



COMP7/8118 M50

# Data Mining

## Decision Tree: II

Xiaofei Zhang

*Slides compiled from Jiawei Han and Raymond C.W. Wong's work*

THE UNIVERSITY OF  
**MEMPHIS**

# Decision Tree Based Classification

- Advantages:
  - Inexpensive to construct
  - Extremely fast at classifying unknown records
  - Easy to interpret for small-sized trees
  - Accuracy is comparable to other classification techniques for many simple data sets

# Decision Tree Examples

- ID3 — Iterative Dichotomiser
  - Impurity Measurement
    - $\text{Gain}(A, T) = \text{Info}(T) - \text{Info}(A, T)$
- C4.5 — Classification
  - Impurity Measurement
    - $\text{Gain}(A, T) = (\text{Info}(T) - \text{Info}(A, T)) / \text{SplitInfo}(A)$ , where  $\text{SplitInfo}(A) = -\sum_{v \in A} p(v) \log p(v)$
- CART — Classification And Regression Trees
  - Impurity Measurement
    - Gini:  $I(P) = 1 - \sum_j p_j^2$

# Entropy

$$\text{Info}(T) = -\frac{1}{2} \log \frac{1}{2} - \frac{1}{2} \log \frac{1}{2} \\ = 1$$

For attribute Gender,

$$\text{Info}(T_{\text{female}}) = -\frac{3}{4} \log \frac{3}{4} - \frac{1}{4} \log \frac{1}{4} = 0.8113$$

$$\text{Info}(T_{\text{male}}) = -\frac{3}{4} \log \frac{3}{4} - \frac{1}{4} \log \frac{1}{4} = 0.8113$$

$$\text{Info}(\text{Gender}, T) = \frac{1}{2} \times \text{Info}(T_{\text{male}}) + \frac{1}{2} \times \text{Info}(T_{\text{female}}) = 0.8113$$

$$\text{Gain}(\text{Gender}, T) = \text{Info}(T) - \text{Info}(\text{Gender}, T) = 1 - 0.8113 = 0.1887$$

For attribute Gender,

$$\text{Gain}(\text{Gender}, T) = 0.1887$$

Gender	Income	Child	Insurance
female	high	no	yes
male	high	yes	yes
male	low	yes	yes
male	low	yes	yes
female	low	no	no
female	low	no	no
female	low	no	no
male	low	no	no

# Entropy

$$\text{Info}(T) = -\frac{1}{2} \log \frac{1}{2} - \frac{1}{2} \log \frac{1}{2} \\ = 1$$

For attribute Income,

$$\text{Info}(T_{\text{high}}) = -1 \log 1 - 0 \log 0 = 0$$

$$\text{Info}(T_{\text{low}}) = -\frac{1}{3} \log \frac{1}{3} - \frac{2}{3} \log \frac{2}{3} = 0.9183$$

$$\text{Info}(\text{Income}, T) = \frac{1}{4} \times \text{Info}(T_{\text{high}}) + \frac{3}{4} \times \text{Info}(T_{\text{low}}) = 0.6887$$

$$\text{Gain}(\text{Income}, T) = \text{Info}(T) - \text{Info}(\text{Income}, T) = 1 - 0.6887 = 0.3113$$

