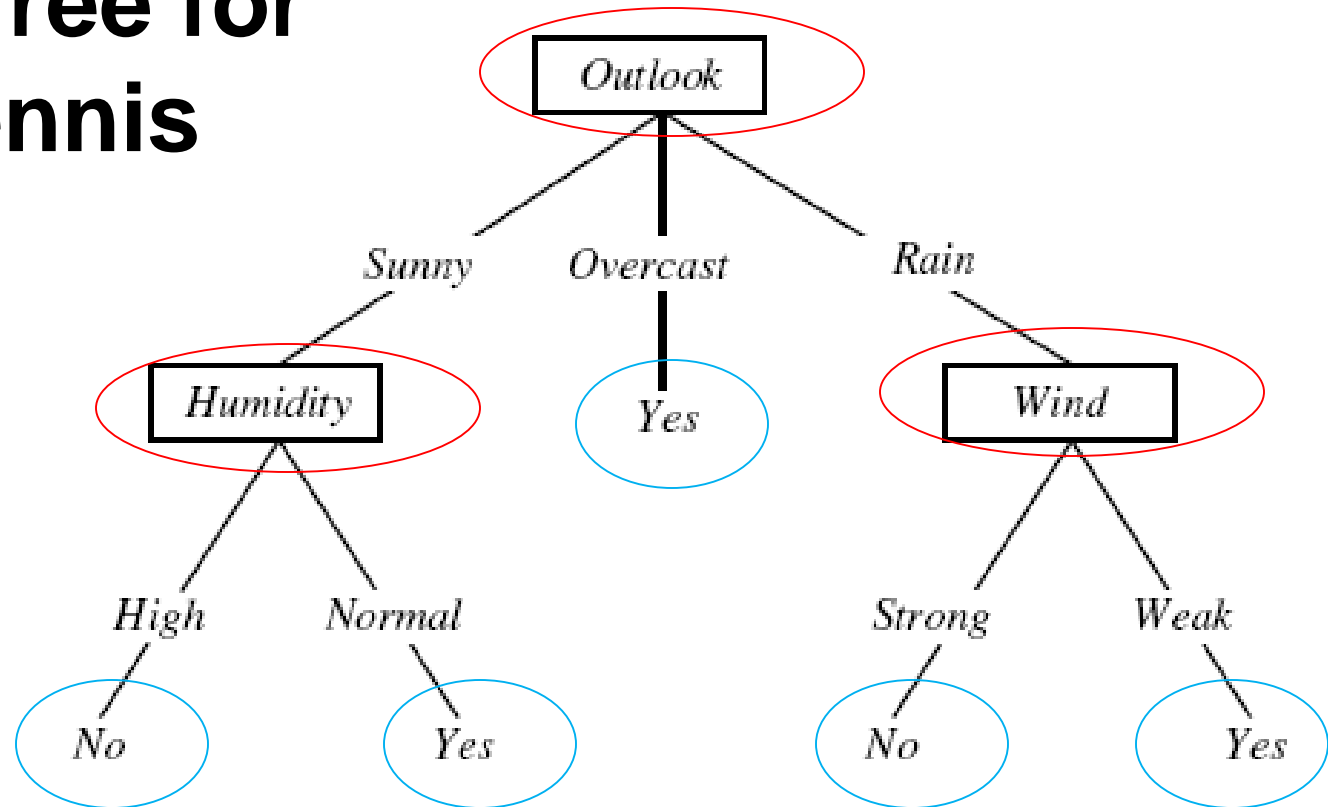


# Decision Tree

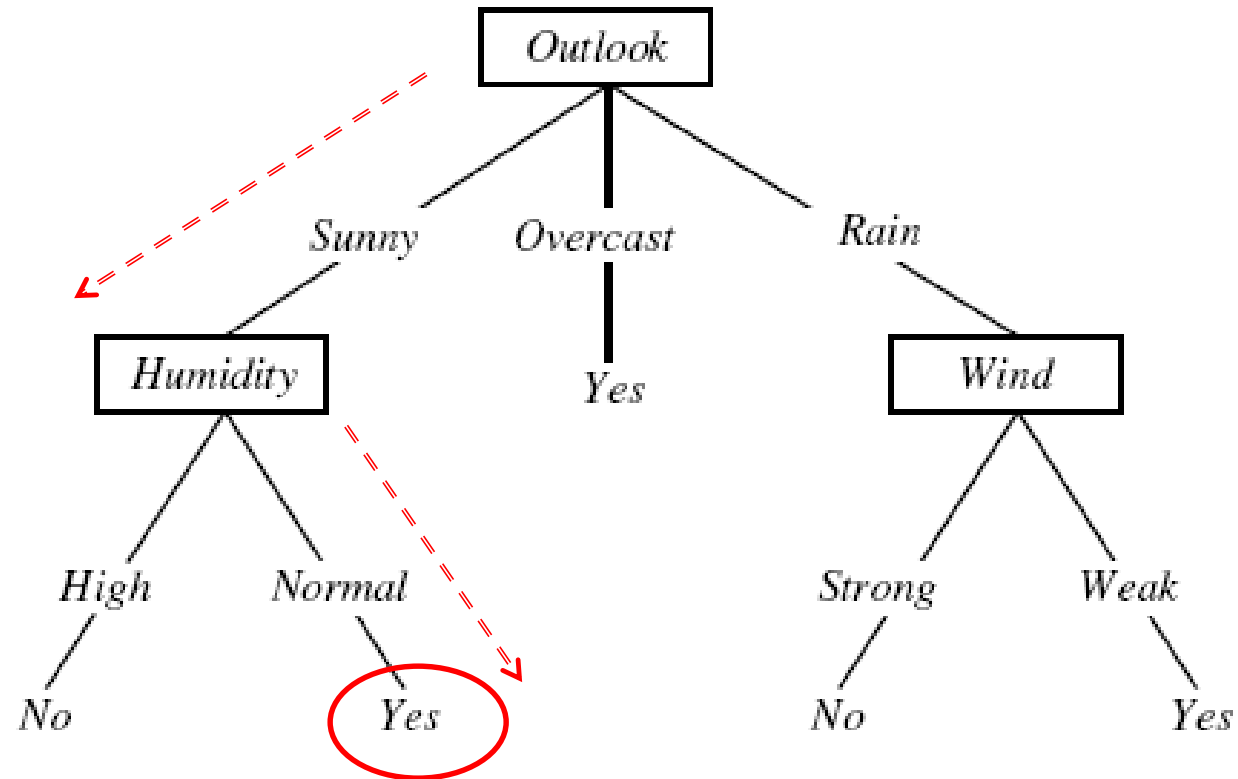


# Decision Tree for Playing Tennis



- Each **internal node** test on an attribute  $x_i$
- Each branch from a node takes a particular value of  $x_i$
- Each leaf node predicts a class label

(outlook=sunny, wind=strong, humidity=normal, ? )



