# CS6375: Machine Learning
## Gautam Kunapuli

# Decision Trees

**THE UNIVERSITY OF TEXAS AT DALLAS**
Erik Jonsson School of Engineering and Computer Science

# Example: Restaurant Recommendation

**Example**: Develop a model to **recommend restaurants** to users depending on their past dining experiences.

Here, the features are **cost** ($x_1$) and the user's **spiciness rating** of the food at the restaurant ($x_2$) and the label is if they liked the food ($y_i = \circ$) or not ($y_i = \blacksquare$).

A data set is **linearly separable** if there exists a hyperplane that separates positive examples from negative examples.
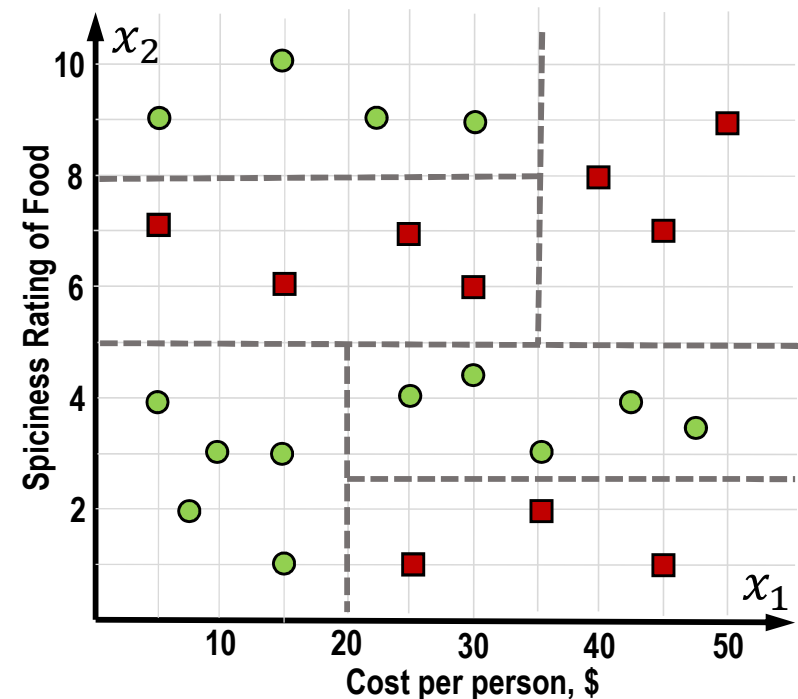- Relatively easy to learn (using standard techniques)
- Easy to visualize and interpret

Many **data sets in real world** are **not linearly separable**!
Two options:
- Use **non-linear features,** and learn a linear classifier in the transformed non-linear feature space
- Use **non-linear classifiers**

Decision Trees can handle nonlinear separable data sets and are one of the **most popular classifiers**
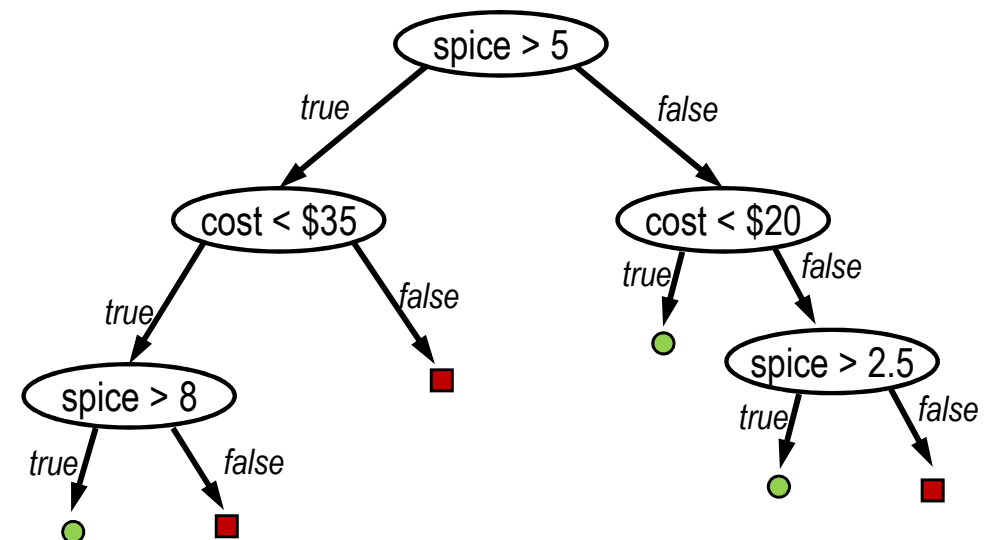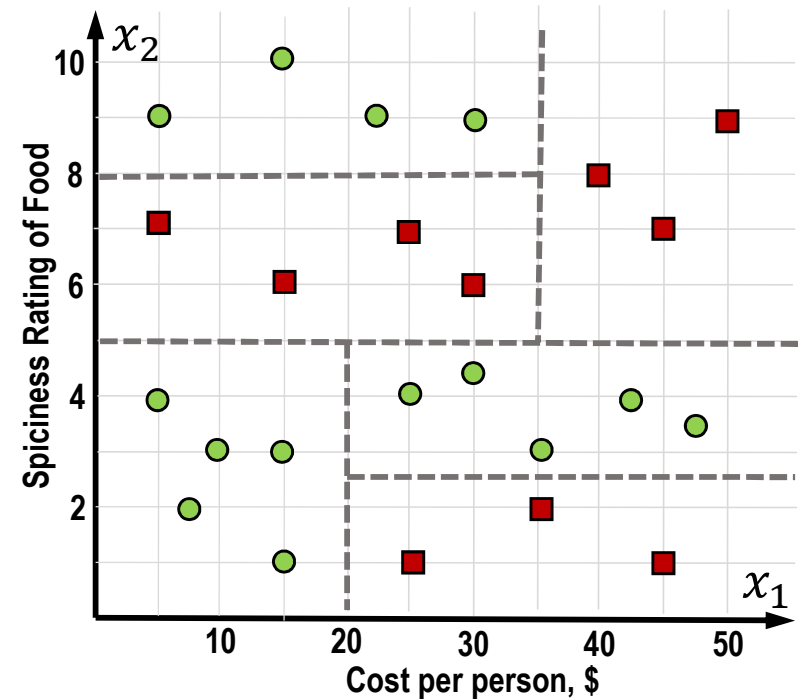
# Decision Trees: Introduction

**Example:** Develop a model to **recommend restaurants** to users depending on their past dining experiences.

Here, the features are **cost** ($x_1$) and the user's **spiciness rating** of the food at the restaurant ($x_2$) and the label is if they liked the food ($y_i = \circ$) or not ($y_i = \blacksquare$).