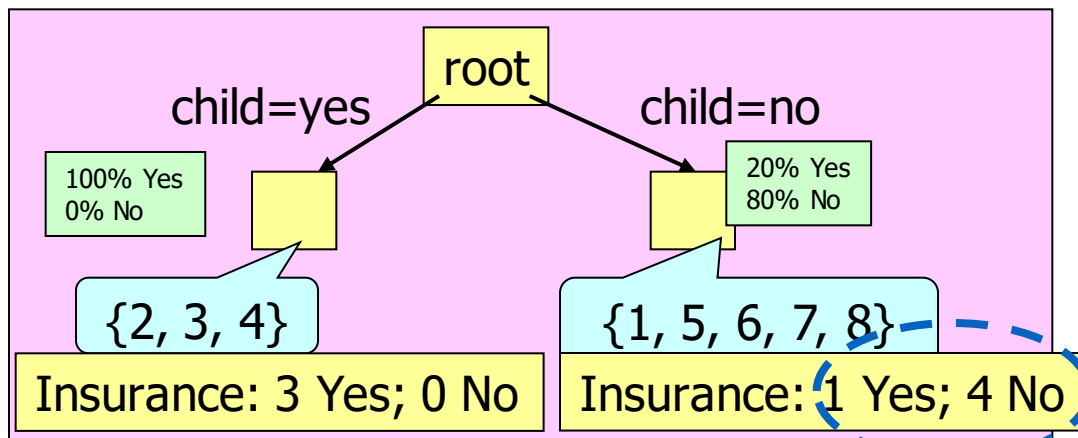
For attribute Gender,

$$\text{Gain}(\text{Gender}, T) = 0.1887$$

For attribute Income,

$$\text{Gain}(\text{Income}, T) = 0.3113$$

Gender	Income	Child	Insurance
female	high	no	yes
male	high	yes	yes
male	low	yes	yes
male	low	yes	yes
female	low	no	no
female	low	no	no
female	low	no	no
male	low	no	no



	Gender	Income	Child	Insurance
1	female	high	no	yes
2	male	high	yes	yes
3	male	low	yes	yes
4	male	low	yes	yes
5	female	low	no	no
6	female	low	no	no
7	female	low	no	no
8	male	low	no	no

$$\text{Info}(T) = -\frac{1}{2} \log \frac{1}{2} - \frac{1}{2} \log \frac{1}{2} = 1$$

For attribute Child,

$$\text{Info}(T_{\text{yes}}) = -1 \log 1 - 0 \log 0 = 0$$

$$\text{Info}(T_{\text{no}}) = -\frac{1}{5} \log \frac{1}{5} - \frac{4}{5} \log \frac{4}{5} = 0.7219$$

$$\text{Info}(\text{Child}, T) = \frac{3}{8} \times \text{Info}(T_{\text{yes}}) + \frac{5}{8} \times \text{Info}(T_{\text{no}}) = 0.4512$$

$$\text{Gain}(\text{Child}, T) = \text{Info}(T) - \text{Info}(\text{Child}, T) = 1 - 0.4512 = 0.5488$$

For attribute Gender,

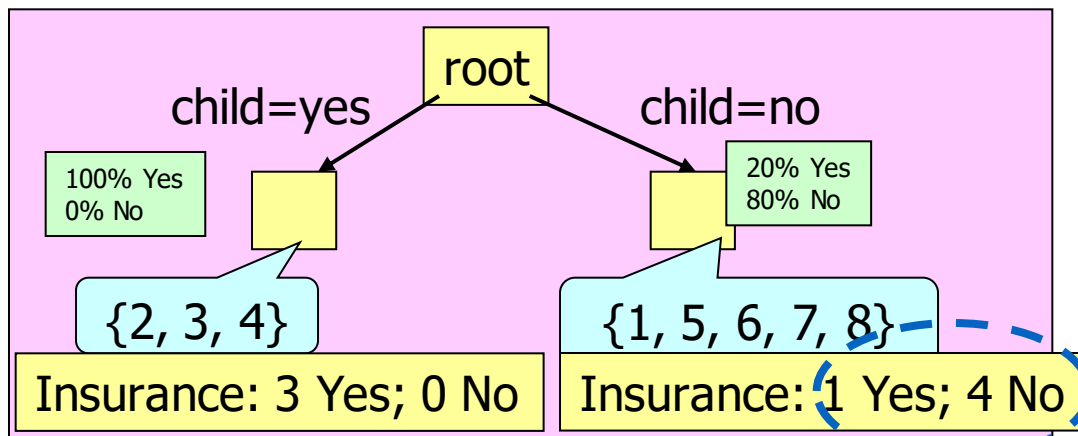
$$\text{Gain}(\text{Gender}, T) = 0.1887$$