# DT for prediction C-section risks

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-

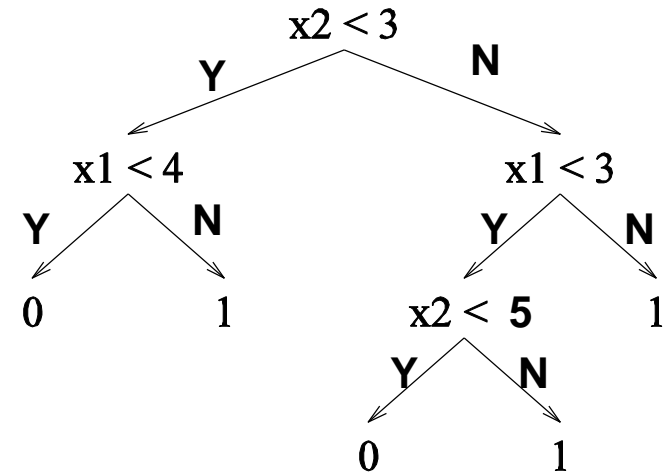
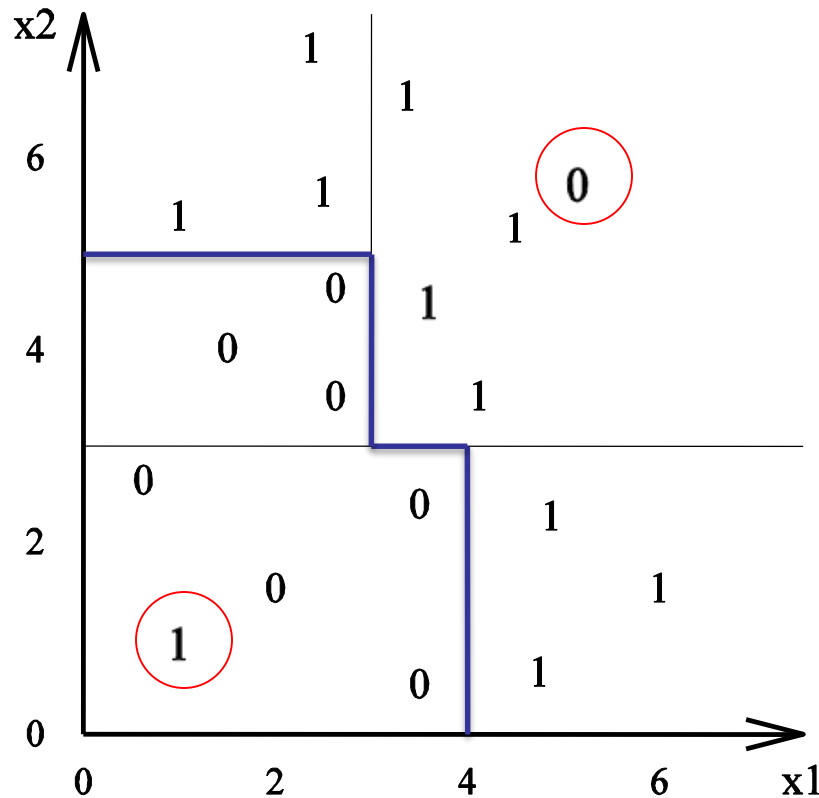
Fetal\_Presentation = 2: [3+,29-] .11+ .89-

Fetal\_Presentation = 3: [8+,22-] .27+ .73-

# Characteristics of Decision Trees

- Decision trees have many appealing properties
  - Similar to human decision process, easy to understand
  - Deal with both **discrete (categorical)** and **continuous** features without the need to transform them
    - For classifiers such as perceptron and logistic regression, to use a categorical feature like “color” (e.g., with value choices: red, white, blue), one will have to create multiple binary features, one for each possible color value.
  - Highly flexible: as the # of nodes (or depth) of the tree increase, decision tree can represent increasingly complex decision boundaries

# DT can represent arbitrarily complex decision boundaries



If needed, the tree can keep on growing until all examples are correctly classified! Although it may not be the best idea

# How to learn decision trees?

- Possible goal: find a decision tree  $h$  that achieves minimum error on training data
  - Trivially achievable – if use a large enough tree
- Another possibility: find the smallest decision tree that achieves the minimum training error
  - NP-hard
- The most commonly used algorithm for decision tree builds a tree incrementally
  - Top down
  - In a greedy fashion

# Greedy Learning For DT

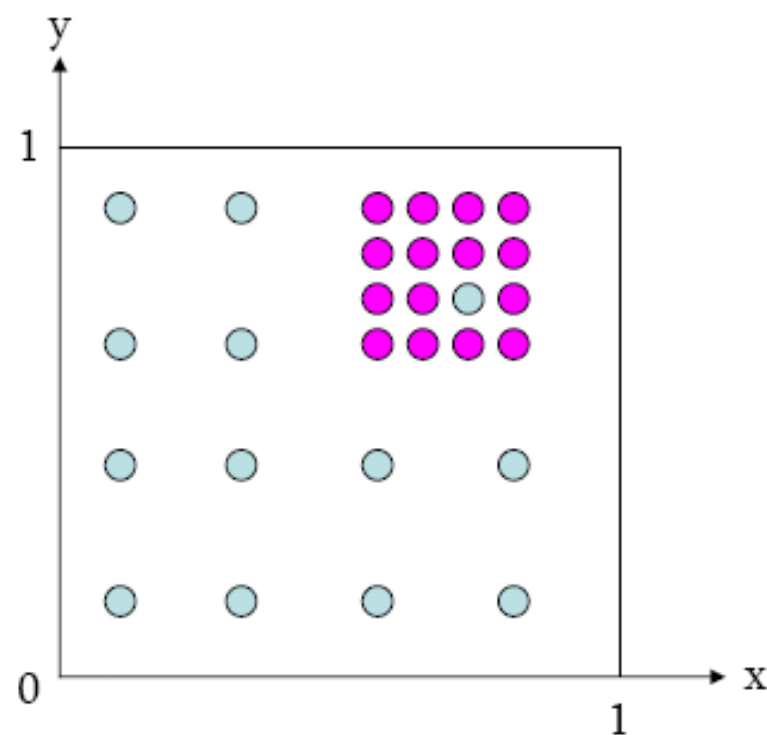
We will study a top-down, greedy search approach. Instead of trying to optimize the whole tree together, we try to build one node (test) at a time.

Basic idea: (assuming discrete features, relax later)

- Choose the best test to be the root of the tree.
- Create a descendant node for each test outcome
- Examples in training set  $S$  are sent to the appropriate descendent node based on the test outcome
- Recursively apply at each descendant node to select the best attribute to test using its subsest of training samples
- If all samples in a node belong to the same class, turn it into a leaf node of that class



