# Classification: Decision Trees

These slides were assembled by Byron Boots, with grateful acknowledgement to Eric Eaton and the many others who made their course materials freely available online. Feel free to reuse or adapt these slides for your own academic purposes, provided that you include proper attribution.

# Function Approximation

**Problem Setting**

- Set of possible instances $\mathcal{X}$

- Set of possible labels $\mathcal{Y}$

- Unknown target function $f : \mathcal{X} \to \mathcal{Y}$

- Set of function hypotheses $H = \{h \mid h : \mathcal{X} \to \mathcal{Y}\}$

**Input**: Training examples of unknown target function f

$$\{\langle \boldsymbol{x}_i, y_i \rangle\}_{i=1}^{n} = \{\langle \boldsymbol{x}_1, y_1 \rangle, \ldots, \langle \boldsymbol{x}_n, y_n \rangle\}$$

**Output**: Hypothesis $h \in H$ that best approximates f
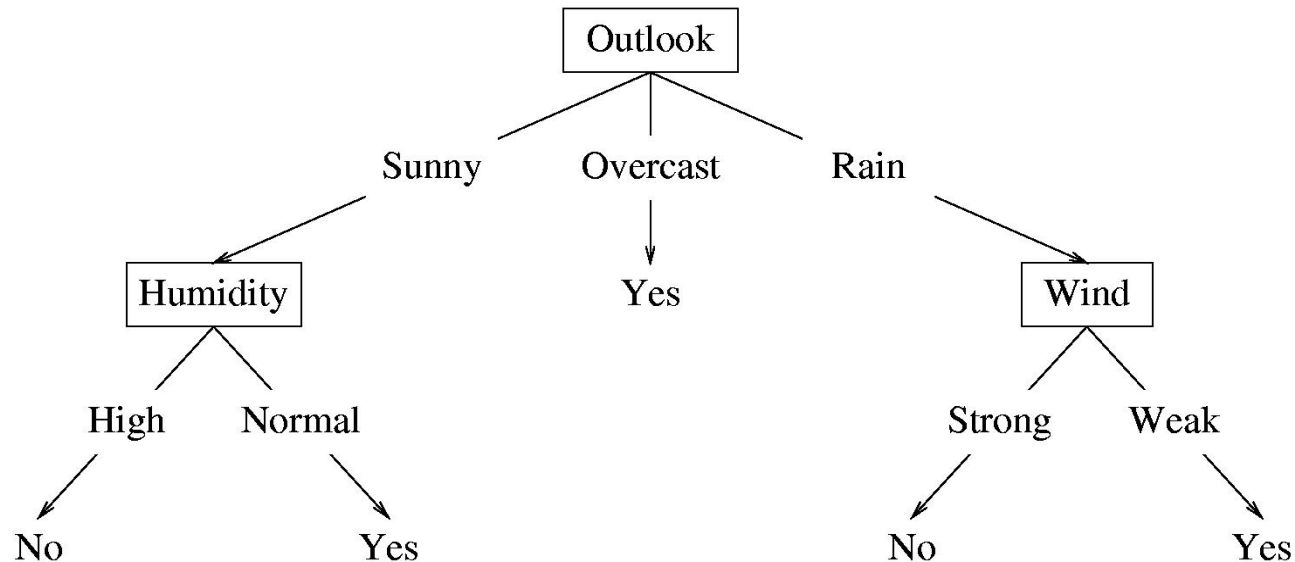
# Sample Dataset (was Tennis Played?)

- Columns denote features $X_i$
- Rows denote labeled instances $\langle \boldsymbol{x}_i, y_i \rangle$
- Class label denotes whether a tennis game was played

$\langle \boldsymbol{x}_i, y_i \rangle$

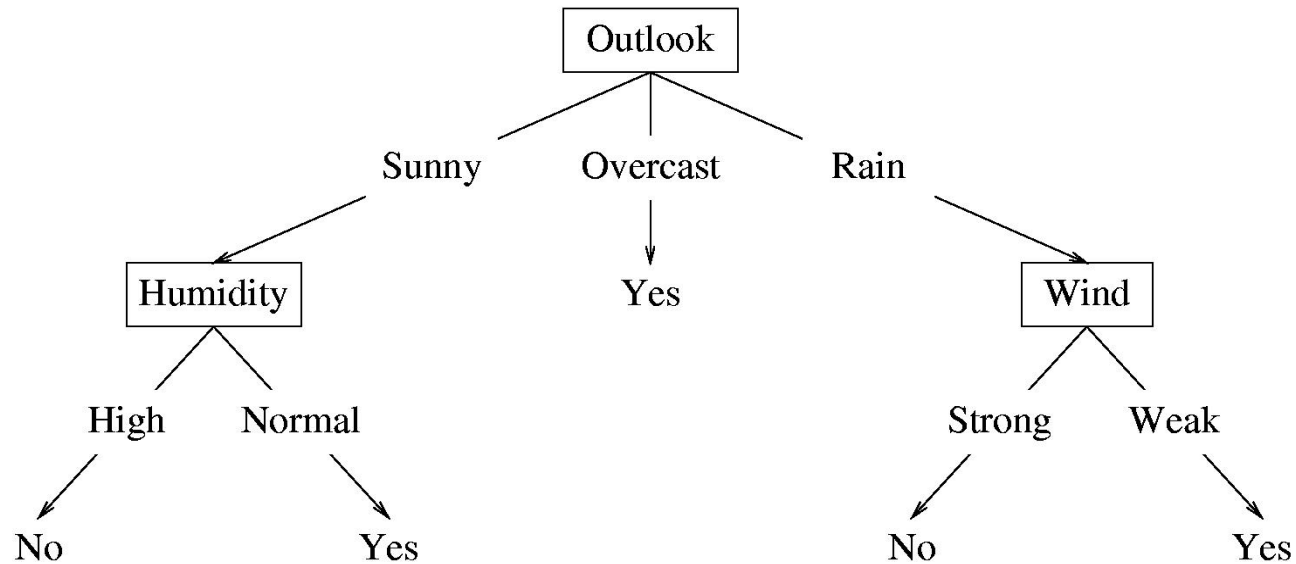| Predictors | | | | Response |
| --- | --- | --- | --- | --- |
| **Outlook** | **Temperature** | **Humidity** | **Wind** | **Class** |
| Sunny | Hot | High | Weak | No |
| Sunny | Hot | High | Strong | No |
| Overcast | Hot | High | Weak | Yes |
| Rain | Mild | High | Weak | Yes |
| Rain | Cool | Normal | Weak | Yes |
| Rain | Cool | Normal | Strong | No |
| Overcast | Cool | Normal | Strong | Yes |
| Sunny | Mild | High | Weak | No |
| Sunny | Cool | Normal | Weak | Yes |
| Rain | Mild | Normal | Weak | Yes |
| Sunny | Mild | Normal | Strong | Yes |
| Overcast | Mild | High | Strong | Yes |
| Overcast | Hot | Normal | Weak | Yes |
| Rain | Mild | High | Strong | No |