For attribute Income,

$$\text{Gain}(\text{Income}, T) = 0.3113$$

For attribute Child,

$$\text{Gain}(\text{Child}, T) = 0.5488$$



	Gender	Income	Child	Insurance
1	female	high	no	yes
2	male	high	yes	yes
3	male	low	yes	yes
4	male	low	yes	yes
5	female	low	no	no
6	female	low	no	no
7	female	low	no	no
8	male	low	no	no

$$\text{Info}(T) = - \frac{1}{5} \log \frac{1}{5} - \frac{4}{5} \log \frac{4}{5} = 0.7219$$

For attribute Gender,

$$\text{Info}(T_{\text{female}}) = - \frac{1}{4} \log \frac{1}{4} - \frac{3}{4} \log \frac{3}{4} = 0.8113$$

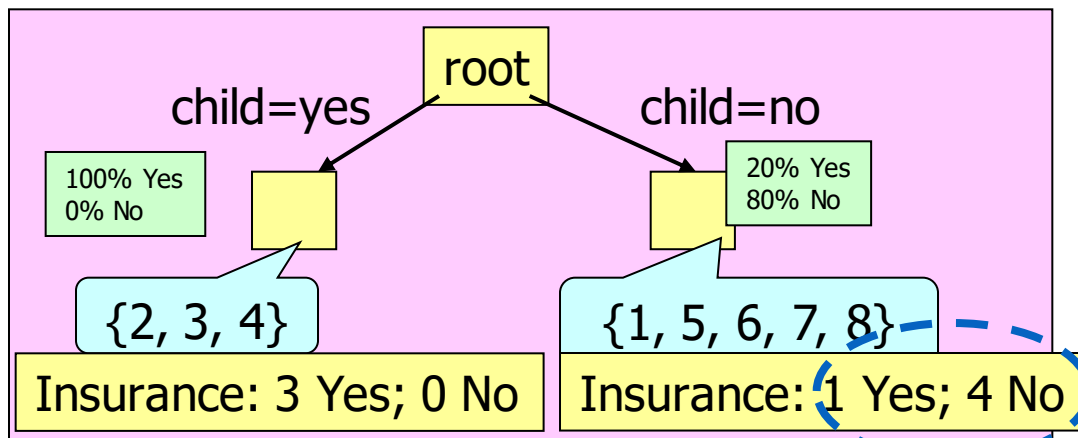
$$\text{Info}(T_{\text{male}}) = - 0 \log 0 - 1 \log 1 = 0$$

$$\text{Info}(\text{Gender}, T) = \frac{4}{5} \times \text{Info}(T_{\text{female}}) + \frac{1}{5} \times \text{Info}(T_{\text{male}}) = 0.6490$$

$$\text{Gain}(\text{Gender}, T) = \text{Info}(T) - \text{Info}(\text{Gender}, T) = 0.7219 - 0.6490 = 0.0729$$

For attribute Gender,

$$\text{Gain}(\text{Gender}, T) = 0.0729$$



	Gender	Income	Child	Insurance
1	female	high	no	yes
2	male	high	yes	yes
3	male	low	yes	yes
4	male	low	yes	yes
5	female	low	no	no
6	female	low	no	no
7	female	low	no	no
8	male	low	no	no

$$\text{Info}(T) = -1/5 \log 1/5 - 4/5 \log 4/5 = 0.7219$$

For attribute Income,

$$\text{Info}(T_{\text{high}}) = -1 \log 1 - 0 \log 0 = 0$$

$$\text{Info}(T_{\text{low}}) = -0 \log 0 - 1 \log 1 = 0$$

$$\text{Info}(\text{Income}, T) = 1/5 \times \text{Info}(T_{\text{high}}) + 4/5 \times \text{Info}(T_{\text{low}}) = 0$$

$$\text{Gain}(\text{Income}, T) = \text{Info}(T) - \text{Info}(\text{Income}, T) = 0.7219 - 0 = 0.7219$$