# Building DT: an example

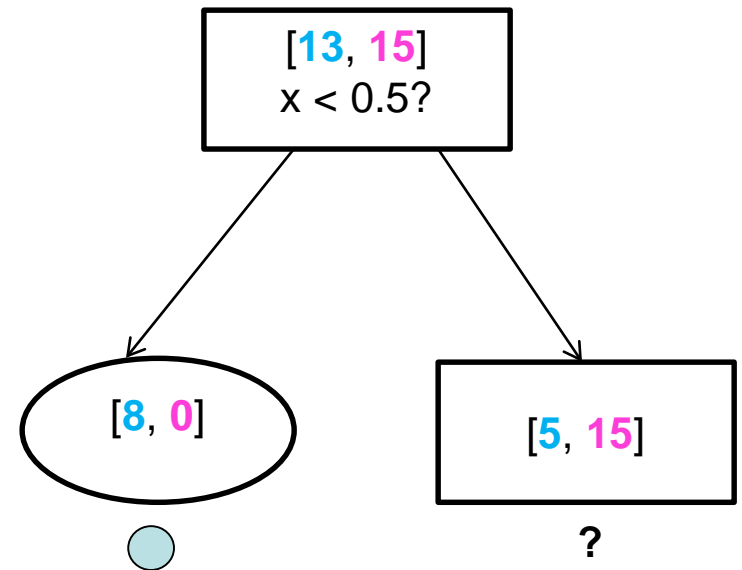
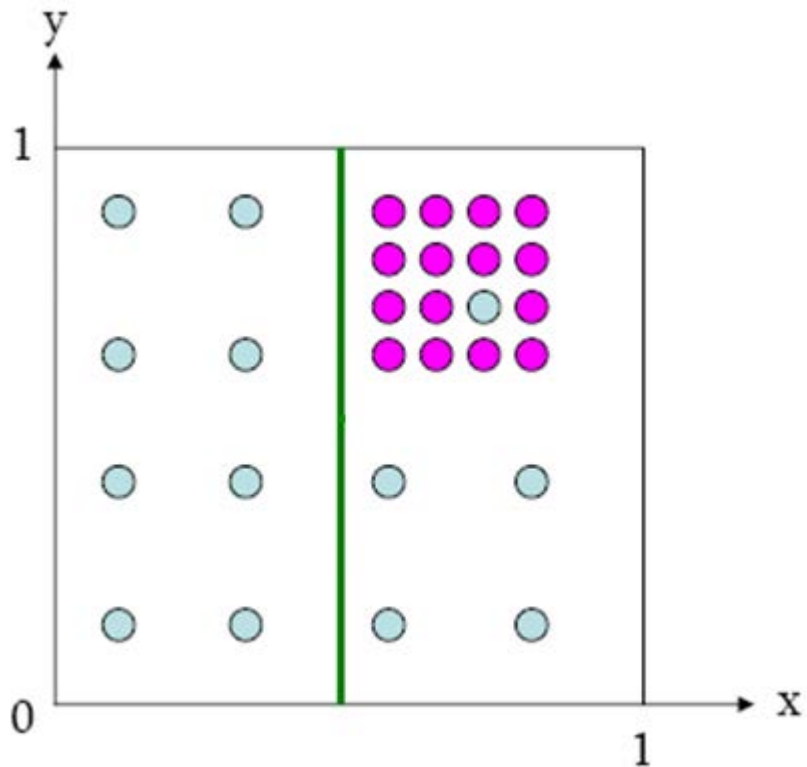


Training data contains

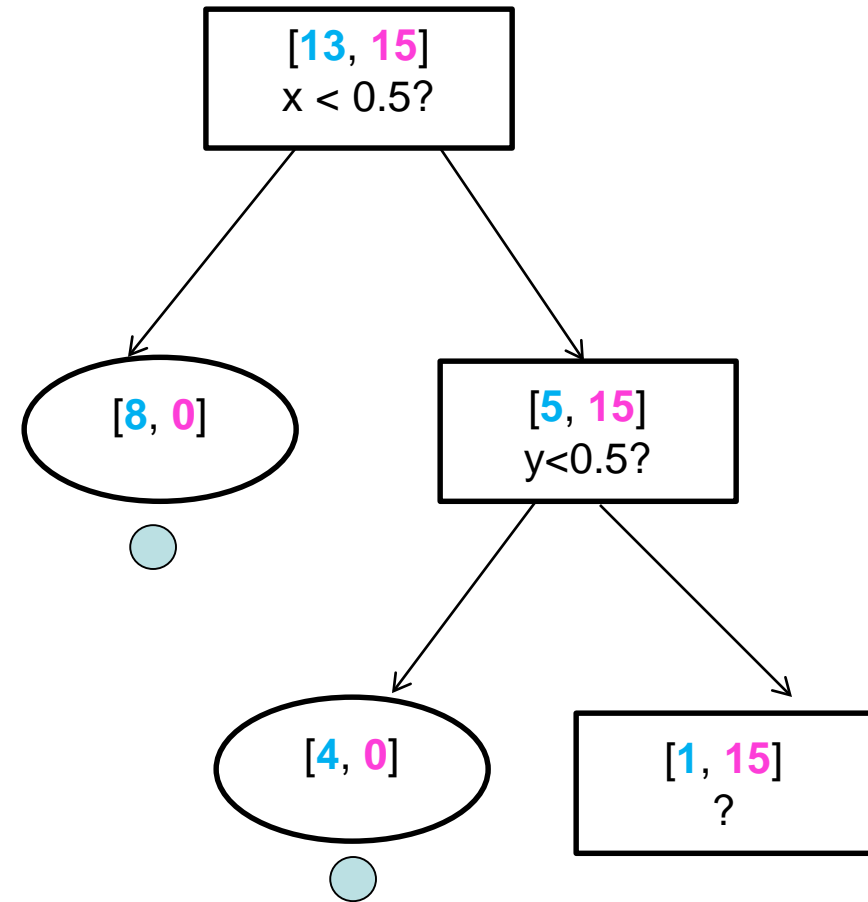
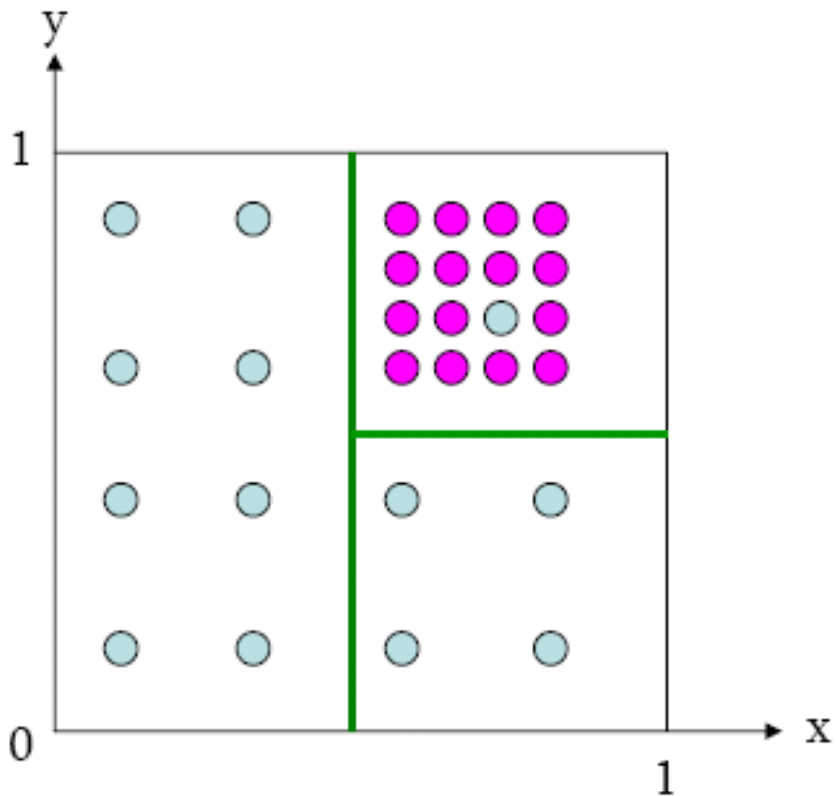
13  15 

13  15 

One possible question: is  $x < 0.5$ ?