Decision Trees represent decision-making as a **checklist of questions**, and visualize it using a tree-structure

Decision Tree **representation**:

- Each **non-leaf node tests an attribute**/feature
- Each **branch corresponds to attribute**/feature value, a decision (to choose a path) as a result of the test
- Each **leaf node assigns a classification**

# Decision Trees: Introduction

**Example:** Develop a model to **recommend restaurants** to users depending on their past dining experiences.

Here, the features are **cost** ($x_1$) and the user's **spiciness rating** of the food at the restaurant ($x_2$) and the label is if they liked the food ($y_i = \bigcirc$) or not ($y_i = \blacksquare$).

- Decision trees divide the feature space into **axis-parallel rectangles**
- Decision Trees can handle **arbitrarily non-linear representations**, given sufficient tree complexity
- Worst-case scenario: the decision tree has an **exponential number of nodes!** (why?)
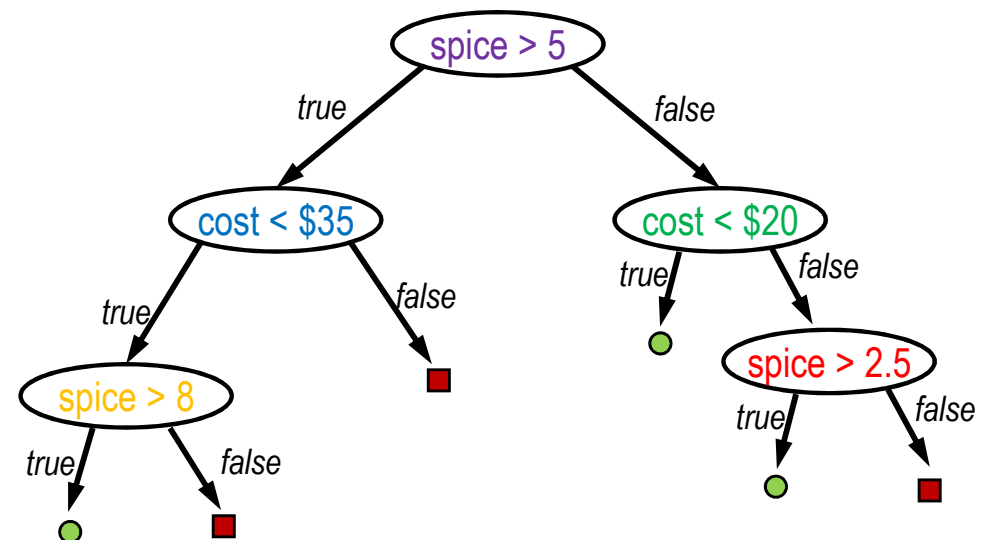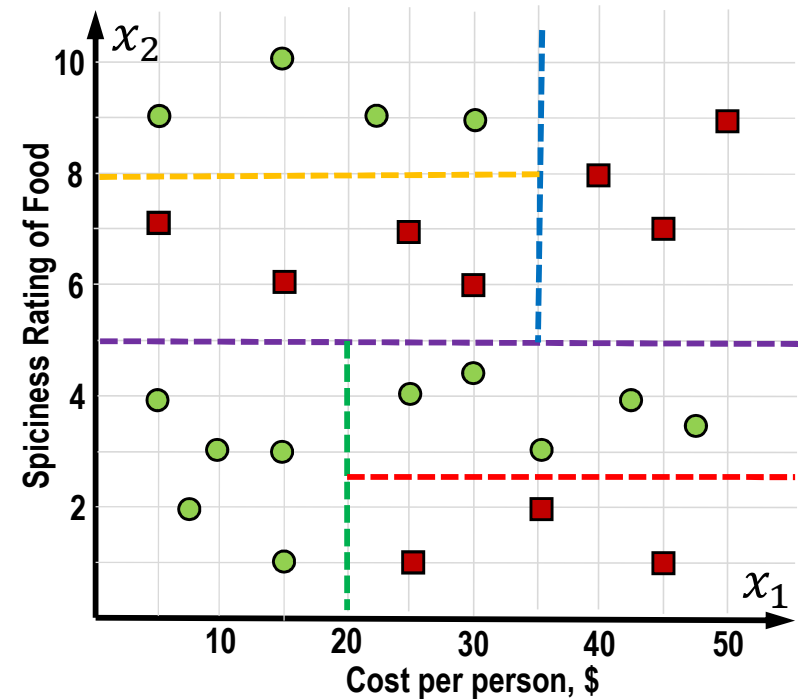
# Decision Trees: Introduction

**Example**: Develop a model to **recommend restaurants** to users depending on their past dining experiences.

Here, the features are **cost** ($x_1$) and the user's **spiciness rating** of the food at the restaurant ($x_2$) and the label is if they liked the food ($y = +1$) or not ($y = +1$).

- Decision trees divide the feature space into **axis-parallel rectangles**
- Decision Trees can handle **arbitrarily non-linear representations**, given sufficient tree complexity
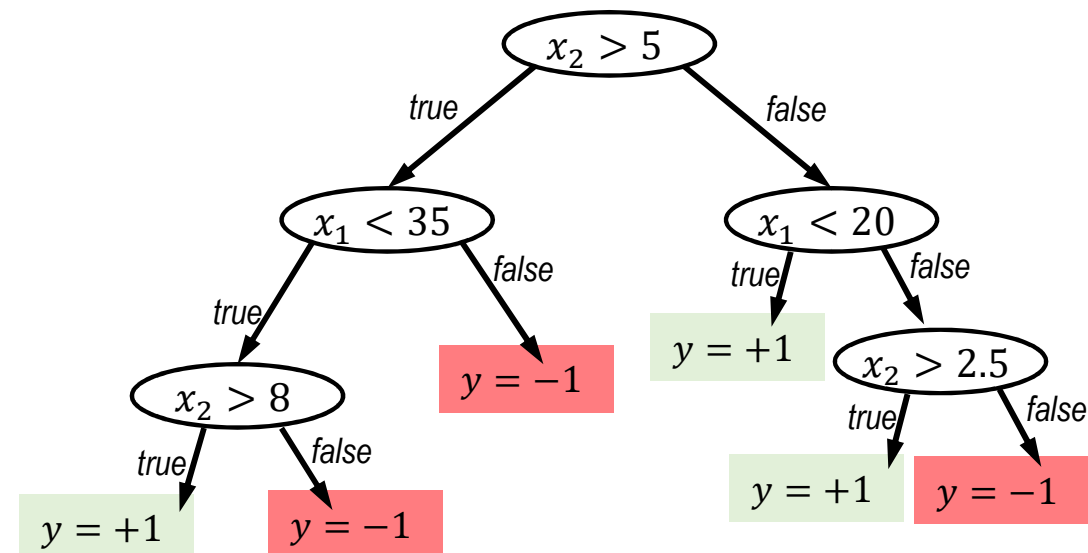- Worst-case scenario: the decision tree has an **exponential number of nodes!**
  - If the target function has $n$ Boolean features, there are $2^n$ possible inputs
  - In the worst case, there is one leaf node for each input (for example: XOR)

