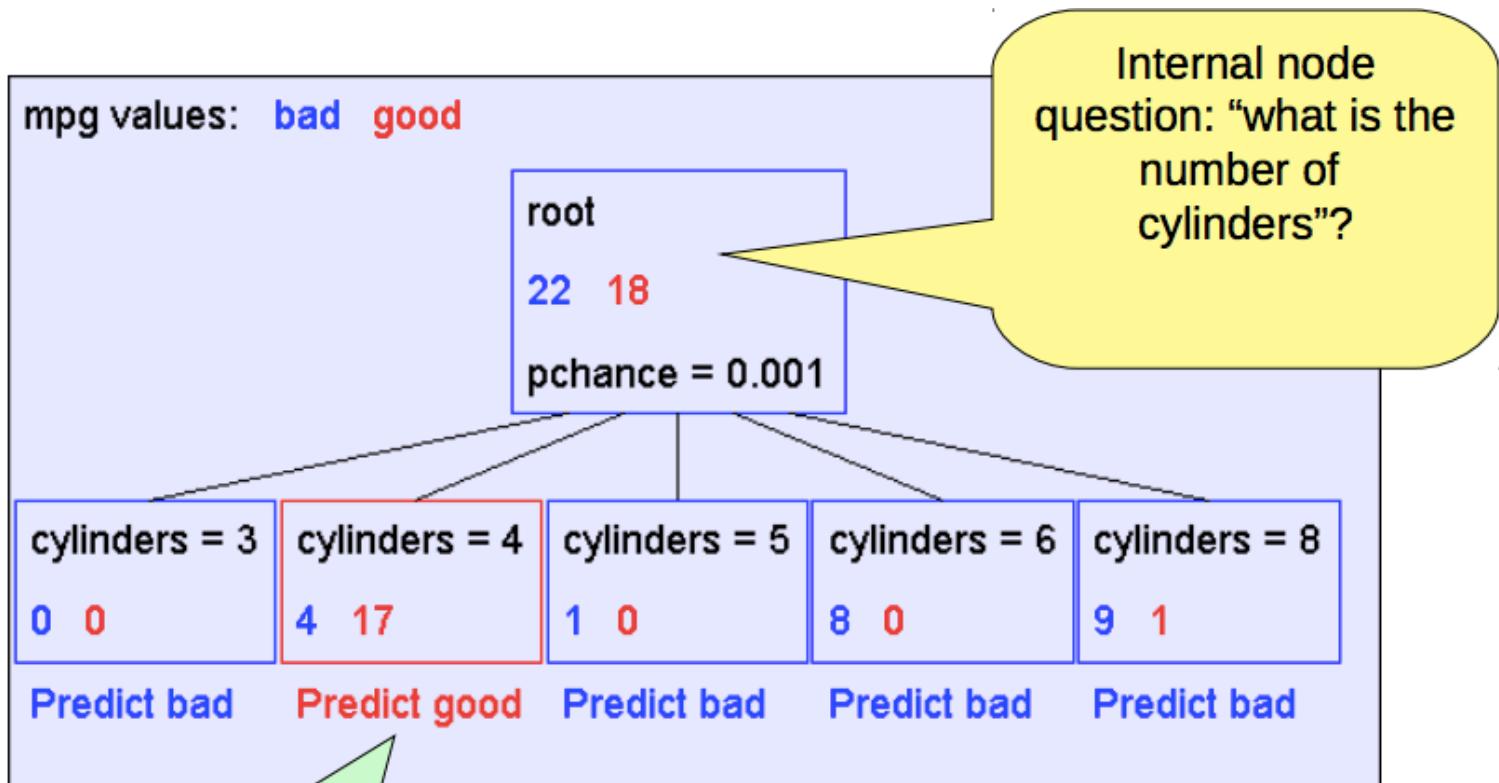
mpg	cylinders	displacement	horsepower	weight	acceleration	modelyear	maker
good	4	low	low	low	high	75to78	asia
bad	6	medium	medium	medium	medium	70to74	america
bad	4	medium	medium	medium	low	75to78	europe
bad	8	high	high	high	low	70to74	america
bad	6	medium	medium	medium	medium	70to74	america
bad	4	low	medium	low	medium	70to74	asia
bad	4	low	medium	low	low	70to74	asia
bad	8	high	high	high	low	75to78	america
:	:	:	:	:	:	:	:
:	:	:	:	:	:	:	:
:	:	:	:	:	:	:	:
bad	8	high	high	high	low	70to74	america
good	8	high	medium	high	high	79to83	america
bad	8	high	high	high	low	75to78	america
good	4	low	low	low	low	79to83	america
bad	6	medium	medium	medium	high	75to78	america
good	4	medium	low	low	low	79to83	america
good	4	low	low	medium	high	79to83	america
bad	8	high	high	high	low	70to74	america
good	4	low	medium	low	medium	75to78	europe
bad	5	medium	medium	medium	medium	75to78	europe

Look at all the information gains...

Suppose we want
to predict MPG.



A Small Decision Tree



Leaves: classify by majority vote

A Growing Decision Tree

mpg values: bad good

question: "what is the value of maker?"

question: "what is the value of horsepower?"

cylinders = 3

0 0

Predict bad

cylinders = 4

4 17

pchance = 0.135

root

22 18

pchance = 0.001

cylinders = 5

1 0

Predict bad

cylinders = 6

8 0

Predict bad

cylinders = 8

9 1

pchance = 0.085

maker = america

0 10

Predict good

maker = asia

2 5

Predict good

maker = europe

2 2

Predict bad

horsepower = low

0 0

Predict bad

horsepower = medium

0 1

Predict good

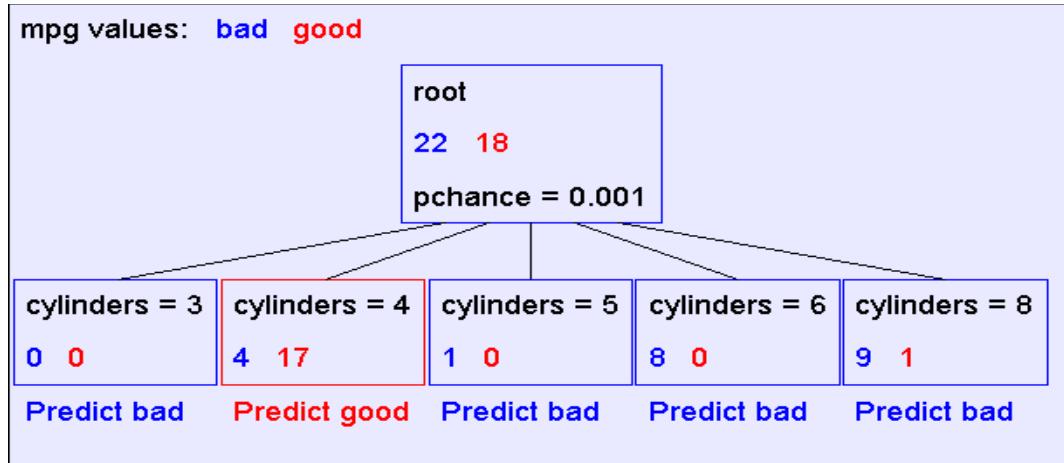
horsepower = high

9 0

Predict bad

Predict "good" is also reasonable by following its parent node instead of the root node.

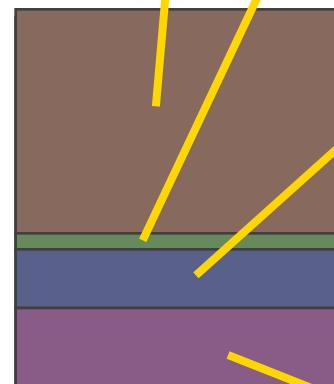
Recursion Step



Take the Original Dataset..