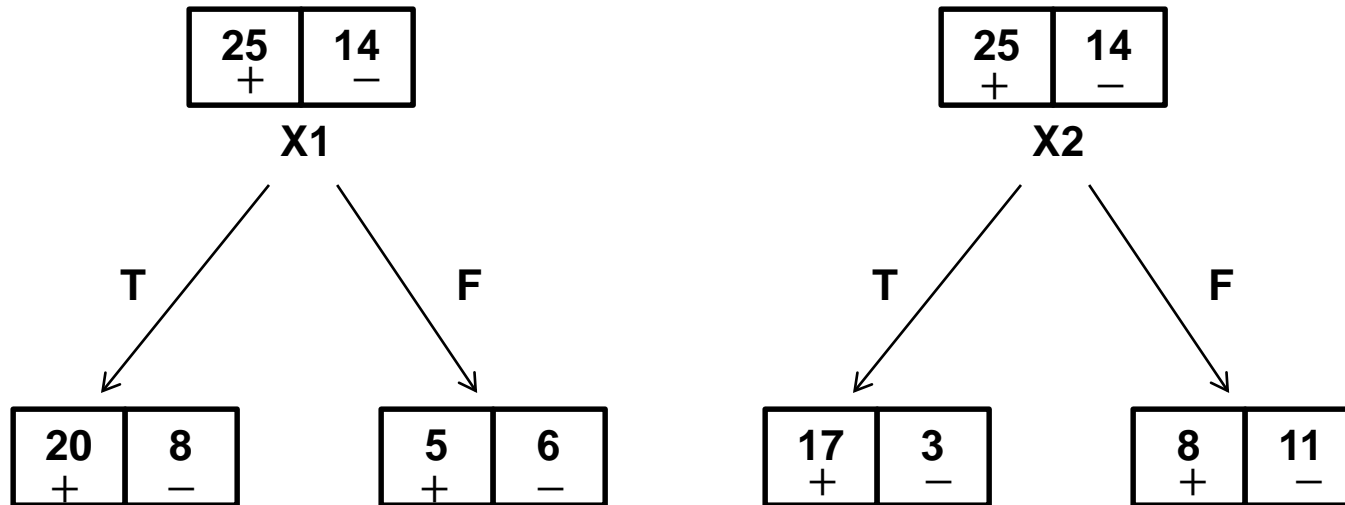


# Continue



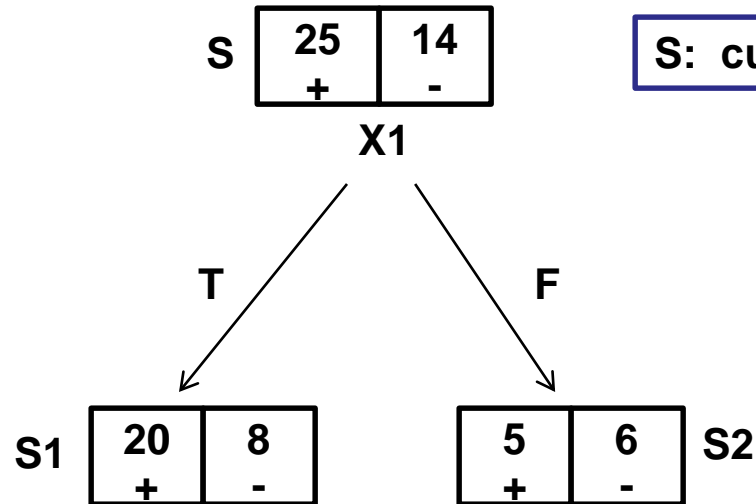
This could keep on going, until all examples are correctly classified.

# How to choose the best test



**Which one is better?**

# Choosing the Best test: A General View



$S$ : current set of training examples

$m$  branches, one for each possible outcome of the test, in this example, only two branches  $T$ , or  $F$

$S_1, S_2, \dots, S_m$ :  $m$  subsets of training examples

$$\text{Benefit of split} = U(S) - \underbrace{\sum_i^m p_i U(S_i)}$$

$U(S)$ : Label Uncertainty in  $S$   
 $p_i$ : The portion of examples in  $S$  that takes branch  $i$   
 $\sum_i^m p_i U(S_i)$ : Total Expected Remaining Label Uncertainty after the test

# Uncertainty Measure: Entropy

- Given a set of training examples  $S$ 
  - Let  $y$  denote the label of a randomly drawn example from  $S$
  - If all examples in  $S$  belong to the same class, then  $y$  will always be the same value – no uncertainty
  - If  $S$  is evenly split between positive and negative, then  $y$  will have a 50-50 chance to be positive or negative – very high uncertainty

**Entropy:** a measure of uncertainty rooted in information theory

# Entropy Definition

Let  $y$  be a categorical random variable that can take  $k$  different values:  $v_1, v_2, \dots, v_k$  and  $p_i = P(y = v_i)$  for  $i = 1, \dots, k$ .

The **entropy** of  $y$ , denoted as  $H(y)$  is defined as:

$$H(y) = - \sum_{i=1}^k p_i \log_2 p_i$$

# A Side Note

- Recall the logistic regression objective:

$$L(\mathbf{w}) = - \sum_{i=1}^n \left( y_i \log \sigma(\mathbf{w}^T \mathbf{x}_i) + (1 - y_i) \log (1 - \sigma(\mathbf{w}^T \mathbf{x}_i)) \right)$$

can be viewed as:

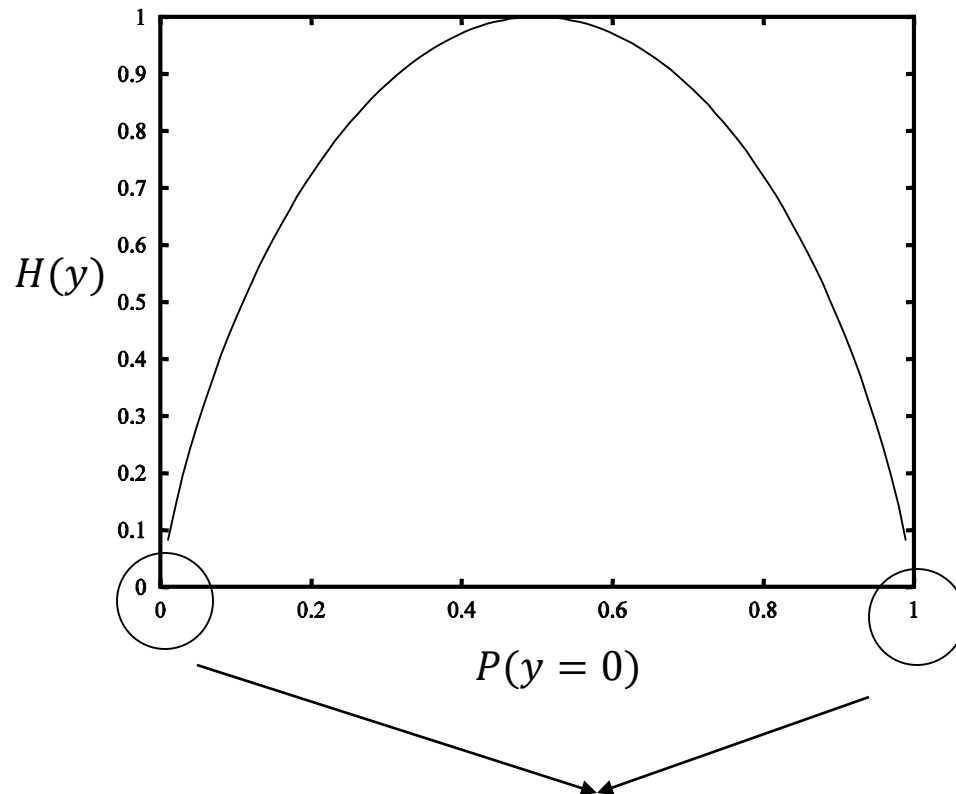
$$\sum_{i=1}^n - (P(y_i = 1) \log P(\hat{y}_i = 1) + P(y_i = 0) \log P(\hat{y}_i = 0))$$