**Decision trees are <u>not</u> unique, and many decision trees can represent the same hypothesis!**
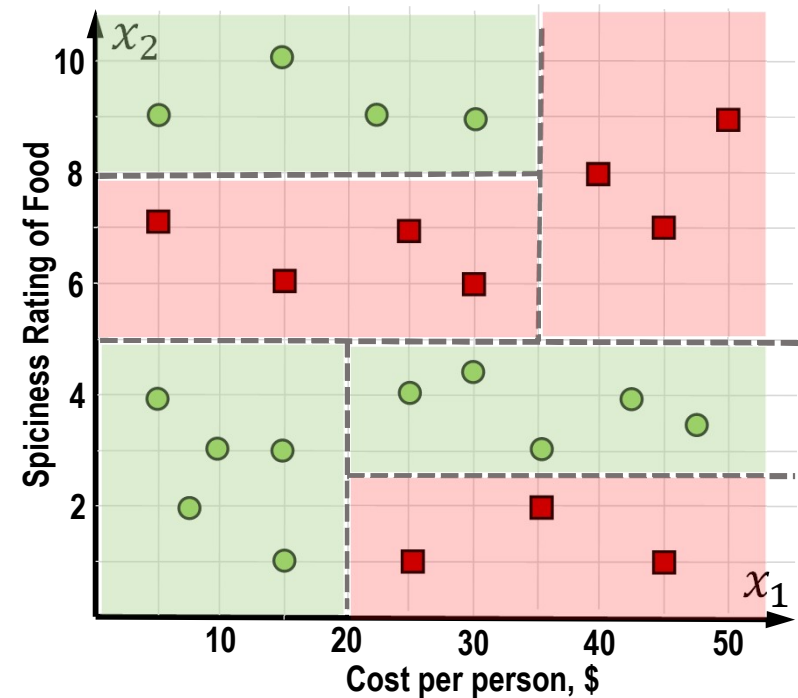
# Decision Trees: Introduction

**Example:** Develop a model to **recommend restaurants** to users depending on their past dining experiences.

Here, the features are **cost** ($x_1$) and the user's **spiciness rating** of the food at the restaurant ($x_2$) and the label is if they liked the food ($y = +1$) or not ($y = +1$).

When do you want Decision Trees?
When instances are **describable by attribute-value pairs**:
- target function is **discrete-valued**
- **disjunctive hypothesis** may be required
- need for **interpretable** model

Examples:
- Equipment or medical diagnosis
- Credit risk analysis
- Modeling calendar scheduling preferences

# Learning Decision Trees

**Problem Formulation:** Find a decision tree **h** that achieves minimum misclassification errors on the training data

- **Solution Approach 1 (Naïve solution):**  Create a decision tree with one path from root to leaf for each training example. *Such a tree would just memorize the training data, and will **not generalize well to new points***.
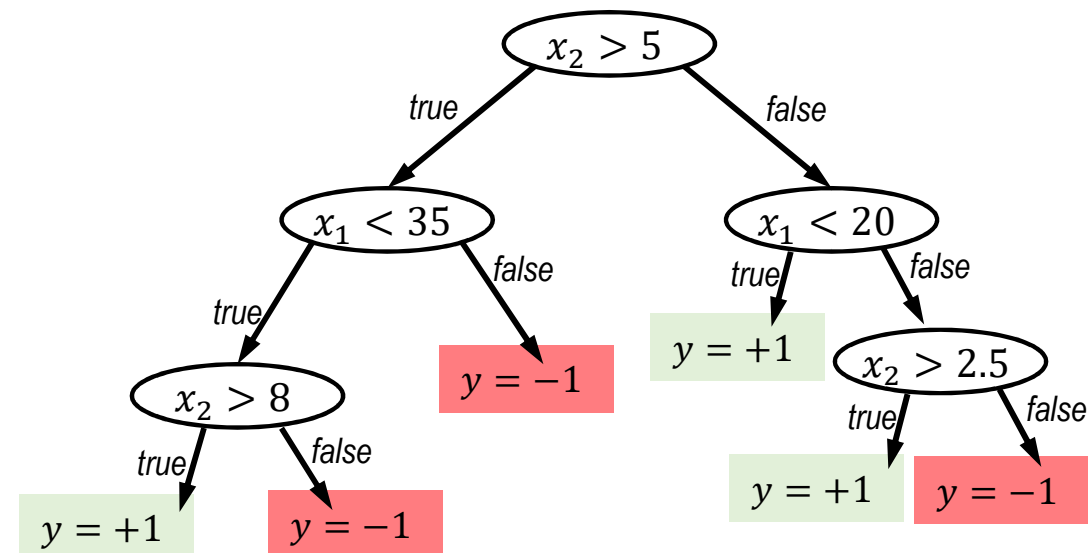
- **Solution Approach 2 (Exact solution):** Find the **smallest** tree that minimizes the classification error. *Finding this solution is **NP-Hard**!*

- **Solution Approach 3 (Heuristic solution)**: Top-down greedy search

**Initialize:** Choose the best feature $f^*$ for the root of the tree
**Function** GrowTree(data, $f^*$)
  [1]Separate data into subsets $\{S_1, S_2, ..., S_k\}$, where each
  subset $S_i$ contains examples that have the **same value for** $f^*$
  [2] for $S_i \in \{S_1, S_2, ..., S_k\}$
    Choose the best feature $f_i^*$ for the next node
    **Recursively** GrowTree($S_i, f_i^*$) until all examples have the
      same class label

# Learning Decision Trees

**Problem Formulation**: Find a decision tree **h** that achieves minimum misclassification errors on the training data

- **Solution Approach 1 (Naïve solution):** Create a decision tree with one path from root to leaf for each training example. *Such a tree would just memorize the training data, and will **not generalize well to new points***.