And partition it according to the value of the attribute we split on



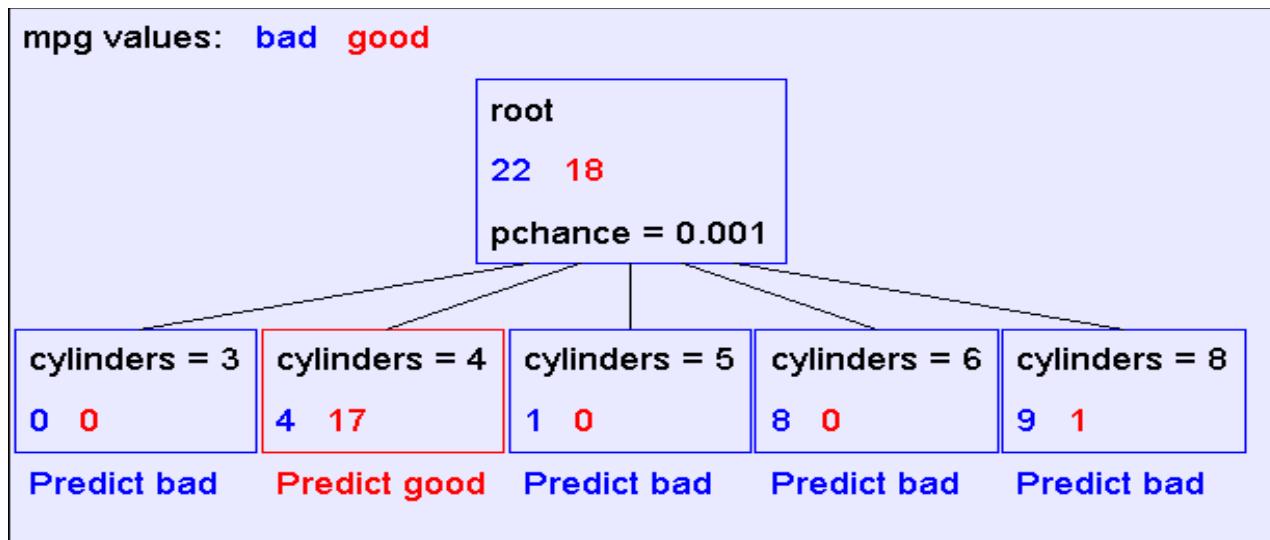
Records in which cylinders = 4

Records in which cylinders = 5

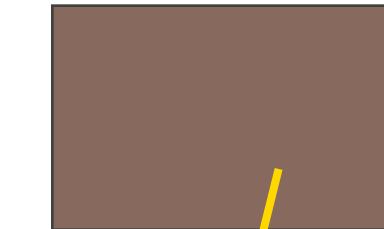
Records in which cylinders = 6

Records in which cylinders = 8

Recursion Step



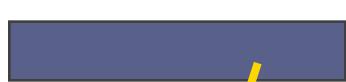
Build tree from
These records..



