



# Classification

Concepts and Techniques





# Classification: Definition

---

Data is a collection of records (**training set**)

Each **record** is characterized by a **tuple**  $(x, y)$ , where  **$x$**  is the **attribute** set and  **$y$**  is the **class** label

**$x$** : attribute, predictor, independent variable → input

**$y$** : class, response, dependent variable → output

Task: Learn a model that maps each attribute set  **$x$**  into one of the predefined class labels

***A classification model is an abstract representation of the relationship between the attribute set and the class label***



# Examples of Classification Task

## Binary classification

problems, in which each data instance can be categorized into one of two classes

## Multiclass classification

problems, in which each data instance can be categorized into one of multiple classes

Task	Attribute set ( $x$ )	Class label ( $y$ )	<i>Classification type</i>
Categorizing email messages	Features extracted from email message header and content	spam or non-spam	Binary
Identifying tumor cells	Features extracted from x-rays or MRI scans	malignant or benign cells	Binary
Cataloging galaxies	Features extracted from telescope images	Elliptical, spiral, or irregular-shaped galaxies	Multiclass

# Classification Models

---



## Roles of Classification Models:

- **Predictive Model:** Classifies previously unlabeled instances.
- **Descriptive Model:** Identifies distinguishing characteristics between different classes.

Example: Medical diagnosis, where justifying predictions is crucial.

## Key Requirements:

- *Accuracy:* Predict outcomes correctly.
- *Efficiency:* fast response times.



# General Framework for Classification

## Key Concepts:

- **Classifier:** A model used to perform classification.
- **Training Set:** A set of instances with attribute values and class labels used to build the model.
- **Learning Algorithm:** The systematic approach used to create the model from the training set. This process is called **Induction** (learning/building a model).
- **Deduction:** Applying the learned model to new, unseen test instances to predict their class labels.

## Steps in Classification:

- **Induction (training):** Learn the model from training data.
- **Deduction (inference):** Apply the model to classify new instances.

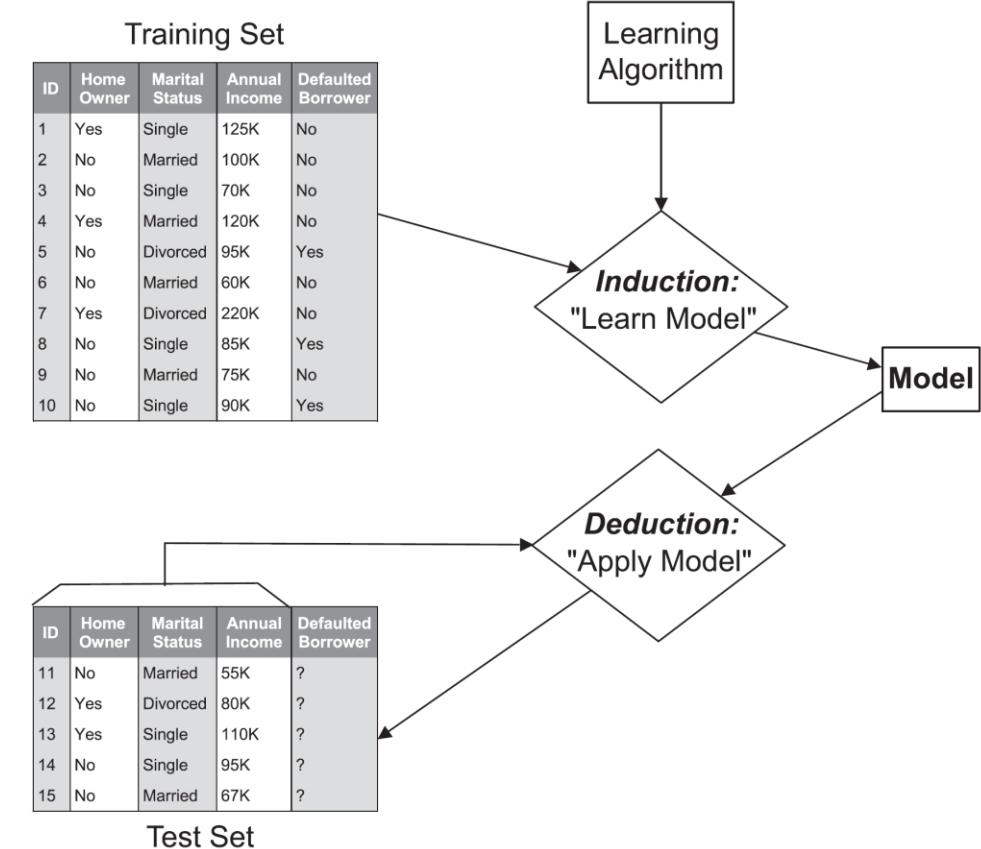


Figure 3.3. General framework for building a classification model.

# Classification performance



## Pivotal concepts

**Independence of Training and Test Sets:** Ensures the model can predict class labels for unseen instances.

**Generalization Performance:** A model with good generalization can accurately predict labels for new, unseen data.

**Model Performance Evaluation:** Compare **predicted labels** against **true labels** to evaluate accuracy.

**Model Performance Evaluation:** Compare **predicted labels** against **true labels** to evaluate accuracy.

- **Confusion Matrix**

		P
R	T <sub>P</sub>	F <sub>N</sub>
	F <sub>P</sub>	T <sub>N</sub>

- **Accuracy:** number of correct over total
- **Error Rate:** number of wrong over total

# Classification Techniques

---



## Base Classifiers

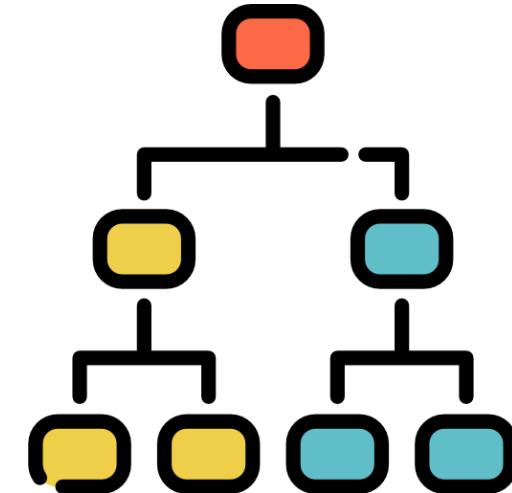
- Decision Tree based Methods
- Nearest-neighbor
- Naïve Bayes and Bayesian Belief Networks
- Support Vector Machines
- Neural Networks, Deep Neural Nets

## Ensemble Classifiers

- Boosting, Bagging, Random Forests



# Decision Tree





# Decision Tree

A **decision tree** is a **supervised** machine learning algorithm used for both **classification** and **regression** tasks. It models decisions and their possible consequences in a tree-like structure, where:

**Root node:** node with no incoming link

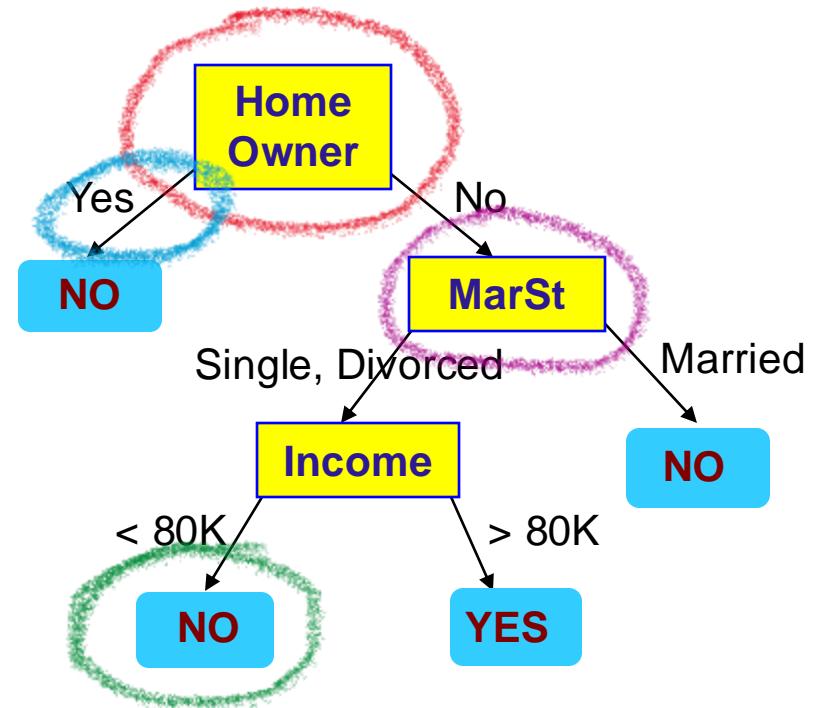
**Nodes:** represent decisions or tests on attributes (features).

**Edges:** represent the outcome of the test, leading to either more nodes or leaf nodes.

**Leaf or terminal nodes:** represent the final outcome or class label in classification tasks.

The tree is built by splitting the data based on feature values that best separate the classes or predict the target variable.

The goal is to create a tree that accurately classifies new instances by following the decision rules in the tree structure.

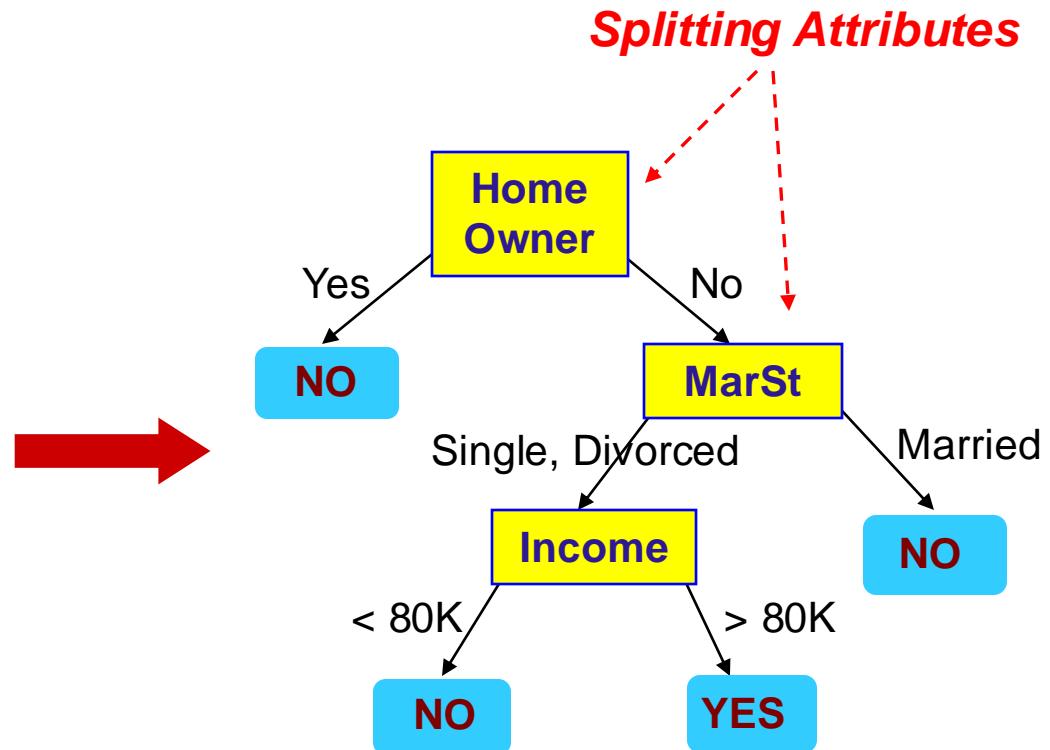




# Example of a Decision Tree

ID	Home Owner	Marital Status	Annual Income	Defaulted Borrower
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes

Training Data



Model: Decision Tree

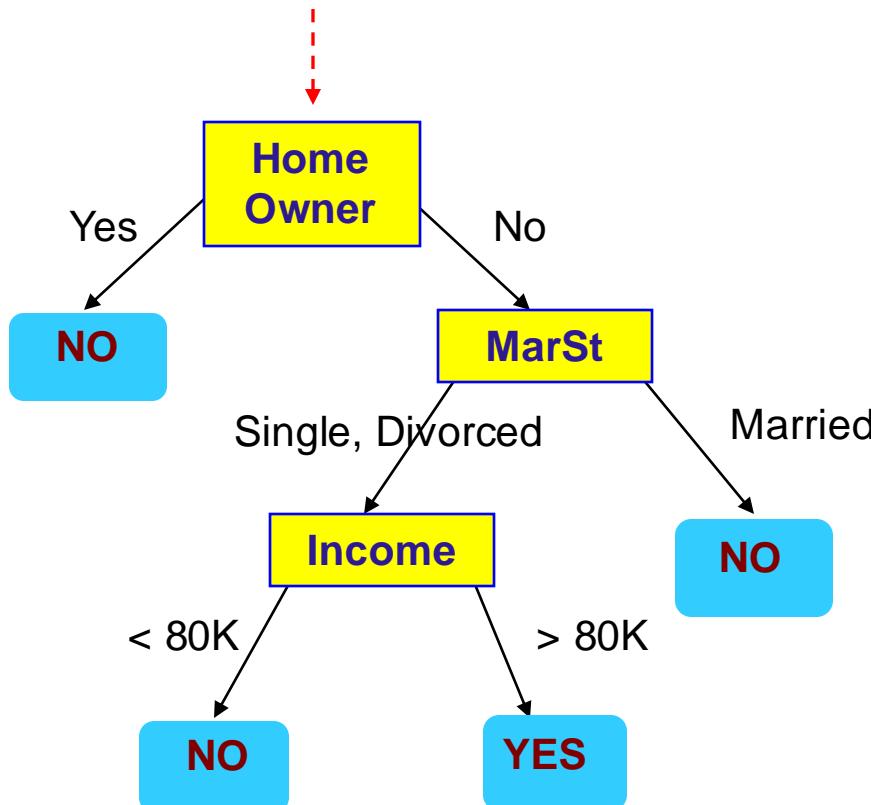


# Apply Model to Test Data

## Test Data

Home Owner	Marital Status	Annual Income	Defaulted Borrower
No	Married	80K	?

Start from the root of tree.

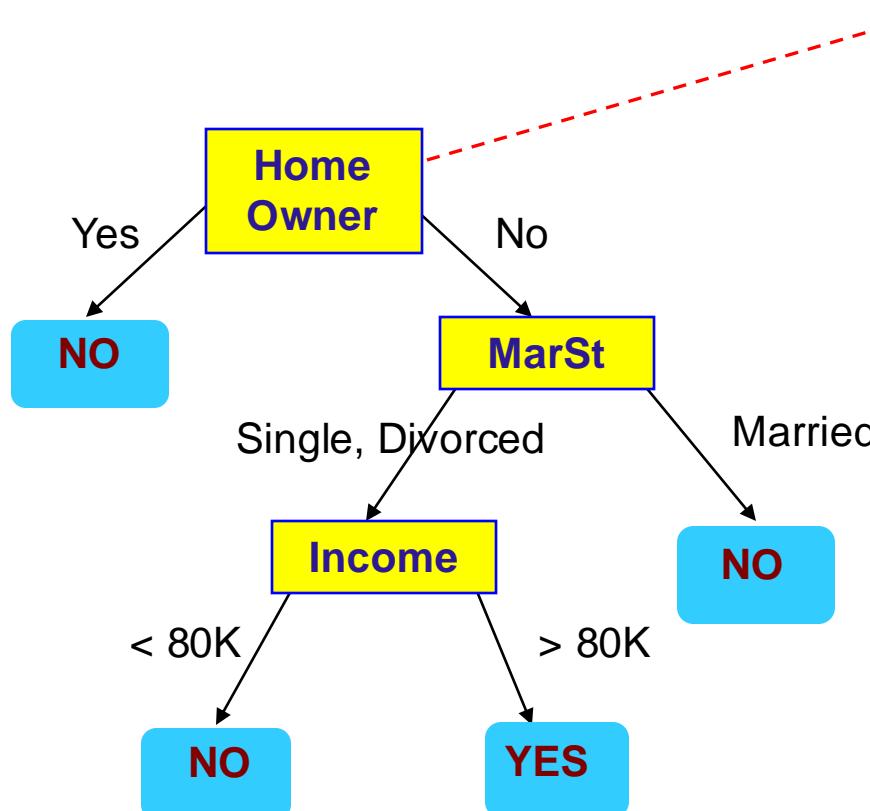


# Apply Model to Test Data



Test Data

Home Owner	Marital Status	Annual Income	Defaulted Borrower
No	Married	80K	?

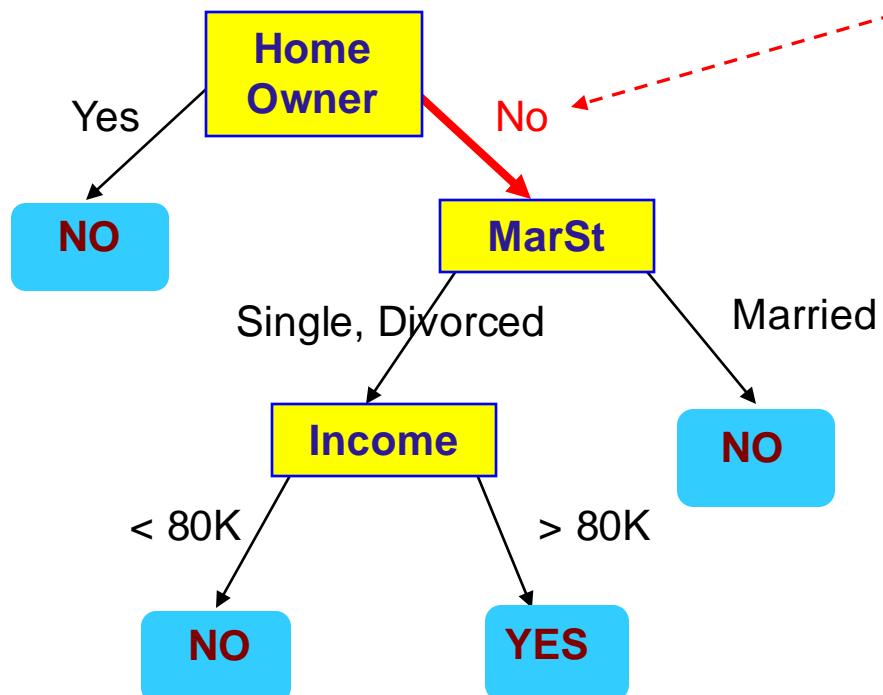


# Apply Model to Test Data



Test Data

Home Owner	Marital Status	Annual Income	Defaulted Borrower
No	Married	80K	?

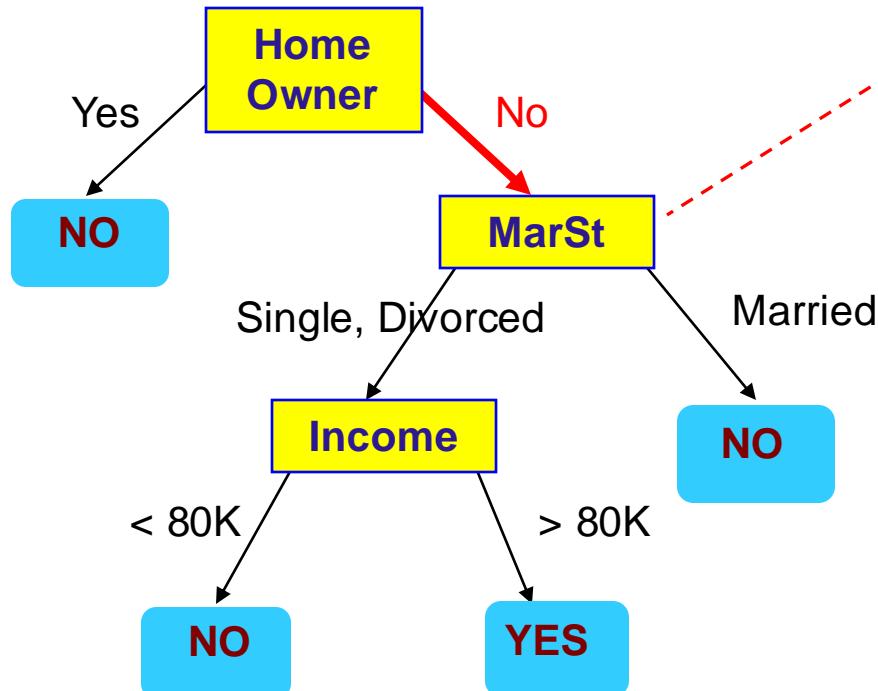


# Apply Model to Test Data



Test Data

Home Owner	Marital Status	Annual Income	Defaulted Borrower
No	Married	80K	?

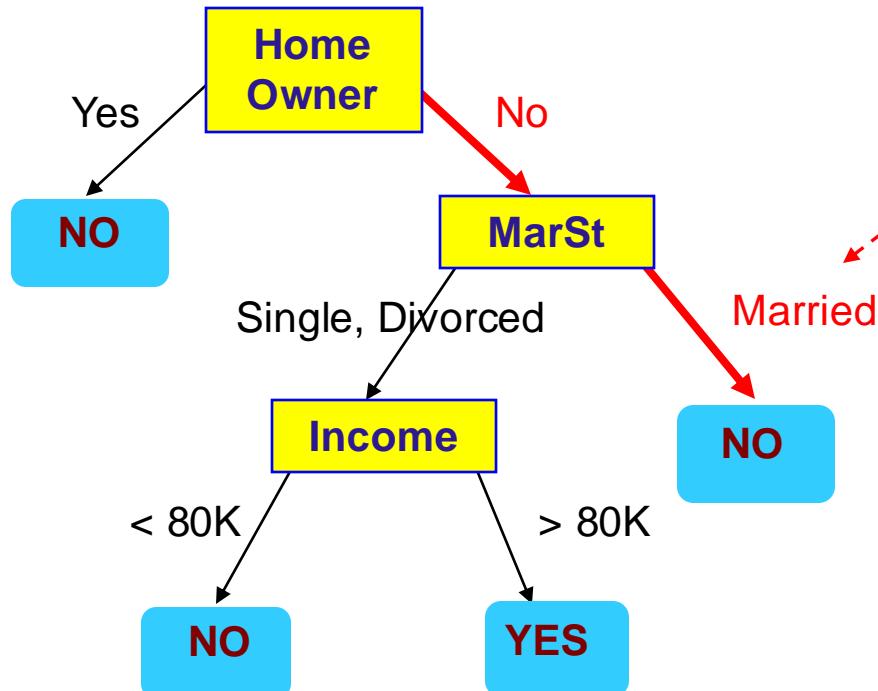




# Apply Model to Test Data

Test Data

Home Owner	Marital Status	Annual Income	Defaulted Borrower
No	Married	80K	?

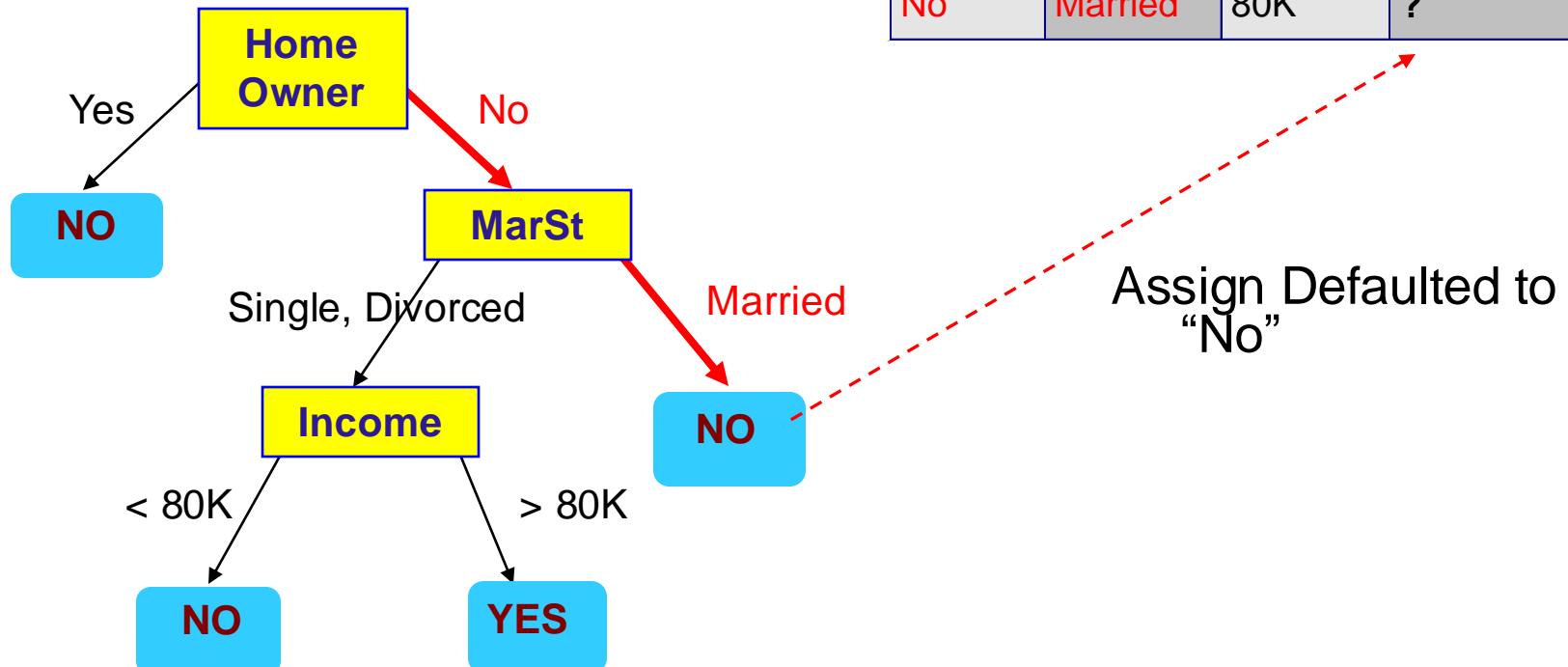


# Apply Model to Test Data



## Test Data

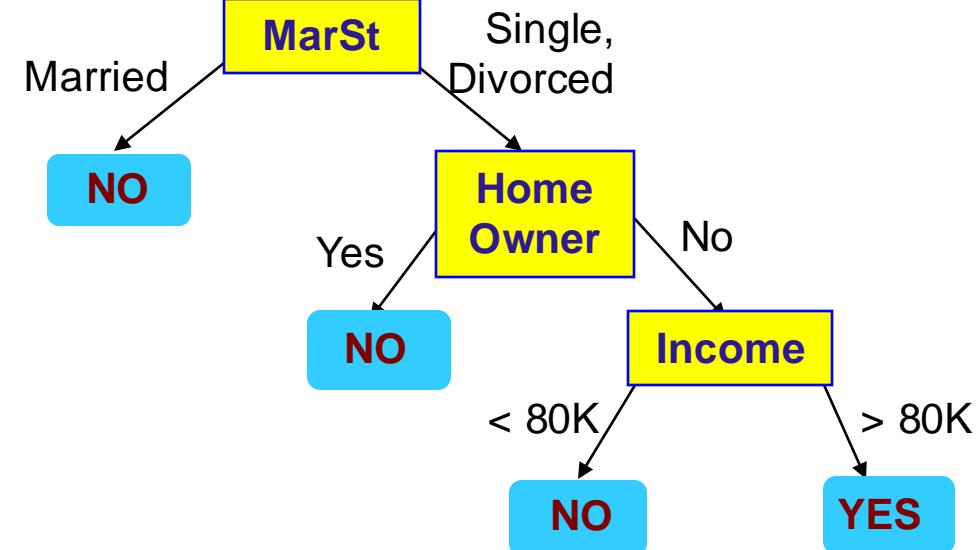
Home Owner	Marital Status	Annual Income	Defaulted Borrower
No	Married	80K	?



# Another Example of Decision Tree



ID	Home Owner	Marital Status	Annual Income	Defaulted Borrower	class
1	Yes	Single	125K	No	
2	No	Married	100K	No	
3	No	Single	70K	No	
4	Yes	Married	120K	No	
5	No	Divorced	95K	Yes	
6	No	Married	60K	No	
7	Yes	Divorced	220K	No	
8	No	Single	85K	Yes	
9	No	Married	75K	No	
10	No	Single	90K	Yes	



There could be more than one tree that fits the same data!

# Decision Tree Induction

---



Many Algorithms:

- Hunt's Algorithm (one of the earliest)
- CART
- ID3, C4.5
- SLIQ, SPRINT



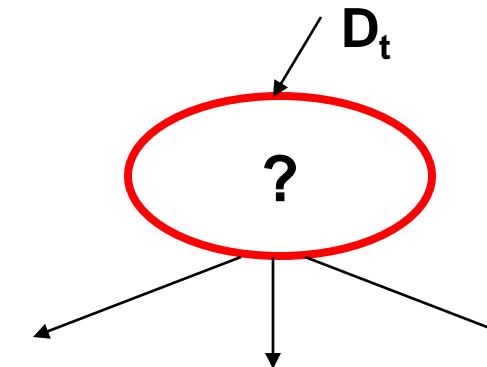
# General Structure of Hunt's Algorithm

Let  $D_t$  be the set of training records that reach a node  $t$

General Procedure:

- If  $D_t$  contains records that belong to the **same class**  $y_t$ , then  $t$  is a leaf node labeled as  $y_t$
- If  $D_t$  contains records that belong to **more than one class**, use an attribute test to split the data into smaller subsets. Recursively apply the procedure to each subset.

ID	Home Owner	Marital Status	Annual Income	Defaulted Borrower
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes



# Hunt's Algorithm



Defaulted

(7,3)  
(no, yes)  
(a)

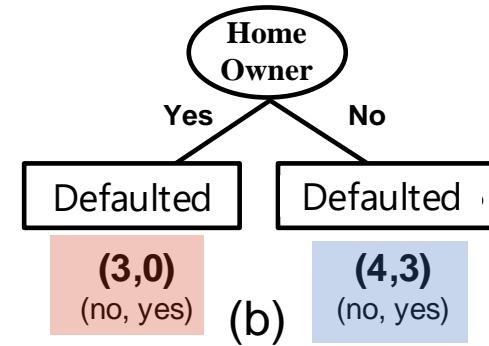
ID	Home Owner	Marital Status	Annual Income	Defaulted Borrower
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes



# Hunt's Algorithm

Defaulted

(7,3)  
(no, yes)  
(a)



ID	Home Owner	Marital Status	Annual Income	Defaulted Borrower
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes



# Hunt's Algorithm

Defaulted

(7,3)  
(no, yes)

(a)

Home  
Owner

Yes      No

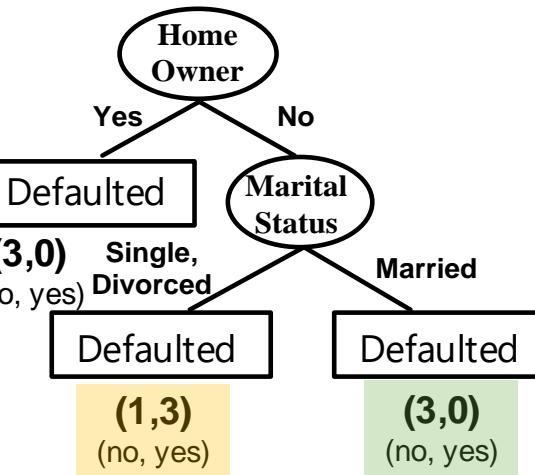
Defaulted

Defaulted

(3,0)  
(no, yes)

(4,3)  
(no, yes)

(b)



(c)

ID	Home Owner	Marital Status	Annual Income	Defaulted Borrower
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes

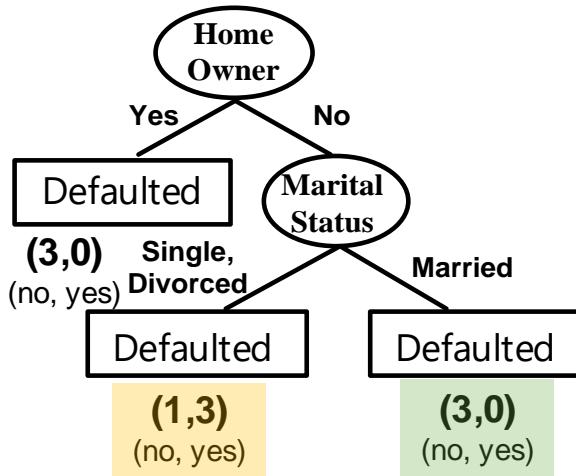


# Hunt's Algorithm

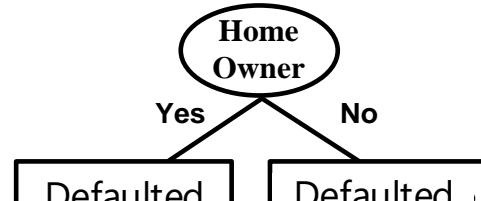
Defaulted

(7,3)  
(no, yes)

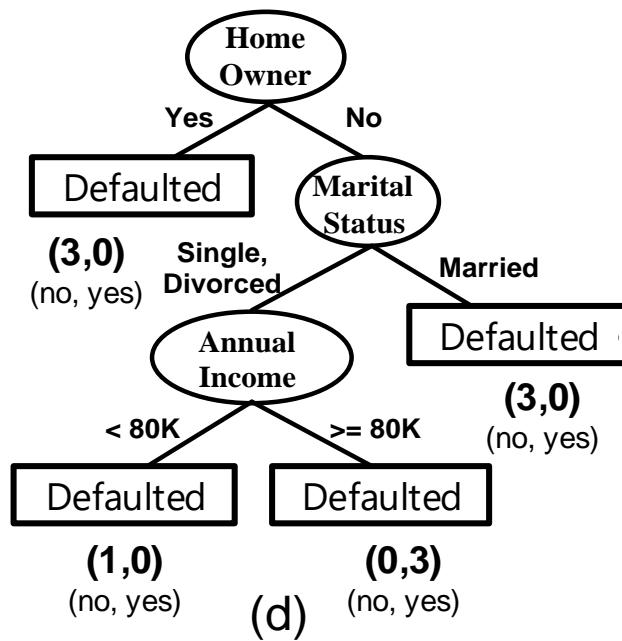
(a)



(c)



(b) (3,0)  
(no, yes) (4,3)  
(no, yes)



(d) (1,0)  
(no, yes) (0,3)  
(no, yes)

ID	Home Owner	Marital Status	Annual Income	Defaulted Borrower
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes

# Design Issues of Decision Tree Induction

---



How should training records be **split**?