# Decision Tree

- A possible decision tree for the data:



- Each internal node: test one attribute $X_i$
- Each branch from a node: selects one value for $X_i$
- Each leaf node: predict $Y$

# Decision Tree

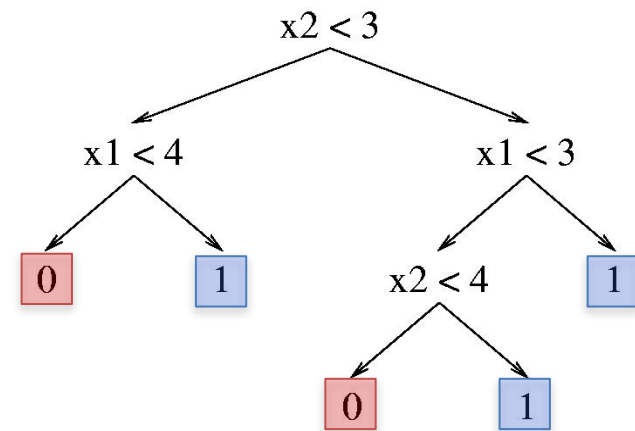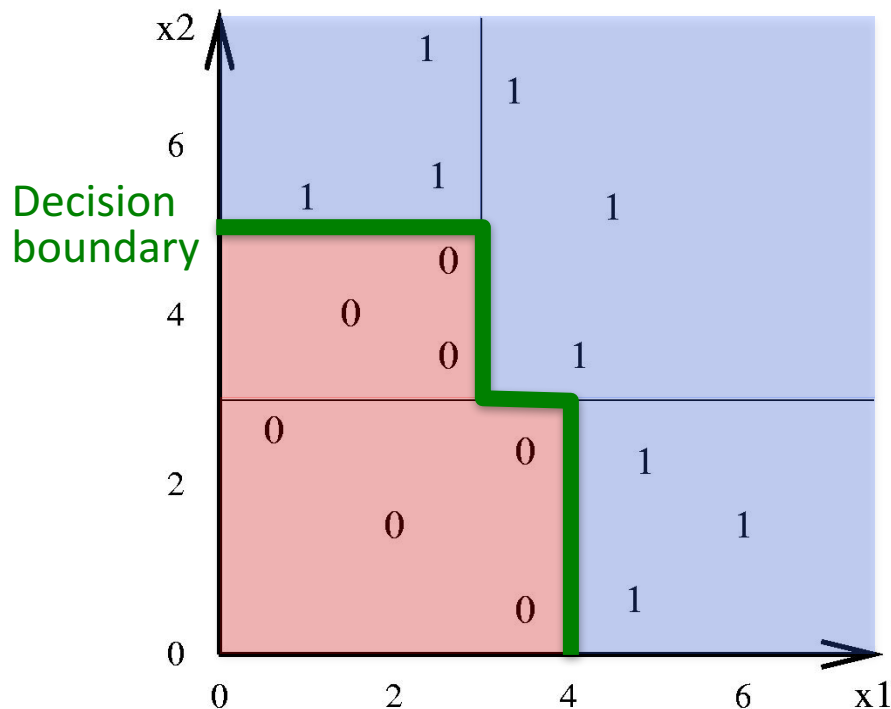- A possible decision tree for the data:



- What prediction would we make for

<outlook=sunny, temperature=hot, humidity=high, wind=weak> ?
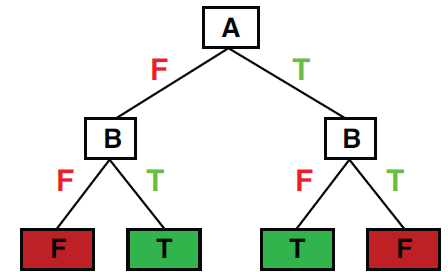
# Decision Tree – Decision Boundary

- Decision trees divide the feature space into axis-parallel (hyper-)rectangles

- Each rectangular region is labeled with one label
  - or a probability distribution over labels

# Expressiveness

- Given a particular space of functions, you may not be able to represent everything

- What **functions** can decision trees represent?