Records in
which cylinders
= 3



Records in
which cylinders
= 4

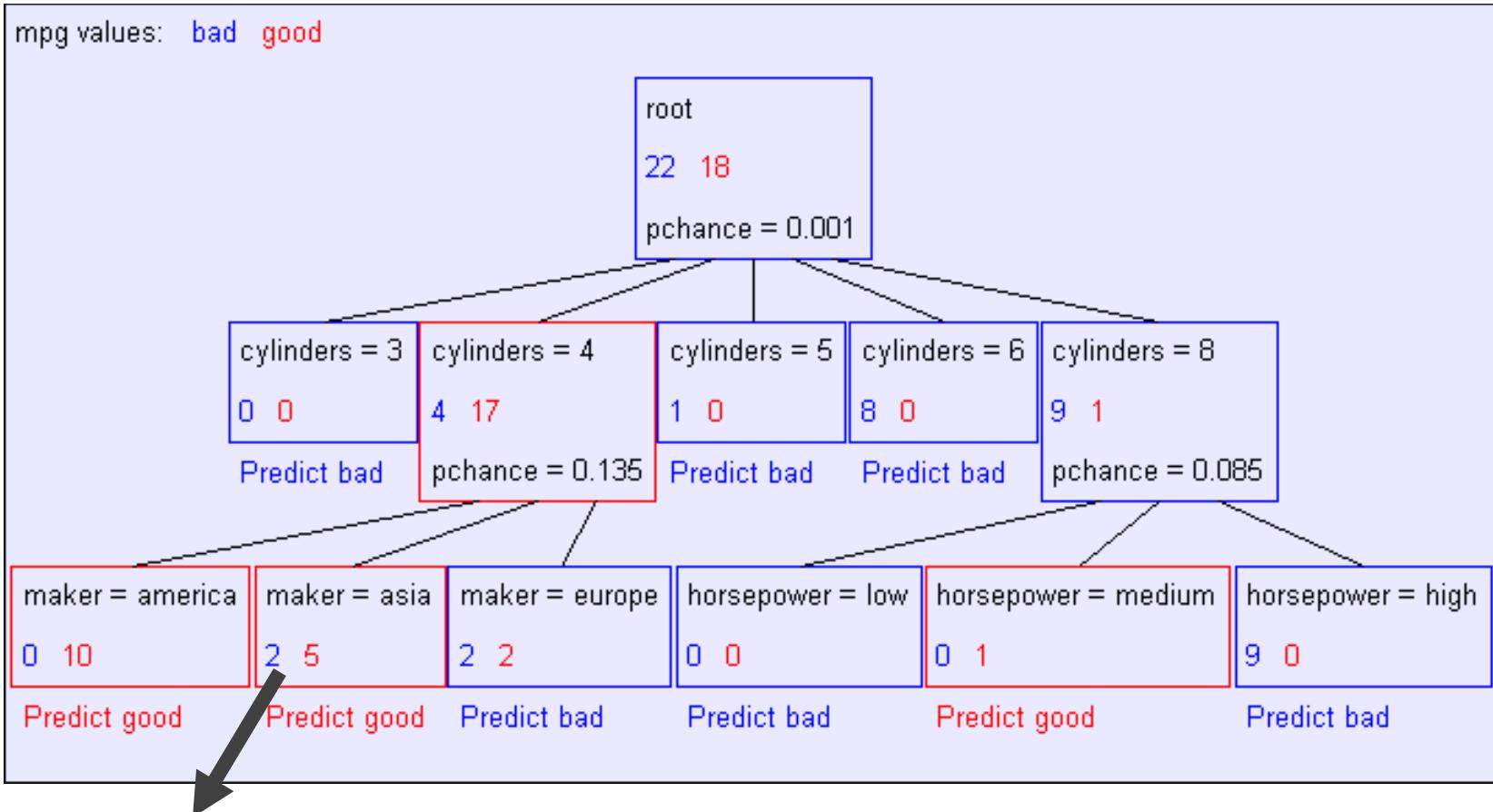


Records in
which cylinders
= 5



Records in
which cylinders
= 8

Second level of tree

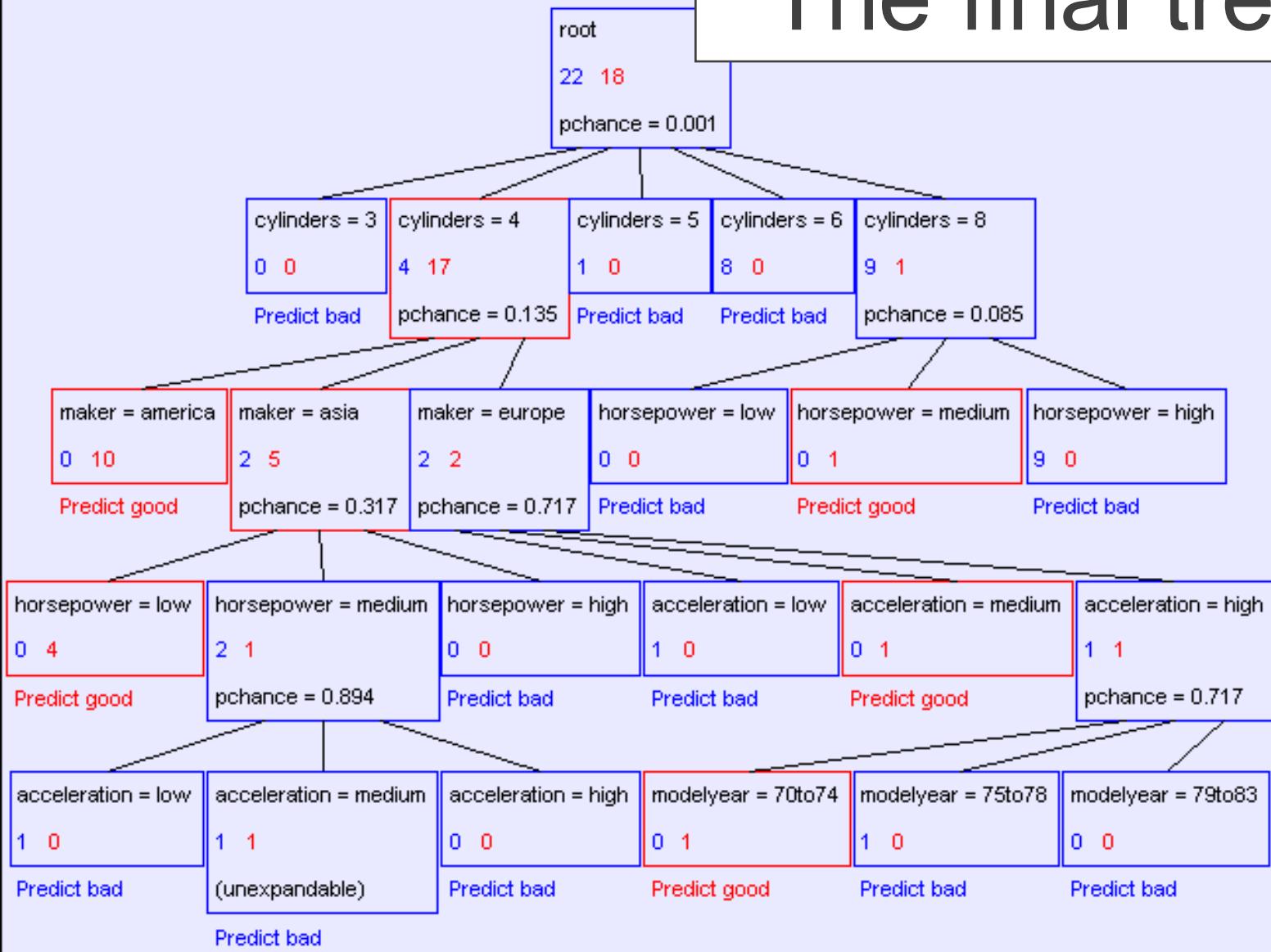


Recursively build a tree from the seven records in which there are four cylinders and the maker was based in Asia

(Similar recursion in the other cases)

mpg values: bad good

The final tree



Significance of Split

Starting with:

- Three cars with 4 cylinders, from Asia, with medium HP
- 2 bad MPG
- 1 good MPG

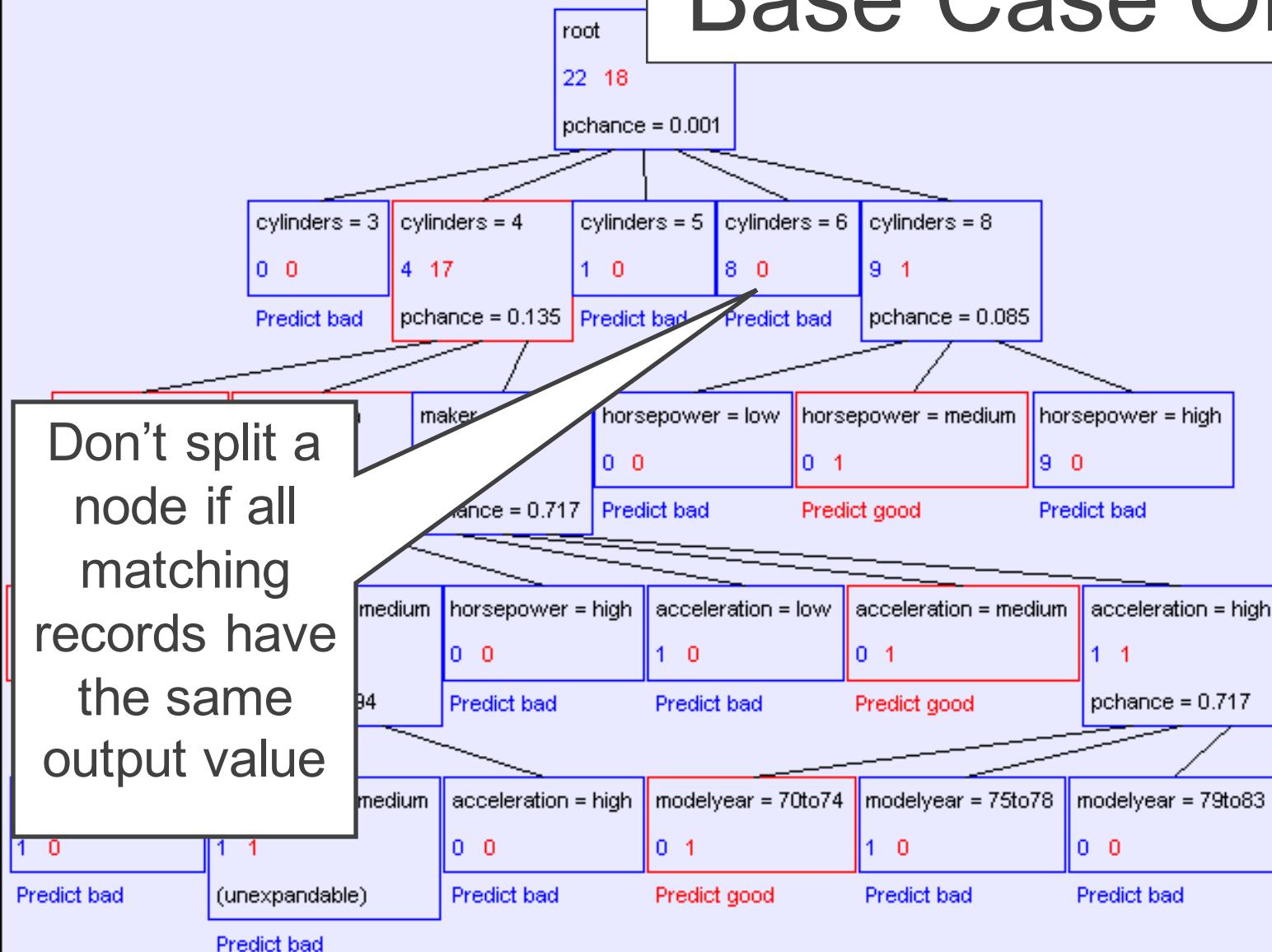
Probably shouldn't split if the counts are so small they could be due to chance A chi-squared test can tell us how likely it is that deviations from a perfect split are due to chance