- Method for expressing **test condition** depending on *attribute types*
- Measure for **evaluating the goodness** of a **test condition**

How should the splitting procedure **stop**?

- Stop splitting if all the records belong to the **same class** or have **identical attribute** values
- Early termination

# Methods for Expressing Test Conditions

---



Depends on attribute types

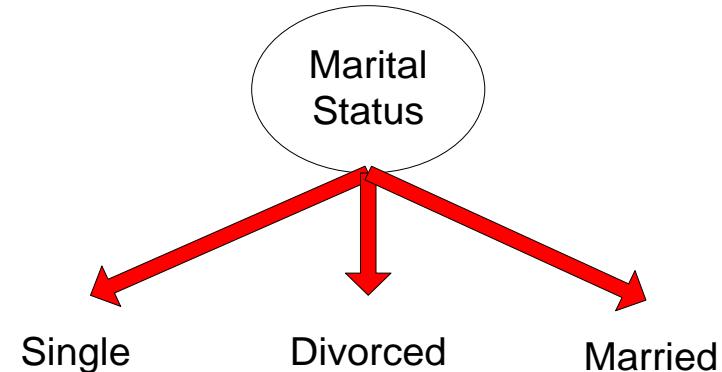
- Binary
- Nominal
- Ordinal
- Continuous

# Test Condition for Nominal Attributes



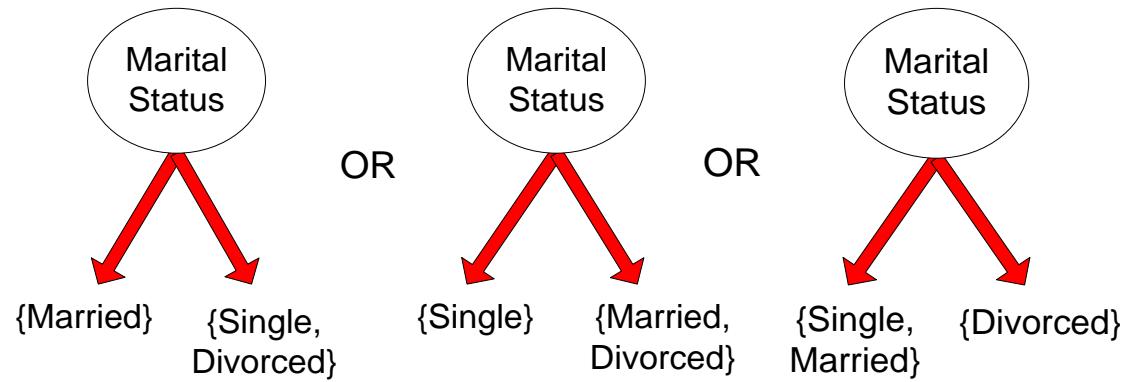
## Multi-way split:

- Use as many partitions as distinct values.



## Binary split:

- Divides values into two subsets

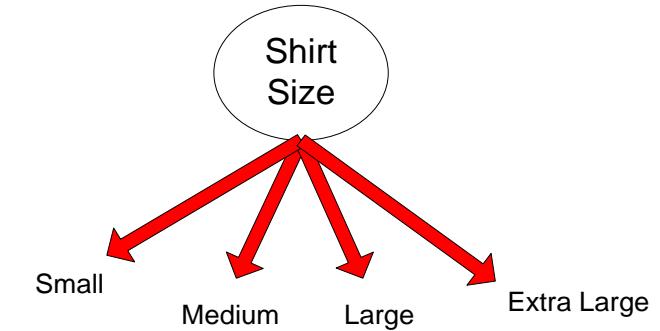


# Test Condition for Ordinal Attributes



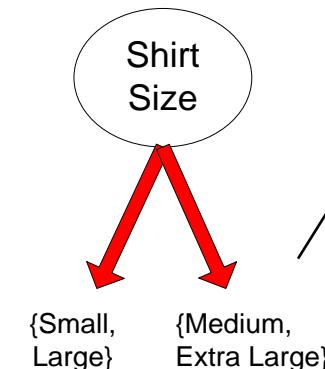
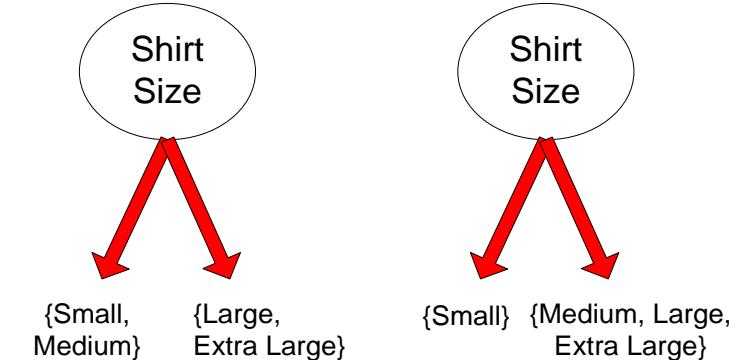
## Multi-way split:

- Use as many partitions as distinct values



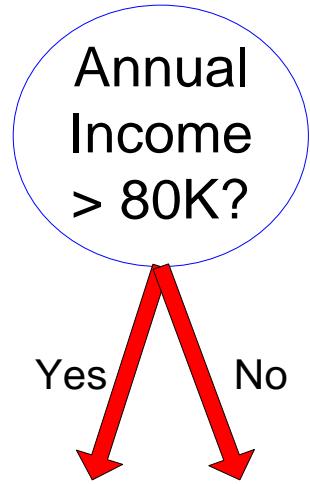
## Binary split:

- Divides values into two subsets
- Preserve order property among attribute values

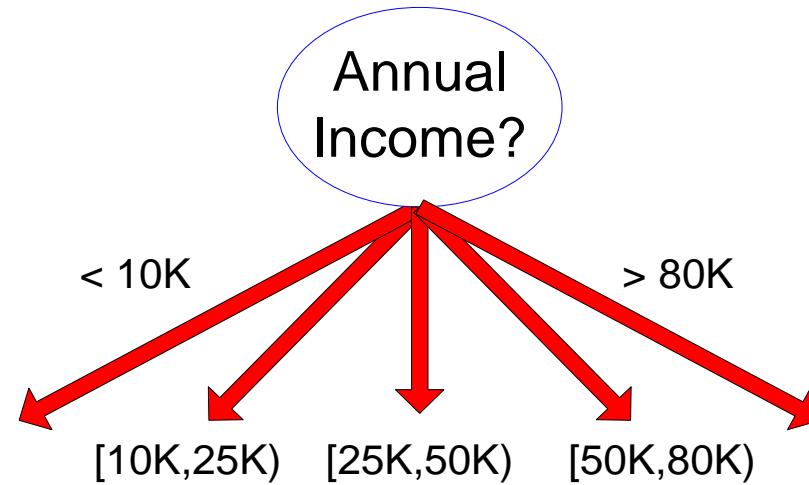


This grouping  
violates order  
property !

# Test Condition for Continuous Attributes



(i) Binary split



(ii) Multi-way split

# Splitting Based on Continuous Attributes



---

Different ways of handling

- **Discretization** to form an ordinal categorical attribute

Ranges can be found by equal interval bucketing, equal frequency bucketing (percentiles), or clustering.

- Static – discretize once at the beginning
- Dynamic – repeat at each node

- **Binary Decision:**  $(A < v)$  or  $(A \geq v)$

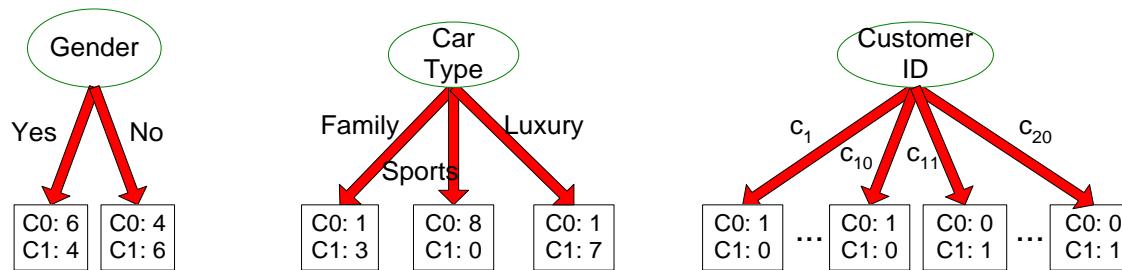
- consider all possible splits and finds the best cut
- can be more compute intensive



# How to determine the Best Split

**Goodness of an attribute test condition:** many measures try to give preference to attribute test conditions that partition the training instances into **purer subsets** in the child nodes, which mostly have the same class labels

**Before Splitting:** 10 records of **class 0**,  
10 records of **class 1**



*Which test condition is the best?*

Customer Id	Gender	Car Type	Shirt Size	Class
1	M	Family	Small	C0
2	M	Sports	Medium	C0
3	M	Sports	Medium	C0
4	M	Sports	Large	C0
5	M	Sports	Extra Large	C0
6	M	Sports	Extra Large	C0
7	F	Sports	Small	C0
8	F	Sports	Small	C0
9	F	Sports	Medium	C0
10	F	Luxury	Large	C0
11	M	Family	Large	C1
12	M	Family	Extra Large	C1
13	M	Family	Medium	C1
14	M	Luxury	Extra Large	C1
15	F	Luxury	Small	C1
16	F	Luxury	Small	C1
17	F	Luxury	Medium	C1
18	F	Luxury	Medium	C1
19	F	Luxury	Medium	C1
20	F	Luxury	Large	C1

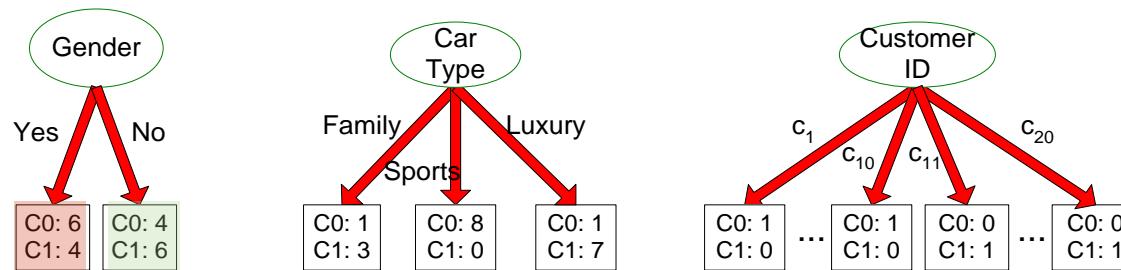
**NOTE:** Larger trees are less desirable as they are more susceptible to model **overfitting**



# How to determine the Best Split

**Goodness of an attribute test condition:** many measures try to give preference to attribute test conditions that partition the training instances into **purer subsets** in the child nodes, which mostly have the same class labels

**Before Splitting:** 10 records of **class 0**,  
10 records of **class 1**



*Which test condition is the best?*

Customer Id	Gender	Car Type	Shirt Size	Class
1	M	Family	Small	C0
2		Sports	Medium	C0
3		Sports	Medium	C0
4		Sports	Large	C0
5		Sports	Extra Large	C0
6		Sports	Extra Large	C0
7	F	Sports	Small	C0
8	F	Sports	Small	C0
9	F	Sports	Medium	C0
10	F	Luxury	Large	C0
11	M	Family	Large	C1
12	M	Family	Extra Large	C1
13	M	Family	Medium	C1
14	M	Luxury	Extra Large	C1
15	F	Luxury	Small	C1
16	F	Luxury	Small	C1
17	F	Luxury	Medium	C1
18	F	Luxury	Medium	C1
19	F	Luxury	Medium	C1
20	F	Luxury	Large	C1

**NOTE:** Larger trees are less desirable as they are more susceptible to model **overfitting**



# How to determine the Best Split

Greedy approach: Nodes with **purer** class distribution are preferred

We need a measure of node impurity !

C0: 5  
C1: 5

**High degree of impurity**

C0: 9  
C1: 1

**Low degree of impurity**

Let's define the **RELATIVE FREQUENCY**  $p_i(t)$

$$C_0 = 5$$

$$P_0(t) = \frac{5}{10}$$

$$C_1 = 5$$

$$P_1(t) = \frac{5}{10}$$

$$C_0 = 1$$

$$P_0(t) = \frac{1}{10}$$

$$C_1 = 9$$

$$P_1(t) = \frac{9}{10}$$



# Measures of Single Node Impurity

## Gini Index

$$Gini\ Index = 1 - \sum_{i=0}^{c-1} p_i(t)^2$$

Where  $p_i(t)$  is the frequency of class  $i$  at node  $t$ , and  $c$  is the total number of classes

## Entropy

$$Entropy = - \sum_{i=0}^{c-1} p_i(t) \log_2 p_i(t)$$

## Misclassification error

$$\begin{aligned} 2^5 &= 32 \\ \log_2(32) &= 5 \end{aligned}$$

$$Classification\ error = 1 - \max[p_i(t)]$$



# Compute Single Node Impurity

C0: 5  
C1: 5

High degree of impurity

$$c_0 = 5 \quad P_0(t) = \frac{5}{10}$$
$$c_1 = 5 \quad P_1(t) = \frac{5}{10}$$

$$GI = 1 - \left( \left(\frac{5}{10}\right)^2 + \left(\frac{5}{10}\right)^2 \right) = 0,5$$

$$E = - \left( \frac{5}{10} \log_2 \frac{5}{10} + \frac{5}{10} \log_2 \frac{5}{10} \right) = 1$$

C0: 9  
C1: 1

Low degree of impurity

$$c_0 = 1 \quad P_0(t) = \frac{1}{10}$$
$$c_1 = 9 \quad P_1(t) = \frac{9}{10}$$

$$GI = 1 - \left( \left(\frac{1}{10}\right)^2 + \left(\frac{9}{10}\right)^2 \right) = 0,18$$

$$E = - \left( \frac{1}{10} \log_2 \frac{1}{10} + \frac{9}{10} \log_2 \frac{9}{10} \right) = 0,469$$



# Compute Collective impurity of Child nodes

Once a decision rule **attribute test condition** have been applied it is possible to compute the collective impurity of child nodes

Consider an **attribute test condition** that splits a node containing N training instances into k children,

$$\{v_1, v_2, \dots, v_k\}$$

Let  $N(v_j)$  be the number of training instances associated with a child node  $v_j$ , whose impurity value is  $I(v_j)$

The **weighted impurity** measure of its **children**

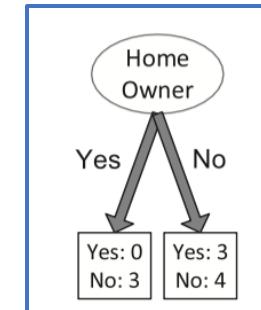
$$I(\text{children}) = \sum_{j=1}^k \frac{N(v_j)}{N} I(v_j),$$

ID	Home Owner	Marital Status	Annual Income	Defaulted Borrower
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes

$$I(\text{Home Owner} = \text{yes}) = -\frac{0}{3} \log_2 \frac{0}{3} - \frac{3}{3} \log_2 \frac{3}{3} = 0$$

$$I(\text{Home Owner} = \text{no}) = -\frac{3}{7} \log_2 \frac{3}{7} - \frac{4}{7} \log_2 \frac{4}{7} = 0.985$$

$$I(\text{Home Owner}) = \frac{3}{10} \times 0 + \frac{7}{10} \times 0.985 = 0.690$$

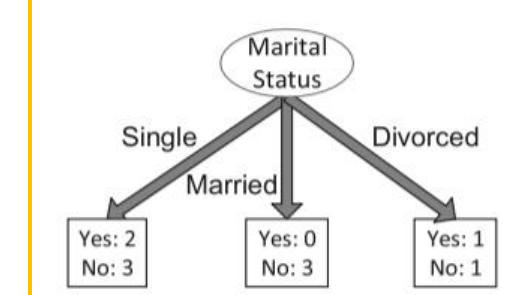


$$I(\text{Marital Status} = \text{Single}) = -\frac{2}{5} \log_2 \frac{2}{5} - \frac{3}{5} \log_2 \frac{3}{5} = 0.971$$

$$I(\text{Marital Status} = \text{Married}) = -\frac{0}{3} \log_2 \frac{0}{3} - \frac{3}{3} \log_2 \frac{3}{3} = 0$$

$$I(\text{Marital Status} = \text{Divorced}) = -\frac{1}{2} \log_2 \frac{1}{2} - \frac{1}{2} \log_2 \frac{1}{2} = 1.000$$

$$I(\text{Marital Status}) = \frac{5}{10} \times 0.971 + \frac{3}{10} \times 0 + \frac{2}{10} \times 1 = 0.686$$





# Finding the Best Split

---

1. Compute impurity measure **I(parent)** before splitting
2. Compute impurity measure **I(children)** after splitting
  - Compute impurity measure of each child node
  - **I(children)** is the weighted impurity of **child** nodes
3. Choose the attribute test condition that produces the highest gain

$$\Delta = I(\text{parent}) - I(\text{children}),$$

*(Handwritten red scribble over the equation)*

or equivalently, lowest impurity measure after splitting



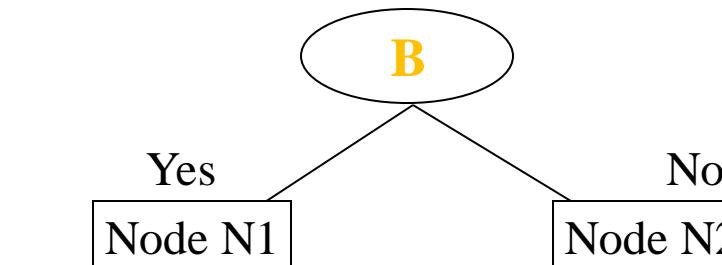
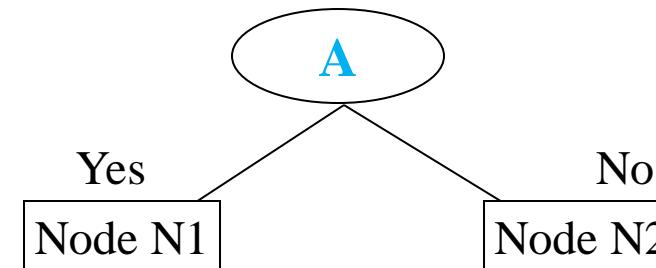
# Finding the Best Split

**Parent**

C0	<b>N00</b>
C1	<b>N01</b>

→ P

**Candidate Split**



**Candidate Children**

C0	<b>N10</b>
C1	<b>N11</b>

C0	<b>N20</b>
C1	<b>N21</b>

C0	<b>N10</b>
C1	<b>N11</b>

C0	<b>N20</b>
C1	<b>N21</b>

$$\text{MA1} = \text{Gini\_Index} (\text{A}, \text{N1})$$

$$\text{MA2} = \text{Gini\_Index} (\text{A}, \text{N2})$$

$$\text{MB1} = \text{Gini\_Index} (\text{B}, \text{N1})$$

$$\text{MB2} = \text{Gini\_Index} (\text{B}, \text{N2})$$

**M A**

**M B**

$$\text{Gain A} = \textcolor{red}{P} - \textcolor{blue}{M A}$$

**vs**

$$\text{Gain B} = \textcolor{red}{P} - \textcolor{blue}{M B}$$



# Finding the Best Split

It can be shown that the gain is non-negative since  $I(\text{parent}) \geq I(\text{children})$  for any reasonable measure

The higher the gain, the purer are the classes in the child nodes relative to the parent node.

The splitting criterion in the decision tree learning algorithm selects the attribute test condition that shows the maximum gain (minimizing the weighted impurity measure of its children)

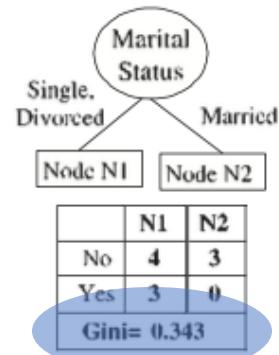
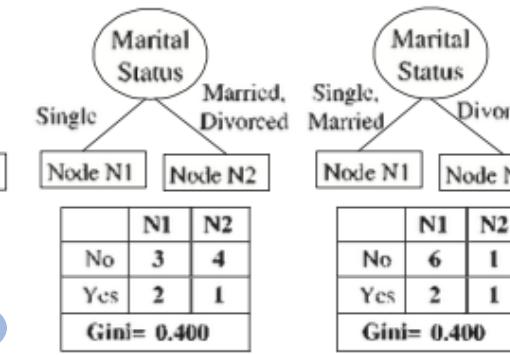
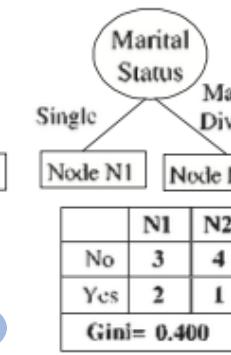
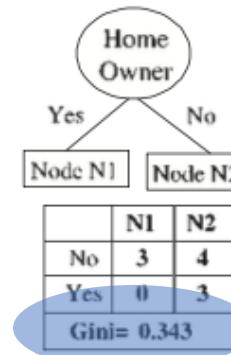
$$\Delta = I(\text{parent}) - I(\text{children}),$$

*(Handwritten note: This is the formula for Gini Index)*

$$I(\text{children}) = \sum_{j=1}^k \frac{N(v_j)}{N} I(v_j),$$

**N(v<sub>j</sub>):** number of training instances associated with the child node v<sub>j</sub>

	Parent
No	7
Yes	3
<b>Gini = 0.420</b>	



*Home Owner and Marital Status (with the highlighted association) shows the minimum weighted impurity*

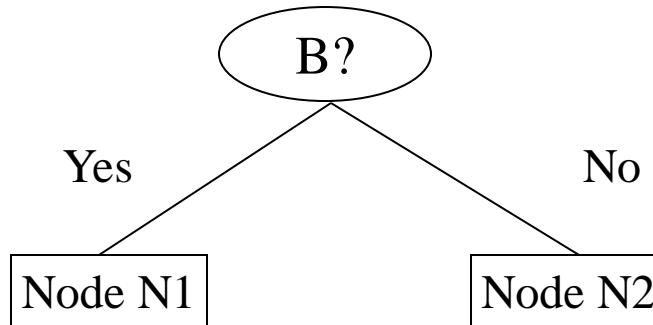
# Binary Attributes: Computing GINI Index



Splits into two partitions (child nodes)

Effect of Weighing partitions:

- Larger and purer partitions are sought



**Gini(N1)**

$$\begin{aligned} &= 1 - (5/6)^2 - (1/6)^2 \\ &= 0.278 \end{aligned}$$

**Gini(N2)**

$$\begin{aligned} &= 1 - (2/6)^2 - (4/6)^2 \\ &= 0.444 \end{aligned}$$

	<b>N1</b>	<b>N2</b>
C1	5	2
C2	1	4
<b>Gini=0.361</b>		

	<b>Parent</b>
C1	7
C2	5
<b>Gini = 0.486</b>	

**Weighted Gini of N1 N2**

$$\begin{aligned} &= 6/12 * 0.278 + \\ &\quad 6/12 * 0.444 \\ &= 0.361 \end{aligned}$$

**Gain = 0.486 – 0.361 = 0.125**

# Categorical Attributes: Computing Gini Index



For each distinct value, gather counts for each class in the dataset

Multi-way split

	CarType		
	Family	Sports	Luxury
C1	1	8	1
C2	3	0	7
Gini	<b>0.163</b>		

Two-way split  
(find best partition of values)

	CarType	
	{Sports, Luxury}	{Family}
C1	9	1
C2	7	3
Gini	<b>0.468</b>	

	CarType	
	{Sports}	{Family, Luxury}
C1	8	2
C2	0	10
Gini	<b>0.167</b>	

Which of these is the best?



# Continuous Attributes: Computing Gini Index

Use Binary Decisions based on one value

Several Choices for the splitting value

- Number of possible splitting values  
= Number of distinct values

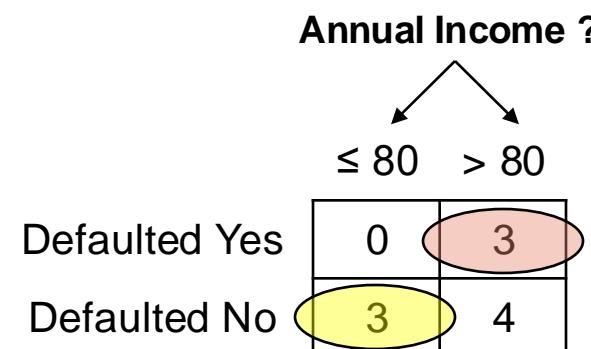
Each splitting value has a count matrix associated with it

- Class counts in each of the partitions,  $A \leq v$  and  $A > v$

Simple method to choose best  $v$

- For each  $v$ , scan the database to gather count matrix and compute its Gini index
- Computationally Inefficient! Repetition of work.

ID	Home Owner	Marital Status	Annual Income	Defaulted
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes





# Continuous Attributes: Computing Gini Index

- ② For efficient computation: for each attribute,
  - Sort the attribute on values
  - Linearly scan these values, each time updating the count matrix and computing gini index
  - Choose the split position that has the least gini index

Cheat	No	No	No	Yes	Yes	Yes	No	No	No	No	
Annual Income											
Sorted Values	→	60	70	75	85	90	95	100	120	125	220

# Continuous Attributes: Computing Gini Index

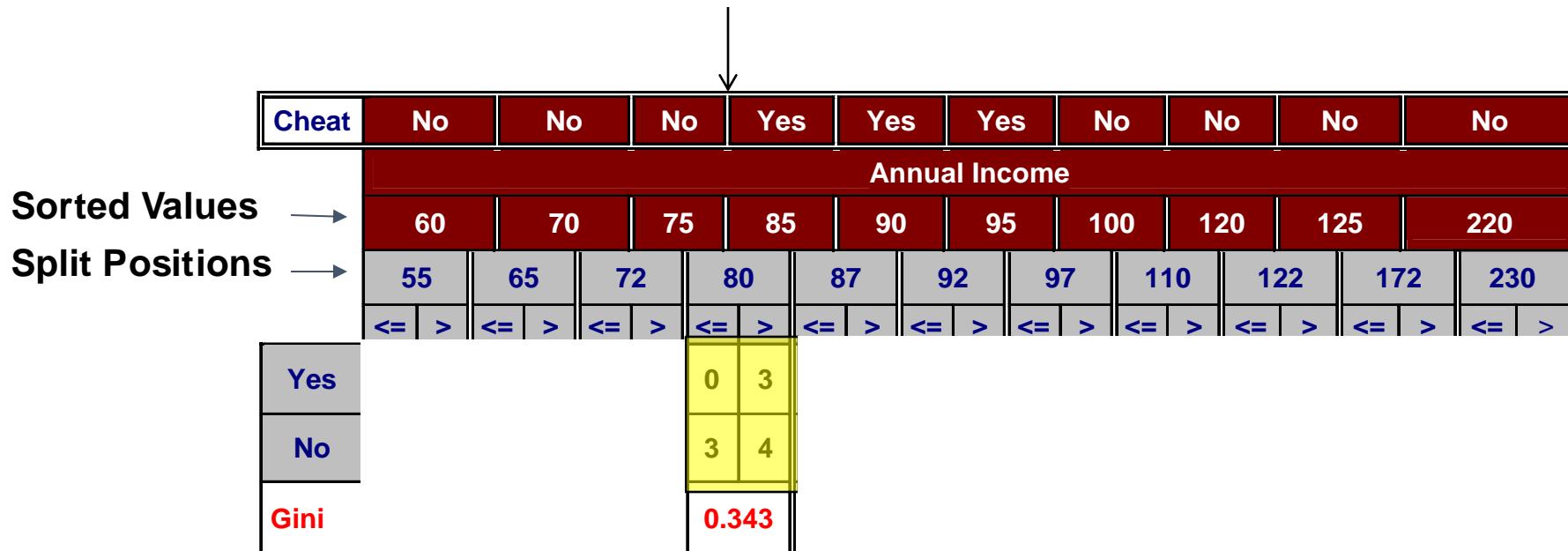


- For efficient computation: for each attribute,
    - Sort the attribute on values
    - Linearly scan these values, each time updating the count matrix and computing gini index
    - Choose the split position that has the least gini index



# Continuous Attributes: Computing Gini Index

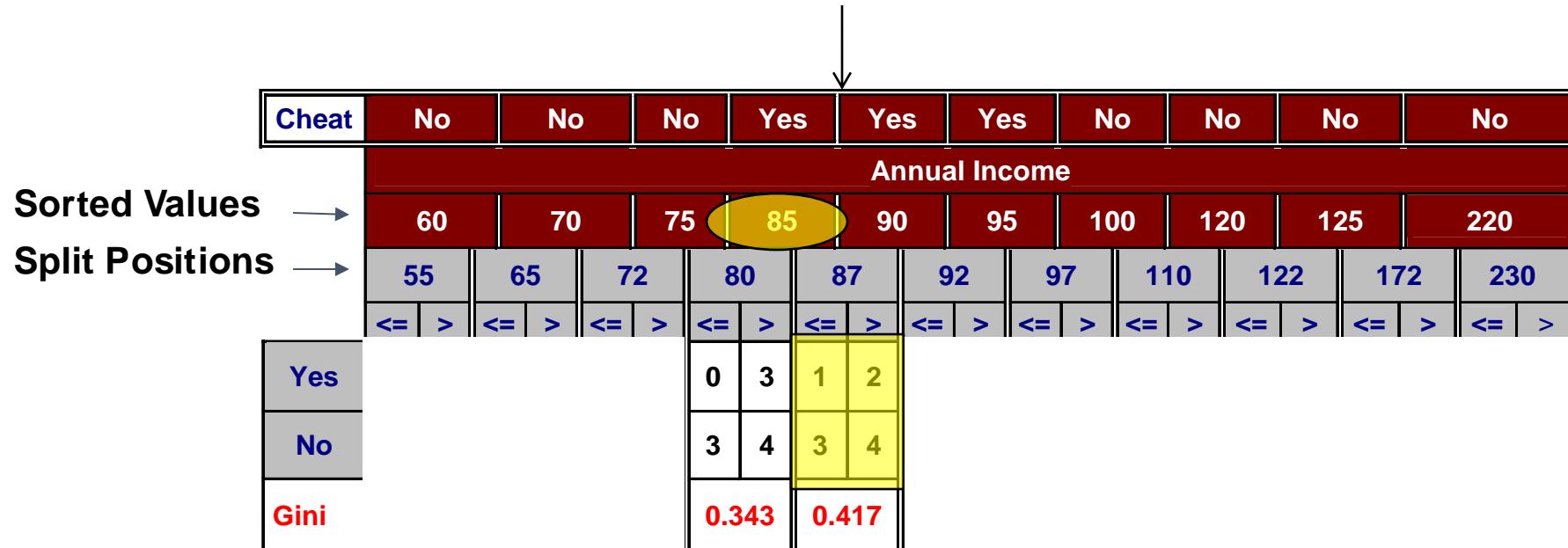
- | For efficient computation: for each attribute,
  - Sort the attribute on values
  - Linearly scan these values, each time updating the count matrix and computing gini index
  - Choose the split position that has the least gini index