- Decision trees can represent any function of the input attributes!
  - Boolean operations (and, or, xor, etc.)?
  - **Yes!**
  - All boolean functions?
  - **Yes!**

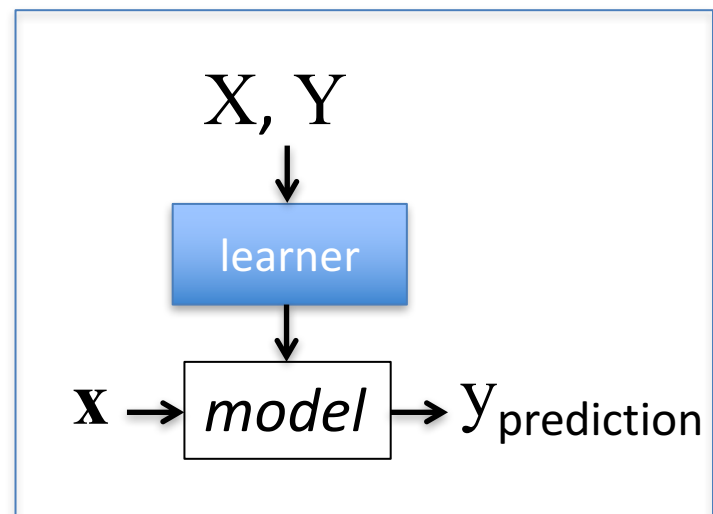| A | B | A xor B |
|---|---|---------|
| F | F | F |
| F | T | T |
| T | F | T |
| T | T | F |

(Figure from Stuart Russell)

# Stages of (Batch) Machine Learning

**Given:** labeled training data $X, Y = \{\langle \boldsymbol{x}_i, y_i \rangle\}_{i=1}^{n}$

- Assumes each $\boldsymbol{x}_i \sim \mathcal{D}(\mathcal{X})$ with $y_i = f_{target}(\boldsymbol{x}_i)$

**Train the model:**

  *model* ← *classifier*.train($X, Y$ )



**Apply the model to new data:**

- Given: new unlabeled instance $\boldsymbol{x} \sim \mathcal{D}(\mathcal{X})$

  $y_{\text{prediction}}$ ← *model*.predict($\mathbf{x}$)

# Basic Algorithm for Top-Down Learning of Decision Trees
[ID3, C4.5 by Quinlan]

*node* = root of decision tree

Main loop:

1.  *A* ← the "best" decision attribute for the next node.
2.  Assign *A* as decision attribute for *node*.
3.  For each value of *A*, create a new descendant of *node*.
4.  Sort training examples to leaf nodes.
5.  If training examples are perfectly classified, stop.  Else, recurse over new leaf nodes.

How do we choose which attribute is best?
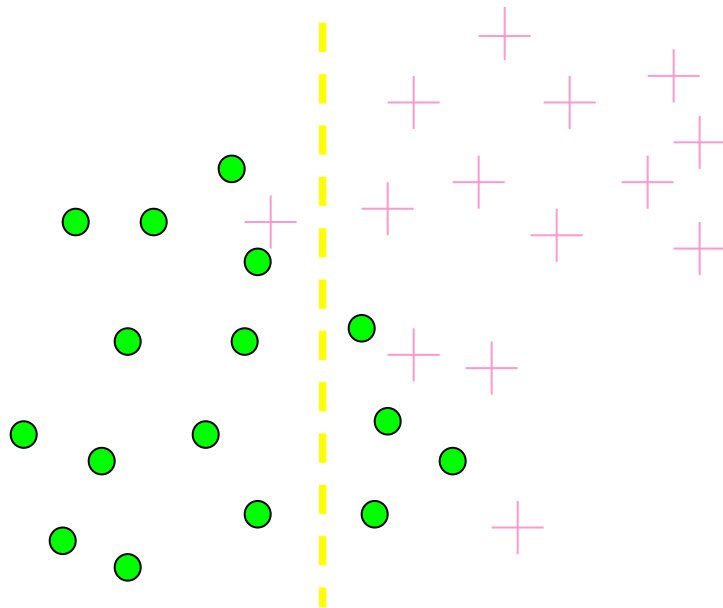
# Choosing the Best Attribute

**Key problem**: choosing which attribute to split a given set of examples

- Some possibilities are:
  - **Random:** Select any attribute at random
  - **Least-Values:** Choose the attribute with the smallest number of possible values
  - **Most-Values:** Choose the attribute with the largest number of possible values
  - **Max-Gain:** Choose the attribute that has the largest expected *information gain*
    - i.e., attribute that results in smallest expected size of subtrees rooted at its children