Each split will have a significance value, pchance

mpg values: bad good

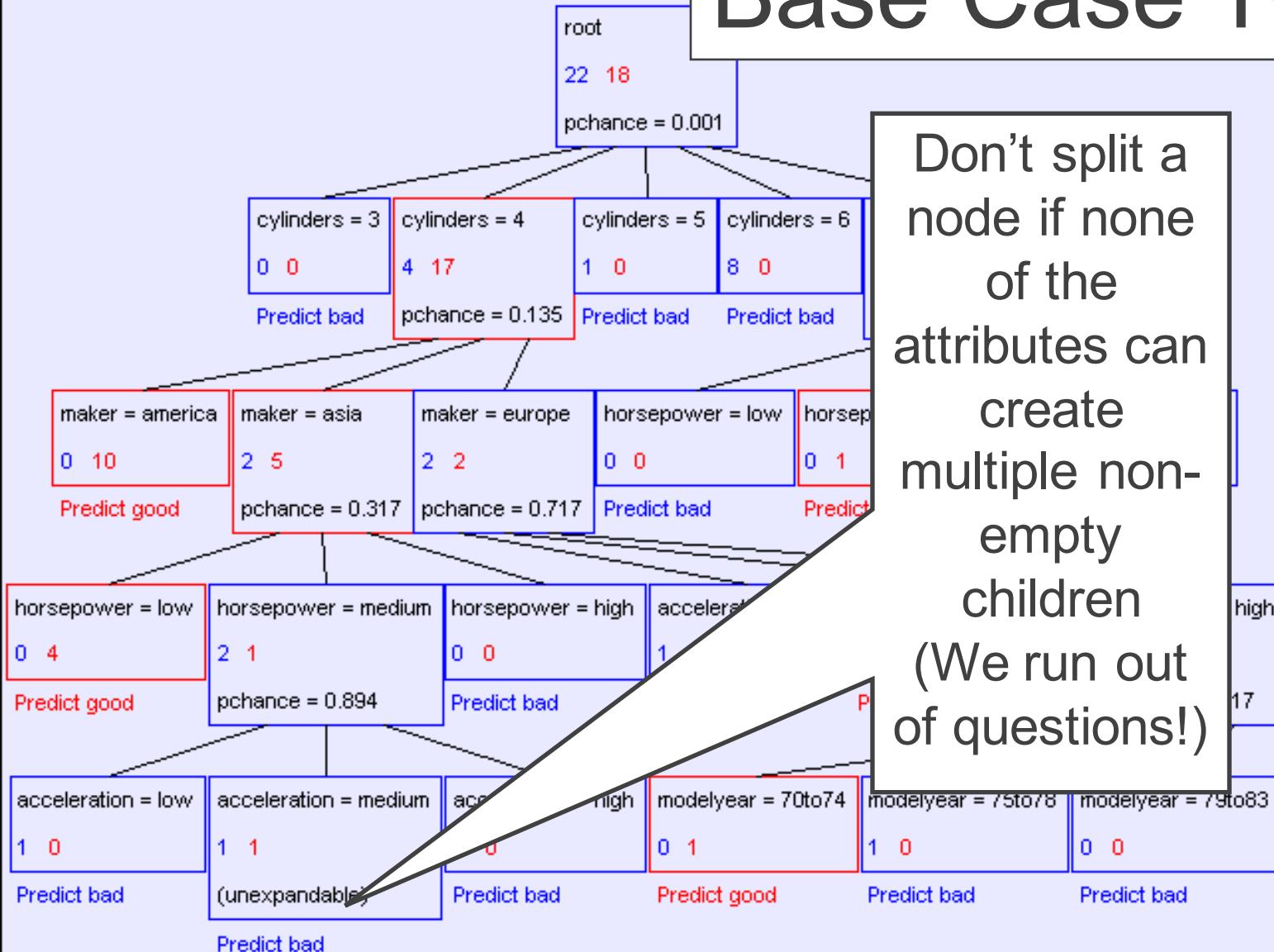
Base Case One

Don't split a node if all matching records have the same output value

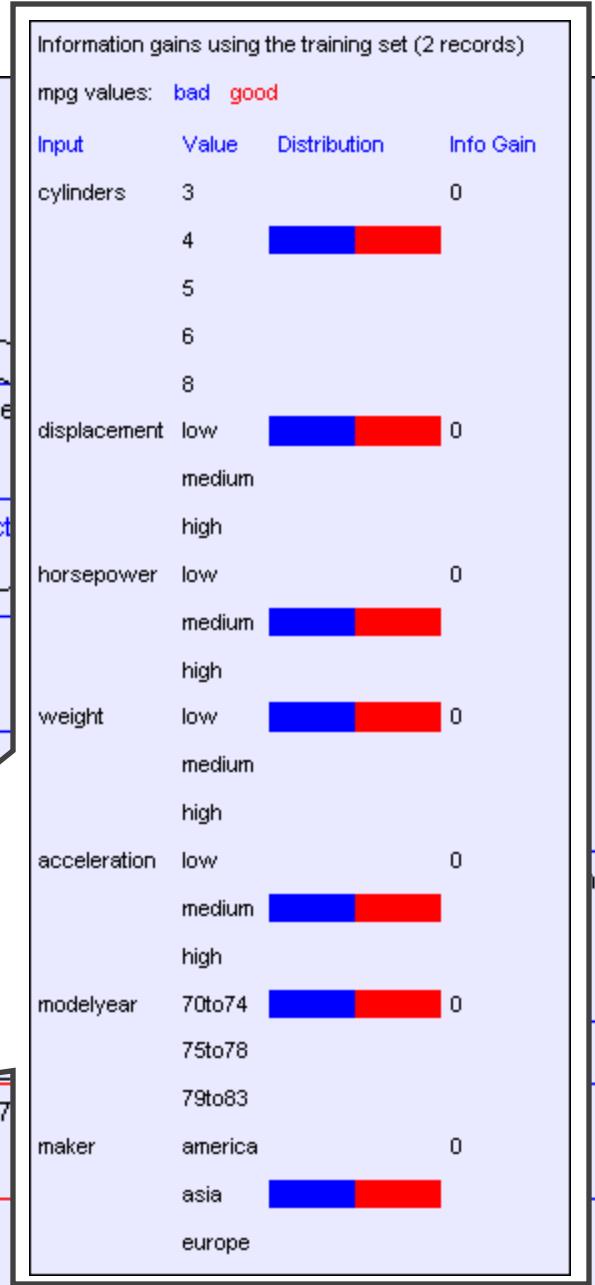
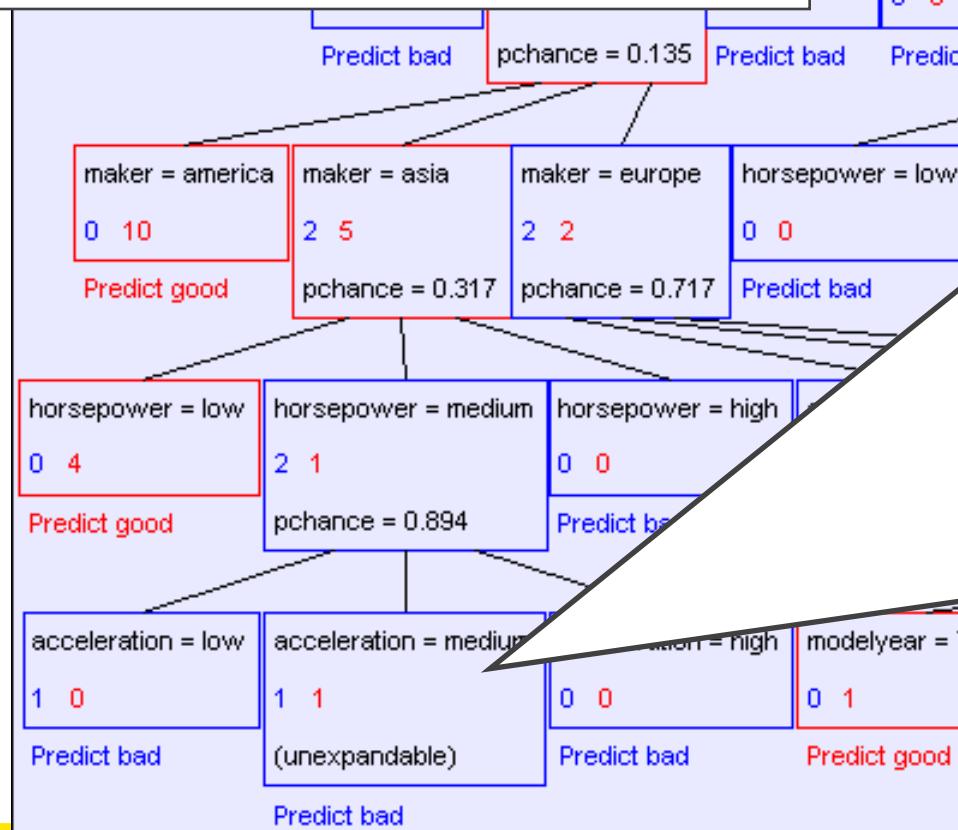