- **Solution Approach 2 (Exact solution):** Find the **smallest** tree that minimizes the classification error. *Finding this solution is **NP-Hard**!*

- **Solution Approach 3 (Heuristic solution)**: Top-down greedy search

**Initialize:** Choose the best feature $f^*$ for the root of the tree
**Function** GrowTree(data, $f^*$)
    ¹Separate data into subsets $\{S_1, S_2, ..., S_k\}$, where each subset $S_i$ contains examples that have the **same value for** $f^*$
    ² for $S_i \in \{S_1, S_2, ..., S_k\}$
        Choose **the best feature** $f_i^*$ for the next node
        **Recursively** GrowTree($S_i, f_i^*$) until all **examples have the same class label**

**How do we pick the best feature?**

**How do we decide when to stop?**

# Learning Decision Trees

**Problem Formulation:** Find a decision tree **h** that achieves minimum misclassification errors on the training data

**Solution Approach 3 (Heuristic solution)**: Top-down greedy search
**Initialize:** Choose the best feature $f^*$ for the root of the tree
**Function** GrowTree(data, $f^*$)
    [1]Separate data into subsets $\{S_1, S_2,..., S_k\}$, where each subset $S_i$ contains examples that have the **same value for** $f^*$
    [2] for $S_i \in \{S_1, S_2,..., S_k\}$
        Choose **the best feature** $f_i^*$ for the next node
        **Recursively** GrowTree($S_i, f_i^*$) until all **examples have the same class label**

**How do we pick the next best feature to place in a decision tree?**
- Random choice
- Largest number of values
- Fewest number of values
- **Lowest classification error**
- Information theoretic measure (Quinlan's approach)



Training examples

| $x_1$ | $x_2$ | $x_3$ | $y$ |
|---|---|---|---|
| 0 | 0 | 0 | 1 |
| 0 | 0 | 1 | 0 |
| 0 | 1 | 0 | 1 |
| 0 | 1 | 1 | 1 |
| 1 | 0 | 0 | 0 |
| 1 | 0 | 1 | 1 |
| 1 | 1 | 0 | 0 |
| 1 | 1 | 1 | 0 |

split on $x_1$
$J = 2$

split on $x_2$
$J = 4$

split on $x_3$
$J = 4$

# Learning Decision Trees

**Problem Formulation:** Find a decision tree **h** that achieves minimum misclassification errors on the training data

**Solution Approach 3 (Heuristic solution):** Top-down greedy search
**Initialize:** Choose the best feature $f^*$ for the root of the tree
**Function** GrowTree(data, $f^*$)
   [1]Separate data into subsets $\{S_1, S_2, ..., S_k\}$, where each subset $S_i$ contains examples that have the **same value for** $f^*$
   [2] for $S_i \in \{S_1, S_2, ..., S_k\}$
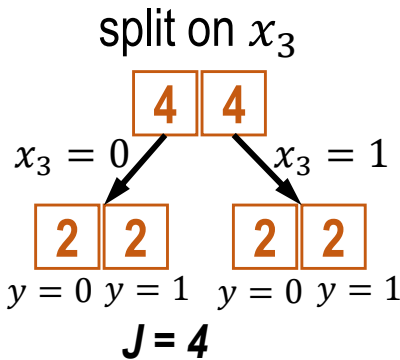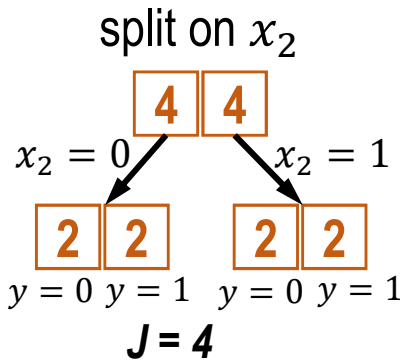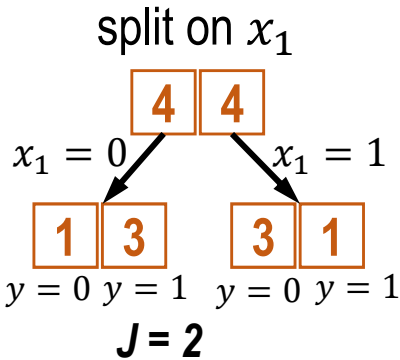      Choose **the best feature** $f_i^*$ for the next node
      **Recursively** GrowTree($S_i, f_i^*$) until all **examples have the same class label**

**How do we pick the next best feature to place in a decision tree?**
- Random choice
- Largest number of values
- Fewest number of values
- **Lowest classification error**
- Information theoretic measure (Quinlan's approach)

| $x_1$ | $x_2$ | $x_3$ | $y$ |
|-------|-------|-------|-----|
| 0 | 0 | 0 | 1 |
| 0 | 0 | 1 | 0 |
| 0 | 1 | 0 | 1 |
| 0 | 1 | 1 | 1 |
| 1 | 0 | 0 | 0 |
| 1 | 0 | 1 | 1 |
| 1 | 1 | 0 | 0 |
| 1 | 1 | 1 | 0 |

**Training examples**

split on $x_1$

$\boxed{4}\ \boxed{4}$
$x_1 = 0$ ↙  ↘ $x_1 = 1$
$\boxed{1}\ \boxed{3}$    $\boxed{3}\ \boxed{1}$
$y=0\ y=1$   $y=0\ y=1$
**J = 2**

$y=0$   $y=1$
$\boxed{1/4}\ \boxed{3/4}$

split on $x_2$

$\boxed{4}\ \boxed{4}$
$x_2 = 0$ ↙  ↘ $x_2 = 1$
$\boxed{2}\ \boxed{2}$    $\boxed{2}\ \boxed{2}$
$y=0\ y=1$   $y=0\ y=1$
**J = 4**

$y=0$   $y=1$
$\boxed{2/4}\ \boxed{2/4}$

split on $x_3$

$\boxed{4}\ \boxed{4}$
$x_3 = 0$ ↙  ↘ $x_3 = 1$
$\boxed{2}\ \boxed{2}$    $\boxed{2}\ \boxed{2}$
$y=0\ y=1$   $y=0\ y=1$
**J = 4**

*Can think of counts as **probability distributions over the labels***

# Learning Decision Trees

**Problem Formulation:** Find a decision tree **h** that achieves minimum misclassification errors on the training data

**Solution Approach 3 (Heuristic solution)**: Top-down greedy search

**Initialize:** Choose the best feature $f^*$ for the root of the tree

**Function** GrowTree(data, $f^*$)

[1]Separate data into subsets $\{S_1, S_2,..., S_k\}$, where each subset $S_i$ contains examples that have the **same value for** $f^*$