- The ID3 algorithm uses the Max-Gain method of selecting the best attribute
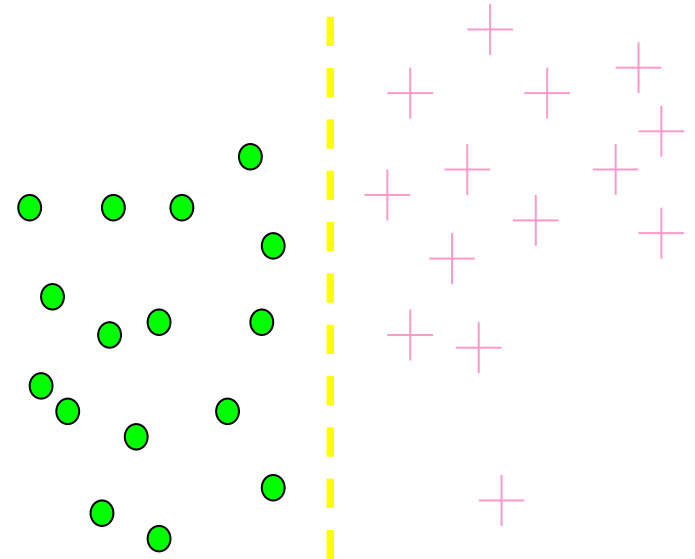
# Information Gain

Which test is more informative?

**Split over whether Balance exceeds 50K**

**Split over whether applicant is employed**
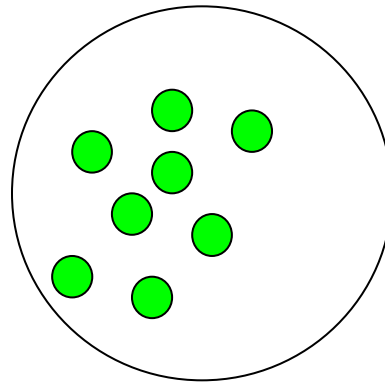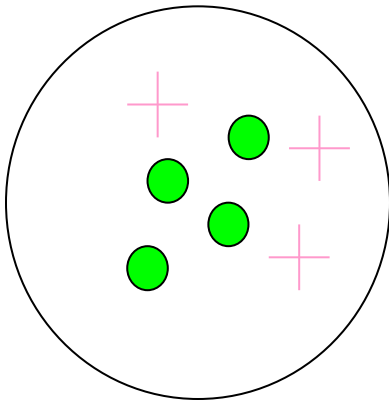
Less or equal 50K    Over 50K
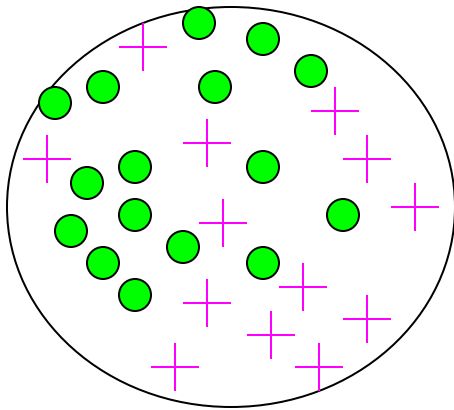
Unemployed    Employed

# Information Gain

**Impurity/Entropy** (informal)

– Measures the level of **impurity** in a group of examples
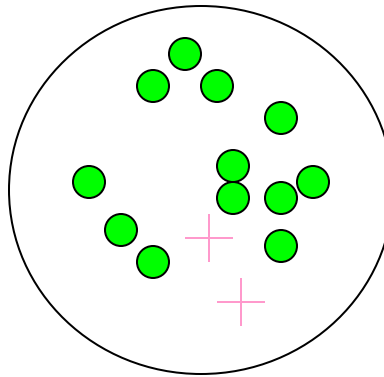
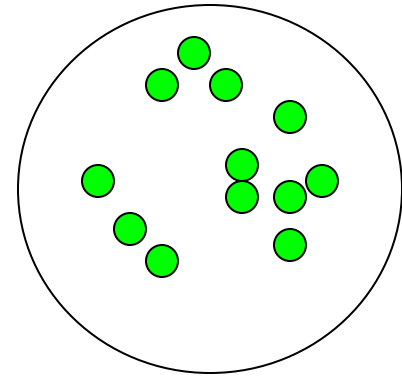# Impurity

**Very impure group**
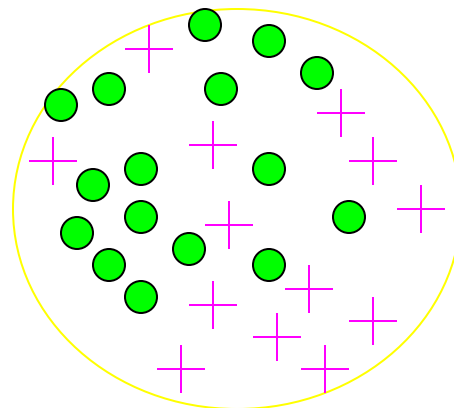
**Less impure**

**Minimum impurity**

# Entropy: a common way to measure impurity

- Entropy = $\sum_i - p_i \log_2 p_i$

  $p_i$ is the probability of class i

  Compute it as the proportion of class i in the set.

- Entropy comes from information theory.  The higher the entropy the more the information content.

  What does that mean for learning from examples?