mpg values: bad good

Base Case Two



Base Case Two: No attributes can distinguish



Base Cases

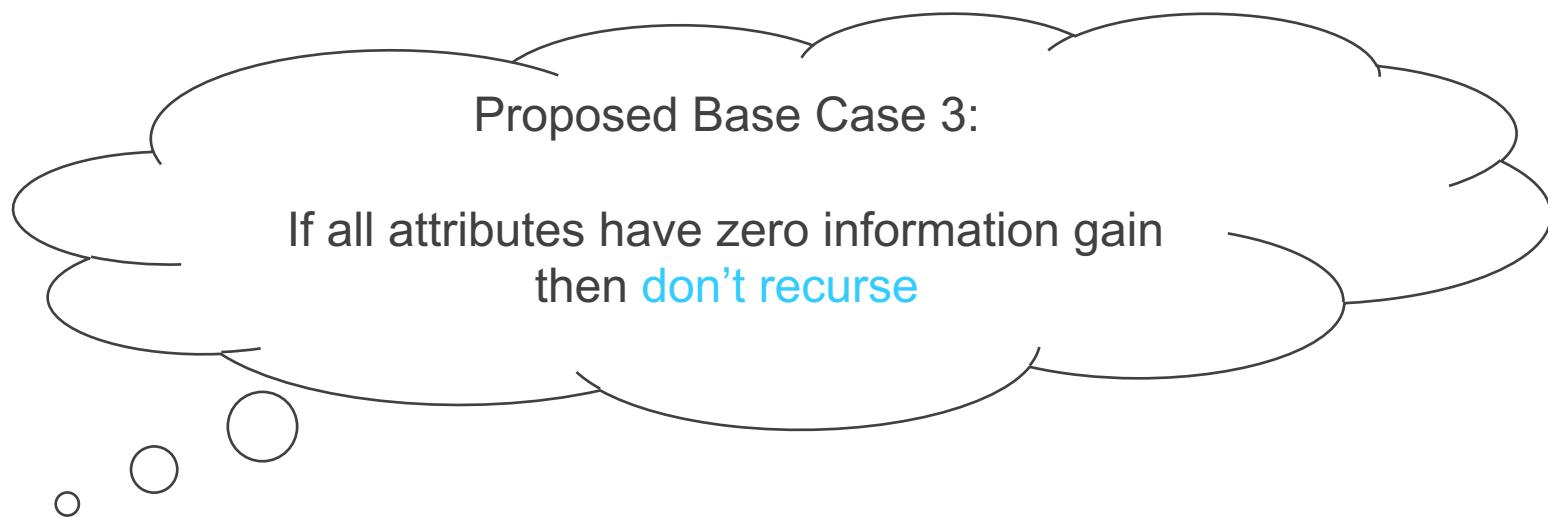
Base Case One: If all records in current data subset have the same output then don't recurse

Base Case Two: If all records have exactly the same set of input attributes then don't recurse

Base Cases: An idea

Base Case One: If all records in current data subset have the same output then don't recurse

Base Case Two: If all records have exactly the same set of input attributes then don't recurse



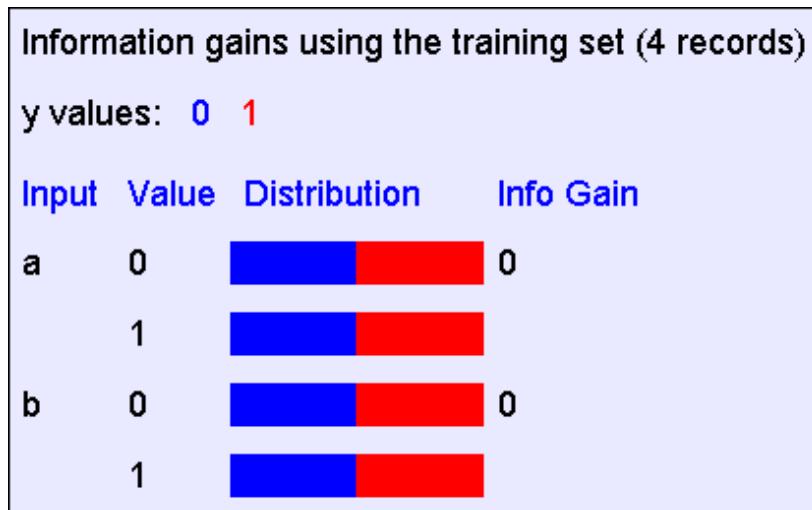
- *Is this a good idea?*

The problem with Base Case 3

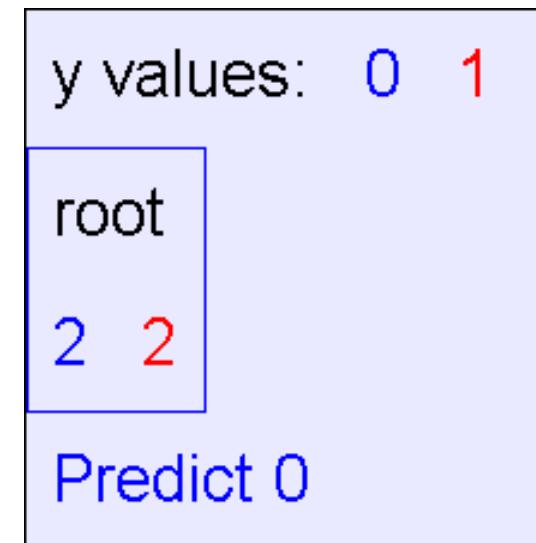
a	b	y
0	0	0
0	1	1
1	0	1
1	1	0

$$y = a \text{ XOR } b$$

The information gains:



The resulting decision tree:

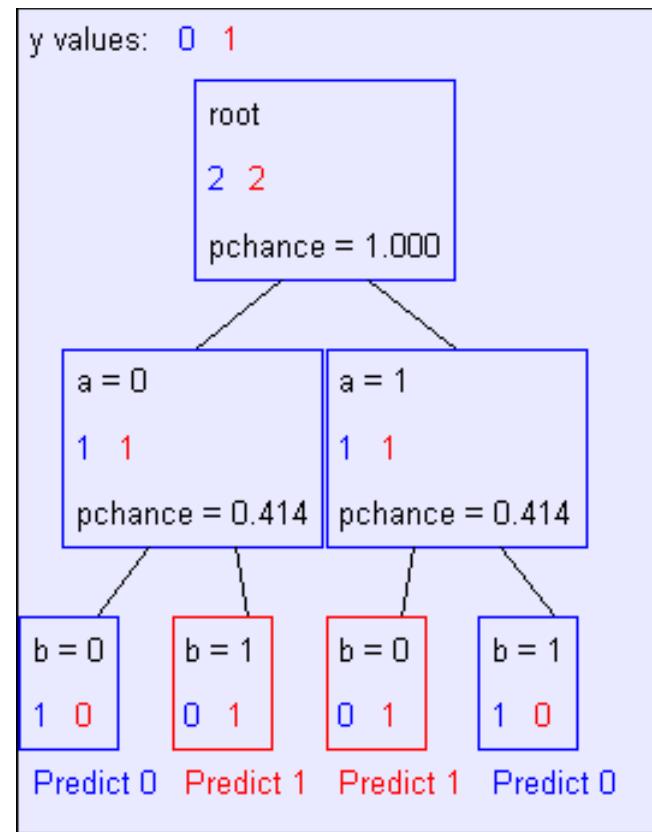


If we omit Base Case 3:

a	b	y
0	0	0
0	1	1
1	0	1
1	1	0

$$y = a \text{ XOR } b$$

The resulting decision tree:



Basic Decision Tree Building Summarized

BuildTree(DataSet,Output)

If all output values are the same in DataSet, return a leaf node that says “predict this unique output”

If all input values are the same, return a leaf node that says “predict the majority output”

Else find attribute X with highest Info Gain

Suppose X has n_X distinct values

- Create and return a non-leaf node with n_X children.
- The i^{th} child should be built by calling
 - BuildTree(DSi,Output)
 - Where DSi built consists of all those records in DataSet for which $X = i^{\text{th}}$ distinct value of X.