[2]for $S_i \in \{S_1, S_2,..., S_k\}$

Choose **the best feature** $f_i^*$ for the next node
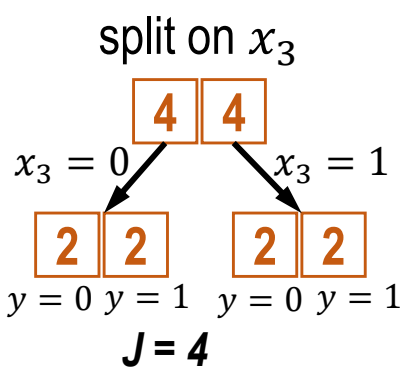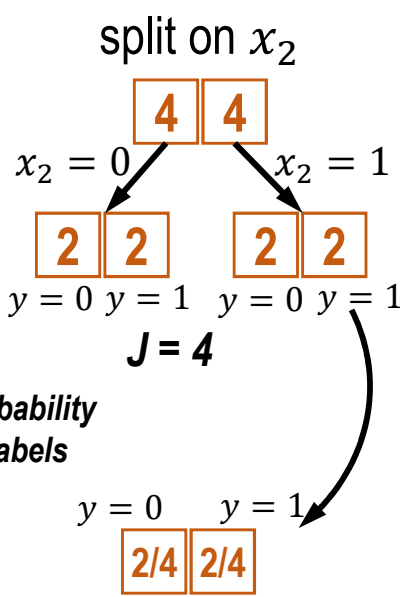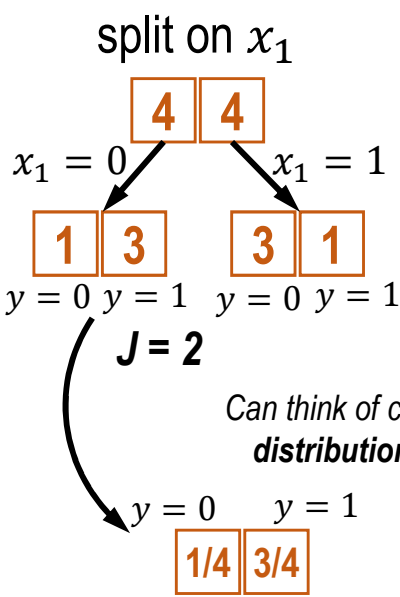
**Recursively** GrowTree($S_i, f_i^*$) until all **examples have the same class label**

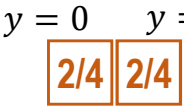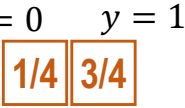### How do we pick the next best feature to place in a decision tree?
- Random choice
- Largest number of values
- Fewest number of values
- **Lowest classification error**
- Information theoretic measure (Quinlan's approach)

The selected attribute is a **good split** if we are **more "certain"** about the classification after the split (compare with the perceptron)

- If each partition with respect to the chosen attribute has a **distinct class label**, we are **completely certain** about the classification

$$y = 0 \quad y = 1$$
$$\boxed{0.0}\ \boxed{1.0}$$

- If **class labels are evenly divided** between partitions, we are **very uncertain** about the classification

$$y = 0 \quad y = 1$$
$$\boxed{0.5}\ \boxed{0.5}$$

### split on $x_1$

$$\boxed{4}\ \boxed{4}$$
$x_1 = 0$    $x_1 = 1$
$$\boxed{1}\ \boxed{3} \quad \boxed{3}\ \boxed{1}$$
$y = 0 \ \ y = 1 \ \ y = 0 \ \ y = 1$

**J = 2**

### split on $x_2$

$$\boxed{4}\ \boxed{4}$$
$x_2 = 0$    $x_2 = 1$
$$\boxed{2}\ \boxed{2} \quad \boxed{2}\ \boxed{2}$$
$y = 0 \ \ y = 1 \ \ y = 0 \ \ y = 1$

**J = 4**

*Can think of counts as **probability distributions over the labels***

$y = 0 \quad y = 1$
$$\boxed{1/4}\ \boxed{3/4}$$

$y = 0 \quad y = 1$
$$\boxed{2/4}\ \boxed{2/4}$$

# Learning Decision Trees

**Problem Formulation:** Find a decision tree **h** that achieves minimum misclassification errors on the training data
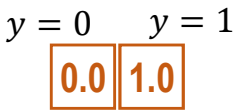
**Solution Approach 3 (Heuristic solution)**: Top-down greedy search
**Initialize:** Choose the best feature $f^*$ for the root of the tree
**Function** GrowTree(data, $f^*$)
   [1]Separate data into subsets $\{S_1, S_2, ..., S_k\}$, where each subset $S_i$ contains examples that have the **same value for** $f^*$
   [2] for $S_i \in \{S_1, S_2, ..., S_k\}$
      Choose **the best feature** $f_i^*$ for the next node
      **Recursively** GrowTree($S_i, f_i^*$) until all **examples have the same class label**

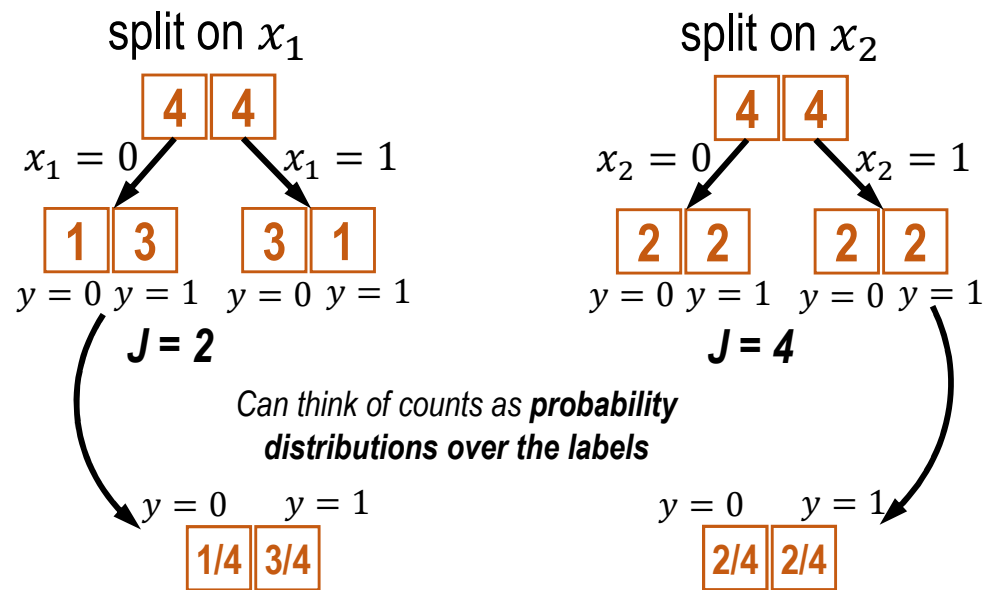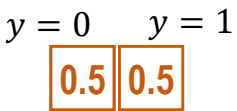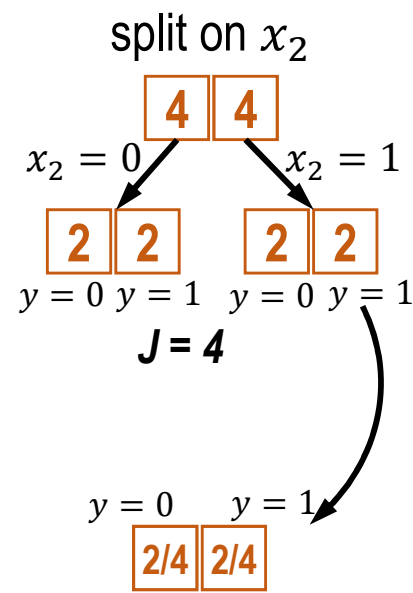How do we pick the next best feature to place in a decision tree?
- Random choice
- Largest number of values
- Fewest number of values
- **Lowest classification error**
- Information theoretic measure (Quinlan's approach)

The selected attribute is a **good split** if we are **more "certain"** about the classification after the split (compare with the perceptron)

We need a better way to resolve the uncertainty!

split on $x_2$



$J = 4$

- If each partition with respect to the chosen attribute has a **distinct class label**, we are **completely certain** about the classification

| $y = 0$ | $y = 1$ |
|---|---|
| 0.0 | 1.0 |

- If **class labels are evenly divided** between partitions, we are **very uncertain** about the classification

| $y = 0$ | $y = 1$ |
|---|---|
| 0.5 | 0.5 |

# Discrete Probability and Information Theory

A **discrete probability distribution** describes the probability
of occurrence of each value of a discrete random variable.