# 2-Class Cases:

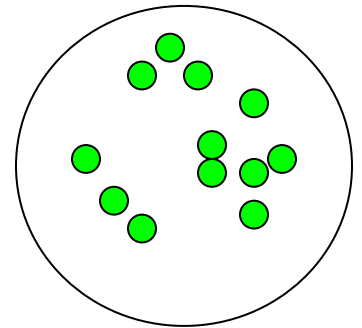$$H(x) = -\sum_{i=1}^{n} P(x = i) \log_2 P(x = i)$$

Entropy

- What is the entropy of a group in which all examples belong to the same class?
  - entropy = - 1 $\log_2$1 = 0

  not a good training set for learning

**Minimum impurity**



- What is the entropy of a group with 50% in either class?
  - entropy = -0.5 $\log_2$0.5 – 0.5 $\log_2$0.5 =1

  good training set for learning

**Maximum impurity**

# 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 $S$

$$H(S) \equiv -p_\oplus \log_2 p_\oplus - p_\ominus \log_2 p_\ominus$$

# Information Gain

- We want to determine which attribute in a given set of training feature vectors is most useful for discriminating between the classes to be learned.

- Information gain tells us how important a given attribute of the feature vectors is.

- We will use it to decide the ordering of attributes in the nodes of a decision tree.

# From Entropy to Information Gain

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

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

Specific conditional entropy $H(X/Y=v)$ of $X$ given $Y=v$ :

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

Conditional entropy $H(X/Y)$ of $X$ given $Y$ :

$$H(X|Y) = \sum_{v \in values(Y)} P(Y=v) H(X|Y=v)$$

Mututal information (aka Information Gain) of $X$ and $Y$ :

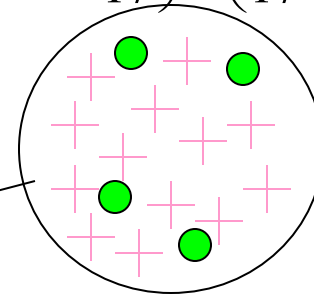$$I(X,Y) = H(X) - H(X|Y) = H(Y) - H(Y|X)$$

# Information Gain

Information Gain is the expected reduction in entropy of target variable Y for data sample S, due to sorting

# Calculating Information Gain

**Information Gain** =    entropy(parent) – [average entropy(children)]

child entropy
$$-\left(\frac{13}{17}\cdot\log_2\frac{13}{17}\right)-\left(\frac{4}{17}\cdot\log_2\frac{4}{17}\right)=0.787$$

Entire population (30 instances)



17 instances

child entropy
$$-\left(\frac{1}{13}\cdot\log_2\frac{1}{13}\right)-\left(\frac{12}{13}\cdot\log_2\frac{12}{13}\right)=0.391$$

parent entropy
$$-\left(\frac{14}{30}\cdot\log_2\frac{14}{30}\right)-\left(\frac{16}{30}\cdot\log_2\frac{16}{30}\right)=0.996$$

13 instances

(Weighted) Average Entropy of Children = $\left(\frac{17}{30}\cdot 0.787\right)+\left(\frac{13}{30}\cdot 0.391\right)=0.615$

**Information Gain= 0.996 - 0.615 = 0.38**

21

Based on slide by Pedro Domingos

# 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?**

S: [9+,5-]
E =0.940

Humidity

High                Normal

[3+,4-]                    [6+,1-]
E =0.985                  E =0.592

Gain (S, Humidity )
  = .940 - (7/14).985 - (7/14).592
  = .151

S: [9+,5-]
E=0.940

Wind

Weak                Strong

[6+,2-]                    [3+,3-]
E=0.811                  E=1.00

Gain (S, Wind)
  = .940 - (8/14).811 - (6/14)1.0
  = .048

{D1, D2, ..., D14}

[9+,5−]

Outlook

Sunny    Overcast    Rain

{D1,D2,D8,D9,D11}        {D3,D7,D12,D13}        {D4,D5,D6,D10,D14}

[2+,3−]        [4+,0−]        [3+,2−]

?        Yes        ?

*Which attribute should be tested here?*

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

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

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

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

Slide by Tom Mitchell

# Which Tree Should We Output?



- ID3 performs heuristic search through space of decision trees

- It stops at smallest acceptable tree. Why?

# Preference bias: Ockham's Razor

- Principle stated by William of Ockham (1285-1347)
  - "*non sunt multiplicanda entia praeter necessitatem*"
  - entities are not to be multiplied beyond necessity
  - AKA Occam's Razor, Law of Economy, or Law of Parsimony

**Idea**: The simplest consistent explanation is the best

- Therefore, the smallest decision tree that correctly classifies all of the training examples is best
  - Finding the provably smallest decision tree is NP-hard
  - …So instead of constructing the absolute smallest tree consistent with the training examples, construct one that is pretty small