Hence the name: cross-entropy



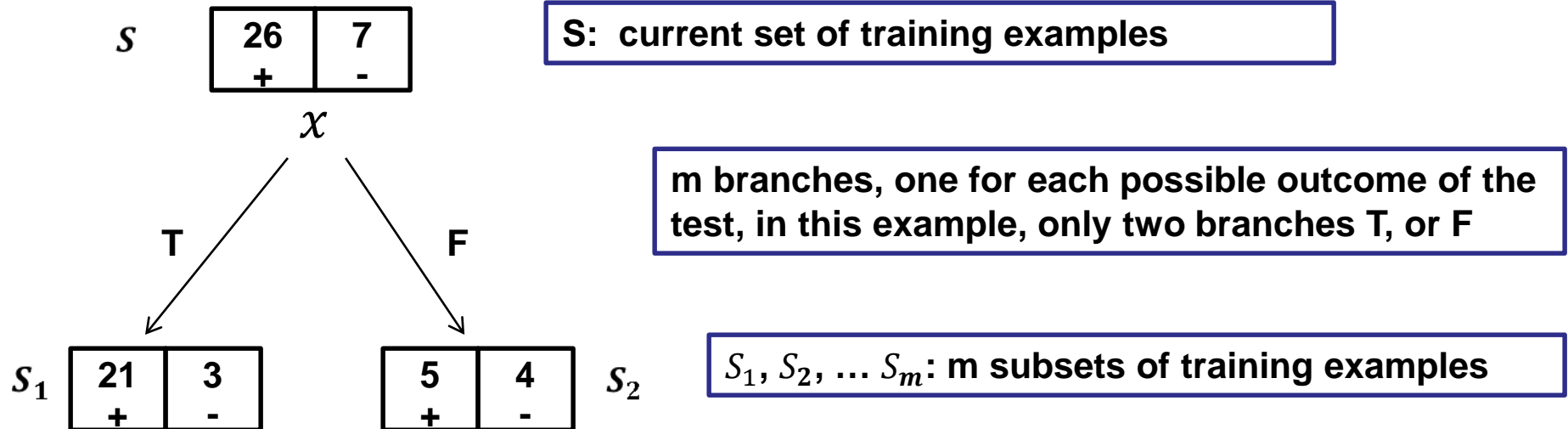
# Entropy of a Binary $y$

- Entropy is a concave function downward



Minimum uncertainty occurs when  $p_0 = 0$  or 1

# Back to building decision tree

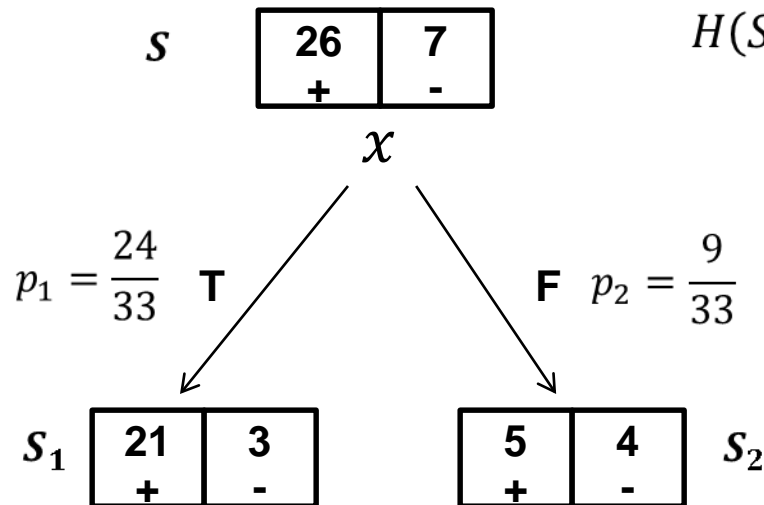


Choose the split that maximizes:

$$\text{Benefit of split} = \underline{U(S)} - \sum_i^m p_i \underline{U(S_i)}$$

Measure the label uncertainty using entropy, referred to as the **Information gain** criterion

# Measuring uncertainty using entropy



$$H(S) = -\frac{26}{33} \log_2 \frac{26}{33} - \frac{7}{33} \log_2 \frac{7}{33} = 0.7455$$

$$H(S_1) = -\frac{21}{24} \log_2 \frac{21}{24} - \frac{3}{24} \log_2 \frac{3}{24} = 0.5436$$

$$H(S_2) = -\frac{5}{9} \log_2 \frac{5}{9} - \frac{4}{9} \log_2 \frac{4}{9} = 0.9911$$

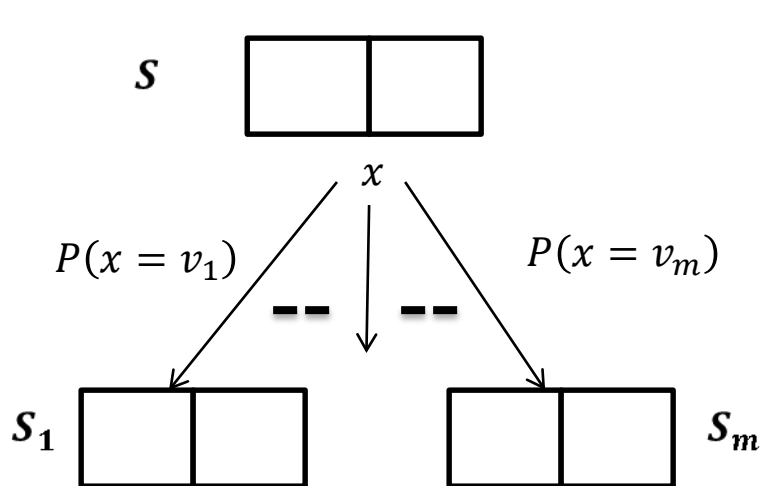
$$H(S) - (p_1 H(S_1) + p_2 H(S_2))$$

$$= 0.7455 - \frac{24}{33} * 0.5436 - \frac{9}{33} * 0.9911 = 0.0799$$

This effectively measures the mutual information between the feature  $x$  and the class label  $y$

Information gain is always non-negative!