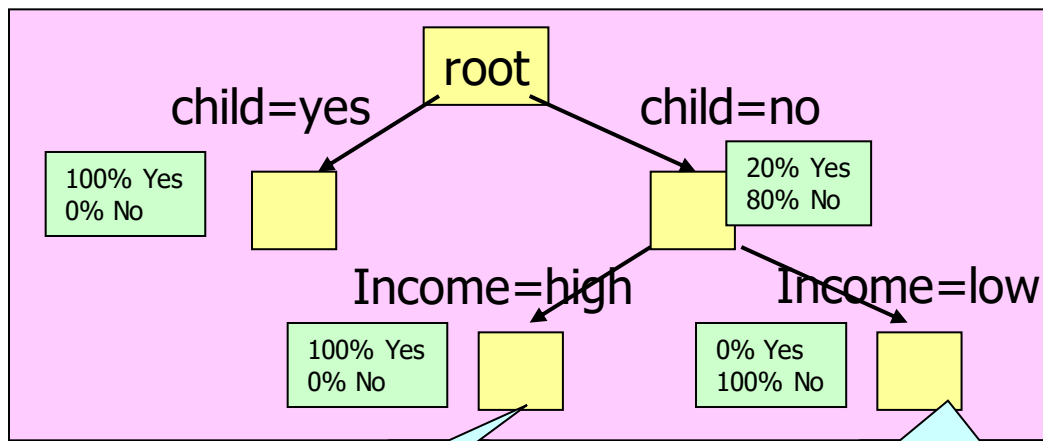
For attribute Gender,

$$\text{Gain}(\text{Gender}, T) = 0.0729$$

For attribute Income,

$$\text{Gain}(\text{Income}, T) = 0.7219$$





	Gender	Income	Child	Insurance
1	female	high	no	yes
2	male	high	yes	yes
3	male	low	yes	yes
4	male	low	yes	yes
5	female	low	no	no
6	female	low	no	no
7	female	low	no	no
8	male	low	no	no

{1}

Insurance: 1 Yes; 0 No

{5, 6, 7, 8}

Insurance: 0 Yes; 4 No

Decision tree

Suppose there is a new person.

Gender	Income	Child	Insurance
male	high	no	?

# Decision Tree Examples

- ID3

- Impurity Measurement

- $\text{Gain}(A, T) = \text{Info}(T) - \text{Info}(A, T)$

- C4.5

- Impurity Measurement

- $\text{Gain}(A, T) = (\text{Info}(T) - \text{Info}(A, T)) / \text{SplitInfo}(A)$ , where  $\text{SplitInfo}(A) = -\sum_{v \in A} p(v) \log p(v)$

- CART

- Impurity Measurement

- Gini:  $I(P) = 1 - \sum_j p_j^2$

# Entropy

$$\text{Info}(T) = -\frac{1}{2} \log \frac{1}{2} - \frac{1}{2} \log \frac{1}{2} \\ = 1$$

For attribute Gender,

$$\text{Info}(T_{\text{female}}) = -\frac{3}{4} \log \frac{3}{4} - \frac{1}{4} \log \frac{1}{4} = 0.8113$$

$$\text{Info}(T_{\text{male}}) = -\frac{3}{4} \log \frac{3}{4} - \frac{1}{4} \log \frac{1}{4} = 0.8113$$

$$\text{Info}(\text{Gender}, T) = \frac{1}{2} \times \text{Info}(T_{\text{female}}) + \frac{1}{2} \times \text{Info}(T_{\text{male}}) = 0.8113$$

$$\text{SplitInfo}(\text{Gender}) = -\frac{1}{2} \log \frac{1}{2} - \frac{1}{2} \log \frac{1}{2} = 1$$

$$\text{Gain}(\text{Gender}, T) = (\text{Info}(T) - \text{Info}(\text{Gender}, T)) / \text{SplitInfo}(\text{Gender}) = (1 - 0.8113) / 1 = 0.1887$$

For attribute Gender,

$$\text{Gain}(\text{Gender}, T) = 0.1887$$

Gender	Income	Child	Insurance
female	high	no	yes
male	high	yes	yes
male	low	yes	yes
male	low	yes	yes
female	low	no	no
female	low	no	no
female	low	no	no
male	low	no	no

# Entropy

$$\text{Info}(T) = -\frac{1}{2} \log \frac{1}{2} - \frac{1}{2} \log \frac{1}{2} = 1$$