# Overfitting in Decision Trees

- Many kinds of "noise" can occur in the examples:
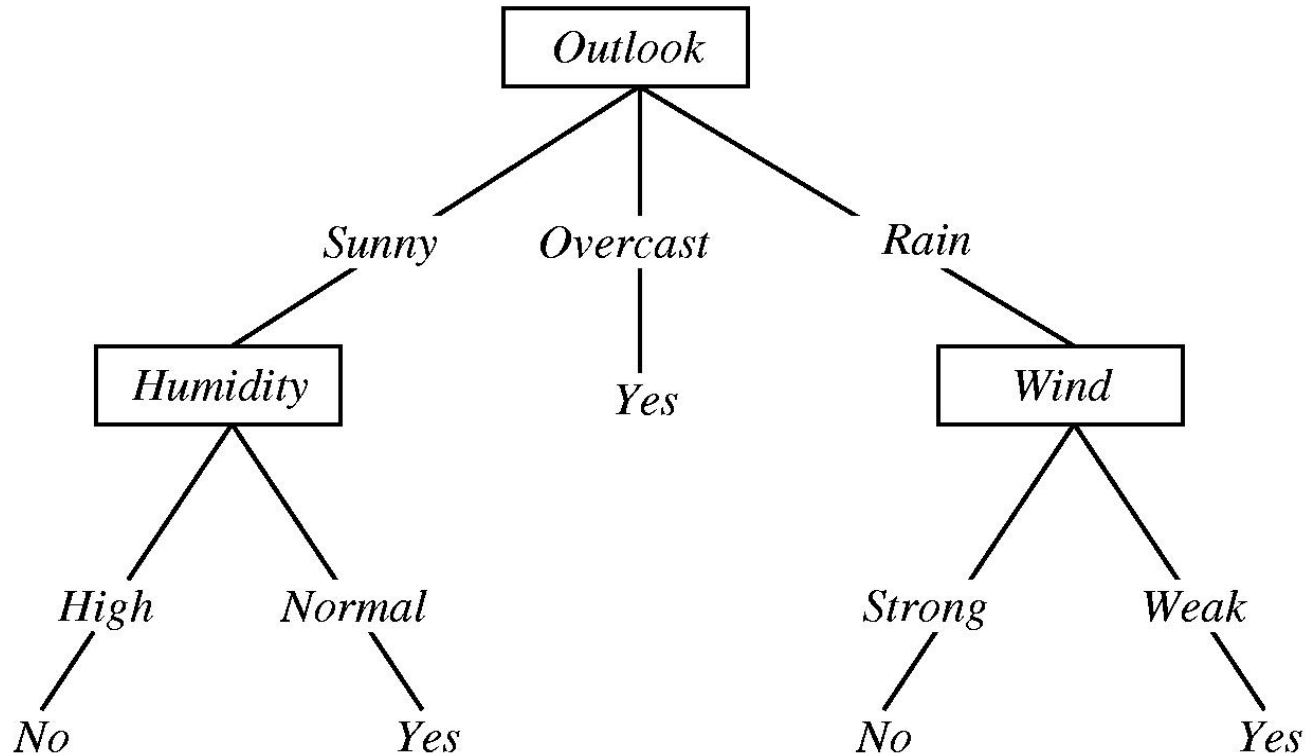  - Two examples have same attribute/value pairs, but different classifications
  - Some values of attributes are incorrect because of errors in the data acquisition process or the preprocessing phase
  - The instance was labeled incorrectly (+ instead of -)

- Also, some attributes are irrelevant to the decision-making process
  - e.g., color of a die is irrelevant to its outcome

# Overfitting in Decision Trees

- Irrelevant attributes can result in *overfitting* the training example data
  - If hypothesis space has many dimensions (large number of attributes), we may find **meaningless regularity** in the data that is irrelevant to the true, important, distinguishing features

- If we have too little training data, even a reasonable hypothesis space will 'overfit'

# Overfitting in Decision Trees

Consider adding a **noisy** training example to the following tree:



What would be the effect of adding:

<outlook=sunny, temperature=hot, humidity=normal, wind=strong, playTennis=No> ?

Based on Slide by Pedro Domingos

# Overfitting in Decision Trees

Consider error of hypothesis $h$ over

- training data: $error_{train}(h)$

- entire distribution $\mathcal{D}$ of data: $error_{\mathcal{D}}(h)$
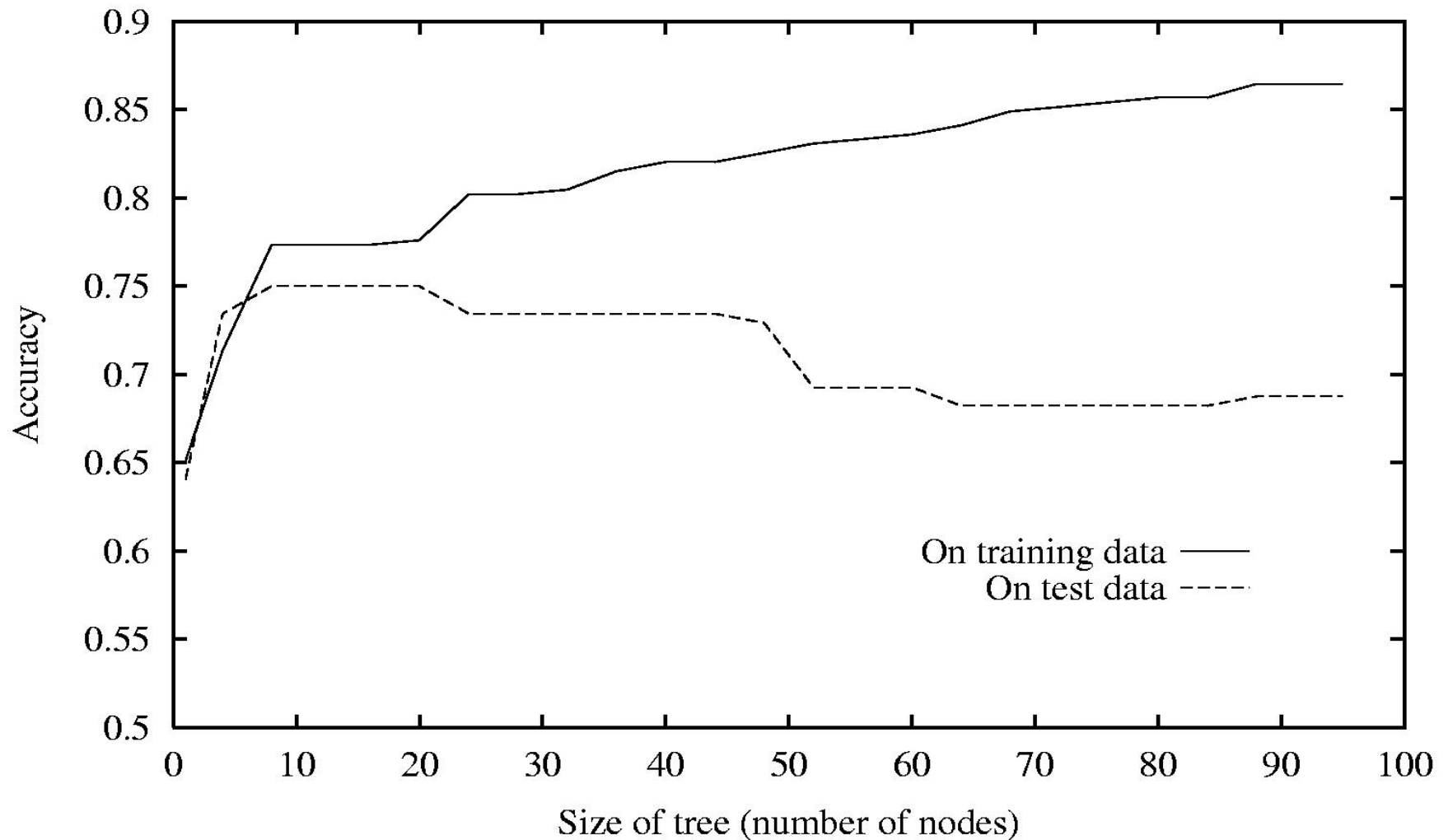
Hypothesis $h \in H$ **overfits** training data if there is an alternative hypothesis $h' \in H$ such that