The **surprise** or **self-information** of each event of $X$ is
defined to be

$$S(X = x) = -\mathbf{log_2} \, \text{Prob}(X = x)$$

- An event with probability 1 has zero surprise; *this is
  because when the content of a message is known
  beforehand with certainty, there is no actual information
  conveyed*
- The **smaller the probability** of event, the **larger the
  quantity of self-information** associated with the
  message that the event occurred
- An event with probability 0 has infinite surprise

- The surprise is the **asymptotic number of bits of
  information** that need to be transmitted to a recipient
  who knows the probabilities of the results. This is also
  called the **description length** of X.

*Random Variable: Number of heads when
tossing a coin 3 times*

| X | 0 | 1 | 2 | 3 |
|---|---|---|---|---|
| Prob(X) | 1/8 | 3/8 | 3/8 | 1/8 |
| $-\log_2 \text{P}(X)$ | 3 | 1.415 | 1.415 | 3 |
| $-\log_e \text{P}(X)$ | 2.079 | 0.980 | 0.980 | 2.079 |
| $-\log_{10} \text{P}(X)$ | 0.903 | 0.426 | 0.426 | 0.903 |

*If the logarithm is base 2, the unit of information is
bits, base e is nats and base 10 hartleys*

# Entropy

A standard way to measure **uncertainty of a random variable** is to use **entropy**

$$H(X) = -\sum_{x} P(X = x)\log_2 P(X = x)$$

• Note that the entropy is computed by **summing over all the events**/outcomes/states of the random variable.

• Entropy is maximized for uniform distributions, where the probability of all outcomes is equal (is this what we want?)

• Entropy is minimized for distributions that place all their probability on a single outcome (or is this what we want?)

The entropy of label distributions can be computed as:

$$H(y) = -P(y = 0)\log_2 P(y = 0) - P(y = 1)\log_2 P(y = 1)$$

split on $x_i$

| 4 | 4 |

$x_1 = 0$      $x_1 = 1$

| 2 | 2 | | 2 | 2 |

$y = 0$ $y = 1$   $y = 0$ $y = 1$

$H(y) = 1$

split on $x_j$

| 4 | 4 |

$x_3 = 0$      $x_3 = 1$

| 1 | 3 | | 3 | 1 |

$y = 0$ $y = 1$   $y = 0$ $y = 1$

$H(y) = 0.8113$

split on $x_k$

| 4 | 4 |

$x_2 = 0$      $x_2 = 1$

| 1 | 0 | | 0 | 1 |

$y = 0$ $y = 1$   $y = 0$ $y = 1$

$H(y) \approx 0$

(use $0 \cdot \log_2 0 = 0$)

# Conditional Entropy and Mutual Information

Entropy can also be computed when conditioned on another variable:

$$H(Y|X) = -\sum_{x} P(X = x) \sum_{y} P(Y = y \mid X = x) \log_2 (Y = y \mid X = x)$$

This is called **conditional entropy** and is the amount of information needed to quantify the random variable $Y$ given the random variable $X$.

The **mutual information** or **information gain** between two random variables is defined as

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

This is the amount of information we learn about $Y$ by knowing the value of $X$ and vice-versa (it is symmetric).

In our case, **larger information gain** corresponds **to less uncertainty about $Y$ (labels) given $X$ (data)**.



$$H(Y) = 0.9183$$

$P(X_1=0)=0.6677 \qquad P(X_1=1)=0.3333$

$H(Y|X_1=0)$
$=-0.6*\log 0.6-0.4*\log 0.4$
$=0.9710$

$H(Y|X_1)=0.8879$
$I(X_1, Y)=0.0304$

$H(Y|X_1=1)$
$=-0.8*\log 0.8-0.2*\log 0.2$
$=0.7219$

# Choosing the Best Feature

split on $x_i$

$P(y = 0)$ | 29 | 35 | $P(y = 1)$

$x_1 = 0$ (48)
$P(x_i = 0)$

$x_1 = 1$ (16)
$P(x_i = 1)$

| 21 | 5 | | 8 | 30 |

$y = 0$ $y = 1$   $y = 0$ $y = 1$

$P(y = 0|x_i = 0)$
$P(y = 1|x_i = 0)$
$P(y = 0|x_i = 1)$
$P(y = 1|x_i = 1)$

*Where are all the probabilities?*

Step 1: Count the various combinations of features and labels

split on $x_1$

| 5 | 3 |

$x_1 = 0$ (4)              $x_1 = 1$ (4)

| 1 | 3 | | 4 | 0 |

$y = 0$ $y = 1$  $y = 0$ $y = 1$

split on $x_2$

| 5 | 3 |

$x_2 = 0$ (4)              $x_2 = 1$ (4)

| 2 | 2 | | 3 | 1 |

$y = 0$ $y = 1$  $y = 0$ $y = 1$

Step 2: Convert to probabilities

split on $x_1$

| 5/8 | 3/8 |

$x_1 = 0$ (4/8)          $x_1 = 1$ (4/8)

| 1/4 | 3/4 | | 4/4 | 0/4 |

$y = 0$ $y = 1$  $y = 0$ $y = 1$

split on $x_2$

| 5/8 | 3/8 |

$x_2 = 0$ (4/8)          $x_2 = 1$ (4/8)

| 2/4 | 2/4 | | 3/4 | 1/4 |

$y = 0$ $y = 1$  $y = 0$ $y = 1$

| $x_1$ | $x_2$ | $y$ |
|-------|-------|-------|
| 1 | 1 | 0 (+) |
| 1 | 0 | 0 (+) |
| 1 | 1 | 0 (+) |
| 1 | 0 | 0 (+) |
| 0 | 1 | 0 (+) |
| 0 | 0 | 1 (-) |
| 0 | 1 | 1 (-) |
| 0 | 0 | 1 (-) |