# Choosing the Best Feature: Summary



$$\text{Benefit of split} = U(S) - \underbrace{\sum_i^m p_i U(S_i)}_{\text{Total Expected Remaining Uncertainty after the test}}$$

Original  
uncertainty

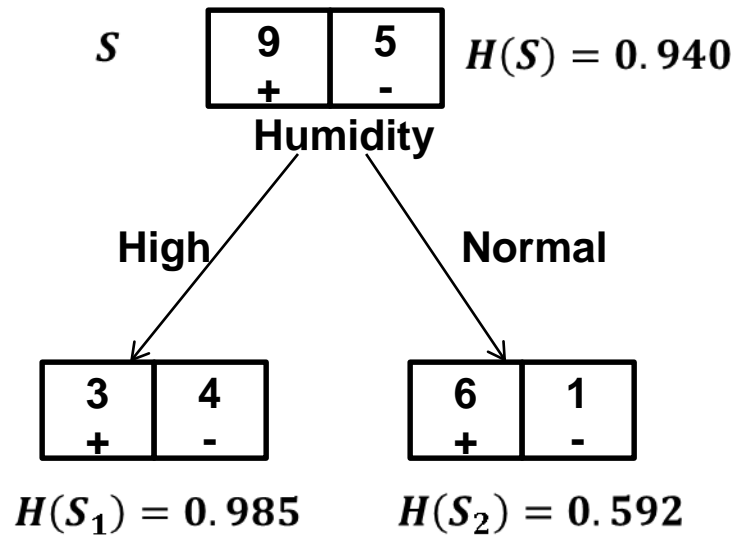
Total Expected  
Remaining Uncertainty  
after the test

Measures of Uncertainty	
Error	$\min(p_+, p_-)$
Entropy	$-p_+ \log_2 p_+ - p_- \log_2 p_-$
Gini Index	$p_+ p_-$

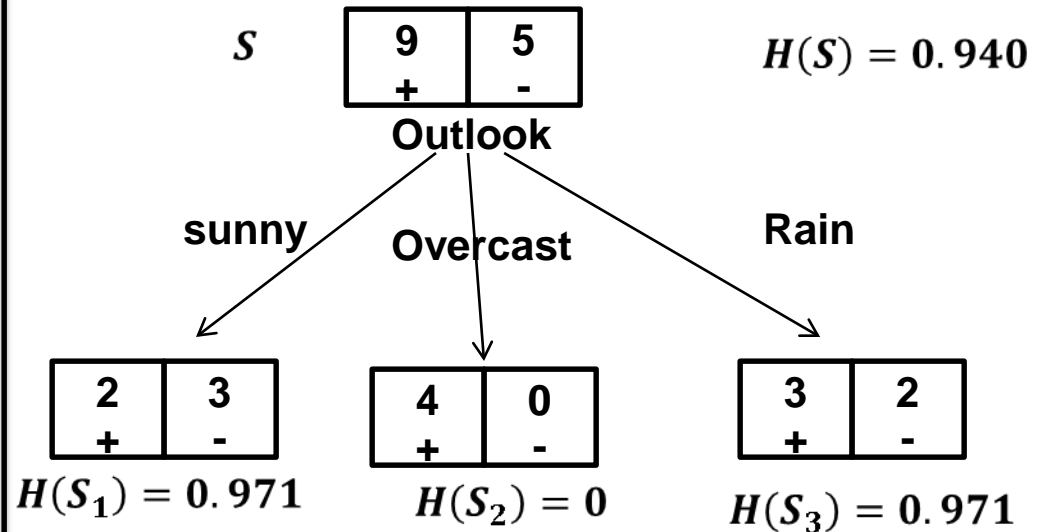
# Example

Day	Outlook	Temperature	Humidity	Wind	PlayTenn
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 root test using information gain

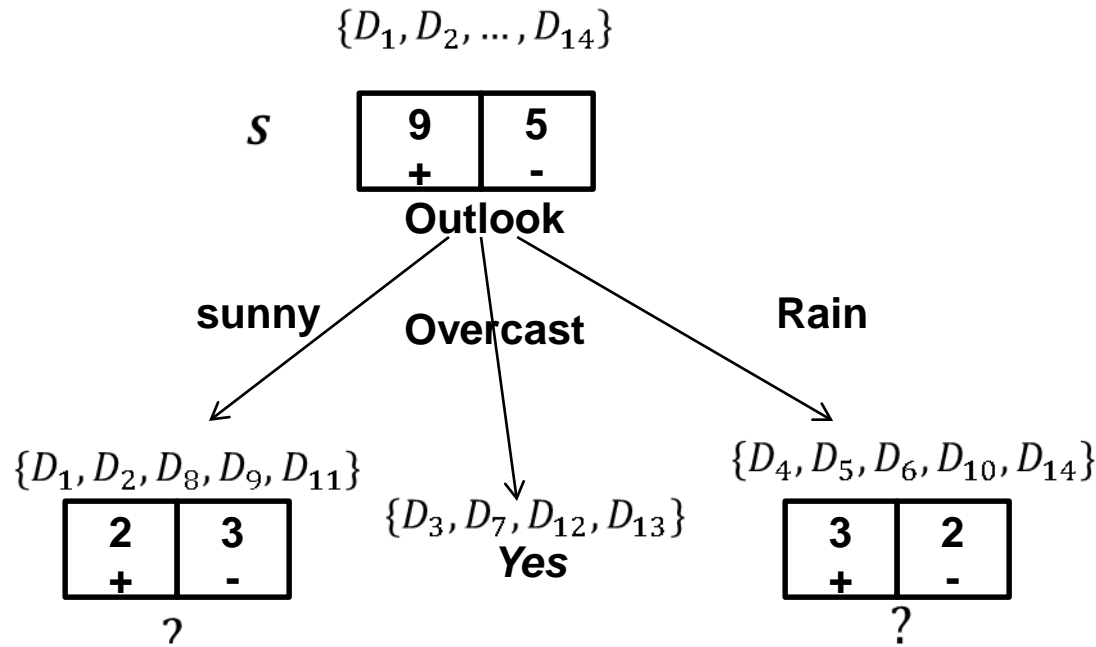


$$\text{Gain}(\text{humidity}) = 0.940 - \frac{1}{2} 0.985 - \frac{1}{2} 0.592 = 0.151$$

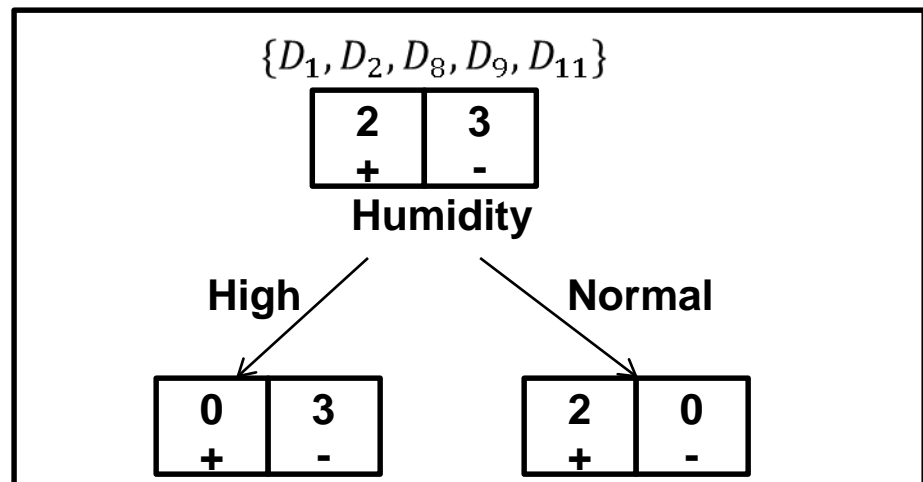


$$\text{Gain}(\text{Outlook}) = 0.940 - \frac{5}{14} 0.971 - \frac{5}{14} 0.971 = 0.2464$$

# Continue building the tree



Which test should be placed here?