For attribute Income,

$$\text{Info}(T_{\text{high}}) = -1 \log 1 - 0 \log 0 = 0$$

$$\text{Info}(T_{\text{low}}) = -\frac{1}{3} \log \frac{1}{3} - \frac{2}{3} \log \frac{2}{3} = 0.9183$$

$$\text{Info}(\text{Income}, T) = \frac{1}{4} \times \text{Info}(T_{\text{high}}) + \frac{3}{4} \times \text{Info}(T_{\text{low}}) = 0.6887$$

$$\text{SplitInfo}(\text{Income}) = -\frac{2}{8} \log \frac{2}{8} - \frac{6}{8} \log \frac{6}{8} = 0.8113$$

$$\text{Gain}(\text{Income}, T) = (\text{Info}(T) - \text{Info}(\text{Income}, T)) / \text{SplitInfo}(\text{Income}) = (1 - 0.6887) / 0.8113 = 0.3837$$

For attribute Gender,

$$\text{Gain}(\text{Gender}, T) = 0.1887$$

For attribute Income,

$$\text{Gain}(\text{Income}, T) = 0.3837$$

For attribute Child,

$$\text{Gain}(\text{Child}, T) = ?$$

Gender	Income	Child	Insurance
female	high	no	yes
male	high	yes	yes
male	low	yes	yes
male	low	yes	yes
female	low	no	no
female	low	no	no
female	low	no	no
male	low	no	no

# Decision Tree Examples

- ID3

- Impurity Measurement

- $\text{Gain}(A, T) = \text{Info}(T) - \text{Info}(A, T)$

- C4.5

- Impurity Measurement

- $\text{Gain}(A, T) = (\text{Info}(T) - \text{Info}(A, T)) / \text{SplitInfo}(A)$ , where  $\text{SplitInfo}(A) = -\sum_{v \in A} p(v) \log p(v)$

- CART

- Impurity Measurement

- Gini:  $I(P) = 1 - \sum_j p_j^2$

# Gini

$$\begin{aligned}\text{Info}(T) &= 1 - \left(\frac{1}{2}\right)^2 - \left(\frac{1}{2}\right)^2 \\ &= \frac{1}{2}\end{aligned}$$

For attribute Gender,

$$\text{Info}(T_{\text{female}}) = 1 - \left(\frac{3}{4}\right)^2 - \left(\frac{1}{4}\right)^2 = 0.375$$

$$\text{Info}(T_{\text{male}}) = 1 - \left(\frac{3}{4}\right)^2 - \left(\frac{1}{4}\right)^2 = 0.375$$

$$\text{Info}(\text{Gender}, T) = \frac{1}{2} \times \text{Info}(T_{\text{female}}) + \frac{1}{2} \times \text{Info}(T_{\text{male}}) = 0.375$$

$$\text{Gain}(\text{Gender}, T) = \text{Info}(T) - \text{Info}(\text{Gender}, T) = \frac{1}{2} - 0.375 = 0.125$$

For attribute Gender,

$$\text{Gain}(\text{Gender}, T) = 0.125$$

Gender	Income	Child	Insurance
female	high	no	yes
male	high	yes	yes
male	low	yes	yes
male	low	yes	yes
female	low	no	no
female	low	no	no
female	low	no	no
male	low	no	no

# Gini

$$\text{Info}(T) = 1 - \left(\frac{1}{2}\right)^2 - \left(\frac{1}{2}\right)^2 \\ = \frac{1}{2}$$

For attribute Income,

$$\text{Info}(T_{\text{high}}) = 1 - 1^2 - 0^2 = 0$$

$$\text{Info}(T_{\text{low}}) = 1 - \left(\frac{1}{3}\right)^2 - \left(\frac{2}{3}\right)^2 = 0.444$$

$$\text{Info}(\text{Income}, T) = \frac{1}{4} \times \text{Info}(T_{\text{high}}) + \frac{3}{4} \times \text{Info}(T_{\text{low}}) = 0.333$$

$$\text{Gain}(\text{Income}, T) = \text{Info}(T) - \text{Info}(\text{Income}, T) = \frac{1}{2} - 0.333 = 0.167$$