# Continuous Attributes: Computing Gini Index

- | For efficient computation: for each attribute,
  - Sort the attribute on values
  - Linearly scan these values, each time updating the count matrix and computing gini index
  - Choose the split position that has the least gini index





# Continuous Attributes: Computing Gini Index

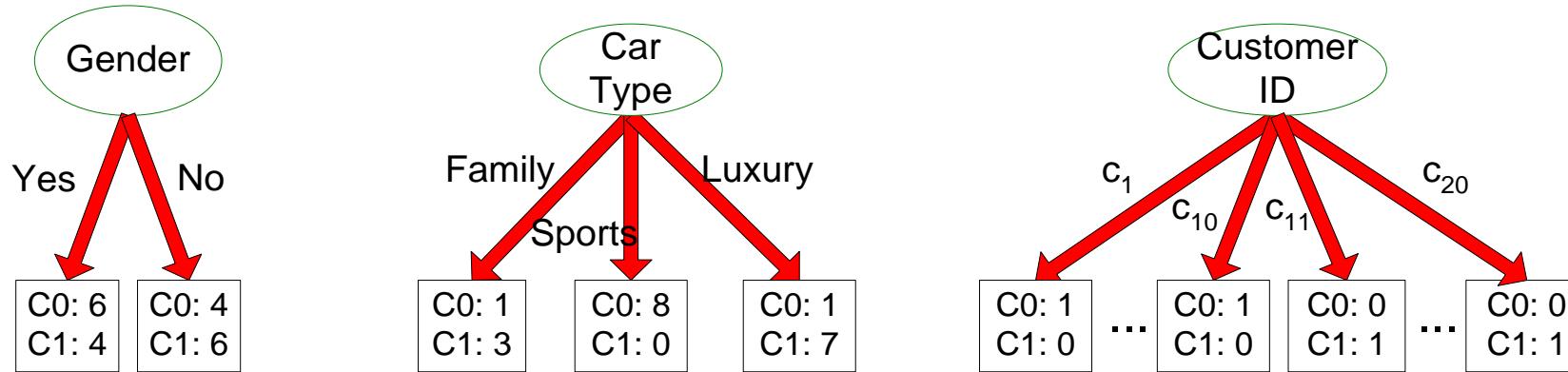
- | For efficient computation: for each attribute,
  - Sort the attribute on values
  - Linearly scan these values, each time updating the count matrix and computing gini index
  - Choose the split position that has the least gini index

Cheat	No	No	No	Yes	Yes	Yes	No	No	No	No
Annual Income										
→	60	70	75	85	90	95	100	120	125	220
→	55	65	72	80	87	92	97	110	122	172
	<= >	<= >	<= >	<= >	<= >	<= >	<= >	<= >	<= >	<= >
Yes	0 3	0 3	0 3	0 3	1 2	2 1	3 0	3 0	3 0	3 0
No	0 7	1 6	2 5	3 4	3 4	3 4	3 4	4 3	5 2	6 1
Gini	0.420	0.400	0.375	0.343	0.417	0.400	0.300	0.343	0.375	0.400



# Problem with large number of partitions

Node **impurity measures** tend to **prefer** splits that result in **large number of partitions**, each being small but pure



Customer Id	Gender	Car Type	Shirt Size	Class
1	M	Family	Small	C0
2	M	Sports	Medium	C0
3	M	Sports	Medium	C0
4	M	Sports	Large	C0
5	M	Sports	Extra Large	C0
6	M	Sports	Extra Large	C0
7	F	Sports	Small	C0
8	F	Sports	Small	C0
9	F	Sports	Medium	C0
10	F	Luxury	Large	C0
11	M	Family	Large	C1
12	M	Family	Extra Large	C1
13	M	Family	Medium	C1
14	M	Luxury	Extra Large	C1
15	F	Luxury	Small	C1
16	F	Luxury	Small	C1
17	F	Luxury	Medium	C1
18	F	Luxury	Medium	C1
19	F	Luxury	Medium	C1
20	F	Luxury	Large	C1

Customer ID has highest information gain because *Gini* or *Entropy* is **0** for all the children

# Gain Ratio



---

The number of children produced by the splitting attribute must be taken into consideration while deciding the best attribute test condition

A measure known as **Gain Ratio** is used to compensate for attributes that produce a large number of child nodes

The **Split Information** measures the entropy of splitting a node into its child nodes and evaluates if the split results in a larger number of equally-sized child nodes or not

$$\underline{\text{Gain Ratio}} = \frac{\Delta}{\underline{\text{Split Info}}}$$

$$\underline{\text{Split Info}} = - \sum_{i=1}^k \frac{n_i}{n} \log_2 \frac{n_i}{n}$$

# Gain Ratio



$$Gain\ Ratio = \frac{\Delta}{Split\ Info}$$

$$Split\ Info = \sum_{i=1}^k \frac{n_i}{n} \log_2 \frac{n_i}{n}$$

	CarType		
	Family	Sports	Luxury
C1	1	8	1
C2	3	0	7
Gini	<b>0.163</b>		

**SplitINFO = 1.52**

	CarType	
	{Sports, Luxury}	{Family}
C1	9	1
C2	7	3
Gini	<b>0.468</b>	

**SplitINFO = 0.72**

	CarType	
	{Sports}	{Family, Luxury}
C1	8	2
C2	0	10
Gini	<b>0.167</b>	

**SplitINFO = 0.97**

Parent Node,  $p$  is split into  $k$  partitions (children)

$n_i$  is number of records in child node  $i$



# Decision Tree Based Classification

---

## Advantages:

- Relatively inexpensive to construct
- Extremely fast at classifying unknown records
- Easy to interpret for small-sized trees
- Robust to noise (especially when methods to avoid overfitting are employed)
- Can easily handle redundant attributes
- Can easily handle irrelevant attributes (unless the attributes are **interacting**)

## Disadvantages: .

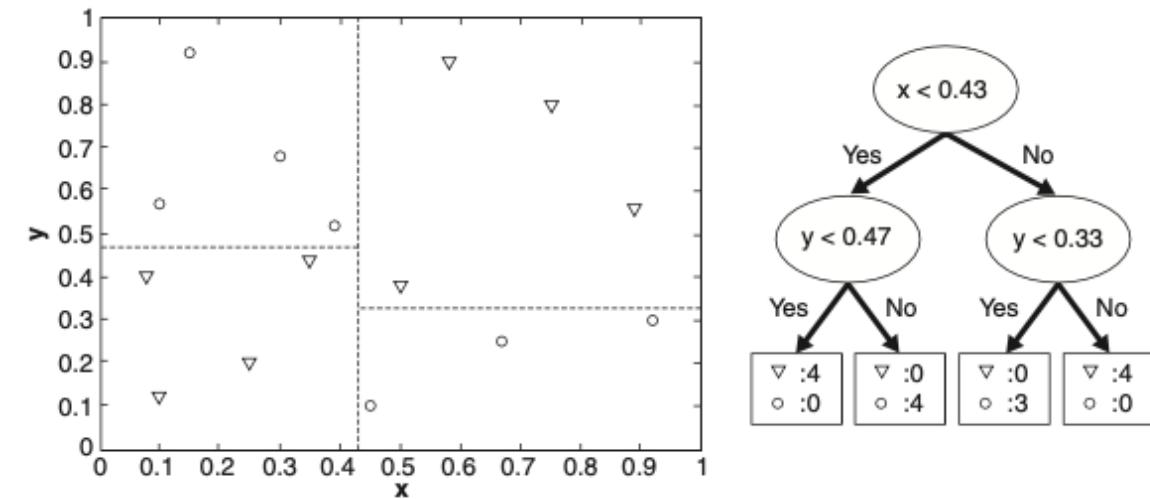
- Due to the greedy nature of splitting criterion, **interacting** attributes (that can distinguish between classes together but not individually) may be passed over in favor of other attributes that are less discriminating.
- Each decision boundary involves only a single attribute



# Handling interactions

## Rectilinear Decision Boundaries:

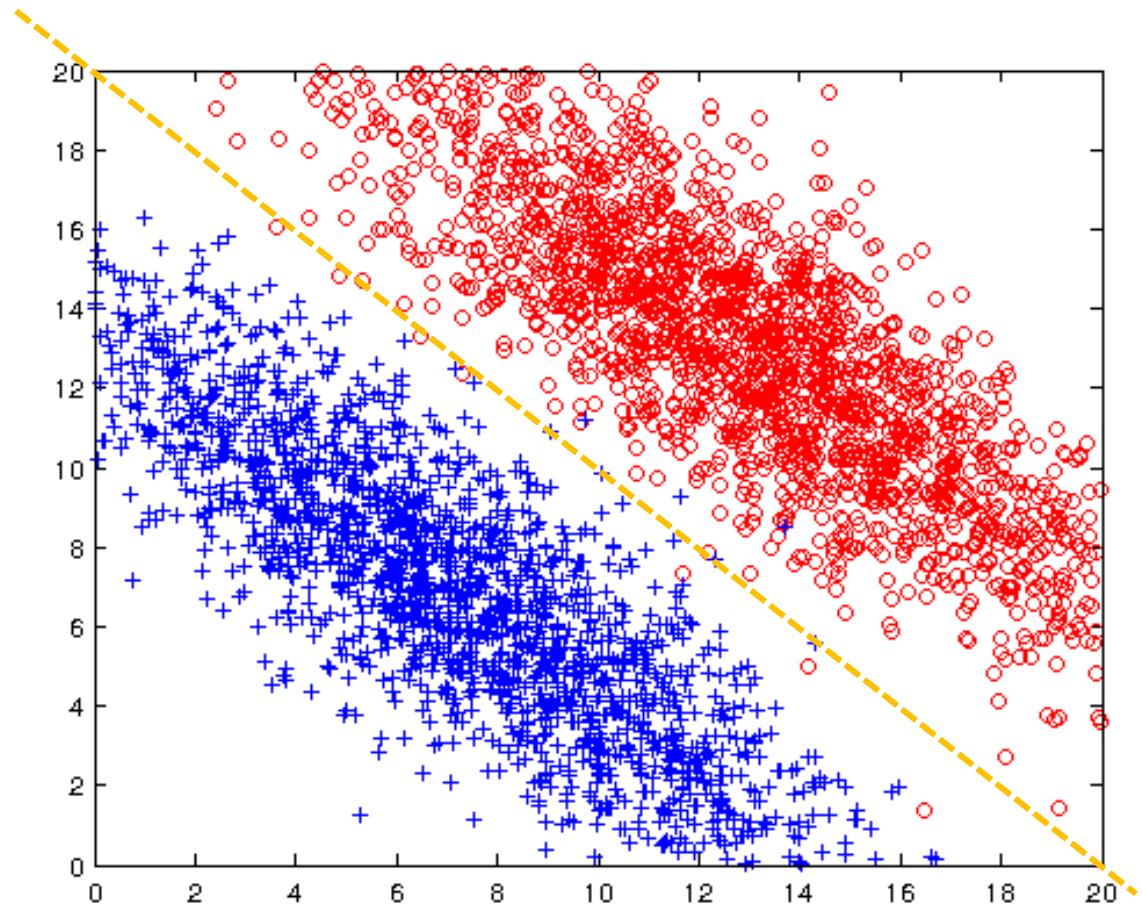
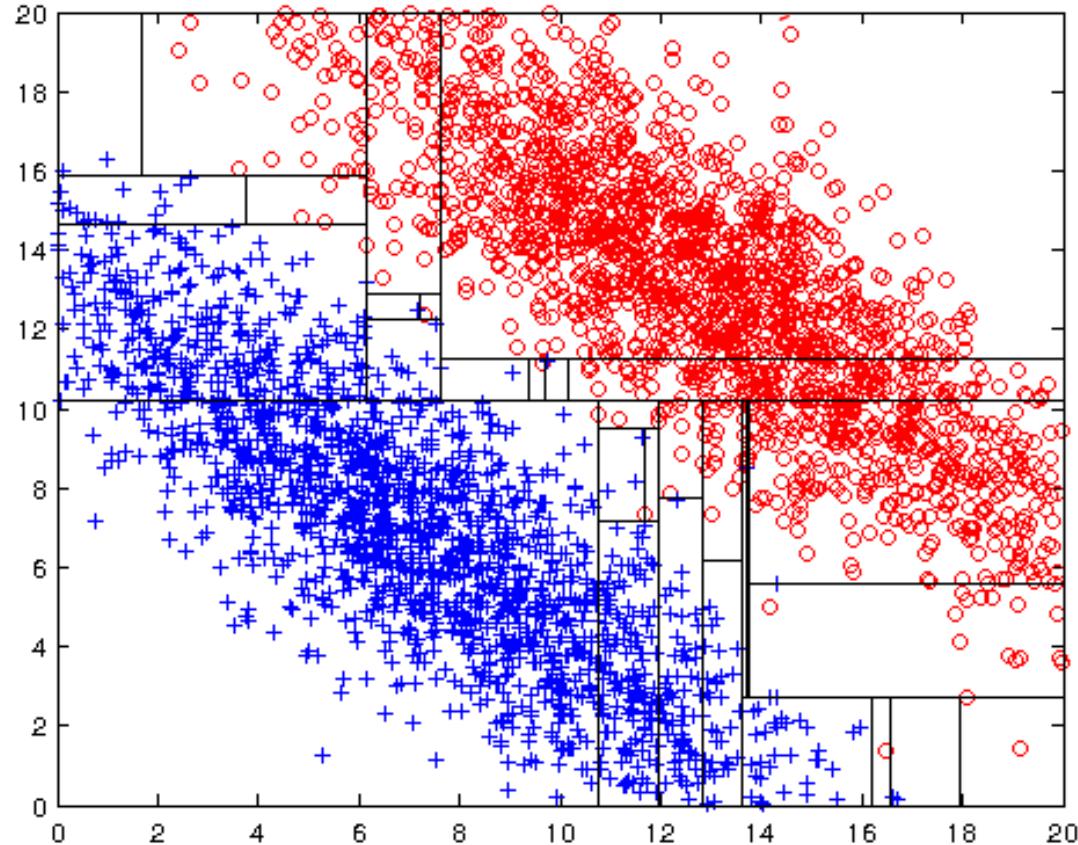
- **Single Attribute Testing:** Trees partition the attribute space based on one attribute at a time.
- **Rectilinear Boundaries:** The decision boundary is parallel to coordinate axes.
- **Limitation:** Restricts decision trees in representing complex boundaries, especially with continuous attributes.



## Oblique Decision Trees:

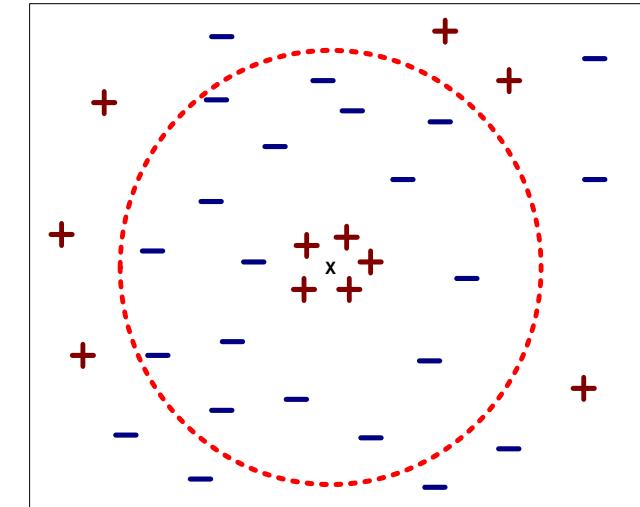
- **Improvement:** Allow tests using multiple attributes (e.g.,  $x + y < 20$ ) for more flexible decision boundaries.
- **Result:** Can better represent complex, non-rectilinear boundaries.

# Limitations of single attribute-based decision boundaries





# Nearest-Neighbor Classifiers



# Nearest-Neighbor Classifiers

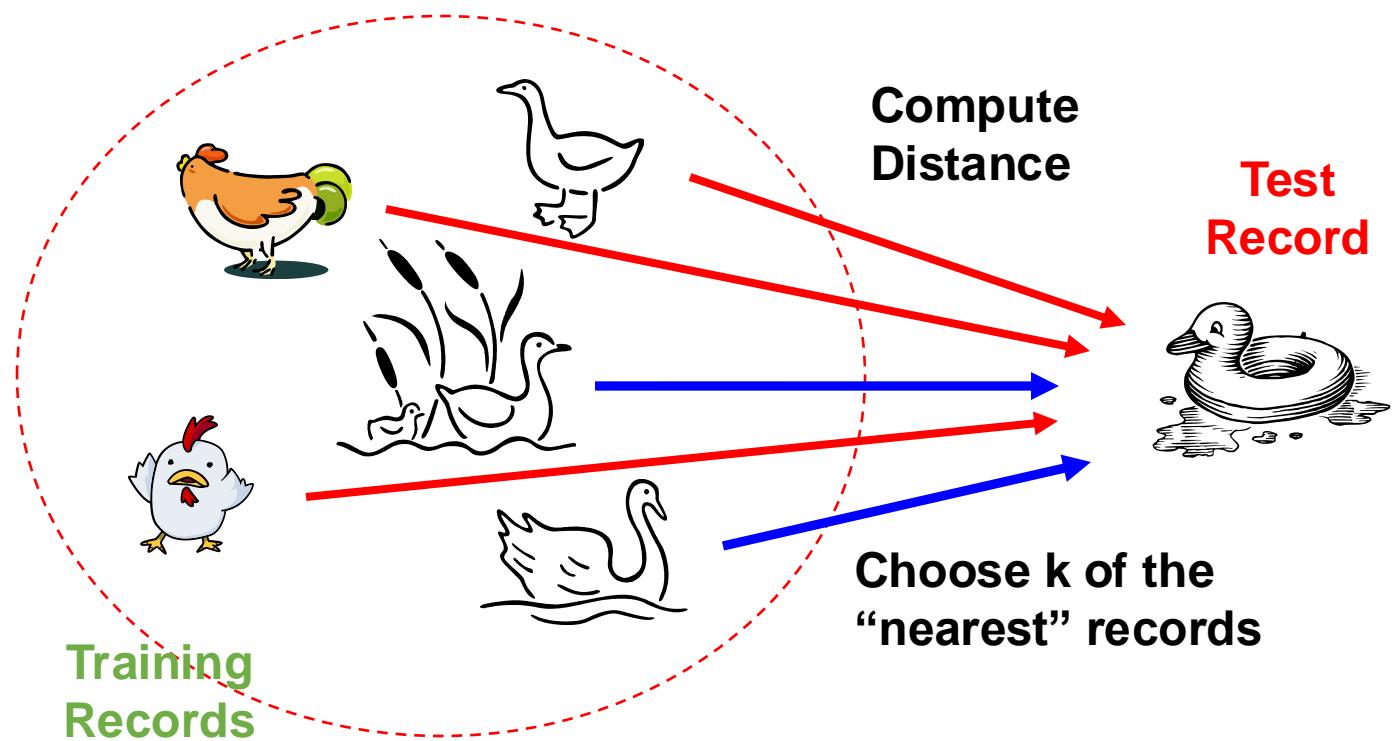


Basic idea: ***If it walks like a duck, quacks like a duck, then it's probably a duck***

A **nearest neighbor classifier** represents each example as a data point in a **d-dimensional space** ( $d$  the number of attributes)

Given a **test instance**, we compute its proximity to the training instances according to one of the proximity measures

The **k-nearest neighbors** of a given **test instance z** refer to the  $k$  training examples that are **closest to z**.

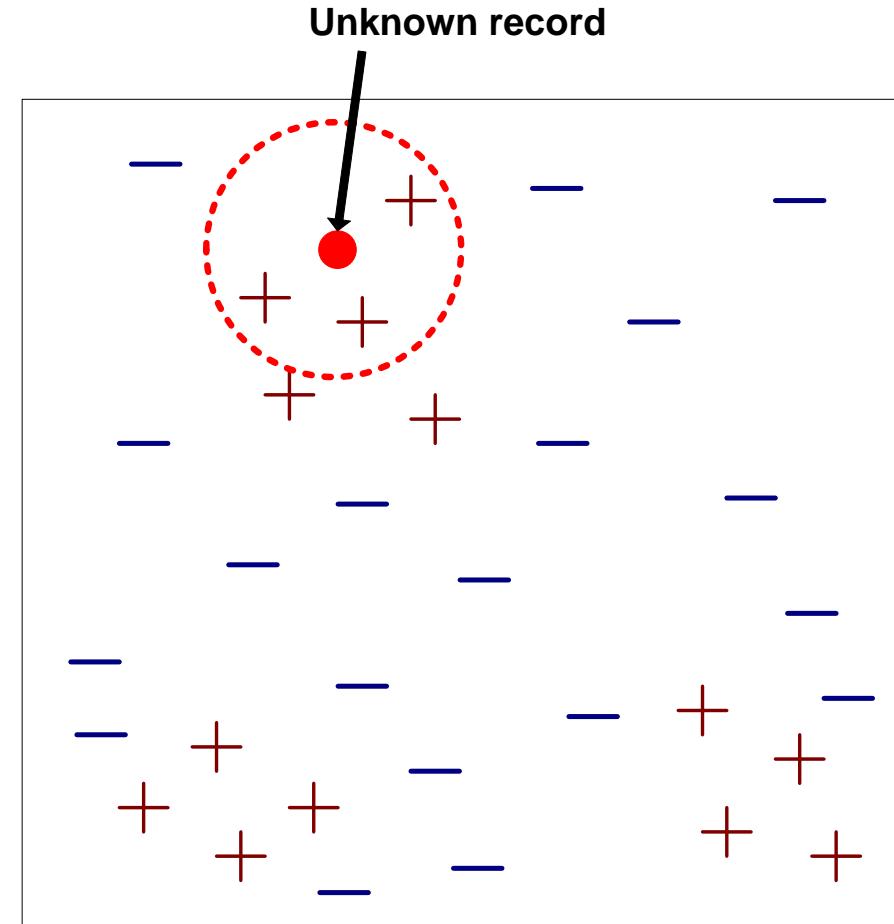




# Nearest-Neighbor Classifiers

Requires the following:

- A set of labeled records
- Proximity metric to compute distance/similarity between a pair of records e.g., Euclidean distance
- The value of  $k$ , the number of nearest neighbors to retrieve
- A method for using class labels of  $K$  nearest neighbors to determine the class label of unknown record (e.g., by taking majority vote)

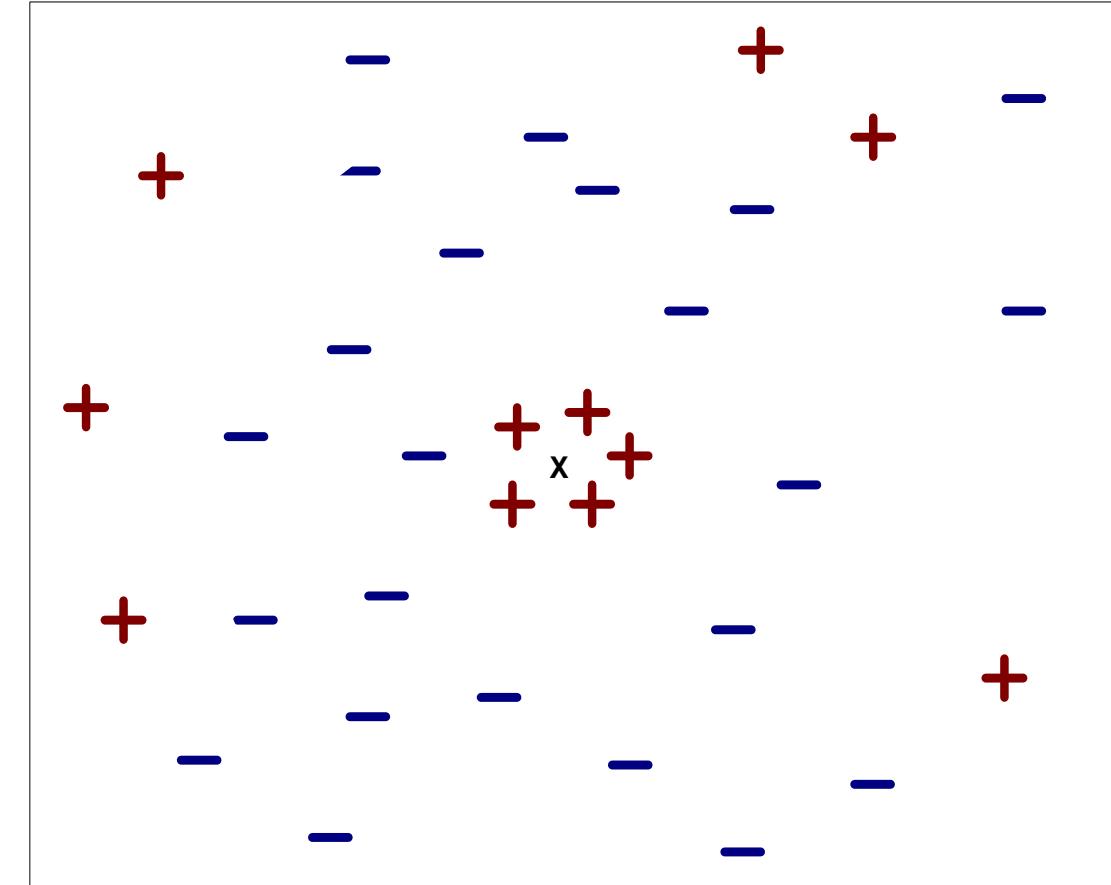




# Nearest-Neighbor Classifiers

Choosing the value of k:

- If k is **too small**, sensitive to noise points (overfitting)
- If k is **too large**, neighborhood may include points from other classes
- An example of **well balanced** k

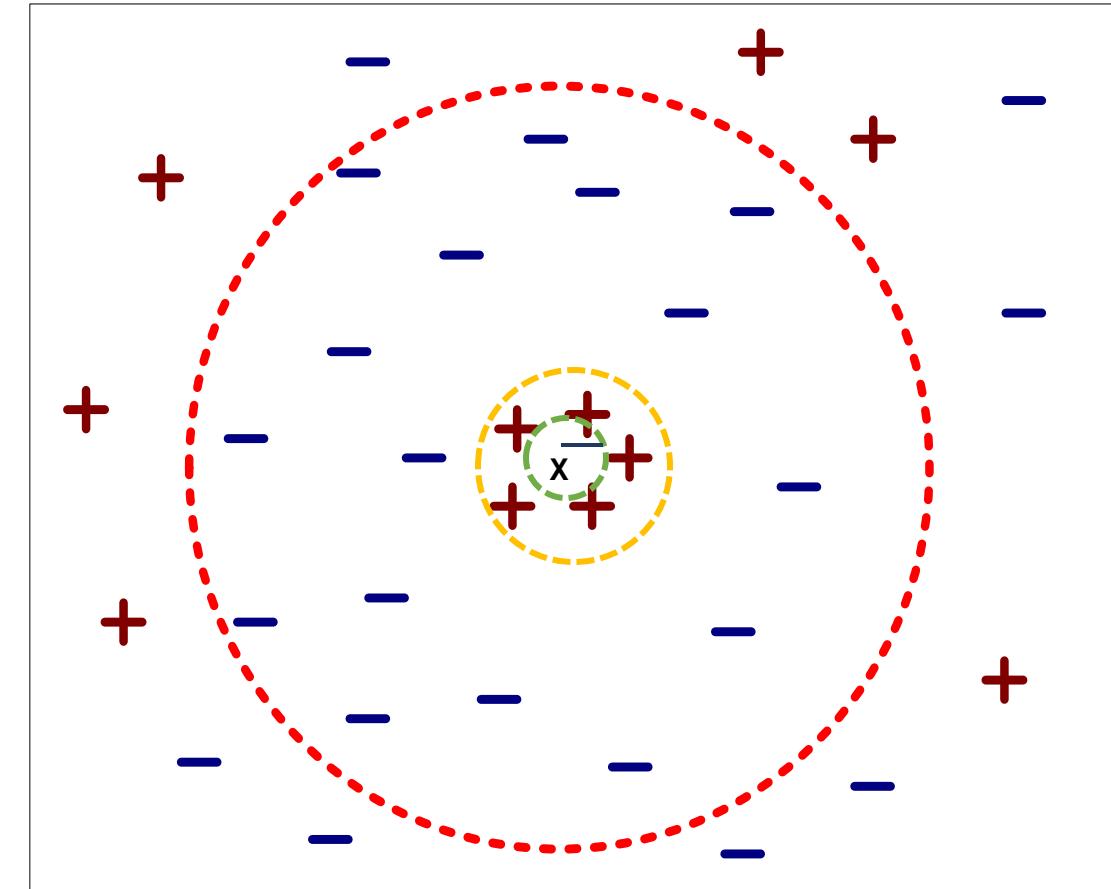




# Nearest-Neighbor Classifiers

Choosing the value of k:

- If k is **too small**, sensitive to noise points (overfitting)
- If k is **too large**, neighborhood may include points from other classes
- An example of **well balanced** k





# Nearest-Neighbor Classifiers

How to Determine the class label of a Test Sample?

- 1: Let  $k$  be the number of nearest neighbors and  $D$  be the set of training examples.
- 2: **for** each test instance  $z = (\mathbf{x}', y')$  **do**
- 3:   Compute  $d(\mathbf{x}', \mathbf{x})$ , the distance between  $z$  and every example,  $(\mathbf{x}, y) \in D$ .
- 4:   Select  $D_z \subseteq D$ , the set of  $k$  closest training examples to  $z$ .
- 5:    $y' = \underset{v}{\operatorname{argmax}} \sum_{(\mathbf{x}_i, y_i) \in D_z} I(v = y_i)$
- 6: **end for**

Take the **majority vote** of class labels among the k-nearest neighbors

Weight the vote according to distance

- weight factor,  $w = 1/d^2$



# Nearest-Neighbor Classifiers

- $v$  is a **class label**
- $y_i$  is the **class label** for **one of the nearest neighbors**
- $I(\cdot)$  is an **indicator function** that returns the value 1 if its argument is true and 0 otherwise

## Majority Voting

$$y' = \operatorname{argmax}_v \sum_{(\mathbf{x}_i, y_i) \in D_z} I(v = y_i),$$

## Distance-Weighted Voting

$$y' = \operatorname{argmax}_v \sum_{(\mathbf{x}_i, y_i) \in D_z} w_i \times I(v = y_i).$$

**weight factor,**  $w = 1/d^2$

# Nearest-Neighbor Classifiers

---



## Data preprocessing is often required

Attributes may have to be **scaled** to prevent *distance measures* from being *dominated* by one of the attributes

- Example:
  - height of a person may vary from 1.5m to 1.8m
  - weight of a person may vary from 90lb to 300lb
  - income of a person may vary from \$10K to \$1M

- It is often useful to **standardize** data to have 0 means a standard deviation of 1

# Nearest-Neighbor Classifiers

---



## How to handle missing values in training and test sets?

- Proximity computations normally require the presence of all attributes
- Some approaches use the subset of attributes present in two instances
  - This may not produce good results since it effectively uses different proximity measures for each pair of instances
  - Thus, proximities are not comparable

## How to handle Irrelevant and Redundant Attributes

- Irrelevant attributes add noise to the proximity measure
- Redundant attributes bias the proximity measure towards certain attributes

# Nearest-Neighbor Classifiers

---



## Characteristics

- Nearest-Neighbor classification is part of a more general technique known as instance-based learning
- **Classifying** a test instance can be quite **expensive** because we need to compute the proximity values individually (lazy learners)
- Nearest neighbor classifiers can produce **decision boundaries of arbitrary shape**
- Nearest neighbor classifiers have **difficulty handling missing values** in both the training and test sets since proximity computations normally require the presence of all attributes
- The presence of irrelevant attributes can distort commonly used proximity measures, especially when the number of irrelevant attributes is large



# Bayesian Classifiers

$$P(A|B)$$
A graphic of a calculator. The display screen is green with a white border and shows the mathematical expression  $P(A|B)$ . Below the screen is a blue rectangular base.