# Choosing the Best Feature

Step 3: Compute information gain for both splits and pick the variable with the biggest gain

split on $x_i$

$$29 \quad 35 \qquad H(y)$$

$x_1 = 0 \qquad\qquad x_1 = 1$
$P(x_i = 0) \qquad\qquad P(x_i = 1)$

$$21 \quad 5 \qquad\quad 8 \quad 30$$

$y = 0 \; y = 1 \quad y = 0 \; y = 1$
$H(y \mid x_i = 0) \quad H(y \mid x_i = 1)$

$H(y \mid x_i)$

*Where are all the entropies?*

$$H(y) = -\frac{5}{8}\log\frac{5}{8} - \frac{3}{8}\log\frac{3}{8}$$

split on $x_1$

$$5/8 \quad 3/8$$

$x_1 = 0 \; (4/8) \qquad x_1 = 1 \; (4/8)$

$$1/4 \quad 3/4 \qquad 4/4 \quad 0/4$$

$y = 0 \; y = 1 \qquad y = 0 \; y = 1$

$H(y|x_1 = 0) = -\frac{1}{4}\log\frac{1}{4} - \frac{3}{4}\log\frac{3}{4}$

$H(y|x_1 = 1) = -1\log 1 - 0\log 0$

$$H(y \mid x_1) = -\frac{4}{8}H(y|x_1 = 0) - \frac{4}{8}H(y|x_1 = 1)$$

$$I(x_1, y) = H(y) - H(y \mid x_1)$$

split on $x_2$

$$5/8 \quad 3/8$$

$x_2 = 0 \; (4/8) \qquad x_2 = 1 \; (4/8)$

$$2/4 \quad 2/4 \qquad 3/4 \quad 1/4$$

$y = 0 \; y = 1 \qquad y = 0 \; y = 1$

$H(y|x_2 = 0) = -\frac{2}{4}\log\frac{2}{4} - \frac{2}{4}\log\frac{2}{4}$

$H(y|x_2 = 1) = -\frac{3}{4}\log\frac{3}{4} - \frac{1}{4}\log\frac{1}{4}$

$$H(y \mid x_2) = -\frac{4}{8}H(y|x_2 = 0) - \frac{4}{8}H(y|x_2 = 1)$$

$$I(x_2, y) = H(y) - H(y \mid x_2)$$

| $x_1$ | $x_2$ | $y$ |
|-------|-------|--------|
| 1 | 1 | 0 (+) |
| 1 | 0 | 0 (+) |
| 1 | 1 | 0 (+) |
| 1 | 0 | 0 (+) |
| 0 | 1 | 0 (+) |
| 0 | 0 | 1 (-) |
| 0 | 1 | 1 (-) |
| 0 | 0 | 1 (-) |

$I(x_1, y) > I(x_2, y) \Rightarrow$ pick feature $x_1$ next

# The ID3 Algorithm

*The ID3 (Iterative Dichotomizer) and its successor, C4.5 were developed by Ross Quinlan in the early to mid 1980s and are widely considered to be a landmark machine learning algorithms, and until at least 2008, were the #1 data mining tool.*

---

ID3($Examples$, $Target\_attribute$, $Attributes$)

*Examples are the training examples. Target_attribute is the attribute whose value is to be predicted by the tree. Attributes is a list of other attributes that may be tested by the learned decision tree. Returns a decision tree that correctly classifies the given Examples.*

- Create a *Root* node for the tree
- If all *Examples* are positive, Return the single-node tree *Root*, with label $= +$
- If all *Examples* are negative, Return the single-node tree *Root*, with label $= -$
- If *Attributes* is empty, Return the single-node tree *Root*, with label $=$ most common value of *Target_attribute* in *Examples*
- Otherwise Begin
    - $A \leftarrow$ the attribute from *Attributes* that best* classifies *Examples*
    - The decision attribute for *Root* $\leftarrow A$
    - For each possible value, $v_i$, of $A$,
        - Add a new tree branch below *Root*, corresponding to the test $A = v_i$
        - Let $Examples_{v_i}$ be the subset of *Examples* that have value $v_i$ for $A$
        - If $Examples_{v_i}$ is empty
            - Then below this new branch add a leaf node with label $=$ most common value of *Target_attribute* in *Examples*
            - Else below this new branch add the subtree
            ID3($Examples_{v_i}$, $Target\_attribute$, $Attributes - \{A\}$))
- End
- Return *Root*

---

# Some Final Details

## When do we terminate?

• If the current set is **"pure"** (i.e., has a single label in the output), stop

• If you **run out of attributes to recurse on**, even if the current data set isn't pure, stop and use a majority vote

• If a partition contains no data points, use the majority vote at its parent in the tree

• If a partition contains no data items, nothing to recurse on

• For fixed depth decision trees, the **final label is determined by majority vote**

## How do we handle real-valued features?

• For continuous attributes, use threshold splits

• Split the tree into $x_k < t$ and $x_k \geq t$

• Can split on the same attribute multiple times on the same path down the tree

## How do we select the splitting threshold?

# Overfitting in Decision Trees

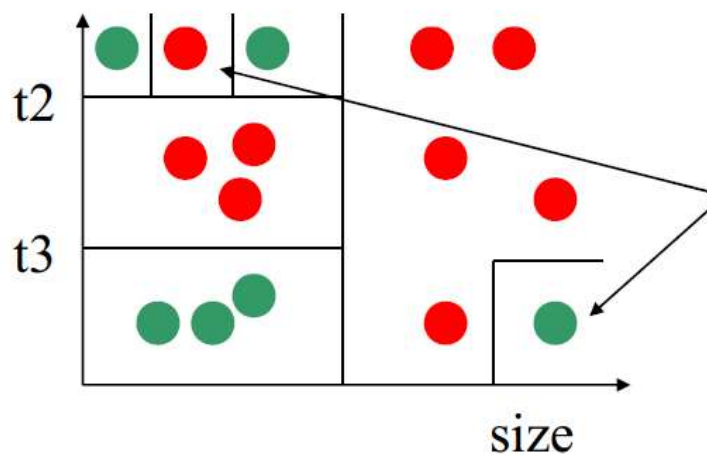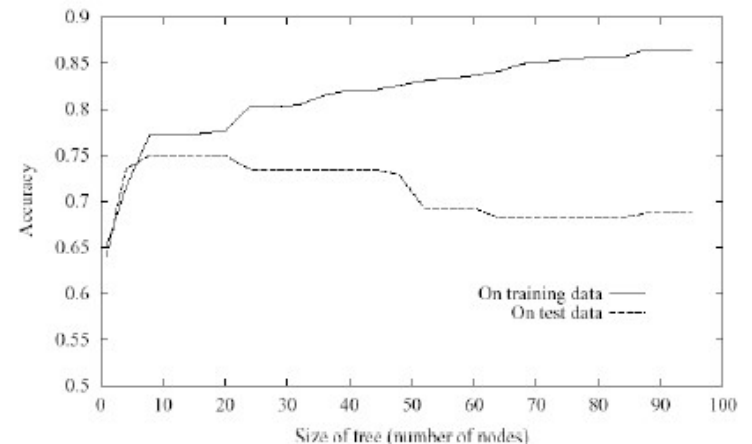**Hypothesis space is complete!** *Target function is surely in there*
**No back tracking;** *Greedy thus local minima*
**Statistics-based search choices;** *Robust to noisy data*
**Inductive bias: heuristically prefer shortest tree**