For attribute Gender,

$$\text{Gain}(\text{Gender}, T) = 0.125$$

For attribute Income,

$$\text{Gain}(\text{Income}, T) = 0.167$$

For attribute Child,

$$\text{Gain}(\text{Child}, T) = ?$$

Gender	Income	Child	Insurance
female	high	no	yes
male	high	yes	yes
male	low	yes	yes
male	low	yes	yes
female	low	no	no
female	low	no	no
female	low	no	no
male	low	no	no

# Other Issues

- Data Fragmentation
- Expressiveness



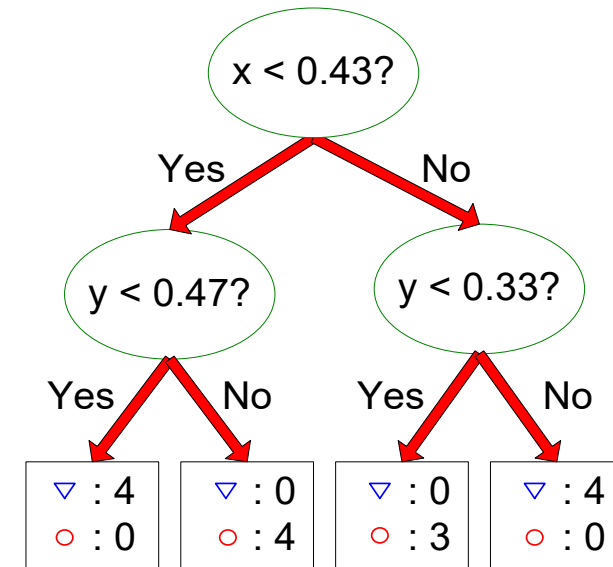
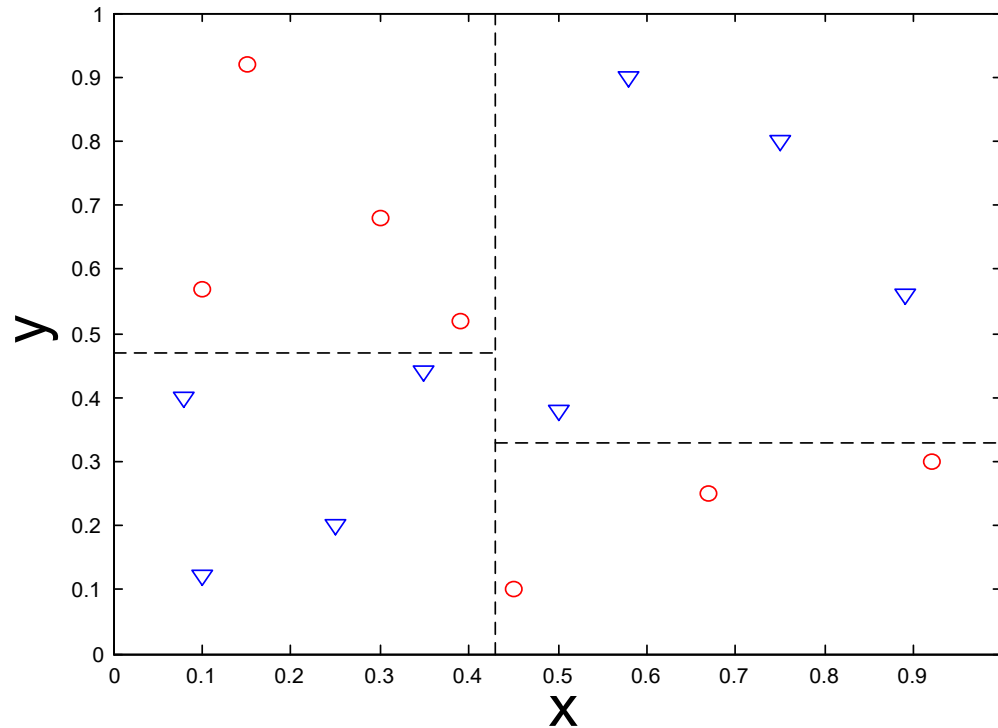
# Data Fragmentation

- Number of instances gets smaller as you traverse down the tree
- Number of instances at the leaf nodes could be **too small** to make any **statistically significant decision**
- You can introduce a lower bound on the number of items per leaf node in the stopping criterion.