# Bayesian Classifiers

---

Many classification problems involve **uncertainty**

- Observed attributes and class labels may be unreliable due to **imperfections** in the measurement process
- Set of attributes may not be fully representative of the target class, resulting in **uncertain** predictions
- Classification model learned over a finite training set may not be able to fully capture the **true relationships** in the overall data

In the presence of **uncertainty** we need provide not only ***predictions*** of class labels but a ***measure of confidence***

**Probability theory** offers a systematic way for ***quantifying*** and ***manipulating uncertainty*** in data

Classification models that make use of ***probability*** are known as **probabilistic classification model**

# Bayesian Classifiers



## Introduction to probability

Consider a **variable  $X$** , which can take any discrete value from the set  $\{x_1, x_2, \dots, x_k\}$

$X$  has the value  $x_i$  for  $n_i$  data objects

Relative frequency of event  $X = x_i$  is  $\frac{n_i}{N}$  where  $N = \sum_{i=1}^k n_i$

The **probability** of an event e, e.g.,  $P(X = x_i)$ , measures how likely it is for the event  $x_i$  to occur

A probability is always a number between 0 and 1

The **sum of probability** values of all possible events (outcomes of a variable  $X$ ) is **equal to 1**

Variables that have probabilities associated with each possible outcome (values) are known as **random variables**

# Bayesian Classifiers



## Using Bayes Theorem for Classification

For the purpose of **classification**, we are interested in computing the probability of observing a *class label*  $y$  for a data instance given its set of attribute  $x$

$P(y|x)$  is the posterior probability

$$P(y|x) = \frac{P(x|y)P(y)}{P(x)}$$

where

$P(x|y)$  is the **class-conditional probability** the likelihood of observing  $x$  from the distribution of instances belonging to  $y$ .

$P(y)$  is the **prior probability** (prior beliefs about the distribution of class labels) [expert knowledge or distribution]

$P(x)$  can be computed as  $P(x) = \sum_{i=1}^k P(x|y_i)P(y_i)$



## Naïve Bayes Assumption for Classification

**Naïve Bayes classifier** assumes that the class-conditional probability  $P(\mathbf{x}|y)$  of all attributes  $\mathbf{x}$  can be factored as a product of class-conditional probabilities  $P(x_i|y)$

$$P(\mathbf{x}|y) = \prod_{i=1}^d P(x_i|y)$$

Where data instance  $\mathbf{x}$  the set  $\{x_1, x_2, \dots, x_d\}$

The basic assumption behind the previous equation is that the attribute values  $x_i$  are **conditionally independent**



### Naïve Bayes Assumption for Classification

The Naïve Bayes classifier computes the posterior probability for a test instance  $\mathbf{x}$  by using the following equation

$$P(y|\mathbf{x}) = P(y) \prod_{i=1}^d P(x_i|y)$$



## Using Bayes Theorem for Classification

- Consider each **attribute** and class label as *random variables*
- Given a record with attributes  $\{x_1, x_2, \dots, x_d\}$  the goal is to predict class  $y$
- Specifically, we want to find the value of  $Y$  that maximizes  $P(y|x_1, x_2, \dots, x_d)$

Can we estimate  $P(y|x_1, x_2, \dots, x_d)$   
directly from data?

Tid	Refund	Marital Status	Taxable Income	Evade
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes

# Bayesian Classifiers: Example



Given a *Test Record*:

$$X = (\text{Refund} = \text{No}, \text{Divorced}, \text{Income} = 120\text{K})$$

$$P(y|x) = P(y) \prod_{i=1}^d P(x_i|y)$$

- We need to estimate

$$P(y = yes|x) \text{ and } P(y = no|x)$$

In the following we will replace  
 $y = Yes$  by Yes, and  
 $y = No$  by No

Tid	Refund	Marital Status	Taxable Income	Evade
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes

# Bayesian Classifiers: Example



Given a *Test Record*:

$$X = (\text{Refund} = \text{No}, \text{Divorced}, \text{Income} = 120\text{K})$$

$$P(\text{yes}|\mathbf{x}) = P(\text{yes}) \prod_{i=1}^d P(x_i|\text{yes})$$

$$P(\text{no}|\mathbf{x}) = P(\text{no}) \prod_{i=1}^d P(x_i|\text{no})$$

Tid	Refund	Marital Status	Taxable Income	Evade
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes

# Bayesian Classifiers: Example



Given a *Test Record*:

$$X = (\text{Refund} = \text{No}, \text{Divorced}, \text{Income} = 120\text{K})$$

$$P(\text{yes}) = 3/10$$

$$P(\text{no}) = 7/10$$

Tid	Refund	Marital Status	Taxable Income	Evade
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes

# Bayesian Classifiers: Example



Given a *Test Record*:

$$X = (\text{Refund} = \text{No}, \text{Divorced}, \text{Income} = 120\text{K})$$

$$\prod_{i=1}^d P(x_i | yes)$$

$$\begin{aligned} P(X | Yes) &= \\ &P(\text{Refund} = \text{No} | Yes) \times \\ &P(\text{Divorced} | Yes) \times \\ &P(\text{Income} = 120\text{K} | Yes) \end{aligned}$$

$$\prod_{i=1}^d P(x_i | no)$$

$$\begin{aligned} P(X | No) &= \\ &P(\text{Refund} = \text{No} | No) \times \\ &P(\text{Divorced} | No) \times \\ &P(\text{Income} = 120\text{K} | No) \end{aligned}$$

Tid	Refund	Marital Status	Taxable Income	Evade
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes

# Bayesian Classifiers: Example



Given a *Test Record*:

$$X = (\text{Refund} = \text{No}, \text{Divorced}, \text{Income} = 120\text{K})$$

$$\prod_{i=1}^d P(x_i | \text{yes})$$

$$P(X | \text{Yes}) =$$

$$P(\text{Refund} = \text{No} | \text{Yes}) \times$$

$$P(\text{Divorced} | \text{Yes}) \times$$

$$P(\text{Income} = 120\text{K} | \text{Yes})$$

$$\prod_{i=1}^d P(x_i | \text{no})$$

$$P(X | \text{No}) =$$

$$P(\text{Refund} = \text{No} | \text{No}) \times$$

$$P(\text{Divorced} | \text{No}) \times$$

$$P(\text{Income} = 120\text{K} | \text{No})$$

Tid	Refund	Marital Status	Taxable Income	Evade
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes

Categorical

Continuous

# Bayesian Classifiers: Example



Given a *Test Record*:

$$X = (\text{Refund} = \text{No}, \text{Divorced}, \text{Income} = 120\text{K})$$

For categorical attributes:

$$P(x_i = c|yes) = n_c/n = n_c/n$$

where  $n_c$  is number of instances having attribute value  $x_i = c$  and belonging to class y

–Examples:

$$P(\text{Status}=\text{Married}|\text{No}) = 4/7$$

$$P(\text{Refund}=\text{Yes}|\text{Yes}) = 0$$

Tid	Refund	Marital Status	Taxable Income	Evade
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes

# Bayesian Classifiers: Example

---



For continuous attributes:

**Discretization:** Partition the range into bins:

Replace continuous value with bin value

Attribute changed from continuous to ordinal

**Probability density estimation:**

Assume attribute follows a normal distribution

Use data to estimate parameters of distribution  
(e.g., mean and standard deviation)

Once probability distribution is known, use it to estimate the conditional probability  
 $P(X_i|Y)$

# Bayesian Classifiers: Example



Given a *Test Record*:

$$X = (\text{Refund} = \text{No}, \text{Divorced}, \text{Income} = 120\text{K})$$

Normal distribution:

$$P(X_i | Y_j) = \frac{1}{\sqrt{2\pi\sigma_{ij}^2}} e^{-\frac{(X_i - \mu_{ij})^2}{2\sigma_{ij}^2}}$$

- One for each  $(X_i, Y_i)$  pair

For (Income, Class=No):

- If Class=No
  - sample mean = 110
  - sample variance = 2975

Tid	Refund	Marital Status	Taxable Income	Evade
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes

$$P(\text{Income} = 120 | \text{No}) = \frac{1}{\sqrt{2\pi}(54.54)} e^{-\frac{(120-110)^2}{2(2975)}} = 0.0072$$

# Bayesian Classifiers: Example



Given a *Test Record*:

$$X = (\text{Refund} = \text{No}, \text{Divorced}, \text{Income} = 120\text{K})$$

## Naïve Bayes Classifier:

$$P(\text{Refund} = \text{Yes} | \text{No}) = 3/7$$

$$P(\text{Refund} = \text{No} | \text{No}) = 4/7$$

$$P(\text{Refund} = \text{Yes} | \text{Yes}) = 0$$

$$P(\text{Refund} = \text{No} | \text{Yes}) = 1$$

$$P(\text{Marital Status} = \text{Single} | \text{No}) = 2/7$$

$$P(\text{Marital Status} = \text{Divorced} | \text{No}) = 1/7$$

$$P(\text{Marital Status} = \text{Married} | \text{No}) = 4/7$$

$$P(\text{Marital Status} = \text{Single} | \text{Yes}) = 2/3$$

$$P(\text{Marital Status} = \text{Divorced} | \text{Yes}) = 1/3$$

$$P(\text{Marital Status} = \text{Married} | \text{Yes}) = 0$$

- $P(X | \text{No}) = P(\text{Refund}=\text{No} | \text{No}) \times P(\text{Divorced} | \text{No}) \times P(\text{Income}=120\text{K} | \text{No}) = 4/7 \times 1/7 \times 0.0072 = 0.0006$
- $P(X | \text{Yes}) = P(\text{Refund}=\text{No} | \text{Yes}) \times P(\text{Divorced} | \text{Yes}) \times P(\text{Income}=120\text{K} | \text{Yes}) = 1 \times 1/3 \times 1.2 \times 10^{-9} = 4 \times 10^{-10}$

For Taxable Income:

If class = No: sample mean = 110

sample variance = 2975

If class = Yes: sample mean = 90

sample variance = 25

Tid	Refund	Marital Status	Taxable Income	Evade
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes



## Bayesian Classifiers: Example

Given a *Test Record*:

$$X = (\text{Refund} = \text{No}, \text{Divorced}, \text{Income} = 120\text{K})$$

$$P(y|x) = P(y) \prod_{i=1}^d P(x_i|y)$$

$$P(\text{no}|x) = P(\text{no}) P(x|\text{no}) = 0.0004$$

$$P(\text{yes}|x) = P(\text{yes}) P(x|\text{yes}) = 1.2 \times 10^{-10}$$

$$\begin{aligned} P(X | \text{No}) &= P(\text{Refund}=\text{No} | \text{No}) \\ &\quad \times P(\text{Divorced} | \text{No}) \\ &\quad \times P(\text{Income}=120\text{K} | \text{No}) \\ &= 4/7 \times 1/7 \times 0.0072 = 0.0006 \end{aligned}$$

$$P(\text{no}) = 7/10$$

$$\begin{aligned} P(X | \text{Yes}) &= P(\text{Refund}=\text{No} | \text{Yes}) \\ &\quad \times P(\text{Divorced} | \text{Yes}) \\ &\quad \times P(\text{Income}=120\text{K} | \text{Yes}) \\ &= 1 \times 1/3 \times 1.2 \times 10^{-9} = 4 \times 10^{-10} \end{aligned}$$

$$P(\text{yes}) = 3/10$$

# Bayesian Classifiers: Issue



Naïve Bayes Classifier can make decisions with partial information about attributes in the test record

## Naïve Bayes Classifier:

$$P(\text{Refund} = \text{Yes} | \text{No}) = 3/7$$

$$P(\text{Refund} = \text{No} | \text{No}) = 4/7$$

$$P(\text{Refund} = \text{Yes} | \text{Yes}) = 0$$

$$P(\text{Refund} = \text{No} | \text{Yes}) = 1$$

$$P(\text{Marital Status} = \text{Single} | \text{No}) = 2/7$$

$$P(\text{Marital Status} = \text{Divorced} | \text{No}) = 1/7$$

$$P(\text{Marital Status} = \text{Married} | \text{No}) = 4/7$$

$$P(\text{Marital Status} = \text{Single} | \text{Yes}) = 2/3$$

$$P(\text{Marital Status} = \text{Divorced} | \text{Yes}) = 1/3$$

$$P(\text{Marital Status} = \text{Married} | \text{Yes}) = 0$$

For Taxable Income:

If class = No: sample mean = 110

sample variance = 2975

If class = Yes: sample mean = 90

sample variance = 25

$$P(\text{Yes}) = 3/10 \quad P(\text{No}) = 7/10$$

**If we only know that marital status is Divorced, then:**

$$P(\text{Yes} | \text{Divorced}) = 1/3 \times 3/10 / P(\text{Divorced})$$

$$P(\text{No} | \text{Divorced}) = 1/7 \times 7/10 / P(\text{Divorced})$$

**If we also know that Refund = No, then**

$$P(\text{Yes} | \text{Refund} = \text{No}, \text{Divorced}) = 1 \times 1/3 \times 3/10 \quad P(\text{Divorced, Refund} = \text{No})$$

$$P(\text{No} | \text{Refund} = \text{No}, \text{Divorced}) = 4/7 \times 1/7 \times 7/10 \quad P(\text{Divorced, Refund} = \text{No})$$

**If we also know that Taxable Income = 120, then**

$$P(\text{Yes} | \text{Refund} = \text{No}, \text{Divorced, Income} = 120) = \\ 1.2 \times 10^{-9} \times 1 \times 1/3 \times 3/10 \quad P(\text{Divorced, Refund} = \text{No, Income} = 120)$$

$$P(\text{No} | \text{Refund} = \text{No}, \text{Divorced, Income} = 120) = \\ 0.0072 \times 4/7 \times 1/7 \times 7/10 \quad P(\text{Divorced, Refund} = \text{No, Income} = 120)$$

# Bayesian Classifiers: Issue



**Given a Test Record:  $X = (\text{Married})$**

**Naïve Bayes Classifier:**

$$P(\text{Refund} = \text{Yes} | \text{No}) = 3/7$$

$$P(\text{Refund} = \text{No} | \text{No}) = 4/7$$

$$P(\text{Refund} = \text{Yes} | \text{Yes}) = 0$$

$$P(\text{Refund} = \text{No} | \text{Yes}) = 1$$

$$P(\text{Marital Status} = \text{Single} | \text{No}) = 2/7$$

$$P(\text{Marital Status} = \text{Divorced} | \text{No}) = 1/7$$

$$P(\text{Marital Status} = \text{Married} | \text{No}) = 4/7$$

$$P(\text{Marital Status} = \text{Single} | \text{Yes}) = 2/3$$

$$P(\text{Marital Status} = \text{Divorced} | \text{Yes}) = 1/3$$

$$P(\text{Marital Status} = \text{Married} | \text{Yes}) = 0$$

For Taxable Income:

If class = No: sample mean = 110

sample variance = 2975

If class = Yes: sample mean = 90

sample variance = 25

$$P(\text{Yes}) = 3/10$$

$$P(\text{No}) = 7/10$$

$$P(\text{Yes} | \text{Married}) = 0 \times 3/10 / P(\text{Married})$$

$$P(\text{No} | \text{Married}) = 4/7 \times 7/10 / P(\text{Married})$$



# Bayesian Classifiers: Issue

Consider the table with record = 7 deleted

Tid	Refund	Marital Status	Taxable Income	Evade
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes

## Naïve Bayes Classifier:

$P(\text{Refund} = \text{Yes} | \text{No}) = 2/6$   
 $P(\text{Refund} = \text{No} | \text{No}) = 4/6$   
→  $P(\text{Refund} = \text{Yes} | \text{Yes}) = 0$   
 $P(\text{Refund} = \text{No} | \text{Yes}) = 1$   
 $P(\text{Marital Status} = \text{Single} | \text{No}) = 2/6$   
→  $P(\text{Marital Status} = \text{Divorced} | \text{No}) = 0$   
 $P(\text{Marital Status} = \text{Married} | \text{No}) = 4/6$   
 $P(\text{Marital Status} = \text{Single} | \text{Yes}) = 2/3$   
 $P(\text{Marital Status} = \text{Divorced} | \text{Yes}) = 1/3$   
 $P(\text{Marital Status} = \text{Married} | \text{Yes}) = 0/3$

For Taxable Income:

If class = No: sample mean = 91  
sample variance = 685  
If class = Yes: sample mean = 90  
sample variance = 25

Given  $X = (\text{Refund} = \text{Yes}, \text{Divorced}, 120\text{K})$

$$P(X | \text{No}) = 2/6 \times 0 \times 0.0083 = 0$$

$$P(X | \text{Yes}) = 0 \times 1/3 \times 1.2 \times 10^{-9} = 0$$

**Naïve Bayes will not be able to  
classify X as Yes or No!**



## Bayesian Classifiers: Issue

If one of the conditional probabilities is zero, then the entire expression becomes zero

Need to use other estimates of conditional probabilities than simple fractions

Probability estimation:

$$\text{original: } P(X_i = c|y) = \frac{n_c}{n}$$

$$\text{Laplace Estimate: } P(X_i = c|y) = \frac{n_c + 1}{n + v}$$

$$\text{m - estimate: } P(X_i = c|y) = \frac{n_c + mp}{n + m}$$

$n$ : number of training instances belonging to class  $y$

$n_c$ : number of instances with  $X_i = c$  and  $Y = y$

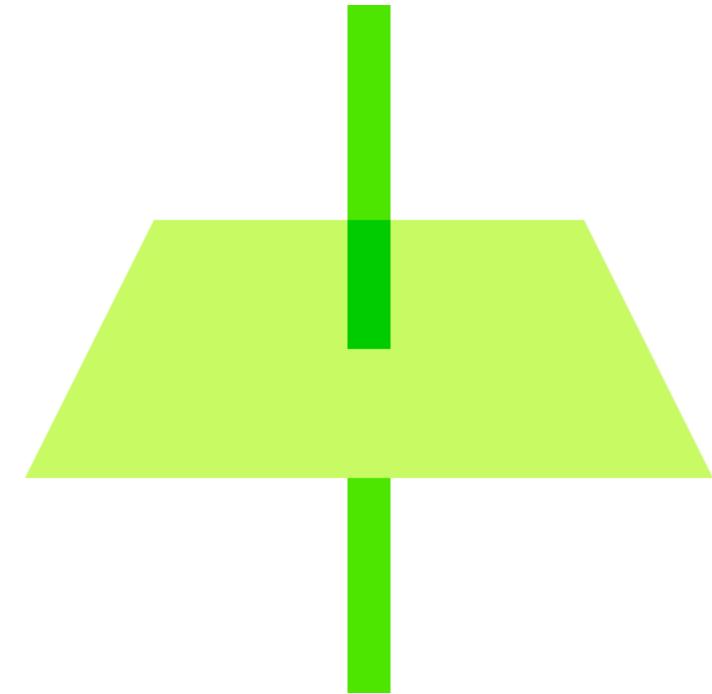
$v$ : total number of attribute values that  $X_i$  can take

$p$ : initial estimate of  $(P(X_i = c|y))$  known apriori

$m$ : hyper-parameter for our confidence in  $p$



# Support Vector Machines





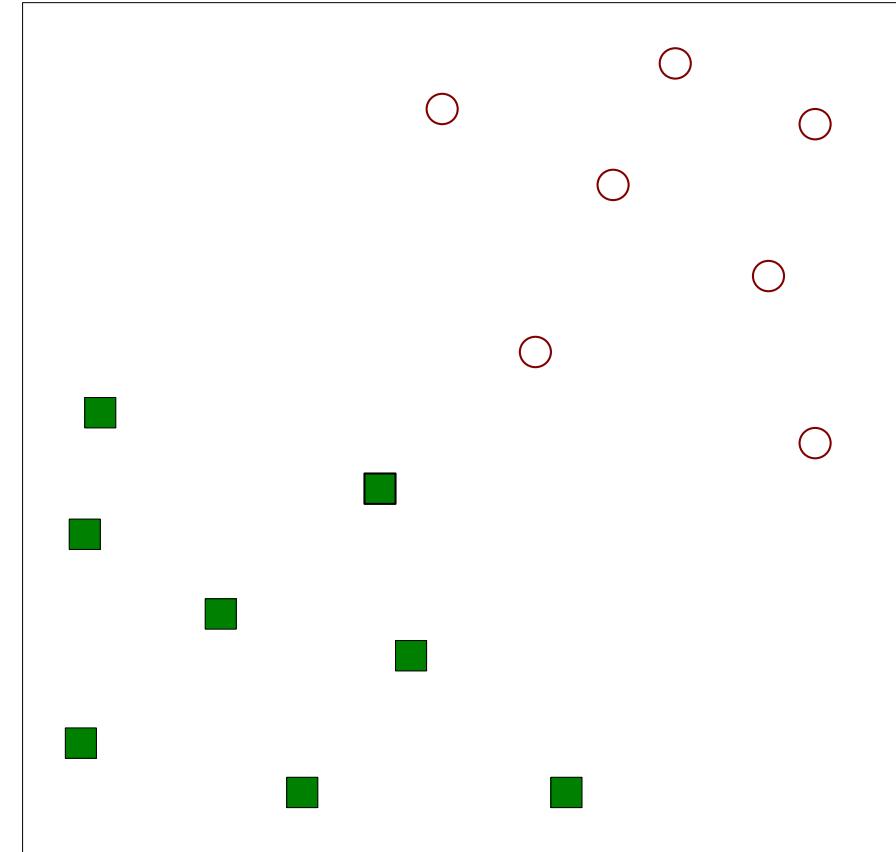
# Support Vector Machines

A support vector machine (SVM) is a discriminative classification model that learns linear or nonlinear decision boundaries in the attribute space to separate the classes.

## SVM

- strong regularization capabilities.
- good generalization performance.
- learn highly expressive models without suffering from overfitting.

**Find a linear hyperplane (decision boundary) that will separate the data**





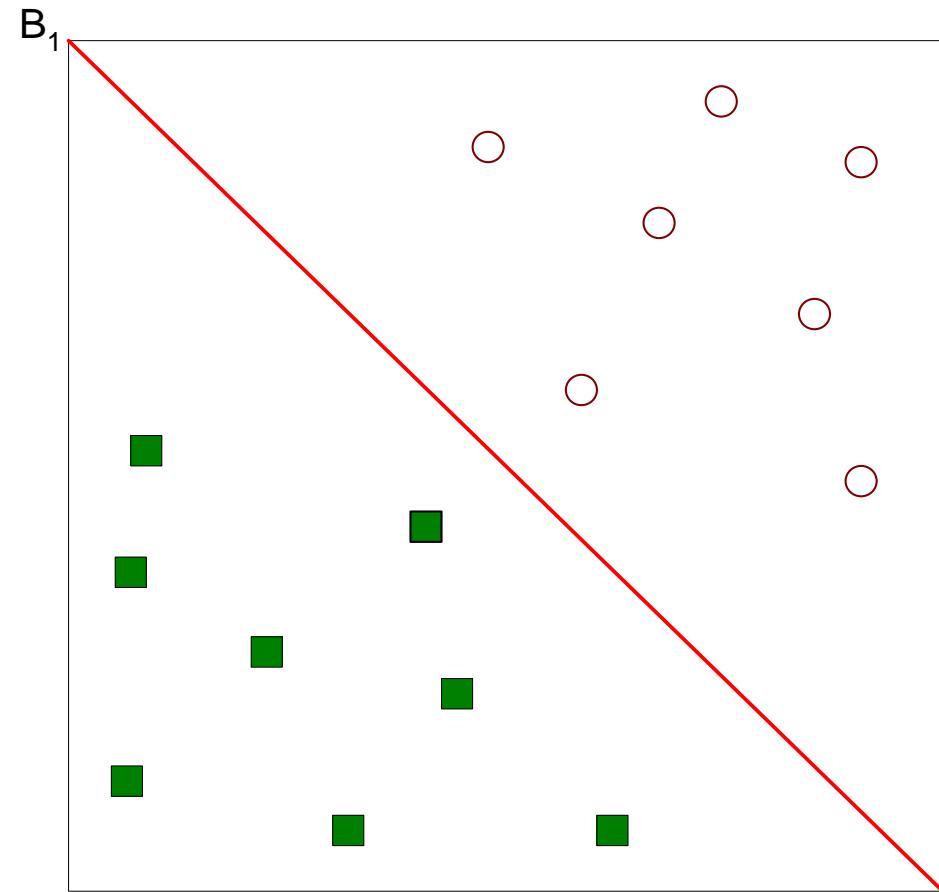
# Support Vector Machines

A support vector machine (SVM) is a discriminative classification model that learns linear or nonlinear decision boundaries in the attribute space to separate the classes.

## SVM

- strong regularization capabilities.
- good generalization performance.
- learn highly expressive models without suffering from overfitting.

One Possible Solution





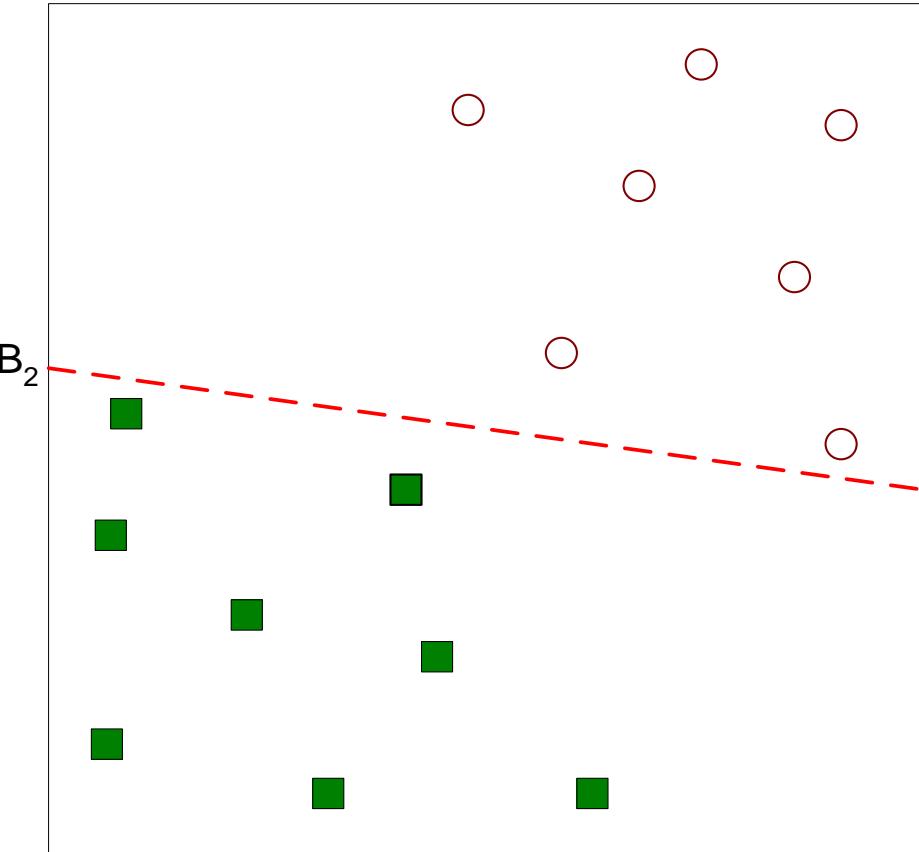
# Support Vector Machines

A support vector machine (SVM) is a discriminative classification model that learns linear or nonlinear decision boundaries in the attribute space to separate the classes.

## SVM

- strong regularization capabilities.
- good generalization performance.
- learn highly expressive models without suffering from overfitting.

Another Possible Solution





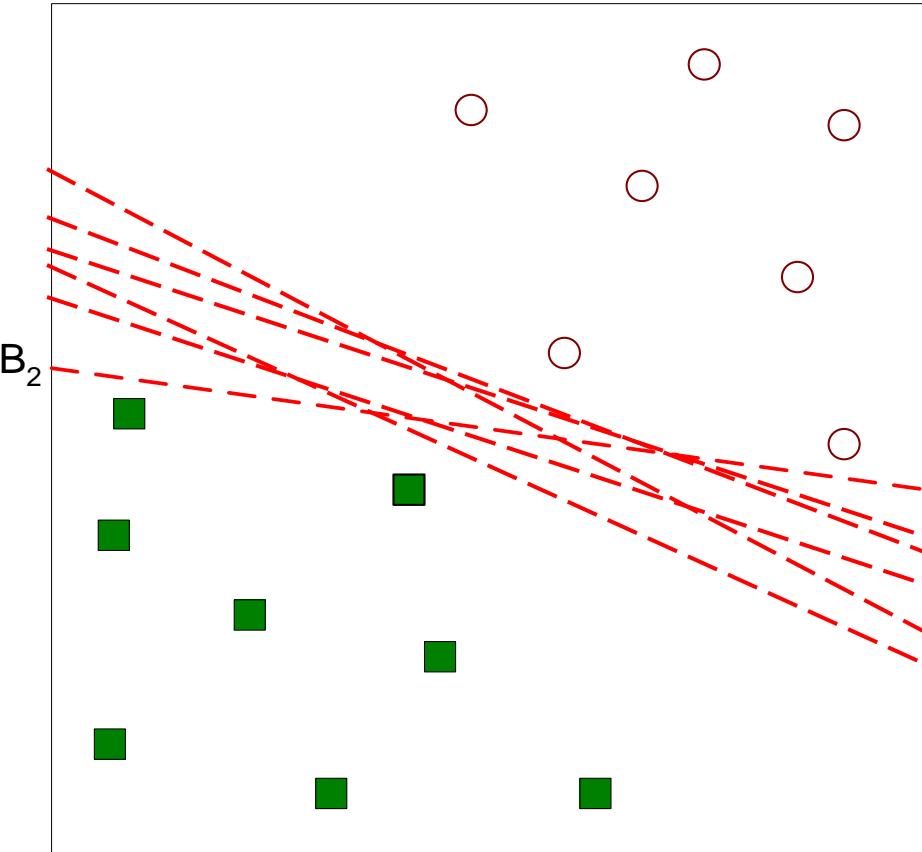
# Support Vector Machines

A support vector machine (SVM) is a discriminative classification model that learns linear or nonlinear decision boundaries in the attribute space to separate the classes.

## SVM

- strong regularization capabilities.
- good generalization performance.
- learn highly expressive models without suffering from overfitting.