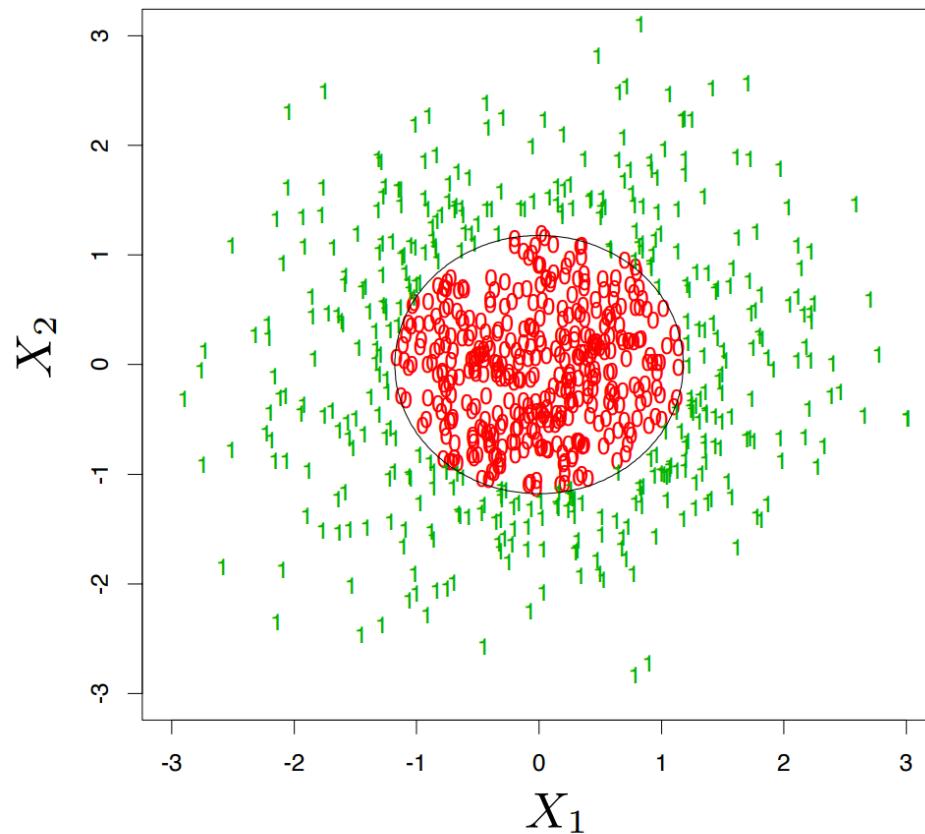
A Quick Summary of Basic Decision Tree

- Can encode any function
- Top-down learning (not perfect!)
- Information gain
- Bottom-up pruning/Cross validation to prevent overfitting

Random Forests

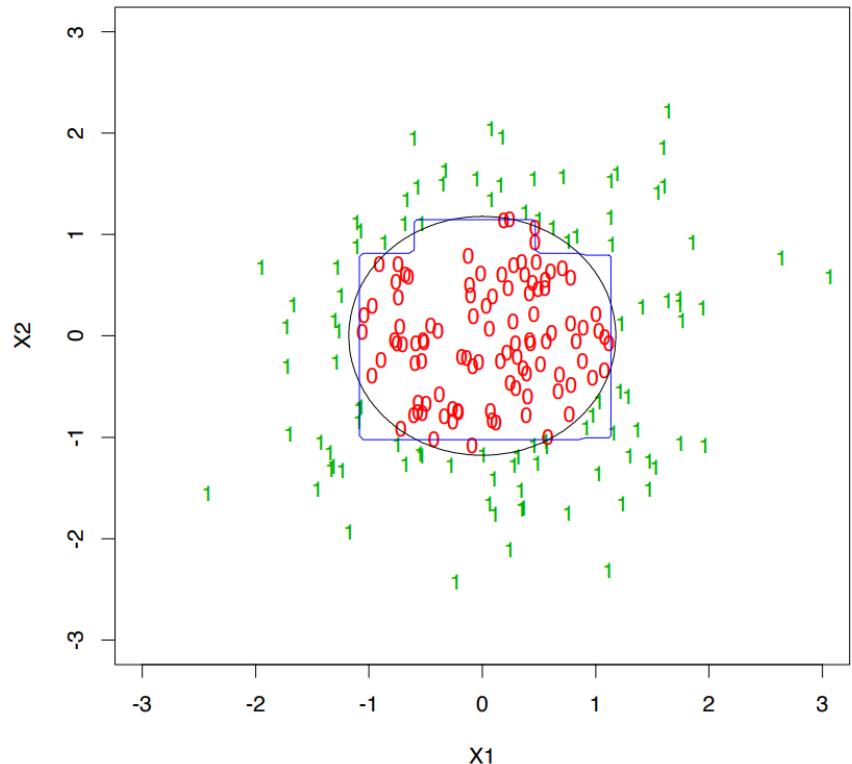
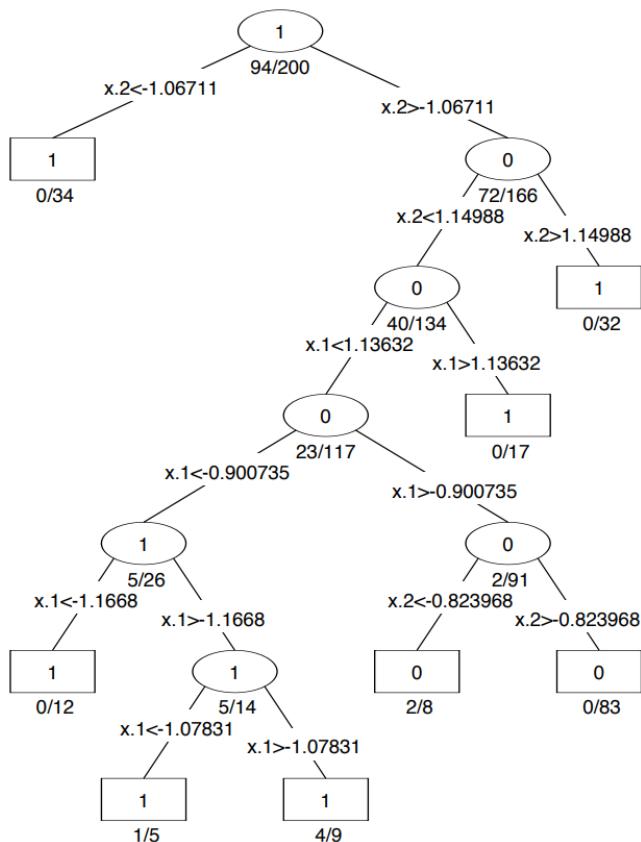
Supervised Learning

Toy Example



Nonlinear separable data.
Optimal decision boundary: $X_1^2 + X_2^2 = 1$

Toy Example



Sample size: 200

7 branching nodes; 6 layers.

Classification error: 7.3% when $d = 2$; > 30% when $d = 10$

Ensemble Learning

Decision trees can be simple, but often overfitting and of high variance

Two minds working together can often achieve better results.

- It is a well-known statistical intuition that averaging measurements can lead to a more stable and reliable

Combination of models: model ensembles

Powerful, but with increased algorithmic and model complexity

Model Averaging

Bagging (Breiman, 1996): Fit many large trees to bootstrap-resampled versions of the training data, and classify by majority vote.

Boosting (Freund & Shapire, 1996): Fit many large or small trees to reweighted versions of the training data. Classify by weighted majority vote.

Random Forests (Breiman 1999): Fancier version of bagging. In general
Boosting>Random Forests>Bagging>Single Tree

Bagging

Bagging, short for ‘bootstrap aggregating’, is a simple but highly effective ensemble method that creates diverse models on different random samples of the original dataset.

These samples are taken uniformly with replacement and are known as bootstrap samples.

Bagging

Bagging is the application of the Bootstrap procedure to a high-variance machine learning algorithm, typically decision trees.

Decision trees are sensitive to the specific data on which they are trained. If the training data is changed (e.g. a tree is trained on a subset of the training data) the resulting decision tree can be quite different and in turn the predictions can be quite different.

Random forest

Definition

Collection of unpruned CARTs

Rule to combine individual tree decisions

Purpose

Improve prediction accuracy

Principle

Encouraging diversity among the tree

Solution: randomness

Bagging

Random decision trees (rCART)

A Forest Model

Parallel, bagged ensemble of the tree models

Strong predictive power, but lower interpretability

Many trees created using subsets of the data and averaged together

More hyperparameters than decisions control the model complexity

- Number of trees
- Sampling rate
- Number of variables to try

Bagging

Advantages:

Technique of ensemble learning...