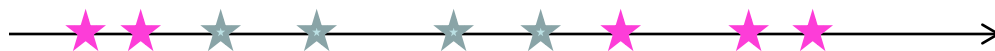
# Issues with Multi-nomial Features

- Multi-nomial features: categorical features with more than 2 possible values
- Consider two features, one is binary, the other has 100 possible values, which one you expect to have higher information gain?
- Conditional entropy of  $y$  given the 100-valued feature will be low – why?
- This bias will prefer multinomial features to binary features
  - Test for one value, e.g., Outlook = sunny?



# Dealing with Continuous Features

- Test against a threshold
- How to compute the best threshold  $\theta_j$  for  $x_j$ ?
  - Sort the examples according to  $x_j$ .
  - Move the threshold  $\theta$  from the smallest to the largest value
  - Select  $\theta$  that gives the best information gain
  - Trick: only need to compute information gain when class label changes



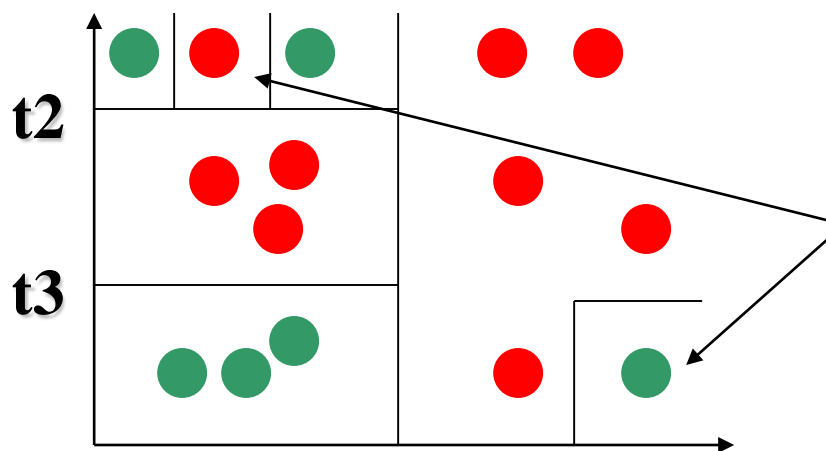
- Note that continuous features can be tested for multiple times on the same path in a DT

# Considering both discrete and continuous features

- If a data set contains both types of features, do we need special handling?
- No, we simply consider all possible splits in every step of the decision tree building process, and choose the one that gives the highest information gain
  - This include all possible (meaningful) thresholds

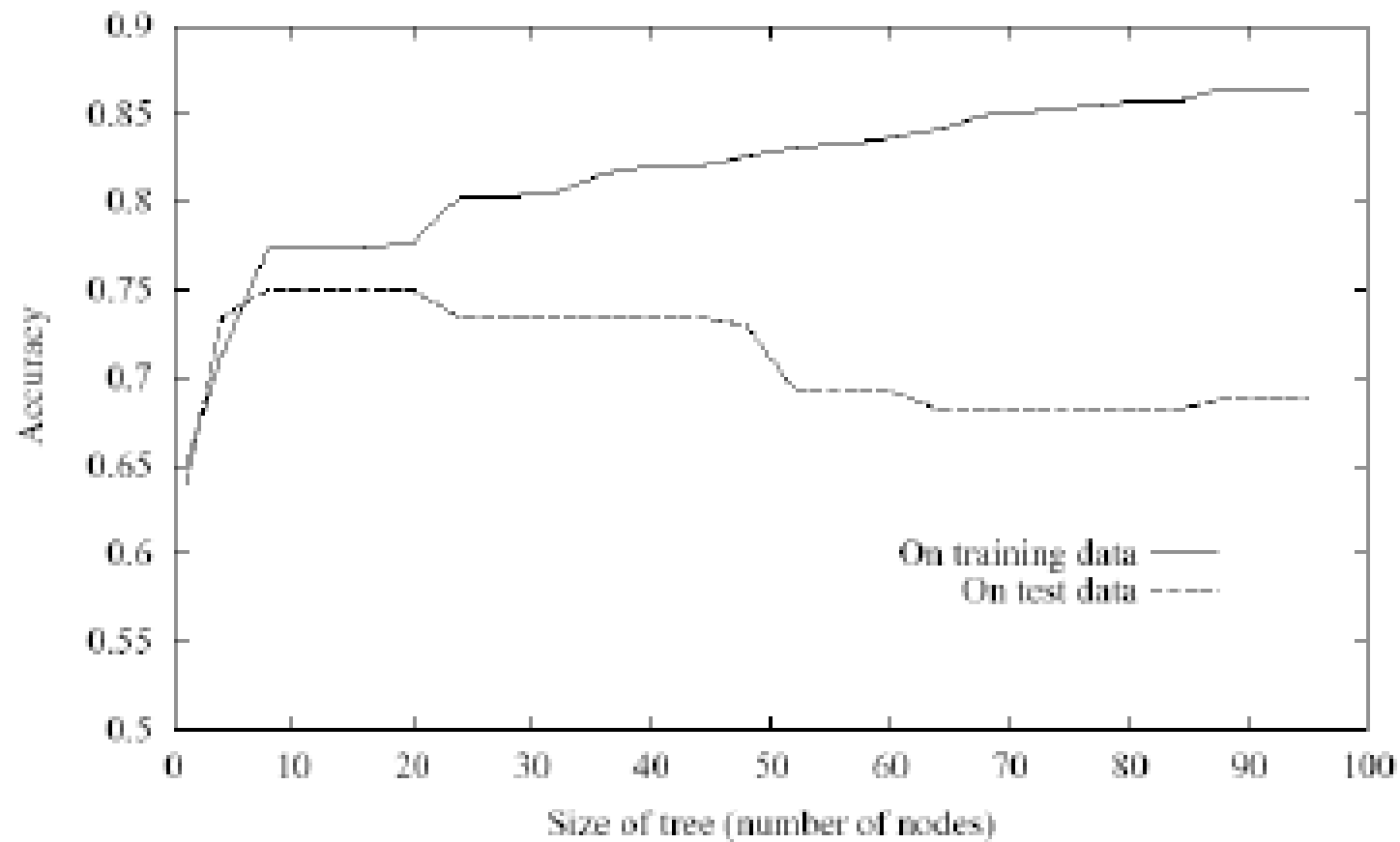
# Issue of Over-fitting

- Decision tree is highly flexible
- As the nodes increase, we can represent arbitrarily complex decision boundaries
- This can lead to over-fitting