$$error_{train}(h) < error_{train}(h')$$

and

$$error_{\mathcal{D}}(h) > error_{\mathcal{D}}(h')$$

# Overfitting in Decision Trees

Slide by Pedro Domingos

# Avoiding Overfitting in Decision Trees

How can we avoid overfitting?

- Stop growing when data split is not statistically significant
- Acquire more training data
- Remove irrelevant attributes  (manual process – not always possible)
- **Grow full tree, then post-prune**


How to select "best" tree:

- Measure performance over training data
- Measure performance over separate validation data set
- Add complexity penalty to performance measure (heuristic: simpler is better)

# Reduced-Error Pruning

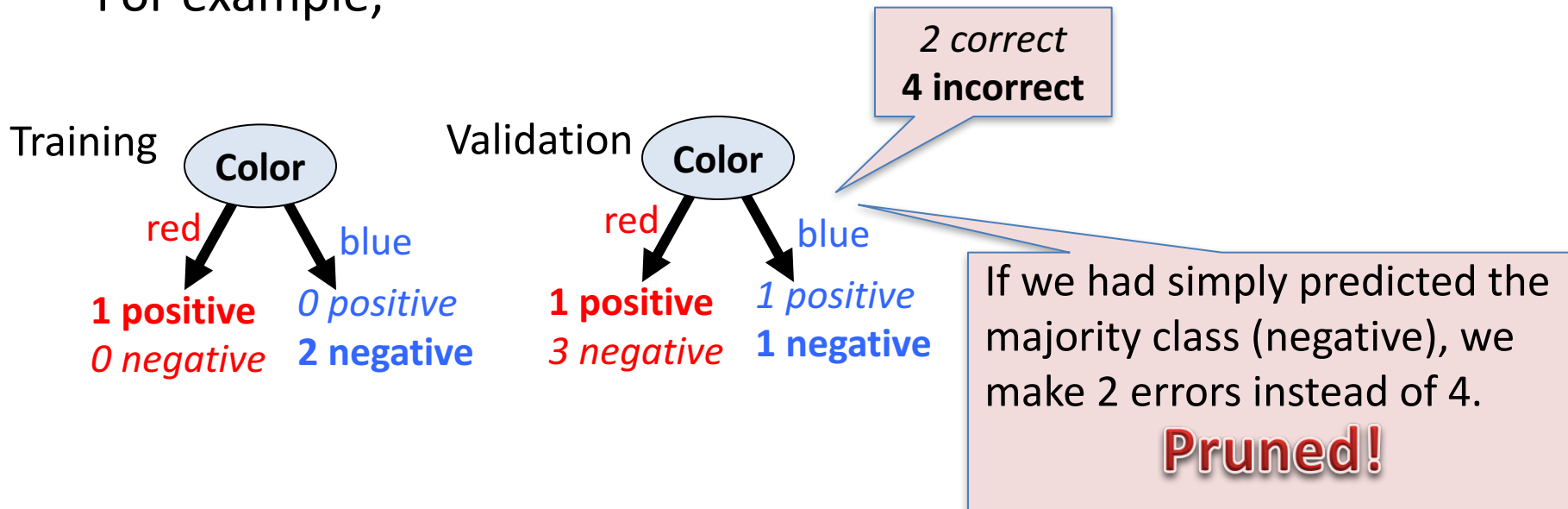Split training data further into *training* and *validation* sets