**Decision trees will always overfit!**
It is always possible to obtain zero training error on the input data
with a deep enough tree (if there is no noise in the labels)



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

# Overfitting in Decision Trees

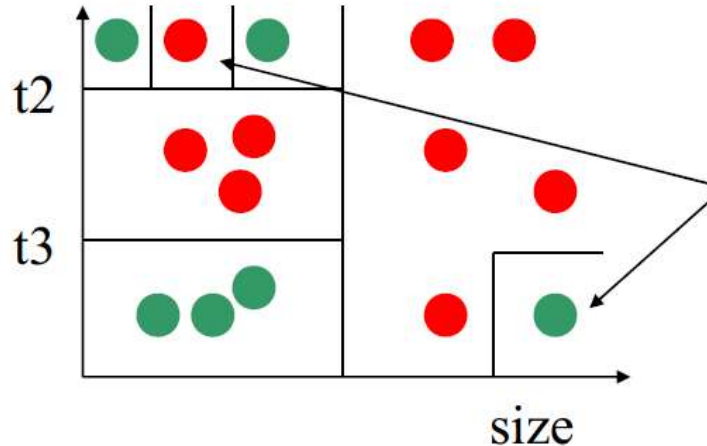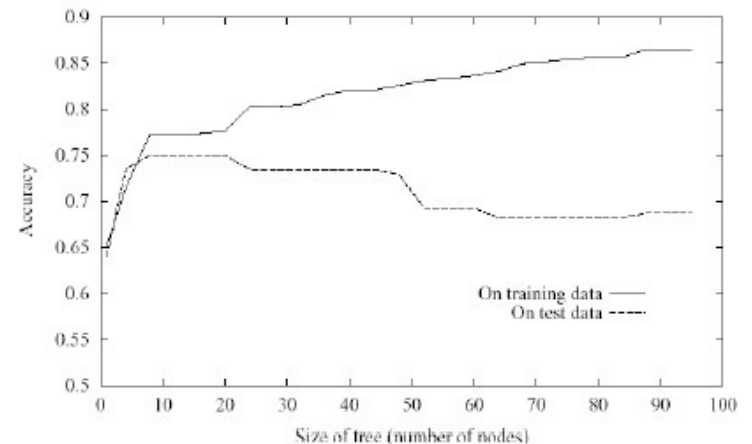**Hypothesis space is complete!** *Target function is surely in there*
**No back tracking;** *Greedy thus local minima*
**Statistics-based search choices;** *Robust to noisy data*
**Inductive bias: heuristically prefer shortest tree**

**Decision trees will always overfit!**

It is always possible to obtain zero training error on the input data
with a deep enough tree (if there is no noise in the labels)



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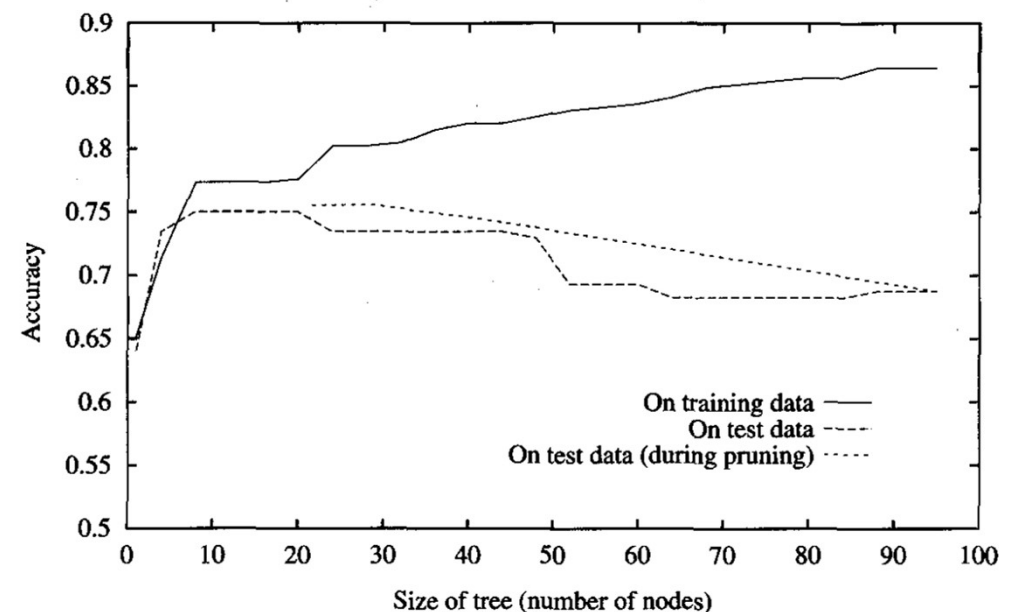the tree is grown larger
to capture these examples

# Avoiding Overfitting in Decision Trees

## Pre-pruning/early stopping

- Typical stopping criterion
  - No error (if all instances belong to same class)
  - IF all the attribute values are same

- More restrictive conditions
  - Stop growing when data split is not statistically significant (example using chi-square test)
  - Stop if the number of instances is less than a predefined threshold
  - Stop if expanding does not significantly improve the measures (information gain)

## Post-pruning after growing a full tree

- Separate data into training and validation sets
- Evaluate impact on validation set **when a node is "pruned"**
- **Greedily remove** node that improves performance the most
- Produces smallest version of most accurate subtree
- Typically use minimum description length (MDL) for post-pruning

# Some Post-pruning Methods

**Reduced-Error Pruning**

- Use a validation set (tuning) to identify errors at every node
- Prune node with highest reduction error
- Repeat until error no longer reduces

**Pessimistic Pruning**

- No necessity of a validation set
- The error estimate at every node is conservative based on the training examples

**Rule-post Pruning**

- Convert tree to equivalent set of rules (how)?
- Prune each rule independently of others
- Sort final rules into desired sequence

# Decision Trees

- **Decision Trees** – popular and a very efficient hypothesis space
    - Variable size: Any Boolean function can be represented
    - Handles discrete and continuous features
    - Handles classification **and regression**
    - Easy to implement
    - Easy to use
    - Computationally cheap

- Constructive **heuristic** search: built top-down by adding nodes
- **Decision trees will overfit!**
    - zero bias classifier (no mistakes) = large variance
    - must use tricks to find simpler trees
        - early stopping, pruning etc.