**Possibly just noise, but  
the tree is grown larger  
to capture these examples**

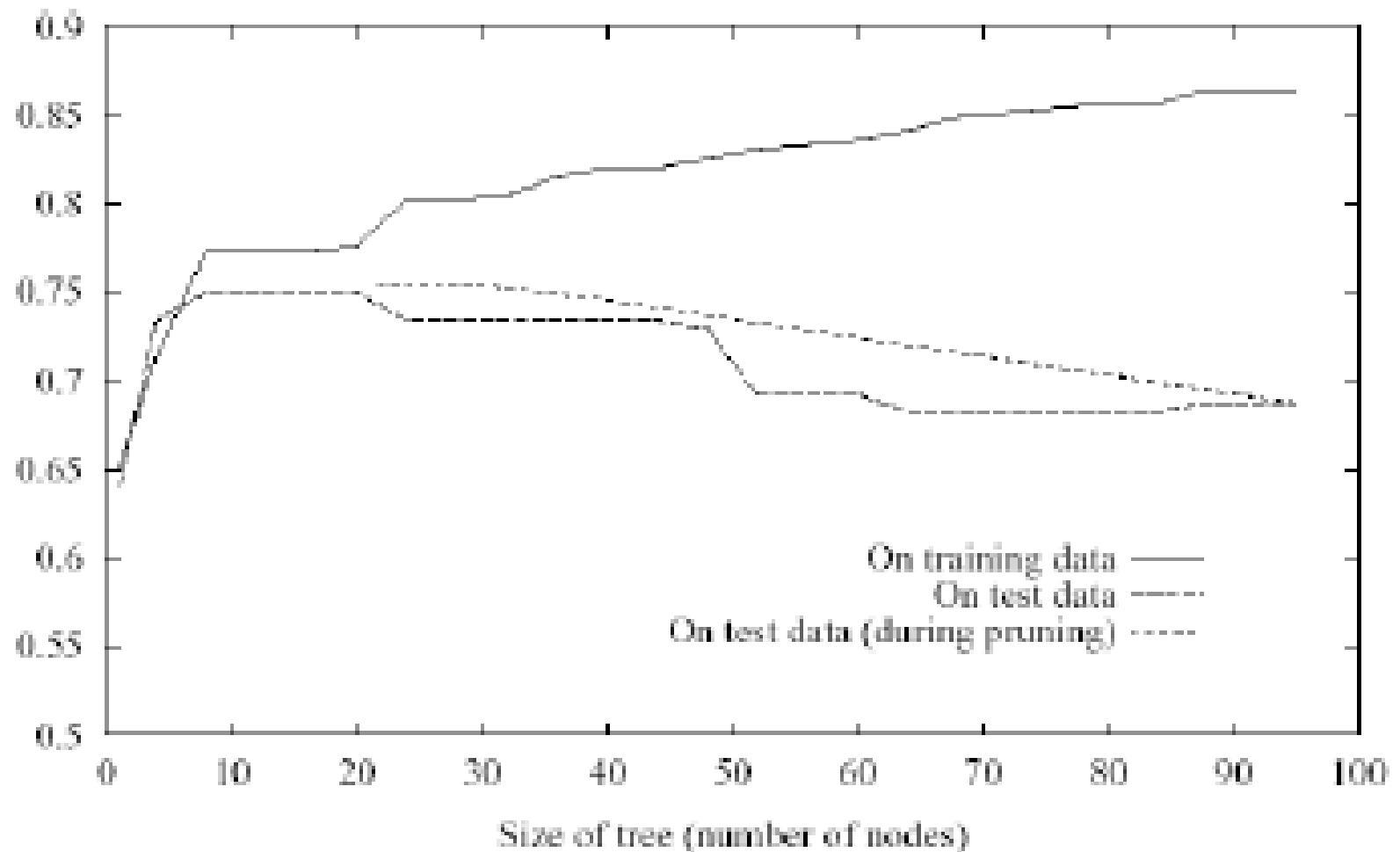
# Over-fitting



# Avoid Overfitting

- Early stop
  - Stop growing the tree when data split does not offer large benefit (e.g., compare information gain to a threshold, or perform statistical testing to decide if the gain is significant)
- Post pruning
  - Separate training data into **training set** and **validating set**
  - Evaluate impact on validation set when pruning each possible node
  - Greedily prune the node that most improves the validation set performance

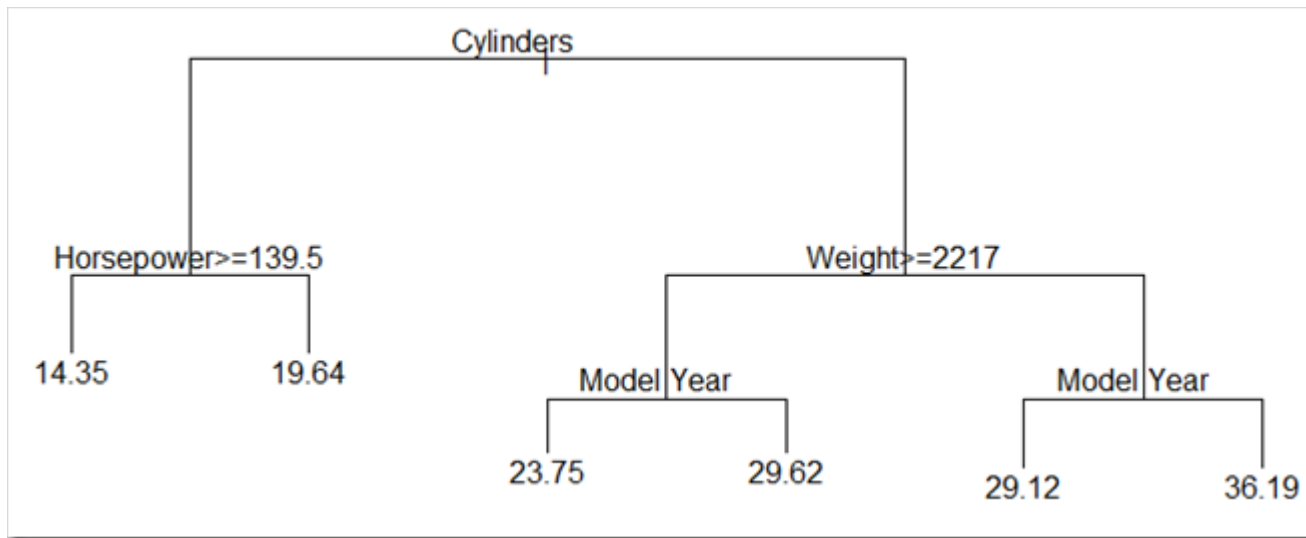
# Effect of Pruning



# Regression Tree

- Similar ideas can be applied for regression problems
- Prediction is computed as the average of the target values of all examples in the leaf node
- Uncertainty is measured by sum of squared errors within the node

# Example Regression Tree



Predicting MPG of a car given its # of cylinders, horsepower, weight, and model year



# Summary

- Decision tree is a very flexible classifier
  - Can model arbitrarily complex decision boundaries
  - By changing the depth of the tree (or # of nodes in the tree), we can increase or decrease the model complexity
  - Handling both continuous and discrete features
- Learning of the decision tree
  - Greedy top-down induction
  - Not guaranteed to find an optimal decision tree
- DT can overfitting to noise and outliers
  - Can be controlled by early stopping or post pruning