Grow tree based on *training set*
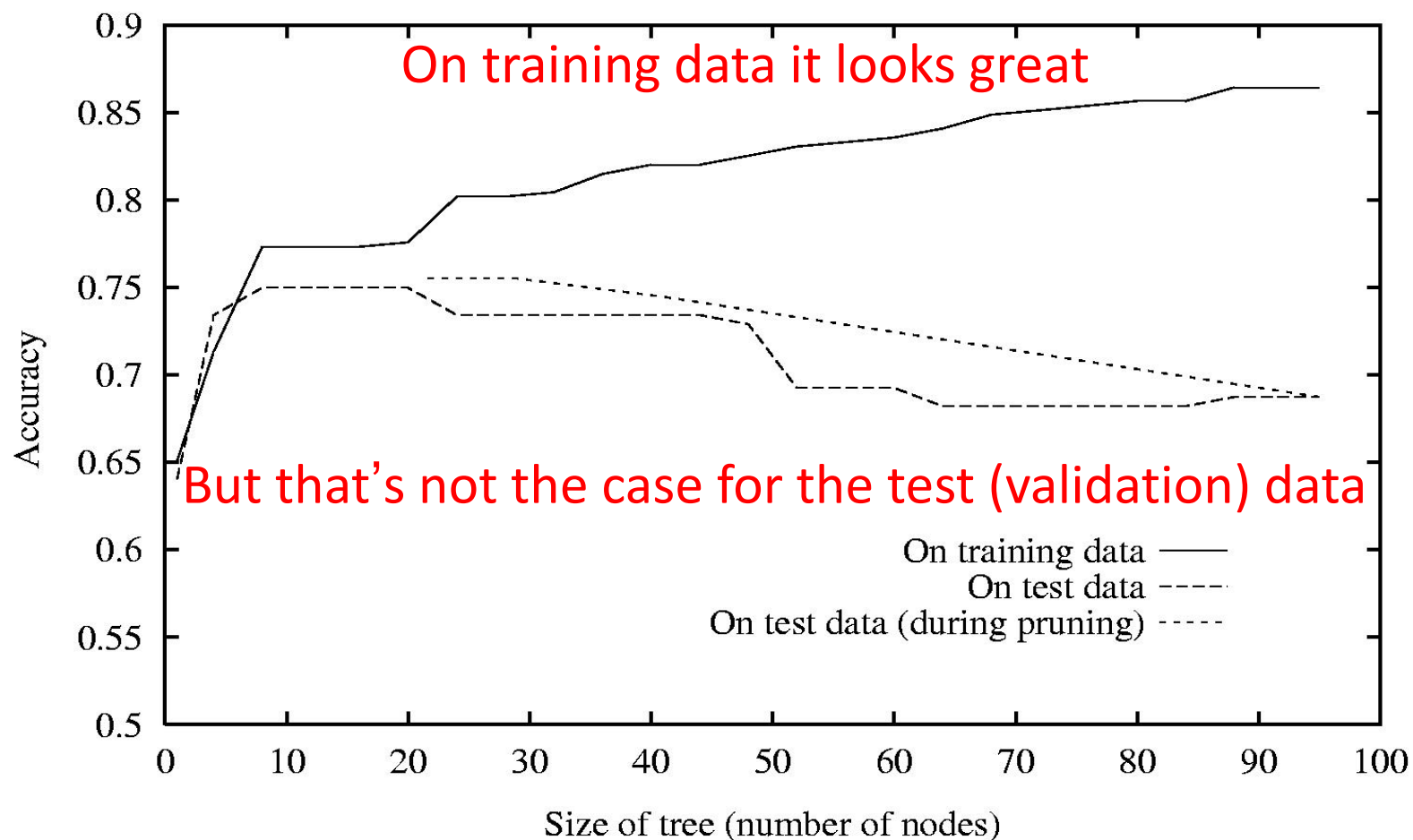
Do until further pruning is harmful:

1. Evaluate impact on validation set of pruning each possible node (plus those below it)

2. Greedily remove the node that most improves *validation set* accuracy

# Pruning Decision Trees

- Pruning of the decision tree is done by replacing a whole subtree by a leaf node.

- The replacement takes place if a decision rule establishes that the expected error rate in the subtree is greater than in the single leaf.

- For example,



*2 correct*
**4 incorrect**

Training   Color

red        blue

**1 positive**   *0 positive*
*0 negative*     **2 negative**

Validation   Color

red        blue

**1 positive**   *1 positive*
*3 negative*     **1 negative**

If we had simply predicted the majority class (negative), we make 2 errors instead of 4.

**Pruned!**

# Effect of Reduced-Error Pruning



On training data it looks great

But that's not the case for the test (validation) data

On training data ——
On test data - - - -
On test data (during pruning) - - - -

# Effect of Reduced-Error Pruning



The tree is pruned back to the red line where it gives more accurate results on the test data

# Summary: Decision Tree Learning

- Widely used in practice

- Strengths include
  - Fast and simple to implement
  - Can convert to rules
  - Handles noisy data

- Weaknesses include
  - Univariate splits/partitioning using only one attribute at a time --- limits types of possible trees
  - Large decision trees may be hard to understand
  - Requires fixed-length feature vectors
  - Non-incremental (i.e., batch method)

# Summary: Decision Tree Learning

- Representation:       decision trees
- Bias:                       prefer small decision trees
- Search algorithm:    greedy
- Heuristic function:  information gain or information
                                content or others
- Overfitting / pruning