Other Possible Solutions





# Support Vector Machines

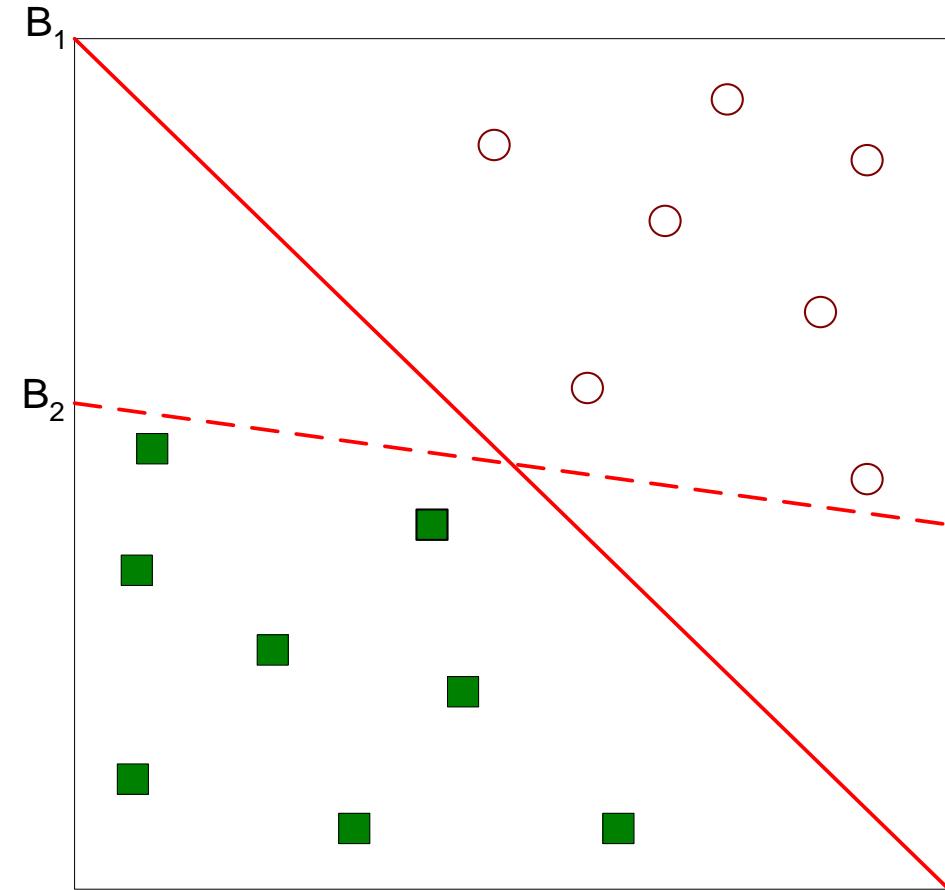
A support vector machine (SVM) is a discriminative classification model that learns linear or nonlinear decision boundaries in the attribute space to separate the classes.

## SVM

- strong regularization capabilities.
- good generalization performance.
- learn highly expressive models without suffering from overfitting.

Which one is better? B1 or B2?

How do you define better?



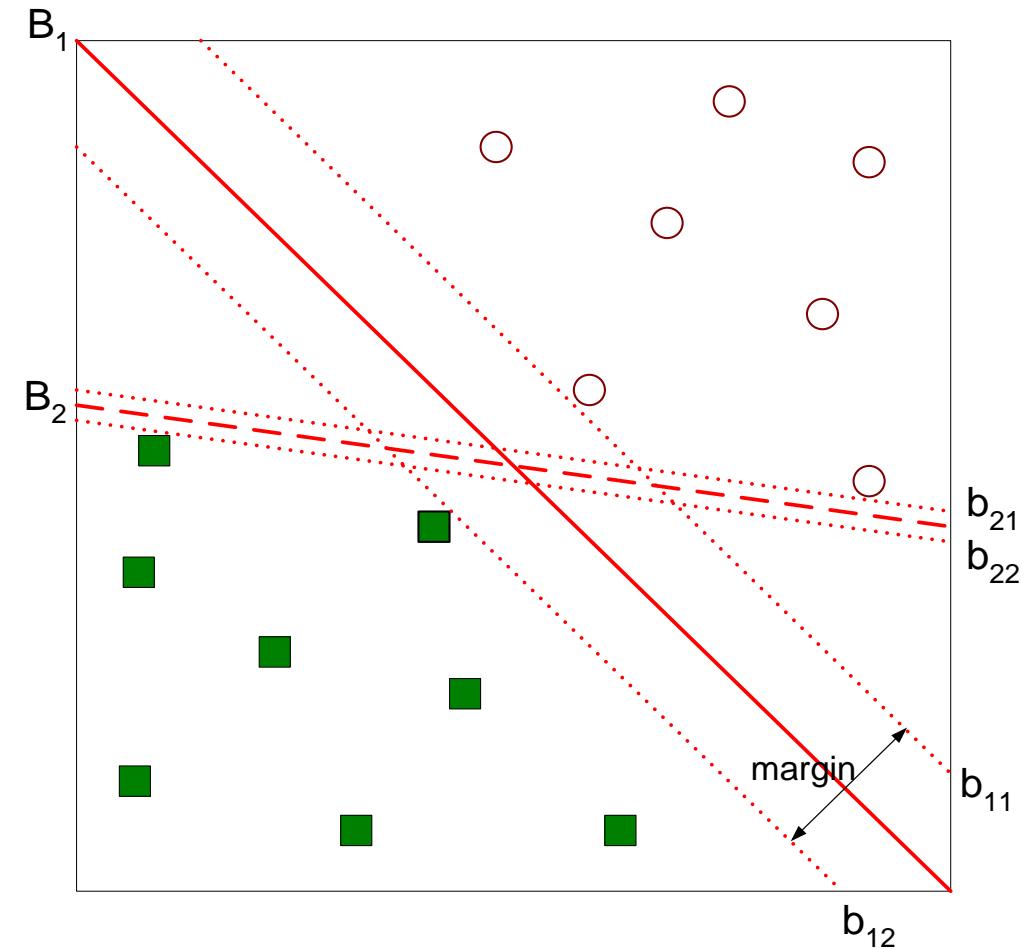


# Support Vector Machines

The point is now about the **margin of a separating hyperplane** and the rationale for choosing such a hyperplane with **maximum margin**.

Find hyperplane **maximizes** the margin

B1 is better than B2





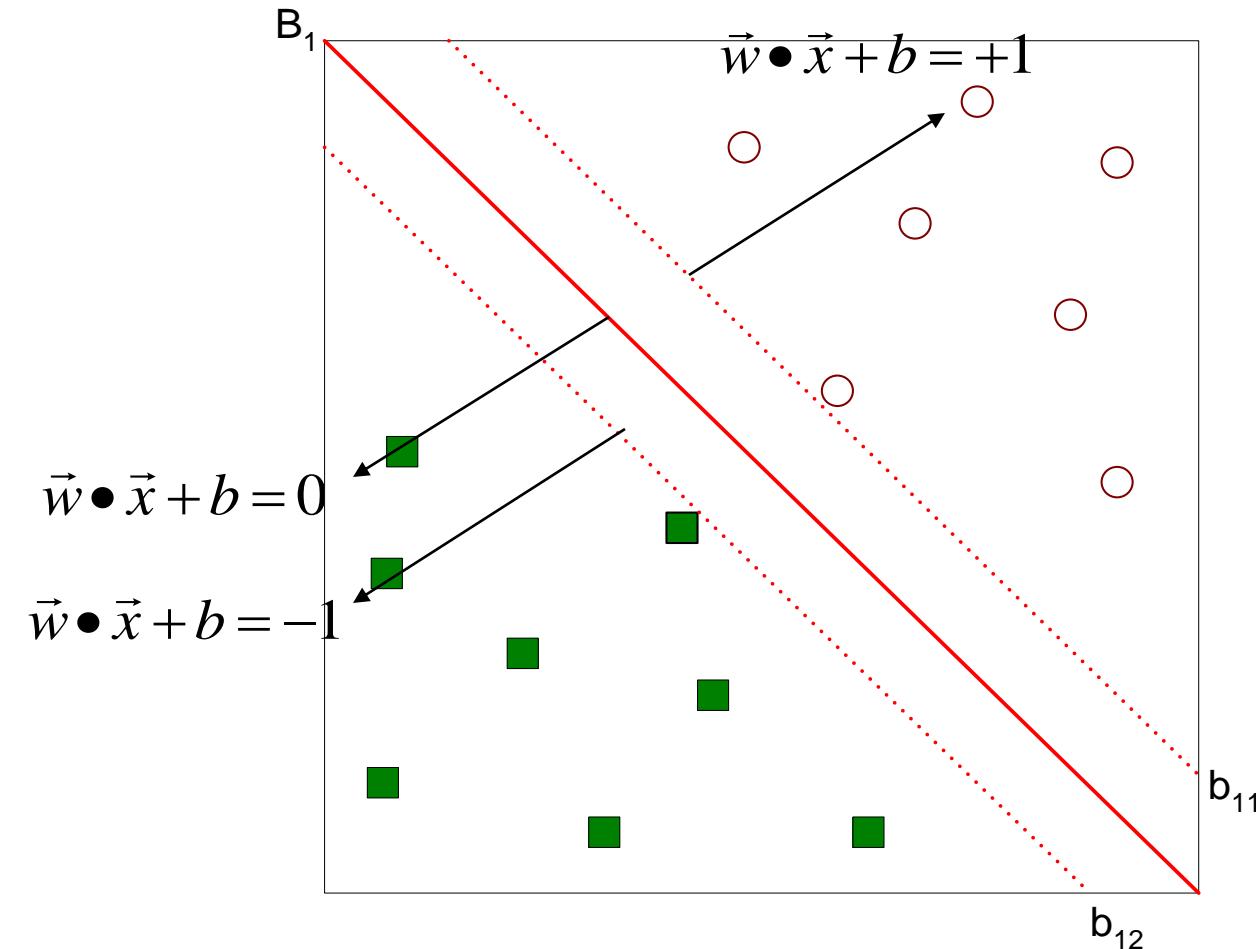
# Support Vector Machines

The generic equation of a separating hyperplane can be written as

$$\vec{w}^T \vec{x} + b = 0$$

where  $x$  represent the attributes and  $(w, b)$  represent the hyperplane.

A data instance  $x_i$  can belong to either side of the hyperplane depending on the sign of  $\vec{w}^T \vec{x} + b$





# Support Vector Machines

The hyperplane parameters  $(w, b)$  working on this condition:

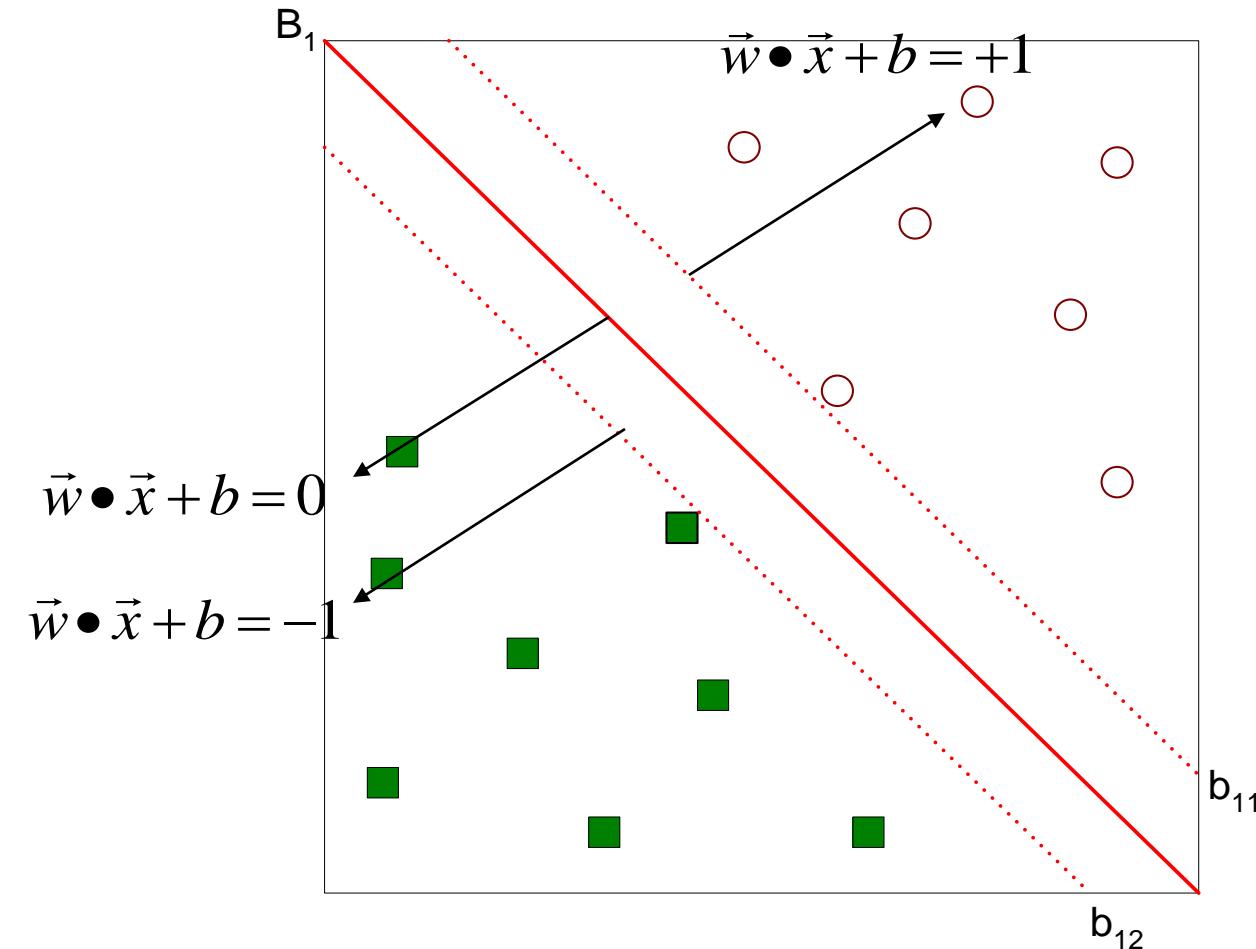
$$w^T x_i + b > 0 \quad \text{if } y_i = 1$$

$$w^T x_i + b < 0 \quad \text{if } y_i = -1$$

That led to:

$$\min_{w,b} \frac{\|w\|^2}{2}$$

subject to  $y_i(w^T x_i + b) > 1$





# Support Vector Machines

SVM to learn linear hyperplanes even in situations where the classes are not linearly separable.

To do this, SVM must consider the trade-off between the width of the margin and the number of training errors committed by the linear hyperplane.

The hyperplane parameters  $(w, b)$  working on this condition:

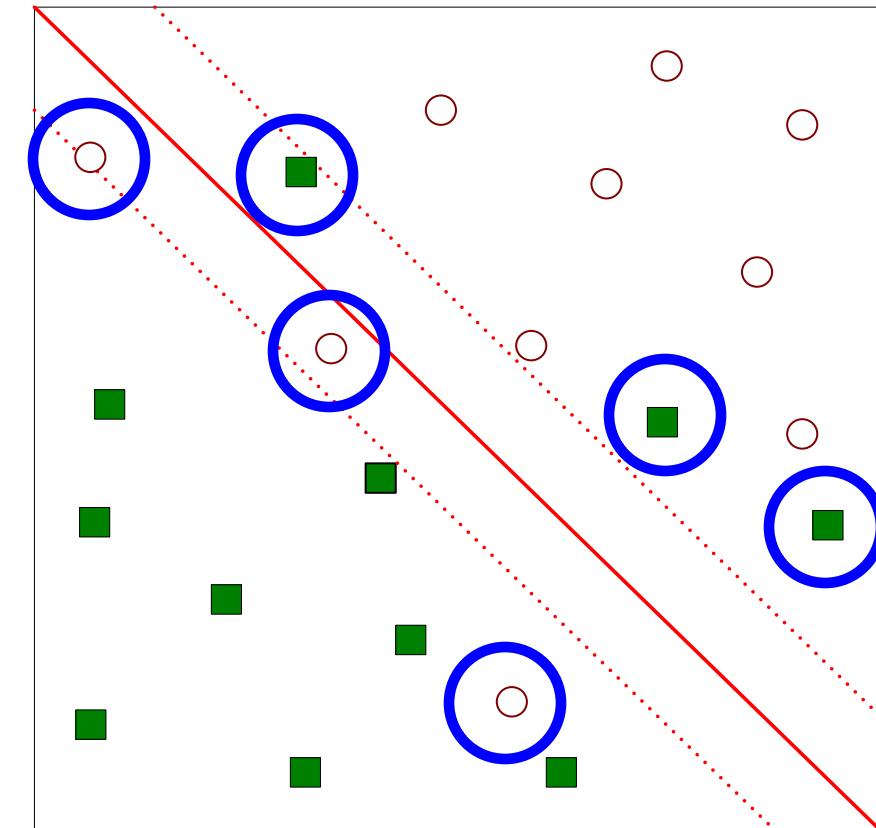
$$w^T x_i + b > 1 - \xi_i$$

That led to:

$$\begin{aligned} & \min_{w, b, \xi_i} \frac{\|w\|^2}{2} + C \sum_{i=1}^n \xi_i \\ \text{subject to } & y_i(w^T x_i + b) > 1 - \xi_i \\ & \xi_i > 0 \end{aligned}$$

C is a hyper-parameter that makes a **trade-off** between **maximizing** the **margin** and **minimizing** the **training error**.

A **large value of C** pays more emphasis on **minimizing the training error** than maximizing the margin

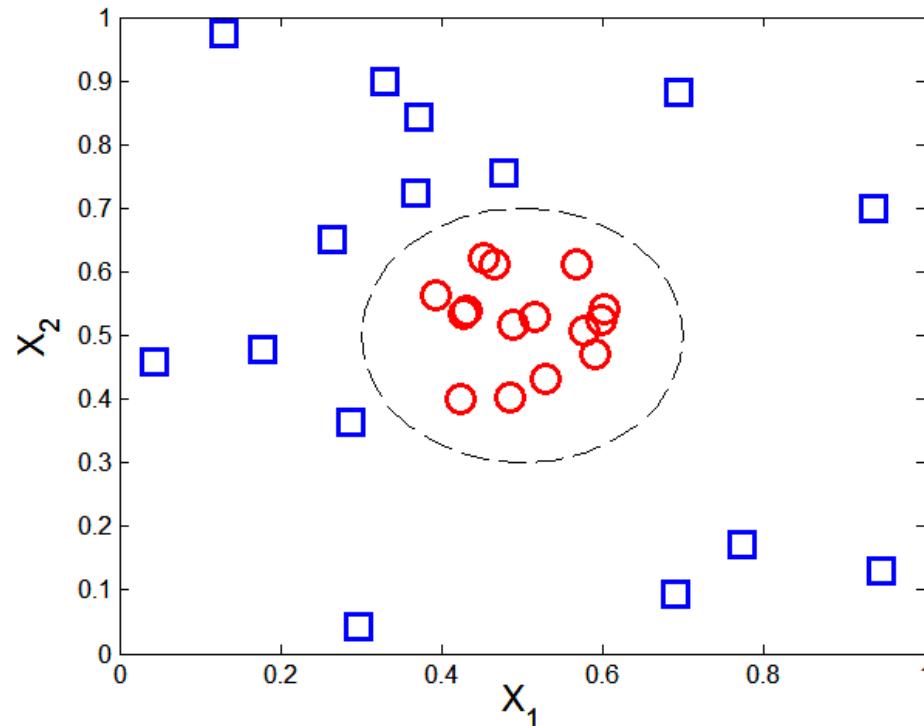




# Nonlinear Support Vector Machines

$$y(x_1, x_2) = \begin{cases} 1 & \text{if } \sqrt{(x_1 - 0.5)^2 + (x_2 - 0.5)^2} > 0.2 \\ -1 & \text{otherwise} \end{cases}$$

*What if decision boundary is not linear?*



A nonlinear transformation  $\varphi$  is needed to map the data from its original attribute space into a new **linear separable** space

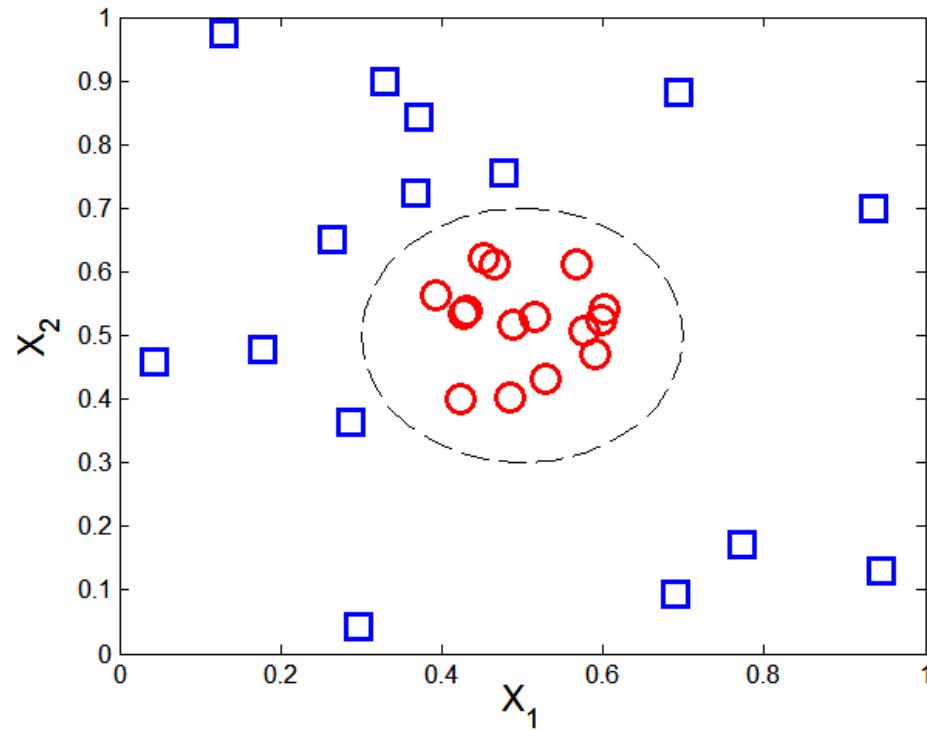
Transform the data from its **original attribute** space in  $x$  into a **new space**  $\varphi(x)$  so that a **linear hyperplane**

**Learned hyperplane** can then be projected back to the original attribute space, resulting in a **nonlinear decision boundary**.



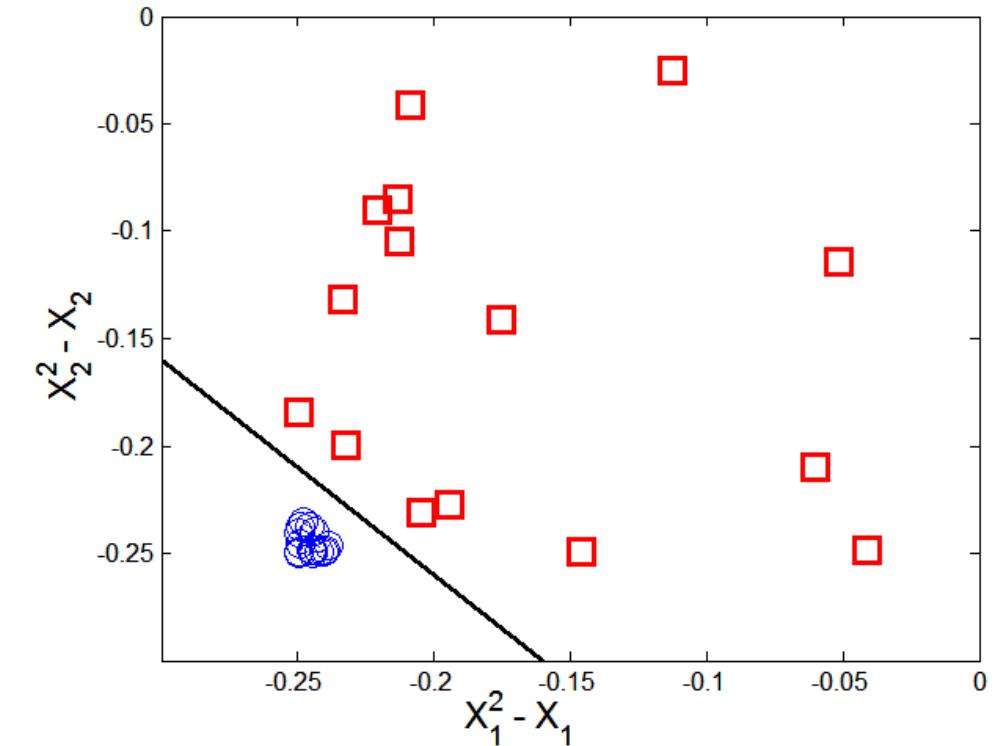
# Nonlinear Support Vector Machines

$$y(x_1, x_2) = \begin{cases} 1 & \text{if } \sqrt{(x_1 - 0.5)^2 + (x_2 - 0.5)^2} > 0.2 \\ -1 & \text{otherwise} \end{cases}$$



The transformation required is:

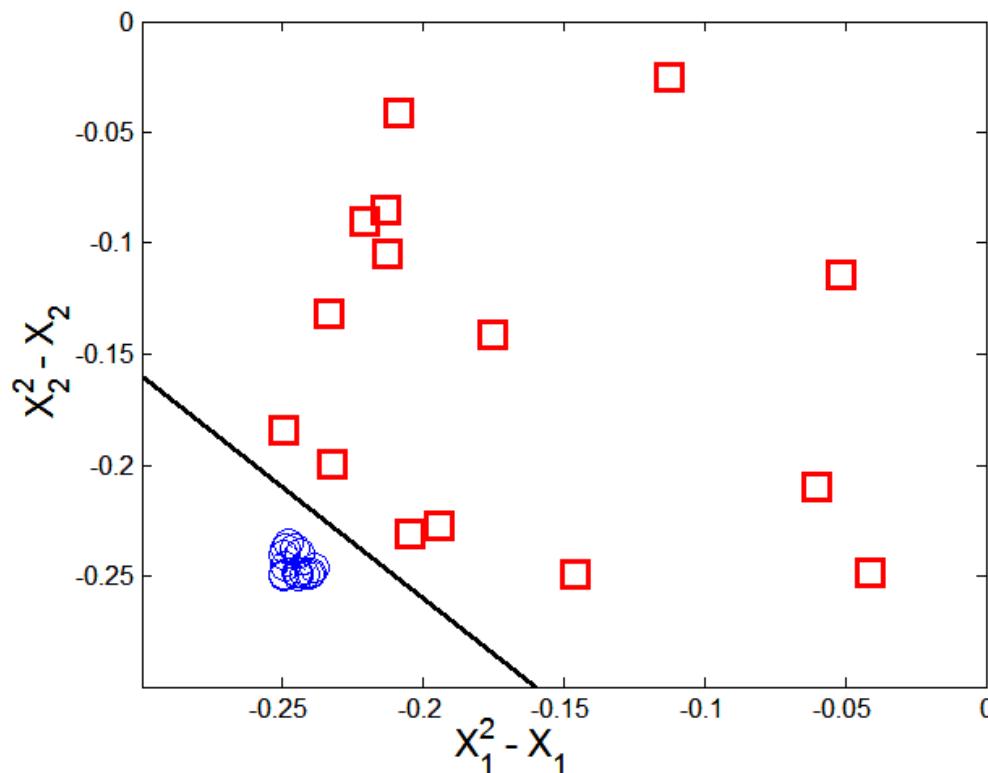
$$\varphi: (x_1, x_2) \rightarrow (x_1^2 - x_1, x_2^2 - x_2)$$





# Nonlinear Support Vector Machines

$$y(x_1, x_2) = \begin{cases} 1 & \text{if } \sqrt{(x_1 - 0.5)^2 + (x_2 - 0.5)^2} > 0.2 \\ -1 & \text{otherwise} \end{cases}$$



*Learning not linear model*

The new decision boundary is:

$$\mathbf{w}^T \varphi(\mathbf{x}_i) + b = 0$$

that results in the new learning

$$\min_{\mathbf{w}, b, \xi_i} \frac{\|\mathbf{w}\|^2}{2} + C \sum_{i=1}^n \xi_i$$

subject to

$$\begin{aligned} y_i (\mathbf{w}^T \varphi(\mathbf{x}_i) + b) &> 1 - \xi_i \\ \xi_i &> 0 \end{aligned}$$

We need only inner products of  $\varphi(\mathbf{x})$

- What type of mapping function  $\varphi$  should be used?
- How to do the computation in high dimensional space?





# Learning Nonlinear SVM

The **kernel trick** is a method for computing this similarity as a *function of the original attribute*

$$K(u, v) = \langle \varphi(u), \varphi(v) \rangle = f(u, v)$$

where  $f()$  is a function that follows certain conditions as stated by the Mercer's Theorem

Polynomial kernel       $K(\mathbf{u}, \mathbf{v}) = (\mathbf{u}^T \mathbf{v} + 1)^p$

Radial Basis Function kernel       $K(\mathbf{u}, \mathbf{v}) = e^{-\|\mathbf{u}-\mathbf{v}\|^2/(2\sigma^2)}$

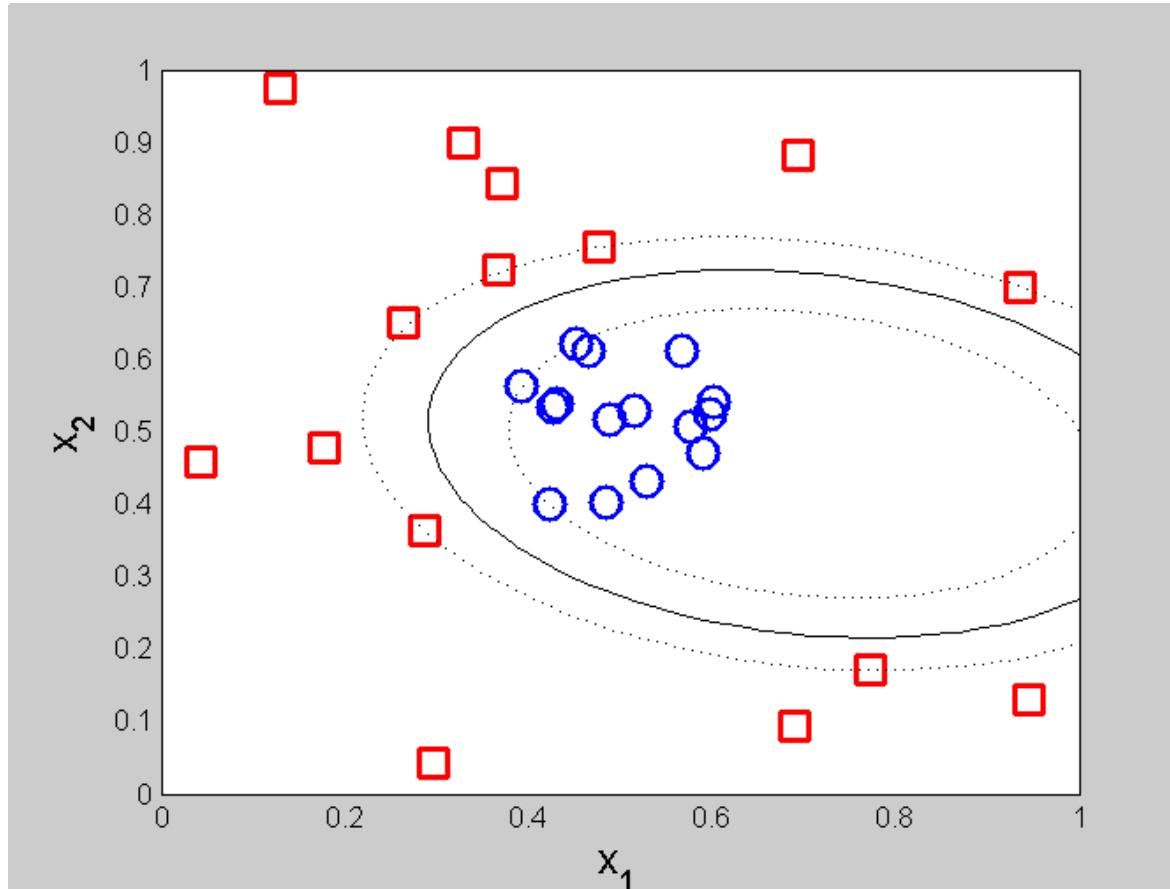
Sigmoid kernel       $K(\mathbf{u}, \mathbf{v}) = \tanh(k\mathbf{u}^T \mathbf{v} - \delta)$

By using a kernel function, we can directly work with inner products in the transformed space without dealing with the **exact forms of the nonlinear transformation function**

**High-dimensional transformations**, performing calculations only in the **original attribute space**



# Example of Nonlinear SVM



**SVM with polynomial degree 2 kernel**



# Learning Nonlinear SVM

---

Advantages of using kernel:

- Don't have to know the mapping function  $\varphi$
- Computing dot product  $\langle \varphi(u), \varphi(v) \rangle$  in the original space avoids curse of dimensionality

Not all functions can be kernels

- Must make sure there is a corresponding  $\Phi$  in some high-dimensional space
- Mercer's theorem (see textbook)

# Characteristics of SVM

---

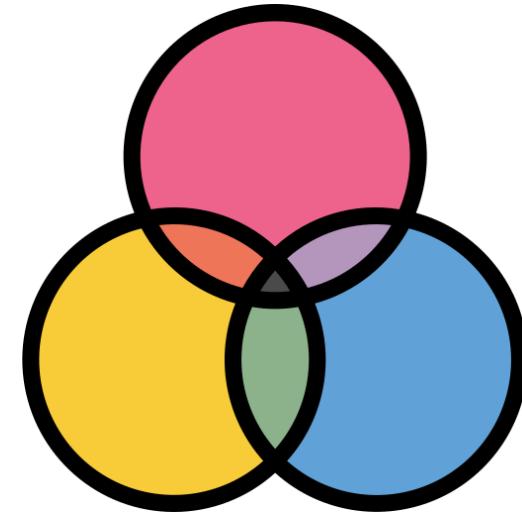


- Robust to noise
- Overfitting is handled by maximizing the margin of the decision boundary,
- SVM can handle irrelevant and redundant attributes better than many other techniques
- The user needs to provide the type of kernel function and cost function
- Difficult to handle missing values



---

# Ensemble Techniques

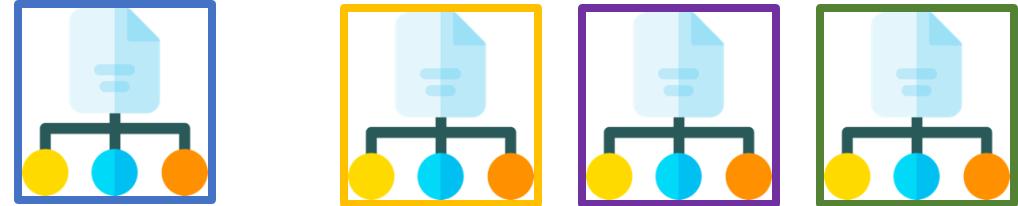


# Ensemble Methods

---



Construct a **set** of **base classifiers** learned from the training data



Predict class label of test records by **combining the predictions** made by multiple classifiers (e.g., by taking majority vote)

Combination rule

Prediction

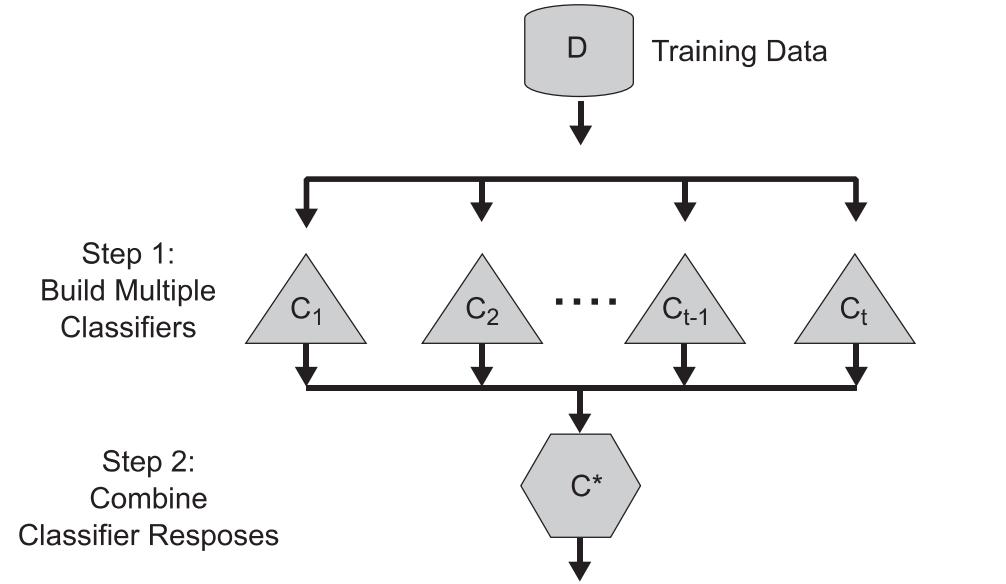


# General Approach of Ensemble Learning

The ensemble of classifiers can be constructed in many ways:

- By manipulating the training set
- By manipulating the input features
- By manipulating the class labels
- By manipulating the learning algorithm

The first three approaches are generic methods that are applicable to any classifier, whereas the fourth approach depends on the type of classifier used.