# Expressiveness

- A classifier defines a **function** that discriminates between two (or more) classes.
- The **expressiveness** of a classifier is the **class of functions** that it can model, and the kind of data that it can **separate**
  - When we have **discrete** (or binary) values, we are interested in the class of **boolean functions** that can be modeled
  - If the data-points are real vectors we talk about the **decision boundary** that the classifier can model

# Decision Boundary

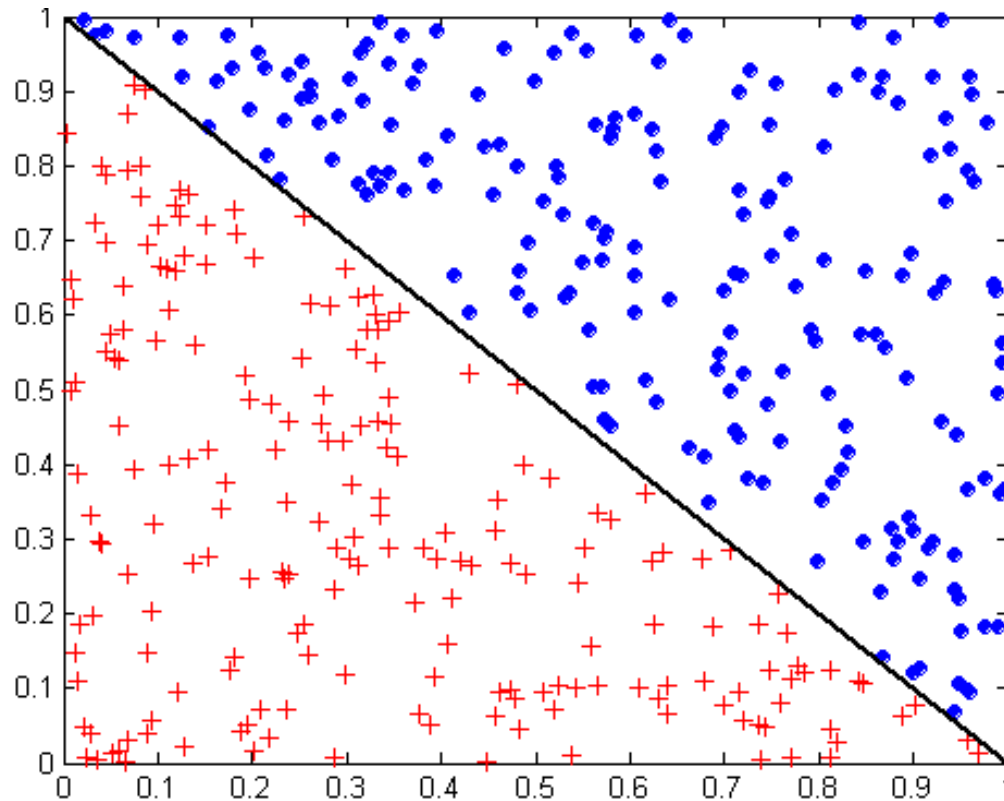


- Border line between two neighboring regions of different classes is known as **decision boundary**
- Decision boundary is **parallel to axes** because test condition involves a single attribute at-a-time

# Expressiveness

- Decision tree provides **expressive** representation for learning discrete-valued function
  - But they do not generalize well to certain types of Boolean functions
    - Example: **parity function**:
      - Class = 1 if there is an **even** number of Boolean attributes with truth value = True
      - Class = 0 if there is an **odd** number of Boolean attributes with truth value = True
    - For accurate modeling, must have a complete tree
- Less expressive for modeling continuous variables
  - Particularly when test condition involves only a single attribute at-a-time

# Oblique Decision Trees



- Test condition may involve multiple attributes
- More expressive representation
- Finding optimal test condition is computationally expensive