- ... to avoid over-fitting

- Important since trees are unpruned

- ... to improve stability and accuracy

Two steps

- Bootstrap sample set

- Aggregation

Bootstrap

L: original learning set composed of p samples

Generate K learning sets L_k...

... composed of q samples, $q \leq p$,...

... obtained by uniform sampling with replacement from L

In consequences, L_k may contain repeated samples

Random forest: $q = p$

Asymptotic proportion of unique samples in L_k = 100 (1 - 1/e) $\sim 63\%$

→ The remaining samples can be used for testing

Aggregation

Learning

For each L_k , one classifier C_k (rCART) is learned

Prediction

S : a new sample

Aggregation = majority vote among the K predictions/votes $C_k(S)$

Bagging/Bootstrap Aggregation

Motivation: Average a given procedure over many samples to reduce the variance.

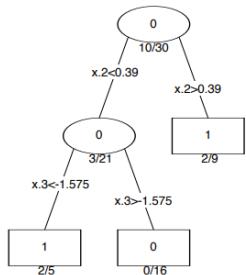
- Given a classifier $C(S, x)$, based on our training data S , producing a predicted class label at input point x .
- To bag C , we draw bootstrap samples S_1, \dots, S_B each of size N with replacement from the training data. Then

$$\hat{C}_{\text{Bag}} = \text{Majority Vote}\{C(S_b, x)\}_{b=1}^B.$$

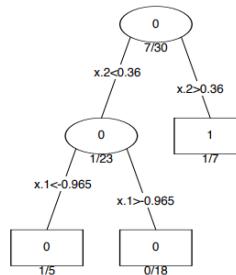
- Bagging can dramatically reduce the variance of unstable procedures (like trees), leading to improved prediction.
- All simple structures in a tree are lost.

Bagging/Bootstrap Aggregation

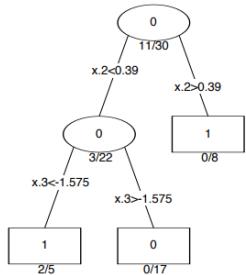
Original Tree



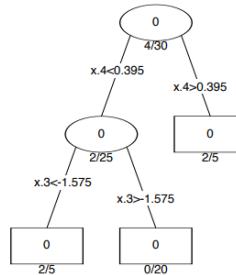
Bootstrap Tree 1



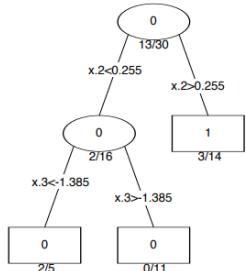
Bootstrap Tree 2



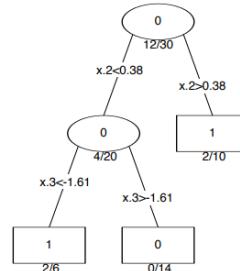
Bootstrap Tree 3



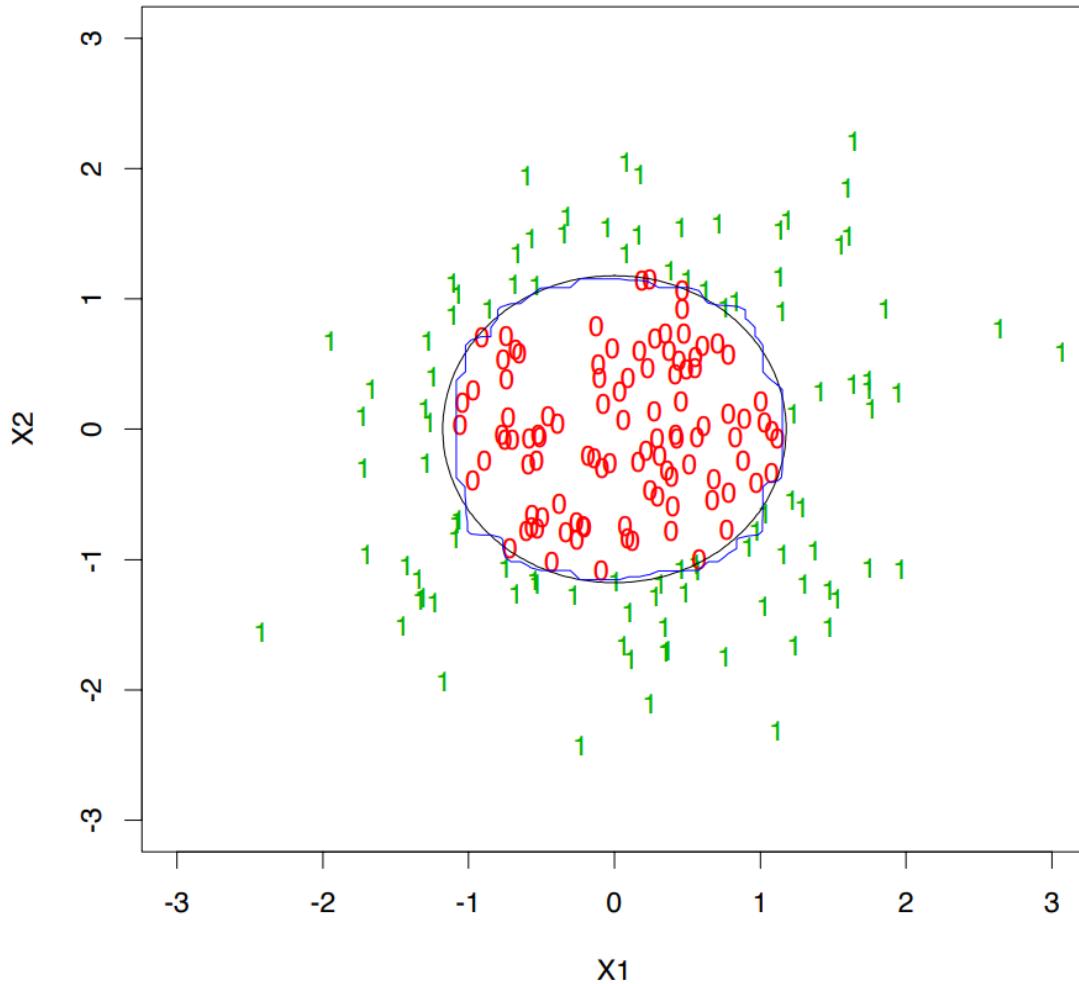
Bootstrap Tree 4



Bootstrap Tree 5



Bagging Decision Trees



A smoother decision boundary.

Classification error: 3.2% (Single deeper tree 7.3%).

Random Forest

A problem with decision trees like CART is that they are greedy. They choose which variable to split on using a greedy algorithm that minimizes error. As such, even with Bagging, the decision trees can have a lot of structural similarities and in turn have high correlation in their predictions.

Combining predictions from multiple models in ensembles works better if the predictions from the sub-models are uncorrelated or at best weakly correlated.

Random forest changes the algorithm for the way that the sub-trees are learned so that the resulting predictions from all of the subtrees have less correlation.

Random Forest Model

Parallel, bagged ensemble of the tree models

Strong predictive power, but lower interpretability

Many trees created using subsets of the data and averaged together

More hyperparameters than decisions control the model complexity

- Number of trees
- Sampling rate
- Number of variables to try

Random Forest

Bagging features and samples simultaneously:

- At each tree split, a random sample of m features is drawn, and only those m features are considered for splitting.
Typically $m = \sqrt{d}$ or $\log_2 d$, where d is the number of features
- For each tree grown on a bootstrap sample, the error rate for observations left out of the bootstrap sample is monitored.
This is called the “out-of-bag” error rate.
- random forests tries to improve on bagging by “de-correlating” the trees. Each tree has the same expectation.

Random decision tree

All labeled samples initially assigned to root node

$N \leftarrow$ root node

With node N do

Find the feature F among a random subset of features + threshold value T ...

... that split the samples assigned to N into 2 subsets S_{left} and S_{right} ...