Using majority vote or weighted majority vote  
(weighted according to their accuracy or relevance)

$$C(x) = f(C_1(x), C_2(x), \dots, C_k(x))$$

# Example: Why Do Ensemble Methods Work?

---



- Suppose there are 25 base classifiers
  - Each classifier has error rate,  $\epsilon = 0.35$
  - Majority vote of classifiers used for classification
  - If all classifiers are identical:
    - ◆ Error rate of ensemble =  $\epsilon$  (0.35)
  - If all classifiers are independent (errors are uncorrelated):
    - ◆ Error rate of ensemble = probability of having more than half of base classifiers being wrong

$$e_{\text{ensemble}} = \sum_{i=13}^{25} \binom{25}{i} \epsilon^i (1 - \epsilon)^{25-i} = 0.06$$

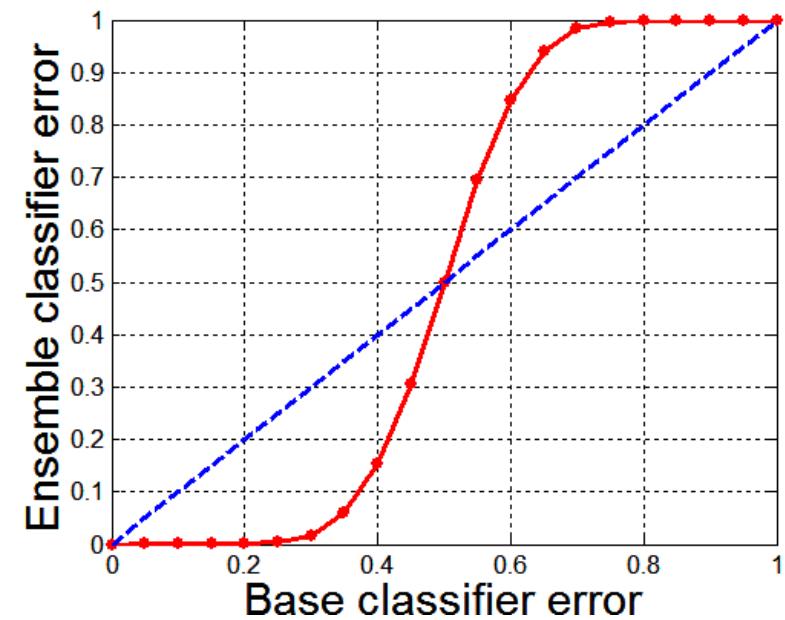
# Necessary Conditions for Ensemble Methods



Ensemble Methods work better than a single base classifier if:

1. All base classifiers are independent of each other
2. All base classifiers perform better than random guessing (error rate  $< 0.5$  for binary classification)

Classification error for an ensemble of 25 base classifiers, assuming their errors are uncorrelated.



# Rationale for Ensemble Learning



---

Ensemble Methods work best with **unstable base classifiers**

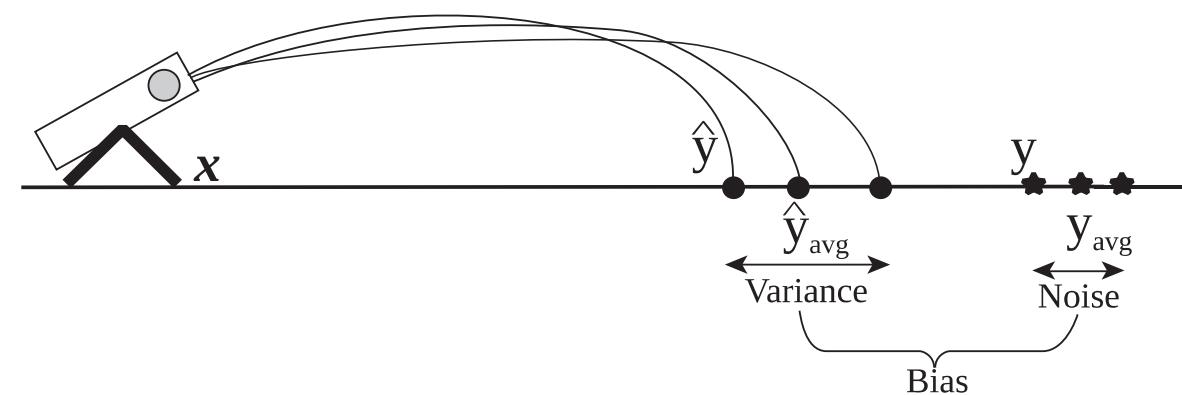
- Classifiers that are sensitive to minor perturbations in training set, due to *high model complexity*
- Examples: Unpruned decision trees, ANNs, ...



# Bias-Variance Decomposition

**Bias-variance decomposition is a formal method for analyzing the generalization error of a predictive model**

- **Noise** term is intrinsic to the target class
- **Bias** and **Variance** terms depend on the choice of the classification model
  - **Bias** of a model represents how close the average prediction of the model is to the average target
  - **Variance** of a model captures the stability of its predictions



For classification, the generalization error of model  $m$  can be given by:

$$\text{gen. error}(m) = c_1 + \text{bias}(m) + c_2 \times \text{variance}(m)$$

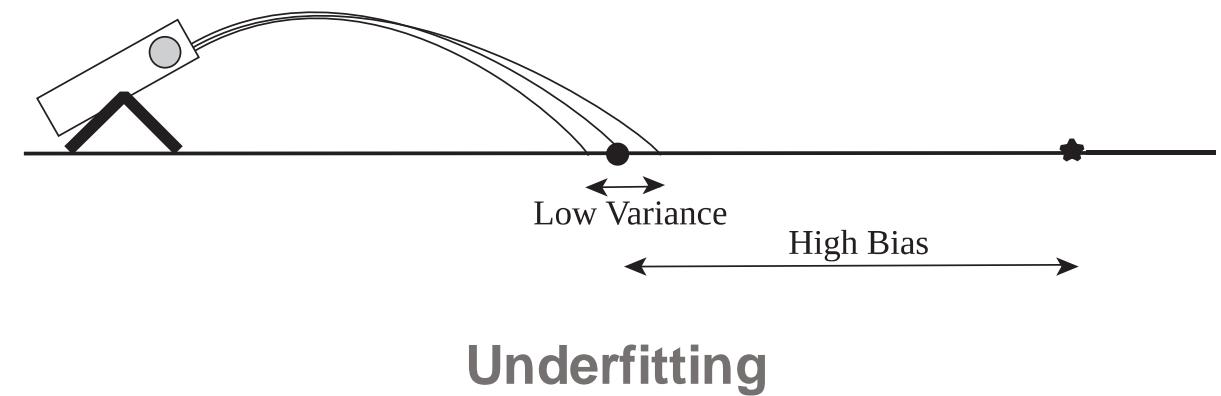
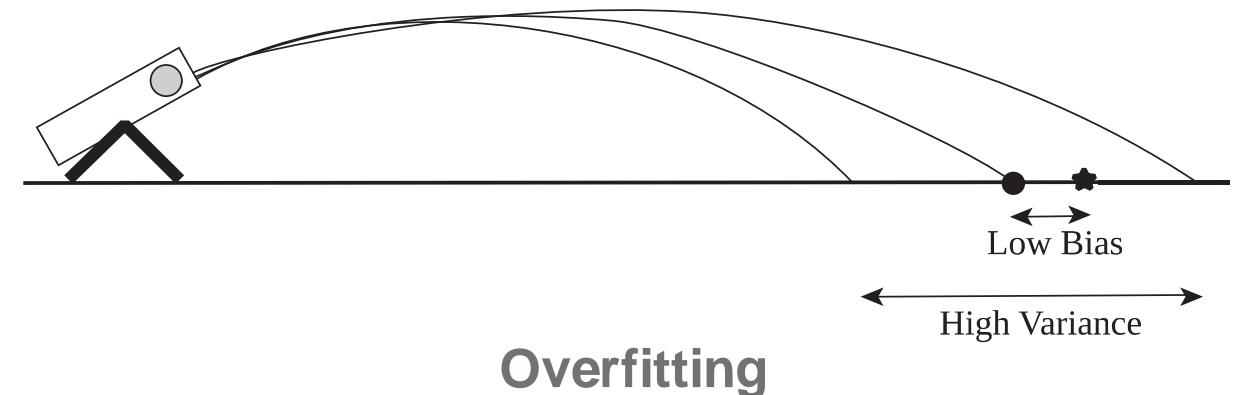
# Bias-Variance Trade-off and Overfitting



Ensemble methods try to **reduce the variance** of complex models (with low bias) by **aggregating** responses of multiple base classifiers (even help in reducing the bias)

Low bias but high variance, it can become susceptible to overfitting,

**Bias-Variance trade-off** can be used to explain why ensemble learning improves the generalization performance of unstable classifiers





# Bagging (Bootstrap AGGREGATING)

**Bagging** (AKA bootstrap aggregating), is a technique that repeatedly samples (with replacement) from a data set according to a uniform probability distribution

- Each bootstrap sample has the **same size** as the **original data**
- Some instances may **appear several times** in the same training set
- Probability of a training instance being selected in a bootstrap sample is:
  - $1 - (1 - 1/n)^n$  (*n: number of training instances*)
  - $\sim 0.632$  when *n* is large

Original Data	1	2	3	4	5	6	7	8	9	10
Bagging (Round 1)	7	8	10	8	2	5	10	10	5	9
Bagging (Round 2)	1	4	9	1	2	3	2	7	3	2
Bagging (Round 3)	1	8	5	10	5	5	9	6	3	7

Build classifier on each bootstrap sample



# Bagging Algorithm

---

---

## Algorithm 4.5 Bagging algorithm.

---

- 1: Let  $k$  be the number of bootstrap samples.
- 2: **for**  $i = 1$  to  $k$  **do**
- 3:   Create a bootstrap sample of size  $N$ ,  $D_i$ .
- 4:   Train a base classifier  $C_i$  on the bootstrap sample  $D_i$ .
- 5: **end for**
- 6:  $C^*(x) = \operatorname{argmax}_y \sum_i \delta(C_i(x) = y).$   
 $\{\delta(\cdot) = 1 \text{ if its argument is true and } 0 \text{ otherwise.}\}$

---

# Bagging Example



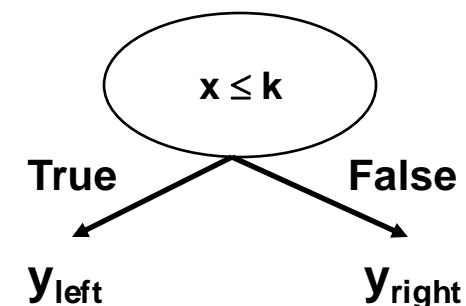
Consider 1-dimensional data set:

Original Data:

x	0.1	0.2	0.3	0.4	0.5	0.6	0.7	0.8	0.9	1
y	1	1	1	-1	-1	-1	-1	1	1	1

Classifier is a ***decision stump*** (decision tree of size 1)

- Decision rule:  $x \leq k$  versus  $x > k$
- Split point  $k$  is chosen based on entropy



# Bagging Example



Bagging Round 1:

x	0.1	0.2	0.2	0.3	0.4	0.4	0.5	0.6	0.9	0.9
y	1	1	1	1	-1	-1	-1	-1	1	1

Bagging Round 2:

x	0.1	0.2	0.3	0.4	0.5	0.5	0.9	1	1	1
y	1	1	1	-1	-1	-1	1	1	1	1

Bagging Round 3:

x	0.1	0.2	0.3	0.4	0.4	0.5	0.7	0.7	0.8	0.9
y	1	1	1	-1	-1	-1	-1	-1	1	1

Bagging Round 4:

x	0.1	0.1	0.2	0.4	0.4	0.5	0.5	0.7	0.8	0.9
y	1	1	1	-1	-1	-1	-1	-1	1	1

Bagging Round 5:

x	0.1	0.1	0.2	0.5	0.6	0.6	0.6	1	1	1
y	1	1	1	-1	-1	-1	-1	1	1	1

$x \leq 0.35 \rightarrow y = 1$   
 $x > 0.35 \rightarrow y = -1$

$x \leq 0.7 \rightarrow y = 1$   
 $x > 0.7 \rightarrow y = 1$

$x \leq 0.35 \rightarrow y = 1$   
 $x > 0.35 \rightarrow y = -1$

$x \leq 0.3 \rightarrow y = 1$   
 $x > 0.3 \rightarrow y = -1$

$x \leq 0.35 \rightarrow y = 1$   
 $x > 0.35 \rightarrow y = -1$

Bagging Round 6:

x	0.2	0.4	0.5	0.6	0.7	0.7	0.7	0.8	0.9	1
y	1	-1	-1	-1	-1	-1	-1	1	1	1

Bagging Round 7:

x	0.1	0.4	0.4	0.6	0.7	0.8	0.9	0.9	0.9	1
y	1	-1	-1	-1	-1	1	1	1	1	1

Bagging Round 8:

x	0.1	0.2	0.5	0.5	0.5	0.7	0.7	0.8	0.9	1
y	1	1	-1	-1	-1	-1	-1	1	1	1

Bagging Round 9:

x	0.1	0.3	0.4	0.4	0.6	0.7	0.7	0.8	1	1
y	1	1	-1	-1	-1	-1	-1	1	1	1

Bagging Round 10:

x	0.1	0.1	0.1	0.1	0.3	0.3	0.8	0.8	0.9	0.9
y	1	1	1	1	1	1	1	1	1	1

$x \leq 0.75 \rightarrow y = -1$   
 $x > 0.75 \rightarrow y = 1$

$x \leq 0.75 \rightarrow y = -1$   
 $x > 0.75 \rightarrow y = 1$

$x \leq 0.75 \rightarrow y = -1$   
 $x > 0.75 \rightarrow y = 1$

$x \leq 0.75 \rightarrow y = -1$   
 $x > 0.75 \rightarrow y = 1$

$x \leq 0.05 \rightarrow y = 1$   
 $x > 0.05 \rightarrow y = 1$



# Bagging Example

Summary of Trained Decision Stumps:

Round	Split Point	Left Class	Right Class
1	0.35	1	-1
2	0.7	1	1
3	0.35	1	-1
4	0.3	1	-1
5	0.35	1	-1
6	0.75	-1	1
7	0.75	-1	1
8	0.75	-1	1
9	0.75	-1	1
10	0.05	1	1



# Bagging Example

Use majority vote (sign of sum of predictions) to determine class of ensemble classifier

Round	x=0.1	x=0.2	x=0.3	x=0.4	x=0.5	x=0.6	x=0.7	x=0.8	x=0.9	x=1.0
1	1	1	1	-1	-1	-1	-1	-1	-1	-1
2	1	1	1	1	1	1	1	1	1	1
3	1	1	1	-1	-1	-1	-1	-1	-1	-1
4	1	1	1	-1	-1	-1	-1	-1	-1	-1
5	1	1	1	-1	-1	-1	-1	-1	-1	-1
6	-1	-1	-1	-1	-1	-1	-1	1	1	1
7	-1	-1	-1	-1	-1	-1	-1	1	1	1
8	-1	-1	-1	-1	-1	-1	-1	1	1	1
9	-1	-1	-1	-1	-1	-1	-1	1	1	1
10	1	1	1	1	1	1	1	1	1	1
Sum	2	2	2	-6	-6	-6	-6	2	2	2
Sign	1	1	1	-1	-1	-1	-1	1	1	1

Predicted Class

Bagging can also increase the complexity (representation capacity) of simple classifiers such as decision stumps



# Bagging Example

Use majority vote (sign of sum of predictions) to determine class of ensemble classifier

Round	x=0.1	x=0.2	x=0.3	x=0.4	x=0.5	x=0.6	x=0.7	x=0.8	x=0.9	x=1.0
1	1	1	1	-1	-1	-1	-1	-1	-1	-1
2	1	1	1	1	1	1	1	1	1	1
3	1	1	1	-1	-1	-1	-1	-1	-1	-1
4	1	1	1	-1	-1	-1	-1	-1	-1	-1
5	1	1	1	-1	-1	-1	-1	-1	-1	-1
6	-1	-1	-1	-1	-1	-1	-1	1	1	1
7	-1	-1	-1	-1	-1	-1	-1	1	1	1
8	-1	-1	-1	-1	-1	-1	-1	1	1	1
9	-1	-1	-1	-1	-1	-1	-1	1	1	1
10	1	1	1	1	1	1	1	1	1	1
Sum	2	2	2	-6	-6	-6	-6	2	2	2
Sign	1	1	1	-1	-1	-1	-1	1	1	1

Predicted Class

Bagging can also increase the complexity (representation capacity) of simple classifiers such as decision stumps



# Bagging Example

Use majority vote (sign of sum of predictions) to determine class of ensemble classifier

Round	x=0.1	x=0.2	x=0.3	x=0.4	x=0.5	x=0.6	x=0.7	x=0.8	x=0.9	x=1.0
1	1	1	1	-1	-1	-1	-1	-1	-1	-1
2	1	1	1	1	1	1	1	1	1	1
3	1	1	1	-1	-1	-1	-1	-1	-1	-1
4	1	1	1	-1	-1	-1	-1	-1	-1	-1
5	1	1	1	-1	-1	-1	-1	-1	-1	-1
6	-1	-1	-1	-1	-1	-1	-1	1	1	1
7	-1	-1	-1	-1	-1	-1	-1	1	1	1
8	-1	-1	-1	-1	-1	-1	-1	1	1	1
9	-1	-1	-1	-1	-1	-1	-1	1	1	1
10	1	1	1	1	1	1	1	1	1	1
Sum	2	2	2	-6	-6	-6	-6	2	2	2
Sign	1	1	1	-1	-1	-1	-1	1	1	1

**Predicted Class**

x	0.1	0.2	0.3	0.4	0.5	0.6	0.7	0.8	0.9	1
y	1	1	1	-1	-1	-1	-1	1	1	1

**Predicted Class**

Bagging can also increase the complexity (representation capacity) of simple classifiers such as decision stumps



# Boosting

---

Boosting is an iterative procedure used to adaptively change the distribution of training examples for learning base classifiers so that they increasingly focus on examples that are hard to classify

Boosting assigns a weight to each training example and may adaptively change the weight at the end of each boosting round

1. They can be used to inform the sampling distribution used to draw a set of bootstrap samples from the original data.
2. They can be used to learn a model that is biased toward examples with higher weight.

Examples that are classified incorrectly will have their weights increased, while those that are classified correctly will have their weights decreased

This forces the classifier to focus on examples that are difficult to classify in subsequent iterations

# AdaBoost



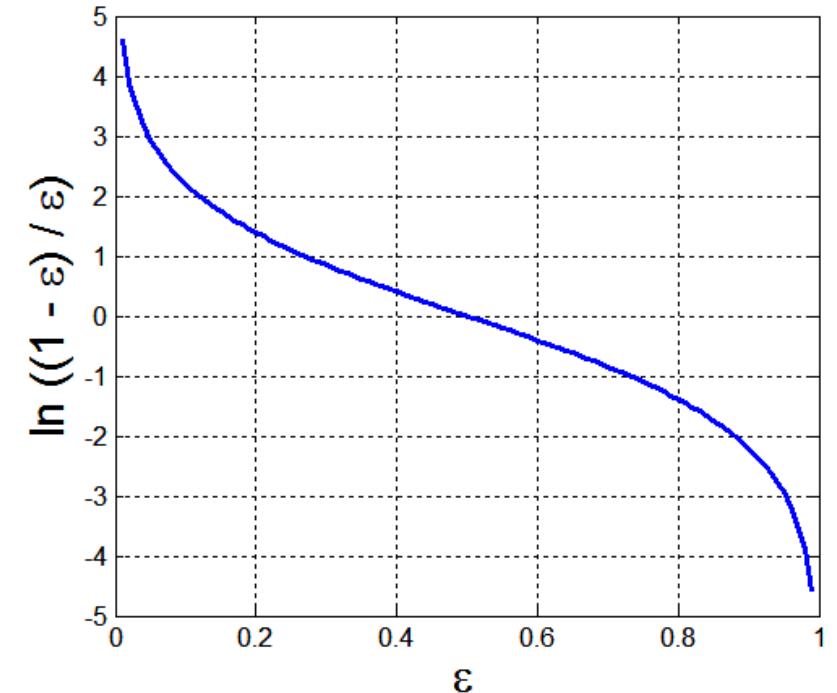
Consider a dataset  $\{(x_i, y_i) | j = 1, 2, \dots, N\}$

In **Adaboost** the **importance** of a classifier  $C_i$  depends on the error rate :

$$\epsilon_i = \frac{1}{N} \sum_{j=1}^N w_j^{(i)} \delta(C_i(x_j) \neq y_j)$$

The **importance** of a base classifier is **given** by

$$\alpha_i = \frac{1}{2} \ln \left( \frac{1 - \epsilon_i}{\epsilon_i} \right)$$





# AdaBoost Algorithm

---

**Weight update:**

$$w_j^{(i+1)} = \frac{w_j^{(i)}}{Z_i} \times \begin{cases} e^{-\alpha_i} & \text{if } C_i(x_j) = y_j \\ e^{\alpha_i} & \text{if } C_i(x_j) \neq y_j \end{cases}$$

Where  $Z_i$  is the normalization factor

If any intermediate rounds produce error rate higher than 50%, the weights are reverted back to  $1/n$  and the resampling procedure is repeated

$$C^*(x) = \arg \max_y \sum_{i=1}^T \alpha_i \delta(C_i(x) = y)$$



# AdaBoost Algorithm

---

**Algorithm 4.6** AdaBoost algorithm.

```
1:  $\mathbf{w} = \{w_j = 1/N \mid j = 1, 2, \dots, N\}$ . {Initialize the weights for all  $N$  examples.}
2: Let  $k$  be the number of boosting rounds.
3: for  $i = 1$  to  $k$  do
4:   Create training set  $D_i$  by sampling (with replacement) from  $D$  according to  $\mathbf{w}$ .
5:   Train a base classifier  $C_i$  on  $D_i$ .
6:   Apply  $C_i$  to all examples in the original training set,  $D$ .
7:    $\epsilon_i = \frac{1}{N} \left[ \sum_j w_j \delta(C_i(x_j) \neq y_j) \right]$  {Calculate the weighted error.}
8:   if  $\epsilon_i > 0.5$  then
9:      $\mathbf{w} = \{w_j = 1/N \mid j = 1, 2, \dots, N\}$ . {Reset the weights for all  $N$  examples.}
10:    Go back to Step 4.
11:   end if
12:    $\alpha_i = \frac{1}{2} \ln \frac{1-\epsilon_i}{\epsilon_i}$ .
13:   Update the weight of each example according to Equation 4.103.
14: end for
15:  $C^*(\mathbf{x}) = \operatorname{argmax}_y \sum_{j=1}^T \alpha_j \delta(C_j(\mathbf{x}) = y)$ .
```

---

# AdaBoost Example



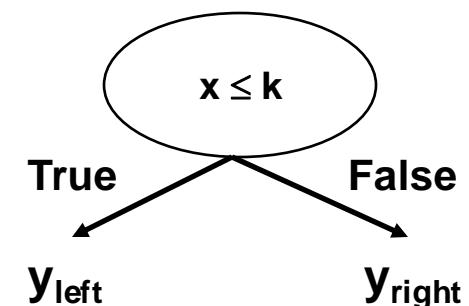
Consider 1-dimensional data set:

Original Data:

x	0.1	0.2	0.3	0.4	0.5	0.6	0.7	0.8	0.9	1
y	1	1	1	-1	-1	-1	-1	1	1	1

Classifier is a ***decision stump*** (decision tree of size 1)

- Decision rule:  $x \leq k$  versus  $x > k$
- Split point  $k$  is chosen based on entropy





# AdaBoost Example

Training sets for the first 3 boosting rounds:

Boosting Round 1:

x	0.1	0.4	0.5	0.6	0.6	0.7	0.7	0.7	0.8	1
y	1	-1	-1	-1	-1	-1	-1	-1	1	1

Boosting Round 2:

x	0.1	0.1	0.2	0.2	0.2	0.2	0.3	0.3	0.3	0.3
y	1	1	1	1	1	1	1	1	1	1

Boosting Round 3:

x	0.2	0.2	0.4	0.4	0.4	0.4	0.5	0.6	0.6	0.7
y	1	1	-1	-1	-1	-1	-1	-1	-1	-1

Weights

Round	x=0.1	x=0.2	x=0.3	x=0.4	x=0.5	x=0.6	x=0.7	x=0.8	x=0.9	x=1.0
1	0.1	0.1	0.1	0.1	0.1	0.1	0.1	0.1	0.1	0.1
2	0.311	0.311	0.311	0.01	0.01	0.01	0.01	0.01	0.01	0.01
3	0.029	0.029	0.029	0.228	0.228	0.228	0.228	0.009	0.009	0.009



# AdaBoost Example

## Alpha

Round	Split Point	Left Class	Right Class	alpha
1	0.75	-1	1	1.738
2	0.05	1	1	2.7784
3	0.3	1	-1	4.1195

## Classification

Predicted  
Class

Round	x=0.1	x=0.2	x=0.3	x=0.4	x=0.5	x=0.6	x=0.7	x=0.8	x=0.9	x=1.0
1	-1	-1	-1	-1	-1	-1	-1	1	1	1
2	1	1	1	1	1	1	1	1	1	1
3	1	1	1	-1	-1	-1	-1	-1	-1	-1
Sum	5.16	5.16	5.16	-3.08	-3.08	-3.08	-3.08	0.397	0.397	0.397
Sign	1	1	1	-1	-1	-1	-1	1	1	1



# Random Forest Algorithm

---

Construct an ensemble of **decision trees** by manipulating training set as well as features

- Use bootstrap sample to train every decision tree (similar to Bagging)
- Use the following tree induction algorithm:
  - At every internal node of decision tree, randomly sample  $p$  attributes for selecting split criterion
  - Repeat this procedure until all leaves are pure (unpruned tree)

# Random Forest Algorithm

---

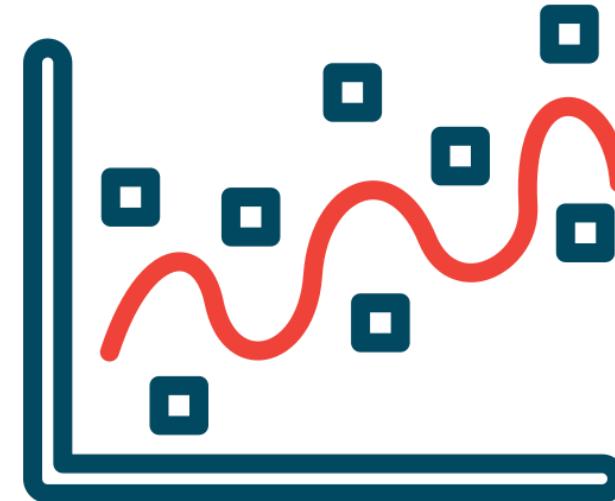


Construct an ensemble of **decision trees** by manipulating training set as well as features

- Base classifier are unpruned trees and hence are unstable classifier
- Base classifier are uncorrelated (due to randomization in training and features)
- Random forest reduce variance of unstable classifier without impacting the bias



# Model Overfitting





# Overfitting

---

If a model fits well over the training data, it can still show poor generalization performance, a phenomenon known as Model **Overfitting**

## Reasons for Model Overfitting

Model overfitting occurs when, in the pursuit of minimizing the training error rate, an overly complex model is selected that captures specific patterns in the training data but fails to learn the true nature of relationships between attributes and class labels in the overall data.

---

The opposite is Model **Underfitting**: the learned model is too simplistic, and thus, incapable of fully representing the true relationship between the attributes and the class labels.



# Classification Errors

**Training errors:** Errors committed on the training set

**Test errors:** Errors committed on the test set

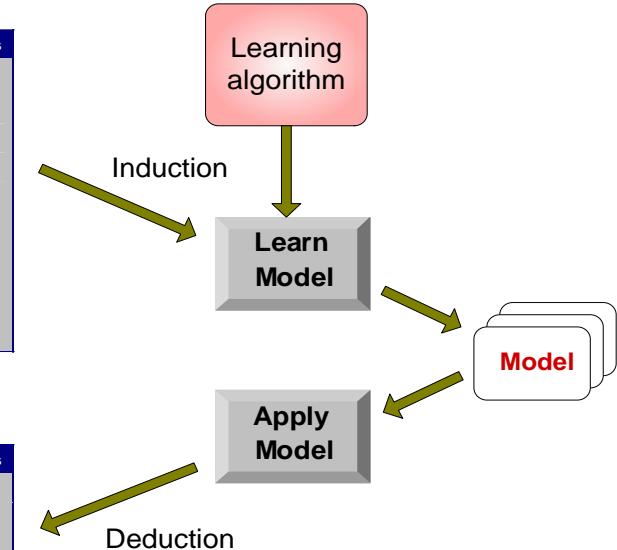
**Generalization errors:** Expected error of a model over random selection of records from same distribution

Tid	Attrib1	Attrib2	Attrib3	Class
1	Yes	Large	125K	No
2	No	Medium	100K	No
3	No	Small	70K	No
4	Yes	Medium	120K	No
5	No	Large	95K	Yes
6	No	Medium	60K	No
7	Yes	Large	220K	No
8	No	Small	85K	Yes
9	No	Medium	75K	No
10	No	Small	90K	Yes

Training Set

Tid	Attrib1	Attrib2	Attrib3	Class
11	No	Small	55K	?
12	Yes	Medium	80K	?
13	Yes	Large	110K	?
14	No	Small	95K	?
15	No	Large	67K	?

Test Set



# Example Data Set: two class problem

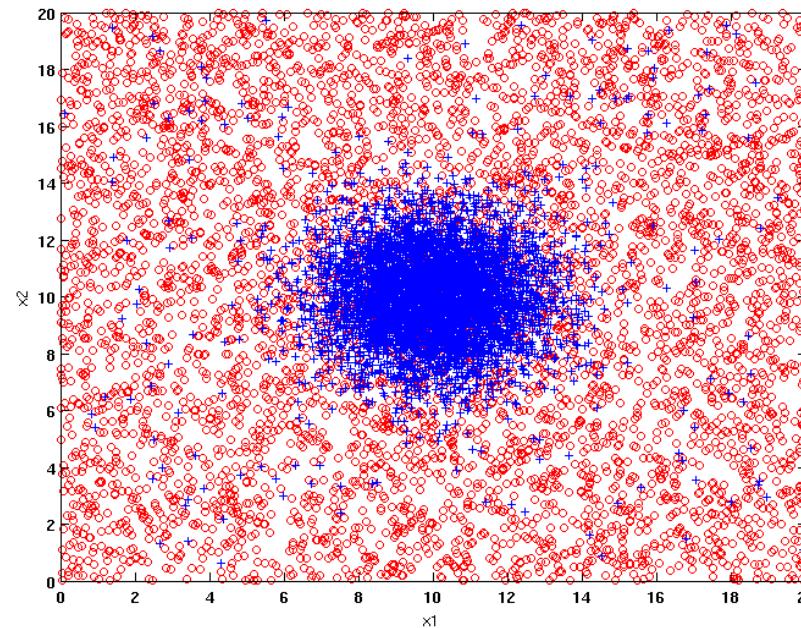


+ : 5400 instances

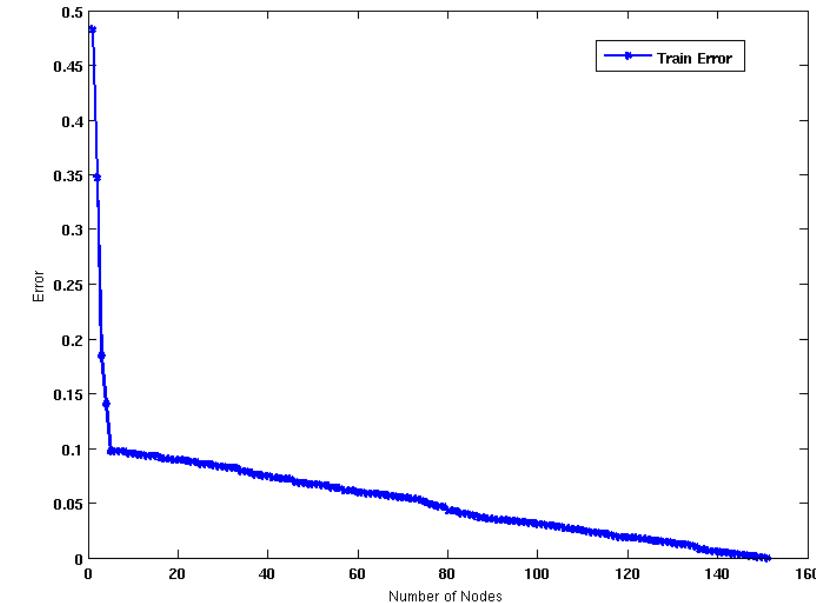
- 5000 instances generated from a Gaussian centered at (10,10)
- 400 noisy instances added

o : 5400 instances

- Generated from a uniform distribution



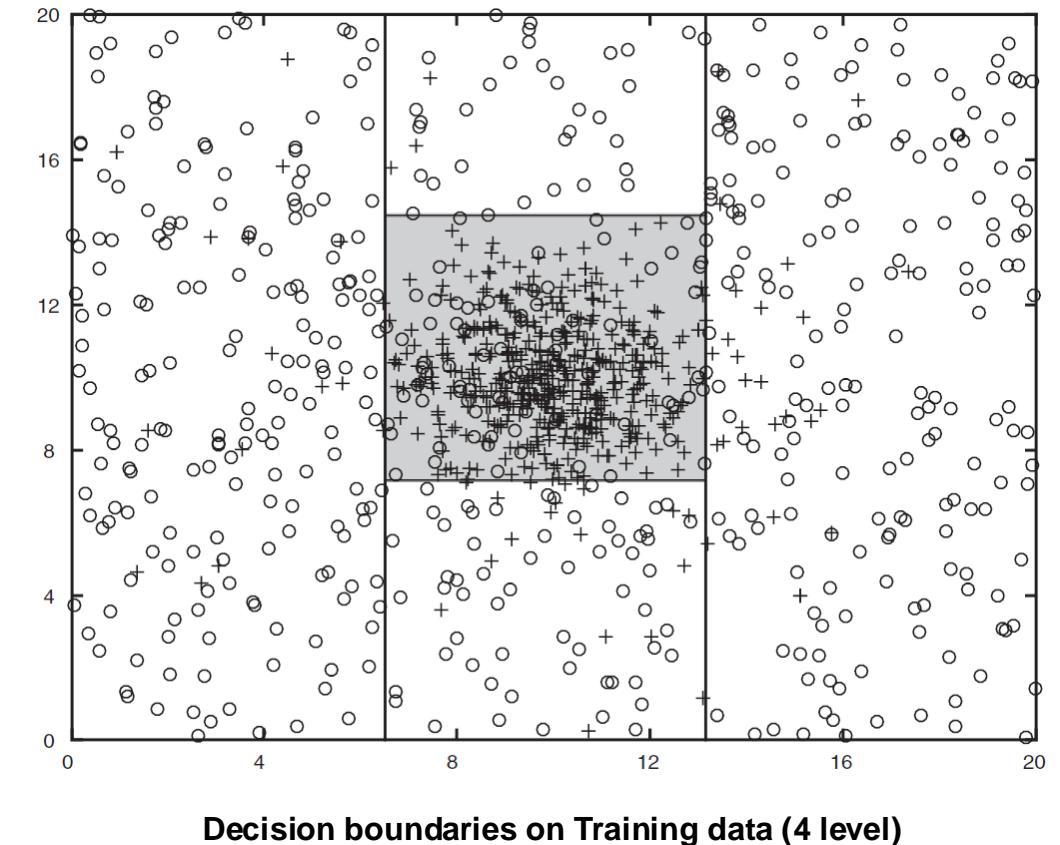
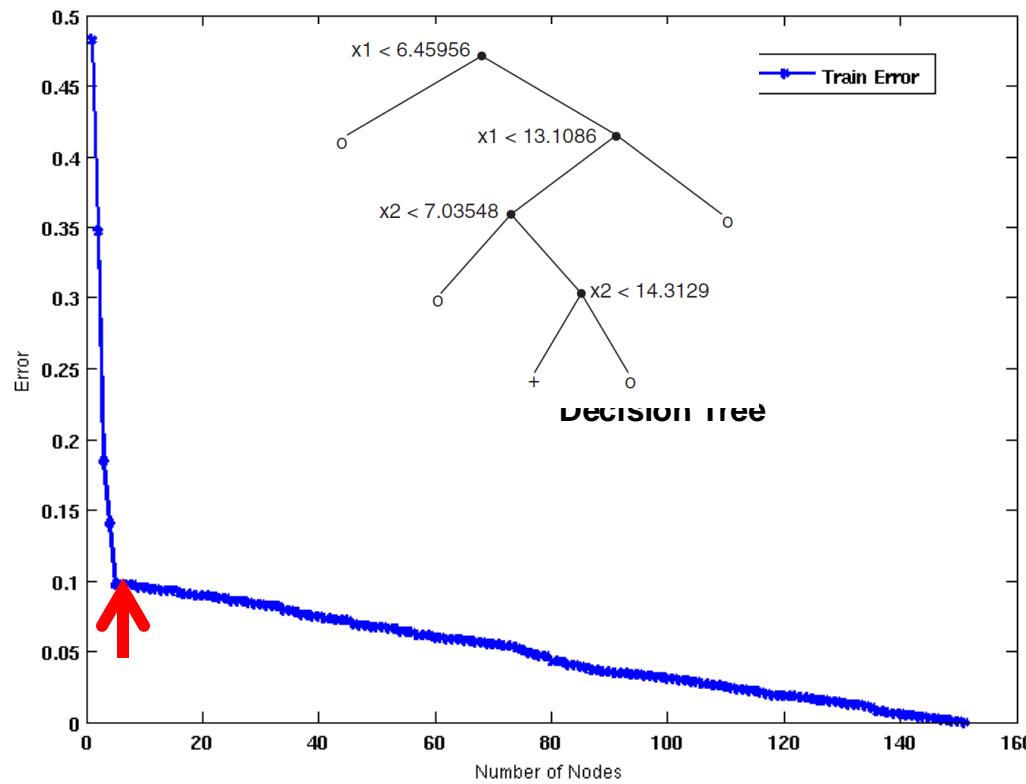
**10 % of the data used for training  
and 90% of the data used for  
testing**





# Decision Tree with 4 nodes

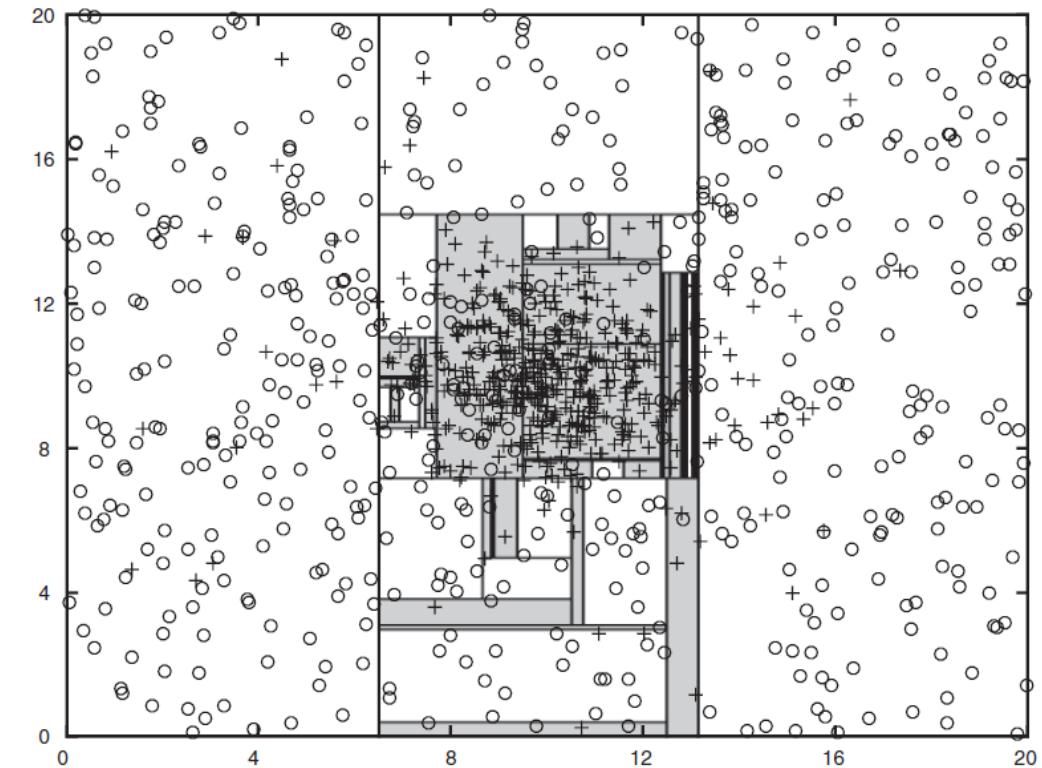
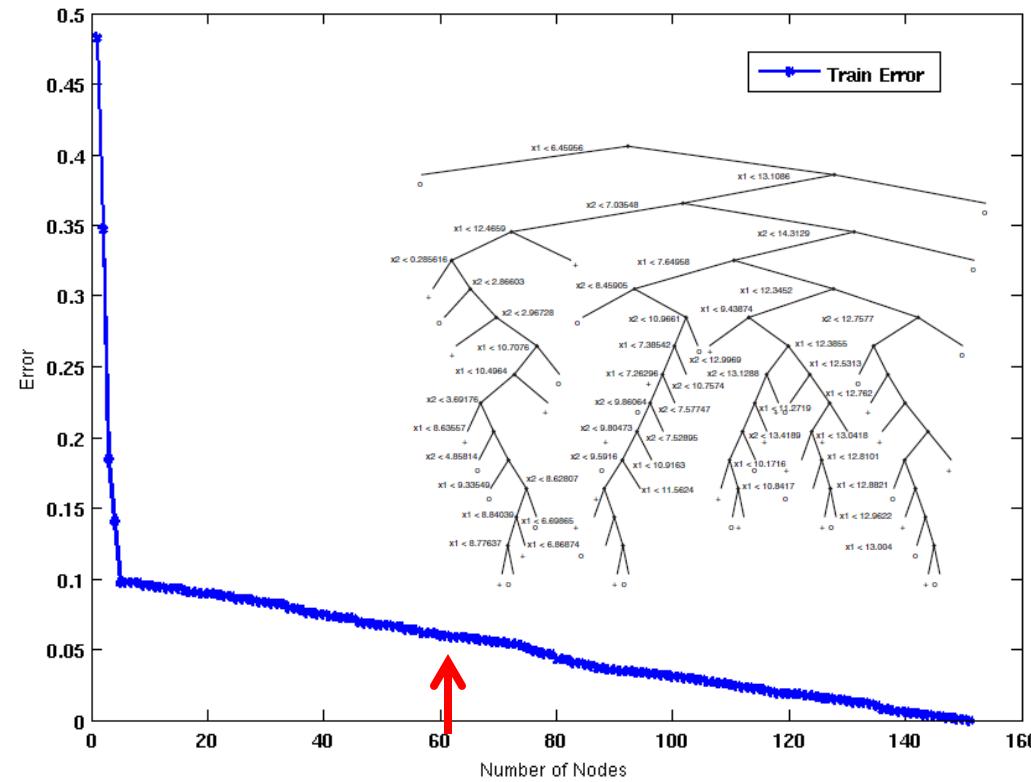
A decision tree made up of **few levels** led to simple decision boundaries





# Decision Tree with 50 nodes

A decision tree made up of **many levels** led to **complex** decision boundaries

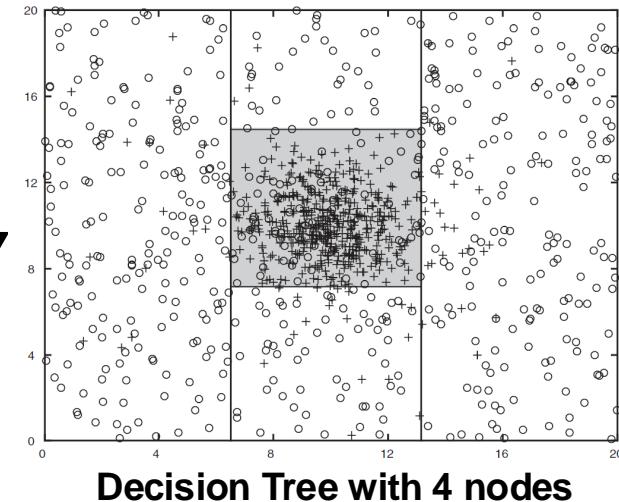
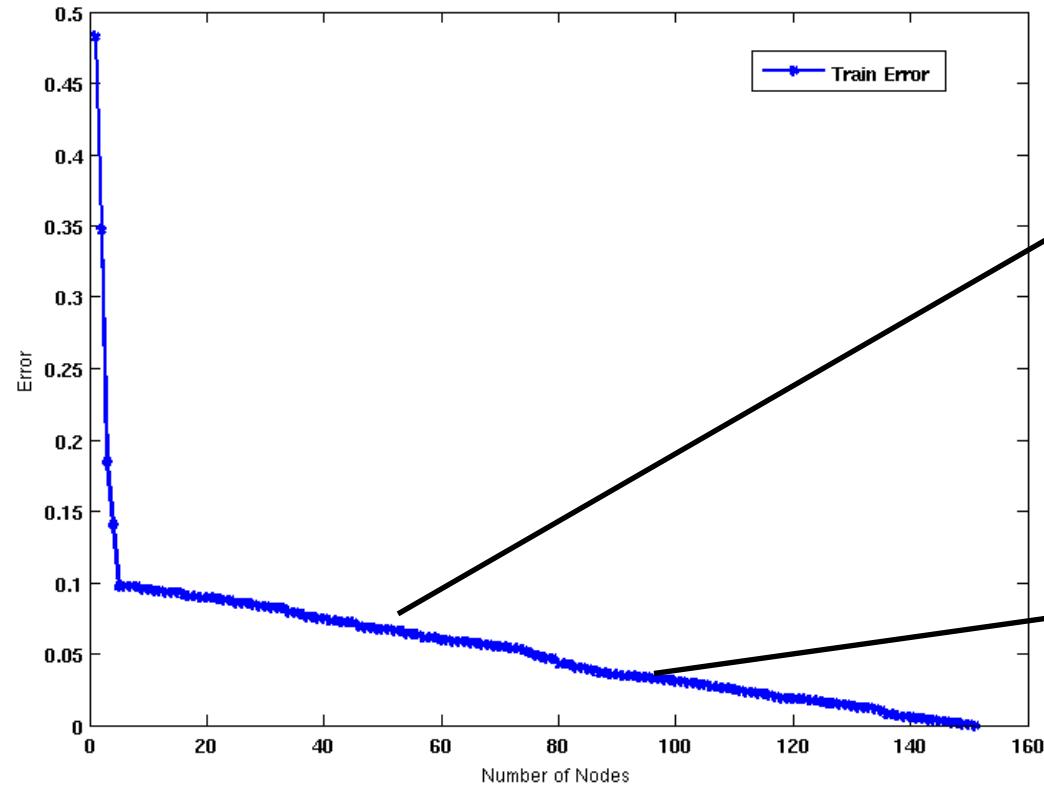


Decision boundaries on Training data (50 level)

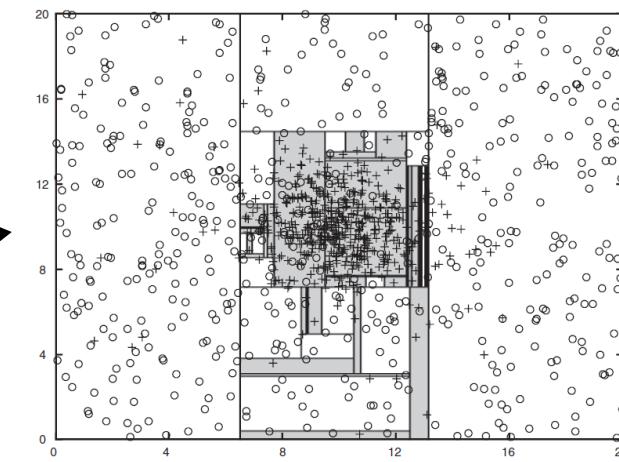


# Which tree is better?

The complex solution involving several layers is better than the simple one?



Decision Tree with 4 nodes



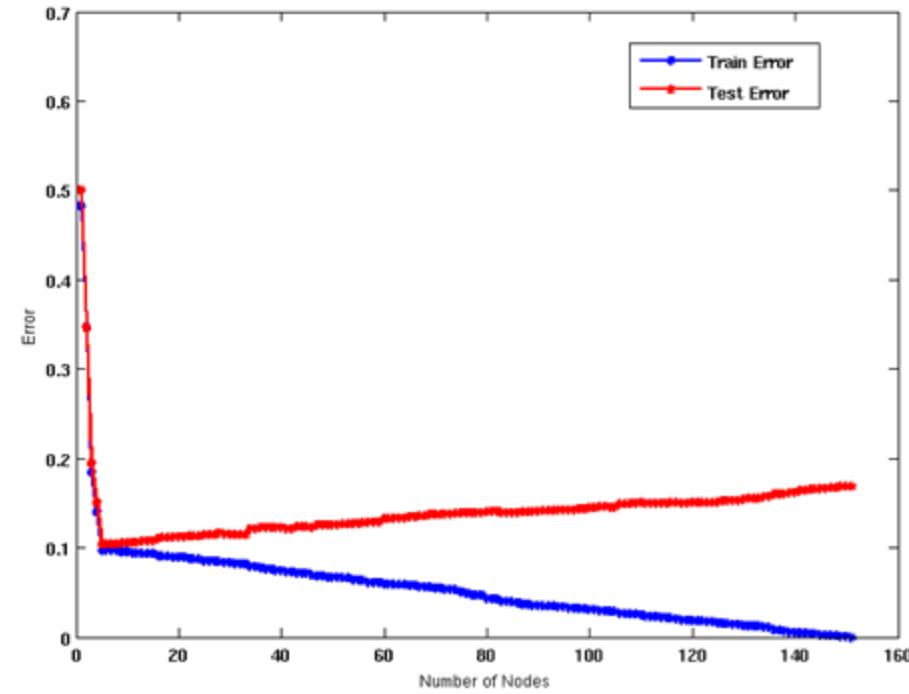
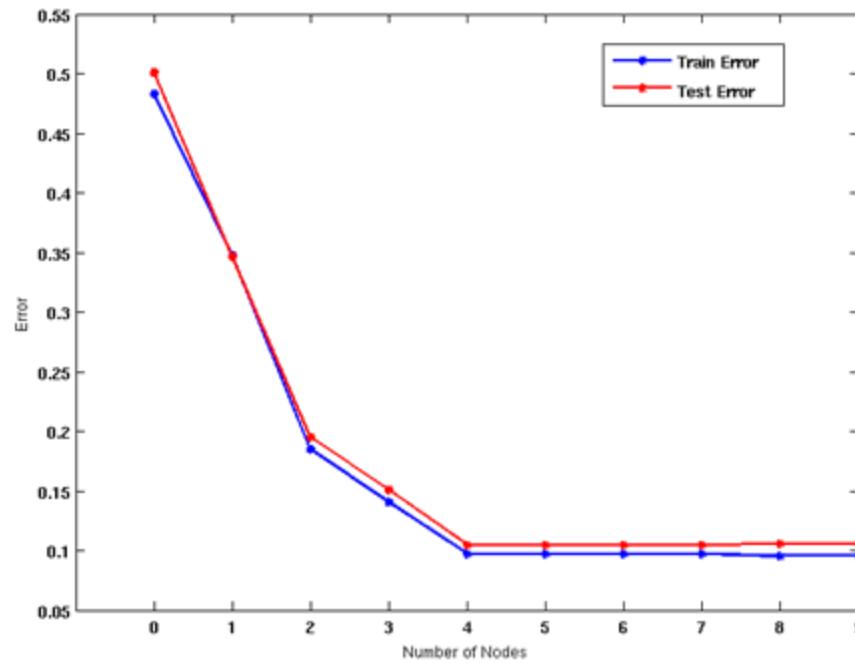
Decision Tree with 50 nodes



# Model Underfitting and Overfitting

Let's involve the **test set**.

As the model becomes more and **more complex**, **test errors** can start **increasing** even though **training error** may be **decreasing**



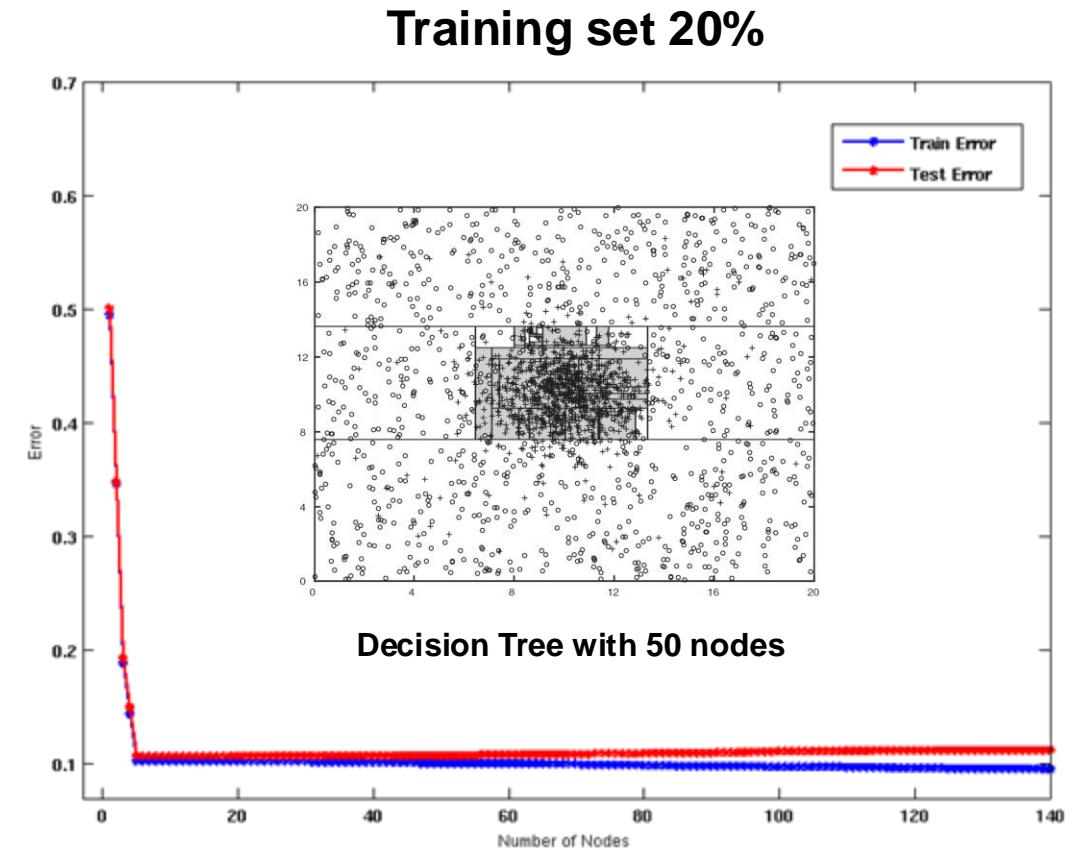
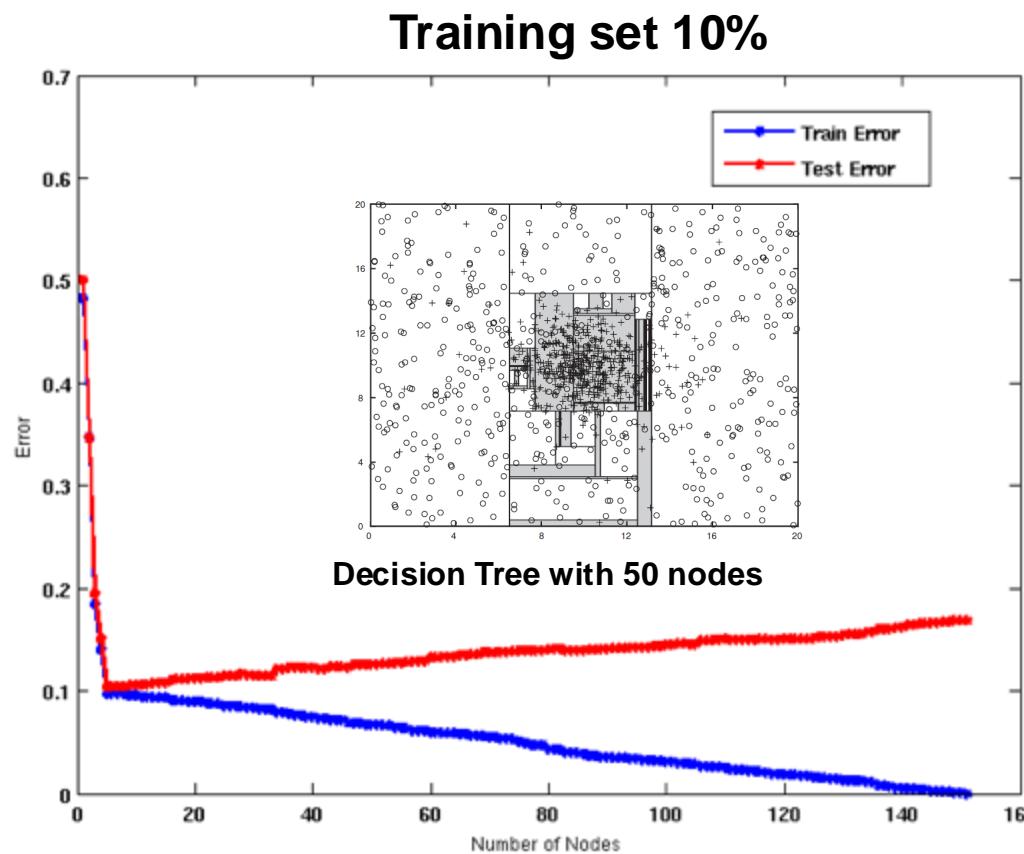
**Underfitting:** when model is too simple, both training and test errors are large

**Overfitting:** when model is too complex, training error is small but test error is large

# Model Overfitting – Impact of Training Data Size



Increasing the size of training data (10% to 20%) reduces the difference between training and testing errors at a given size of model





# Reasons for Model Overfitting

---

## Limited Training Size

A training set consisting of a **finite number of instances** can only provide a **limited representation** of the overall data.

It is possible that the patterns learned from a training set do not fully represent the true patterns in the overall data

## High Model Complexity

A more **complex** model can better represent **complex patterns** in the data.

An **overly complex model** also tends to **learn specific patterns** in the training set that do **not generalize** well over unseen instances

With a limited training set and a complex model, there is a high chance of picking a **spurious combination of parameters** that maximizes the evaluation metric just by random chance.



---

# Model Evaluation

(and Selection)



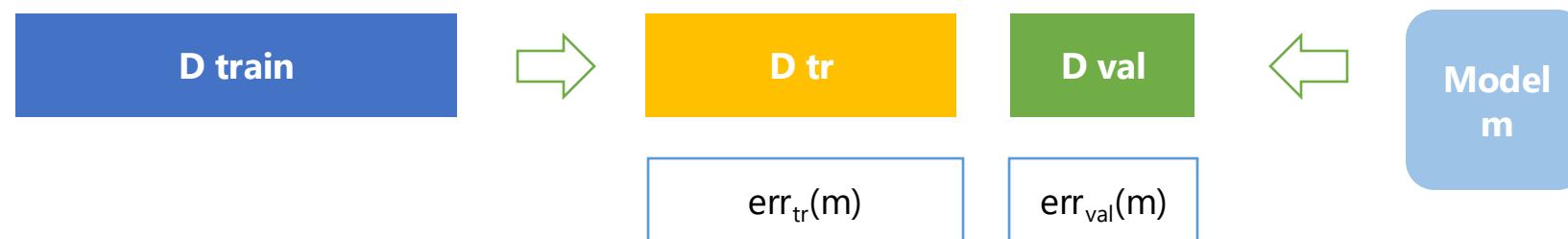


# Model Selection

The process of *selecting a model* with the right *level of complexity*, which is expected to *generalize well* over unseen test instances, is known as **model selection**

We can estimate the **generalization error rate** of a model by using “out-of-sample” estimates on **VALIDATION SET**

The validation error rate is a good indicator of generalization performance since the **validation set** has not been used for **training** the model



Small Training Set → Poor Classification model

Small Val Set → Unreliable result for selecting the model

# Model Selection

---



**Occam's razor** principle (principle of parsimony) suggests that given **two models** with the **same errors**, the **simpler model** is preferred over the more complex model.