... so as to maximize the label purity within these subsets

Assign (F, T) to N

If S_{left} and S_{right} too small to be splitted

Attach child leaf nodes L_{left} and L_{right} to N

Tag the leaves with the most present label in S_{left} and S_{right} , resp.

else

Attach child nodes N_{left} and N_{right} to N

Assign S_{left} and S_{right} to them, resp.

Repeat procedure for $N = N_{\text{left}}$ and $N = N_{\text{right}}$

Random subset of features

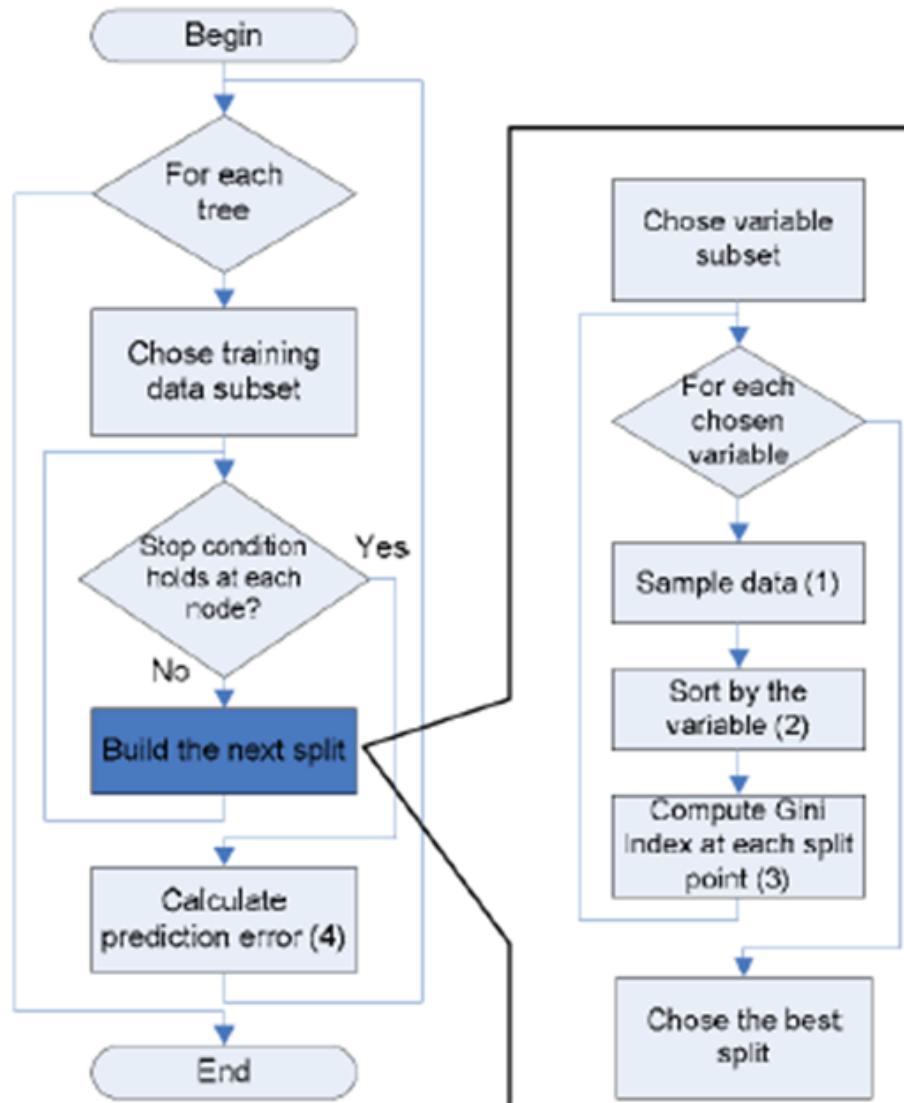
Random drawing repeated at each node

For D -dimensional samples, typical subset size = $\text{round}(\sqrt{D})$ (also $\text{round}(\log_2(x))$)

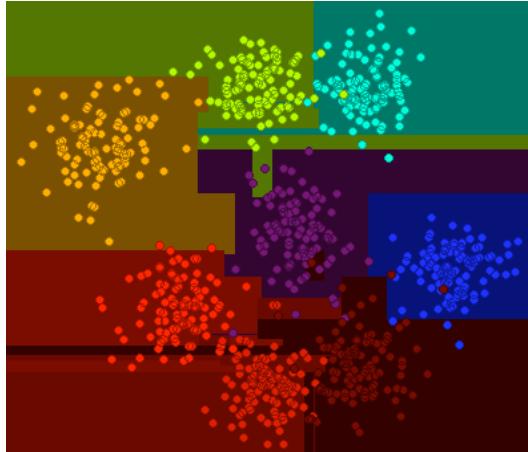
→ Increases diversity among the rCARTs + reduces computational load

Typical purity: Gini index

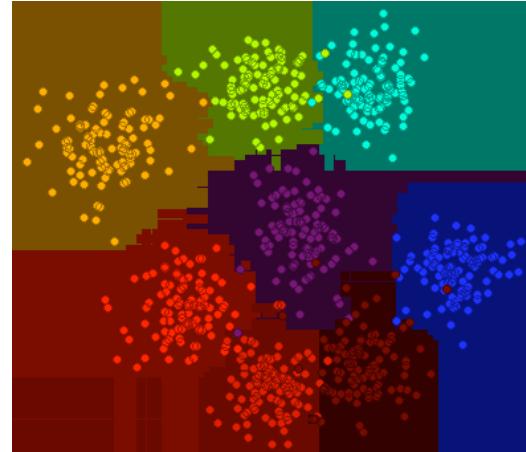
Random decision tree



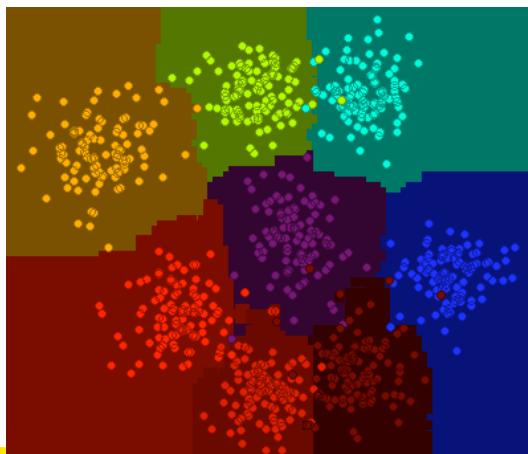
Illustration



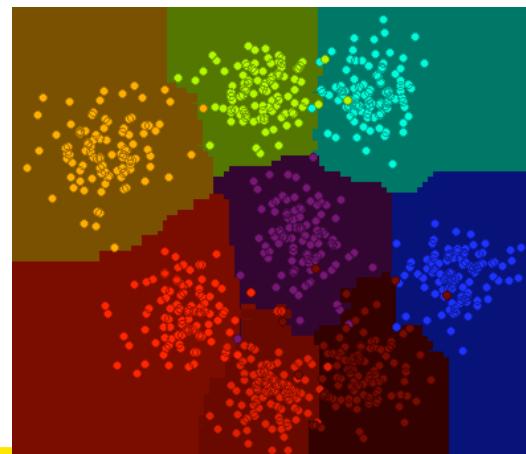
1 rCART



10 rCARTs



100 rCARTs



500 rCARTs

Features and Advantages

- It is one of the most accurate learning algorithms available.
For many data sets, it produces a highly accurate classifier.
- It runs efficiently on large databases.
- It can handle thousands of input variables without variable deletion.

Features and Advantages

It gives estimates of what variables are important in the classification.

It generates an internal unbiased estimate of the generalization error as the forest building progresses.

It has an effective method for estimating missing data and maintains accuracy when a large proportion of the data are missing.

Features and Advantages

It has methods for balancing error in class population unbalanced data sets.

Generated forests can be saved for future use on other data.

Prototypes are computed that give information about the relation between the variables and the classification.

Features and Advantages

It computes proximities between pairs of cases that can be used in clustering, locating outliers, or (by scaling) give interesting views of the data.

The capabilities of the above can be extended to unlabeled data, leading to unsupervised clustering, data views and outlier detection.

It offers an experimental method for detecting variable interactions.