We represent ***gen.error(m)*** as a function of not just the **training error** rate but also the **model complexity**

$$gen.error(m) = train.error(m, D.train) + \alpha \times complexity(M)$$

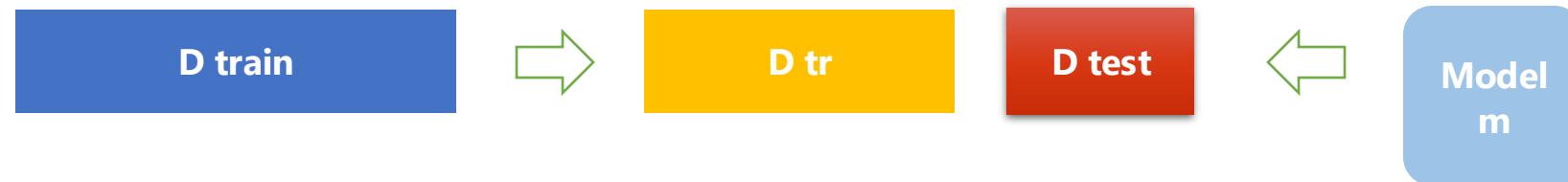


# Model Evaluation

**Model evaluation** estimates its generalization performance, i.e., its performance on unseen instances outside of train set

Holdout Method → training set ***D.train*** and the test set ***D.test***

A correct approach for **model evaluation** would be to *assess the performance* of a learned model on a labeled test set has not been used at any stage of **model selection**



Small Tr Set → Poor Classification model

Small Test Set → Unreliable result for selecting the model

# Model Evaluation: Cross Validation



**Cross-validation** is a model evaluation method that aims to make effective use of all labeled instances in D for both **training** and **testing**.

**Cross-Validation** is also referred as **k-fold** Cross Validation

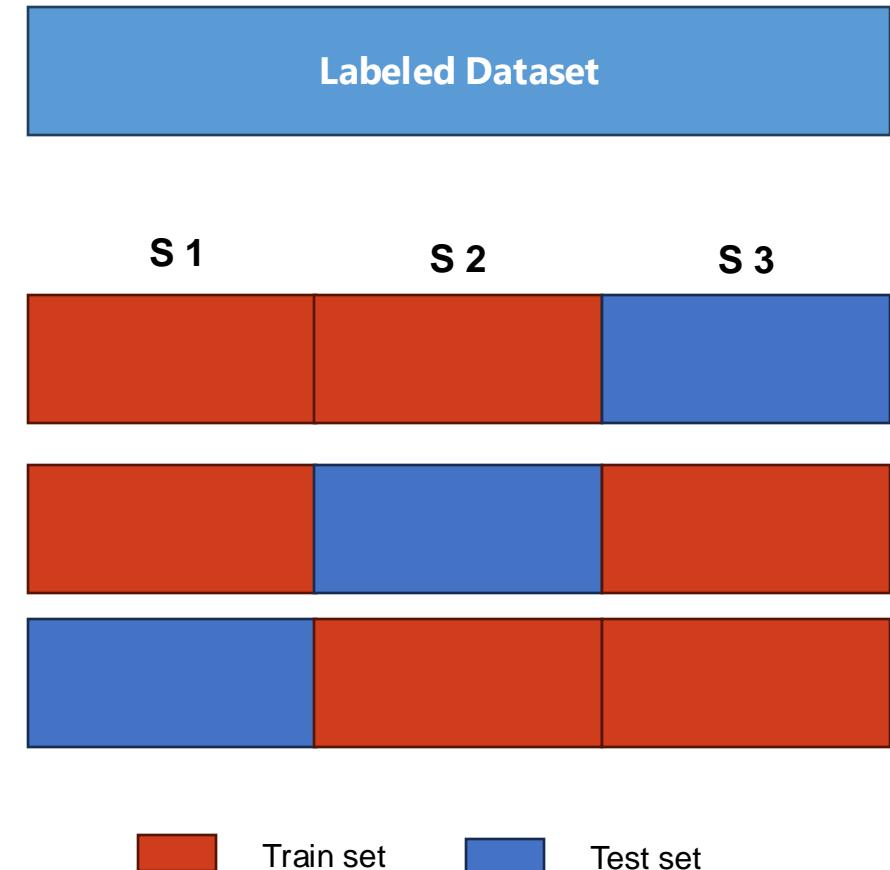
Let's partition the Labeled Dataset into k equal sized subset S1, S2 and S3.

Run 1 → Train (S1, S2), Test (S3)

Run 2 → Train (S1, S3), Test (S2)

Run 3 → Train (S2, S3), Test (S1)

Example with k=3





# Model Evaluation: k-fold

---

Procedure:

- Segment the labeled data D (of size N) into k equal-sized partitions (or folds)
- During the i-th run
  - The i-th partitions of D is chosen as  $D.\text{test}(i)$
  - The rest of the partitions are used as  $D.\text{train}(i)$

The total test error rate,  $\text{err}_{\text{test}}$ , is then computed as

$$\text{err}_{\text{test}} = \frac{\sum_{i=1}^k \text{err}_{\text{test}}(i)}{N}$$



# Model Evaluation: k-fold

---

**Small** value of k will result in **a smaller training** set and **larger test** set at every run

**Large** value of k will result in **a larger training** set and **small test** set at every run

A special case is the **Leave-one-out** approach ( $k=N$ ) → **Computationally Expensive**

This approach of estimating the generalization error rate and its variance is known as the **complete cross-validation** approach

A more practical solution is to **repeat the cross-validation** approach multiple times, using a different **random partitioning**

Usually, k is chosen equal to 5 or 10

**Stratified cross-validation** → Guarantee the equality of positive and negative instances in every partition of the data



# Model Evaluation: k-fold

k-fold cross-validation → A different model at every run →  $err_{test}(i)$

$err_{test}$  does not reflect the generalization error rate of any of the k models

$err_{test}$  represent the expected generalization error rate of the model selection approach

Applied on a training set of the same size as one of the training folds ( $N(k - 1)/k$ )

When the size of the training folds is closer to the size of the overall data (when k is large),

$err_{test}$  resembles the expected performance of a model learned over a data set of the same size as D

# Hyper-parameter Selection

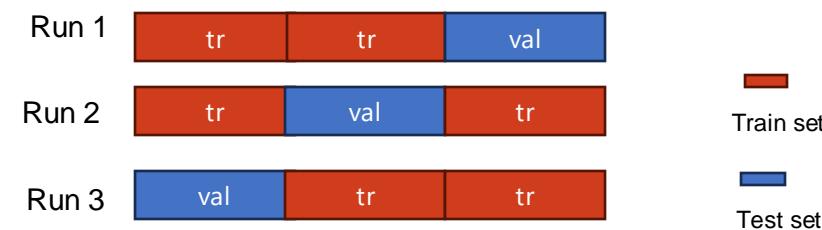


**Hyper-parameters** are parameters of learning algorithms that need to be determined before learning the classification model

Hyper-parameters do not appear in the final model

Regular model parameters, such as the test conditions in the decision tree, appear in the final model

The **hyper-parameter selection** must be performed **during model selection** and considered in **model evaluation**



# Hyper-parameter Selection

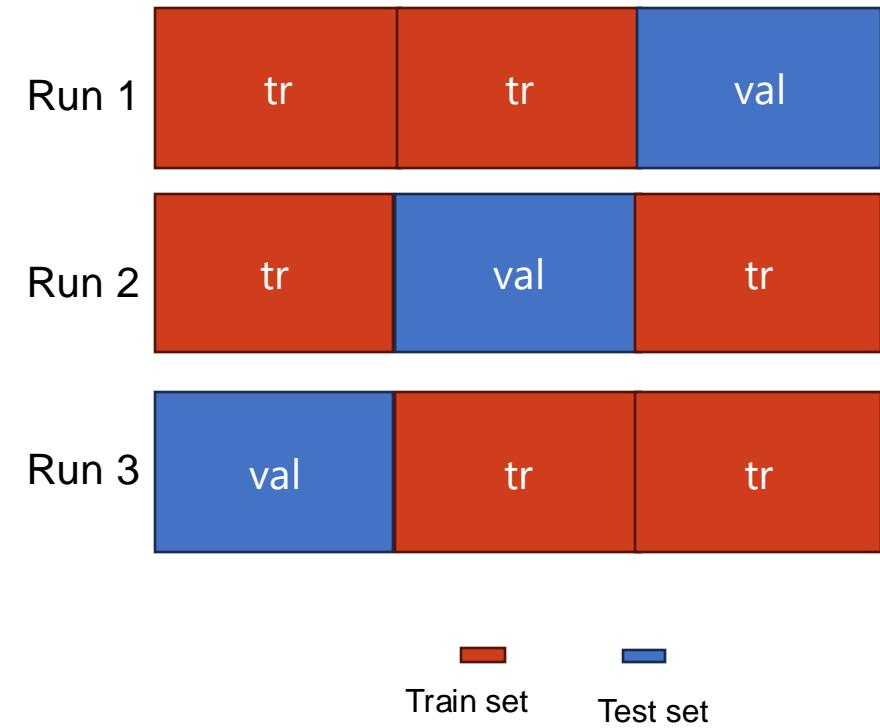


---

**Algorithm 3.2** Procedure **model-select**( $k, \mathcal{P}, D.\text{train}$ )

```
1:  $N_{\text{train}} = |D.\text{train}|$  {Size of  $D.\text{train}$ .}
2: Divide  $D.\text{train}$  into  $k$  partitions,  $D.\text{train}_1$  to  $D.\text{train}_k$ .
3: for each run  $i = 1$  to  $k$  do
4:    $D.\text{val}(i) = D.\text{train}_i$ . {Partition used for validation.}
5:    $D.\text{tr}(i) = D.\text{train} \setminus D.\text{train}_i$ . {Partitions used for training.}
6:   for each parameter  $p \in \mathcal{P}$  do
7:      $m = \text{model-train}(p, D.\text{tr}(i))$ . {Train model}
8:      $\text{err}_{\text{sum}}(p, i) = \text{model-test}(m, D.\text{val}(i))$ . {Sum of validation errors.}
9:   end for
10:  end for
11:   $\text{err}_{\text{val}}(p) = \sum_i^k \text{err}_{\text{sum}}(p, i)/N_{\text{train}}$ . {Compute validation error rate.}
12:   $p^* = \operatorname{argmin}_p \text{err}_{\text{val}}(p)$ . {Select best hyper-parameter value.}
13:   $m^* = \text{model-train}(p^*, D.\text{train})$ . {Learn final model on  $D.\text{train}$ }
14:  return  $(p^*, m^*)$ .
```

---



# Nested Cross-Validation



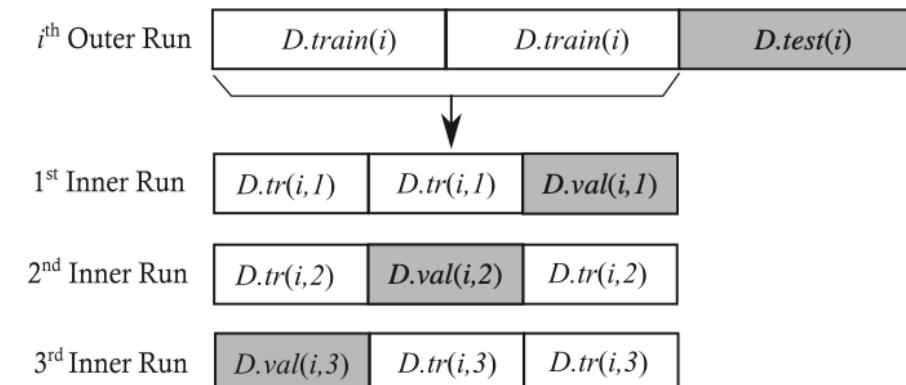
Algorithm 3.2 returns the **final classification model  $m^*$**  (learned with parameter  $p^*$ ) but not an estimate of its generalization performance  $err_{test}$

The validation error rates used in Algorithm 3.2 cannot be used as estimates of generalization performance

## Nested cross-validation or double cross-validation

### Algorithm 3.3 The nested cross-validation approach for computing $err_{test}$ .

```
1: Divide  $D$  into  $k$  partitions,  $D_1$  to  $D_k$ .  
2: for each outer run  $i = 1$  to  $k$  do  
3:    $D.test(i) = D_i$ . {Partition used for testing.}  
4:    $D.train(i) = D \setminus D_i$ . {Partitions used for model selection.}  
5:    $(p^*(i), m^*(i)) = \text{model-select}(k, \mathcal{P}, D.train(i))$ . {Inner cross-validation.}  
6:    $err_{sum}(i) = \text{model-test}(m^*(i), D.test(i))$ . {Sum of test errors.}  
7: end for  
8:  $err_{test} = \sum_i^k err_{sum}(i)/N$ . {Compute test error rate.}
```



# Class Imbalance Problem

---



In many datasets (**skewed** or **imbalanced**), there are a disproportionate number of instances that belong to different classes

This is frequent, especially in critical applications (health care, fraud)

The degree of imbalance varies across different applications or different datasets of the same application.

Accuracy is not well suited to evaluate models in the presence of class imbalance

We need of:

- methods for **building classifiers** when there is class imbalance in the training set
- methods for **adapting classification** decisions in the presence of a skewed test set.
- methods for **evaluating classification** in the presence of a skewed test set.

# Class Imbalance Problem: Undersampling/Oversampling

---



**Undersampling:** the frequency of the majority class is **reduced** to match the frequency of the minority class

- N random samples of the majority class are selected (where N is the number of instances of the minority class)
- Some useful examples could be lost
- Reducing the training set could result in low model capabilities

**Oversampling:** artificial examples of the minority class are **created** to make them equal to the number of negative instances

- Duplicate every positive instance  $n_{\text{majority}}/n_{\text{minority}}$
- Generate synthetic positive instances in the neighborhood of existing positive instances
- Duplicating a positive instance is analogous to doubling its weight during the training stage
- It does not help improve the representation of the positive class outside the boundary of existing positive instances.



# Class Imbalance Problem: Assigning Scores to Test Instances

---

If a classifier returns an ordinal score  $s(x)$  for a test instance  $x$ :

- Higher score → greater likelihood of  $x$  belonging to the positive class
- Lower score → greater likelihood of  $x$  belonging to the negative class

Moving a threshold  $s_T$  we can create a new binary classifier where a test instance  $x$  is classified **positive** only if  $s(x) > s_T$

This way, we can compensate bias introduced by the poor number of instances in the minority class.

This approach works if, in general,  $s(x_1) > s(x_2)$  if  $P(y = 1|x_1) > P(y = 1|x_2)$



# Confusion Matrix

**True positive (TP)** or  $f_{++}$ , which corresponds to the number of positive examples correctly predicted by the classifier.

**False positive (FP)** or  $f_{-+}$  (also known as Type I error), which corresponds to the number of negative examples wrongly predicted as positive by the classifier.

**False negative (FN)** or  $f_{+-}$  (also known as Type II error), which corresponds to the number of positive examples wrongly predicted as negative by the classifier.

**True negative (TN)** or  $f_{--}$ , which corresponds to the number of negative examples correctly predicted by the classifier.

		PREDICTED CLASS	
GT		Class=Yes	Class>No
	Class=Yes	TP	FN
	Class>No	FP	TN

$$\text{Accuracy} = \frac{TP + TN}{TP + TN + FP + FN}$$

# Confusion Matrix



Consider a 2-class problem

- Number of Class NO examples = 990
- Number of Class YES examples = 10

If a model predicts everything to be class NO, accuracy is  $990/1000 = 99\%$

- This is misleading because this trivial model does not detect any class YES example
- Detecting the rare class is usually more interesting (e.g., frauds, intrusions, defects, etc)

		PREDICTED CLASS	
		Class=Yes	Class=No
GT	Class=Yes	TP=0	FN=10
	Class=No	FP=0	TN=990



$$\text{Accuracy} = \frac{TP+TN}{TP+TN+FP+FN} = 99\%$$



# Confusion Matrix

- **True positive rate (TPR), Recall, Sensitivity:** fraction of positive test instances correctly predicted
- **True negative rate (TNR), Specificity:** fraction of negative test instances correctly predicted
- **False positive rate (FPR):** fraction of negative test instances wrongly predicted as positive
- **False negative rate (FNR):** fraction of positive test instances wrongly predicted as negative

$$TPR = \frac{TP}{TP + FN}$$

$$TNR = \frac{TN}{FP + TN}$$

$$FPR = \frac{FP}{FP + TN}$$

$$FNR = \frac{FN}{FN + TP}$$

		PREDICTED CLASS	
		Class=Yes	Class>No
GT	Class=Yes	TP=0	FN=10
	Class>No	FP=0	TN=990

$$Accuracy = \frac{TP+TN}{TP+TN+FP+FN} = 0.99$$

$$\boxed{Balanced Accuracy = \frac{TPR+TNN}{2} = \frac{0+1}{2} = 0.5}$$



# Which model is better?

		PREDICTED CLASS	
		Class=Yes	Class>No
GT	Class=Yes	TP=0	FN=10
	Class>No	FP=0	TN=990

		PREDICTED CLASS	
		Class=Yes	Class>No
GT	Class=Yes	TP=10	FN=0
	Class>No	FP=500	TN=490



# Which model is better?

		PREDICTED CLASS	
		Class=Yes	Class>No
GT	Class=Yes	TP=0	FN=10
	Class>No	FP=0	TN=990

$$\text{Balanced Accuracy} = \frac{TPR+TNN}{2} = \frac{0+1}{2} = 0.5$$

		PREDICTED CLASS	
		Class=Yes	Class>No
GT	Class=Yes	TP=10	FN=0
	Class>No	FP=500	TN=490

$$\text{Balanced Accuracy} = \frac{TPR+TNN}{2} = \frac{1+0.5}{2} = 0.75$$



# Alternative Measures

**Precision** (also called positive predictive value) A classifier that has a high precision is likely to have most of its positive predictions correct.

**Recall** is the fraction of positive test instances correctly predicted

Low Precision → many a **false positive**

Low Recall → Many **false negative**

**Precision VS Recall:** may be valuable to prioritize one metric over the other in cases where the outcome of a false positive or false negative is costly

In medical diagnosis, the **cost of a false negative** is high, and **recall** may be a valuable metric

In fraud detection, the **cost of a false positive** is high and **precision** may be a more valuable metric

	PREDICTED CLASS		
GT		Class=Yes	Class>No
	Class=Yes	TP	FN
	Class>No	FP	TN

$$Precision(p) = \frac{TP}{TP + FP}$$

$$Recall(r) = \frac{TP}{TP + FN} = TPR$$

$$F1\ score = \frac{2rp}{r + p} = \frac{2\ TP}{2\ TP + FP + FN}$$

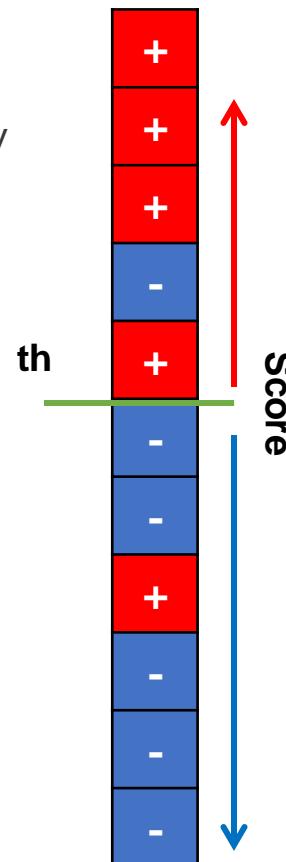
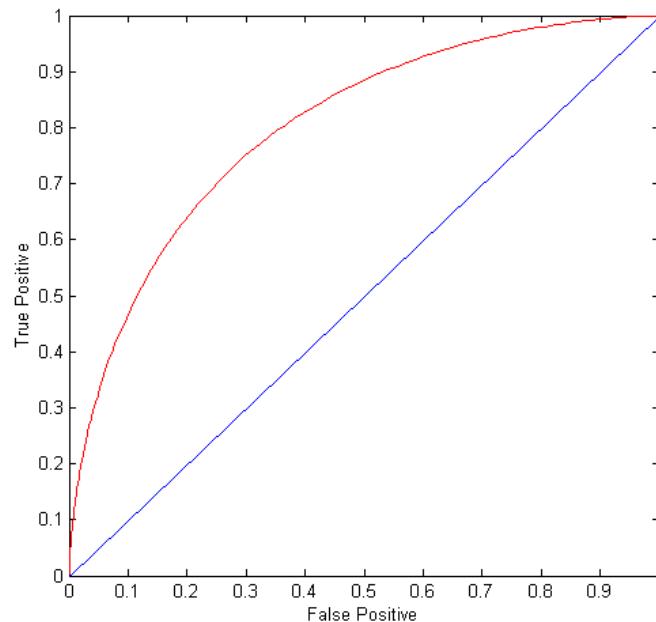
# ROC (Receiver Operating Characteristic)



A graphical approach for displaying **trade-off** between **TPR** (detection rate) and **FPR** (false alarm rate)

Developed in 1950s for signal detection theory to analyze noisy signals

ROC curve plots TPR against FPR



Let's say we have a classifier providing a score  $s$  where:

$$s(x_1) > s(x_2) \quad \text{if.} \quad P(y = 1|x_1) > P(y = 1|x_2)$$

- Higher score → greater likelihood of  $x$  belonging to the positive class
- Lower score → greater likelihood of  $x$  belonging to the negative class

A threshold  $th$  can divide positive and negative instances

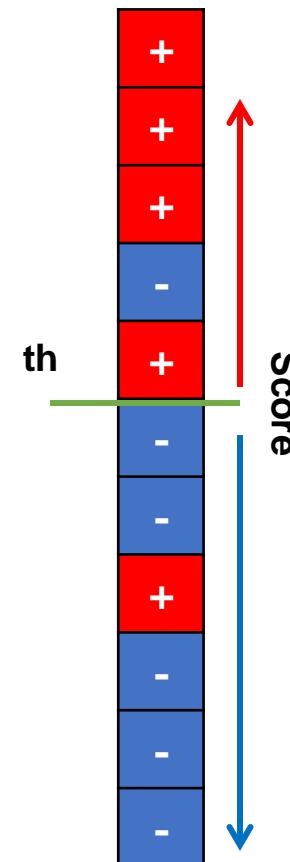
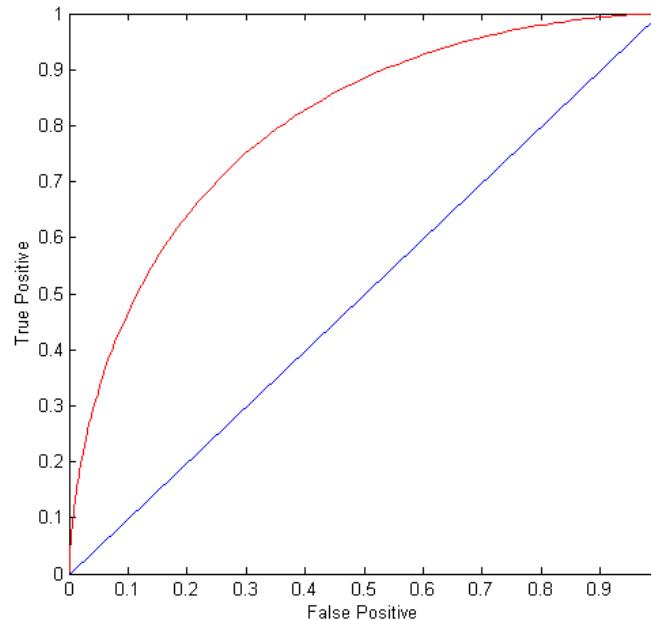
- Higher  $th$  → greater risk of misclassify positive as negative
- Lower  $th$  → greater risk of misclassify negative as positive

# ROC (Receiver Operating Characteristic)



ROC curve plots TPR against FPR

By using **different thresholds** on this value, we can create different **variations of the classifier** with TPR/FPR tradeoffs



Let's say we have a classifier providing a score  $s$  where:

$$s(x_1) > s(x_2) \text{ if. } P(y = 1|x_1) > P(y = 1|x_2)$$

- Higher score → greater likelihood of  $x$  belonging to the positive class
- Lower score → greater likelihood of  $x$  belonging to the negative class

A threshold  $th$  can divide positive and negative instances

- Higher  $th$  → greater risk of misclassify positive as negative
- Lower  $th$  → greater risk of misclassify negative as positive

# ROC (Receiver Operating Characteristic)

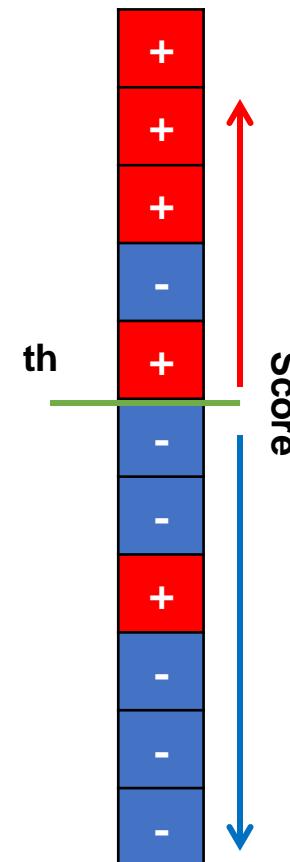
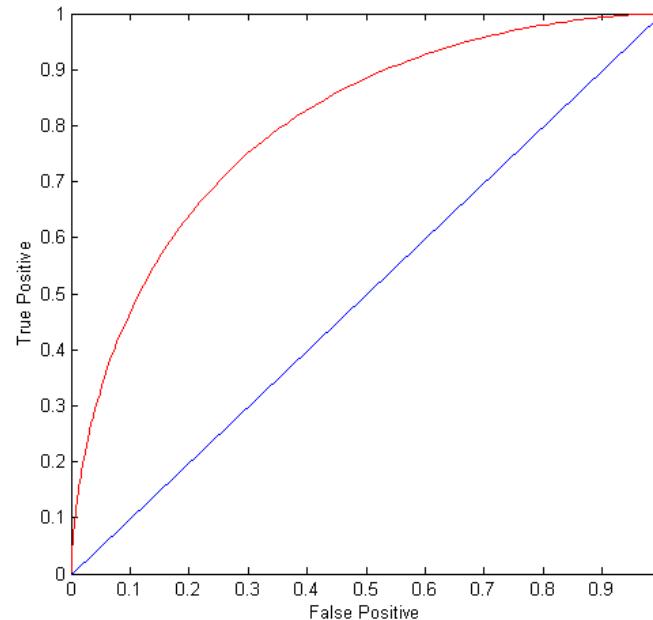


ROC curve plots TPR against FPR

(0,0): declare everything to be negative class

(1,1): declare everything to be positive class

(1,0): ideal



Let's say we have a classifier providing a score  $s$  where:

$$s(x_1) > s(x_2) \quad \text{if.} \quad P(y = 1|x_1) > P(y = 1|x_2)$$

- Higher score → greater likelihood of  $x$  belonging to the positive class
- Lower score → greater likelihood of  $x$  belonging to the negative class

A threshold  $th$  can divide positive and negative instances

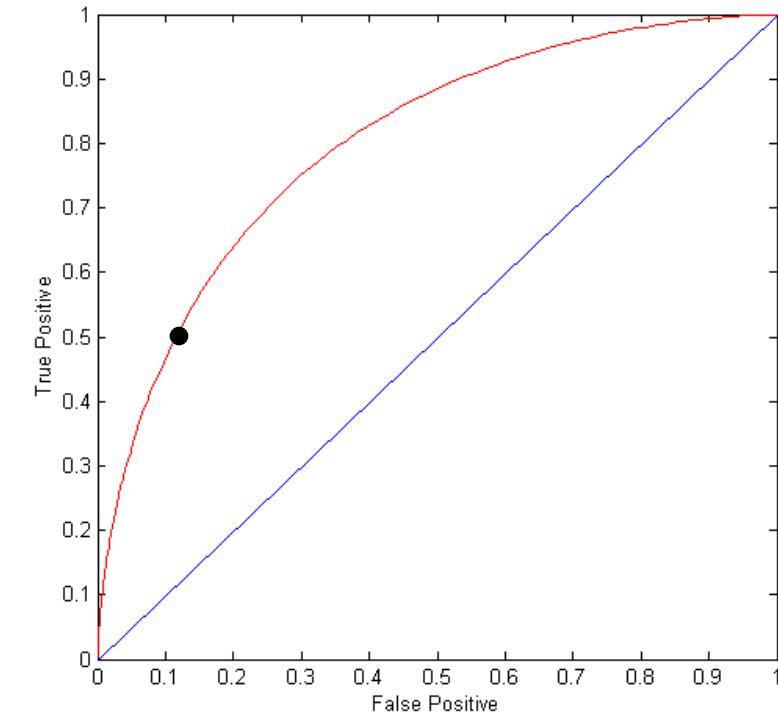
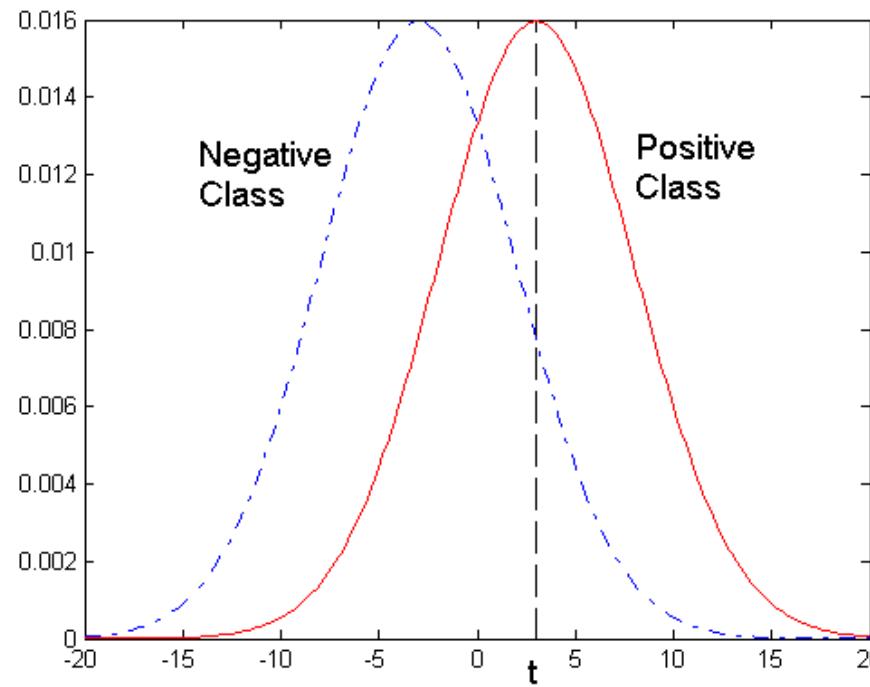
- Higher  $th$  → greater risk of misclassify positive as negative
- Lower  $th$  → greater risk of misclassify negative as positive



# ROC Curve Example

1-dimensional data set containing 2 classes (positive and negative)

Any points located at  $x > t$  is classified as positive



At threshold  $t$ : TPR=0.5, FNR=0.5, FPR=0.12, TNR=0.88

# How to Construct an ROC curve



Instance	Score	True Class
1	0.95	+
2	0.93	+
3	0.87	-
4	0.85	-
5	0.85	-
6	0.85	+
7	0.76	-
8	0.53	+
9	0.43	-
10	0.25	+

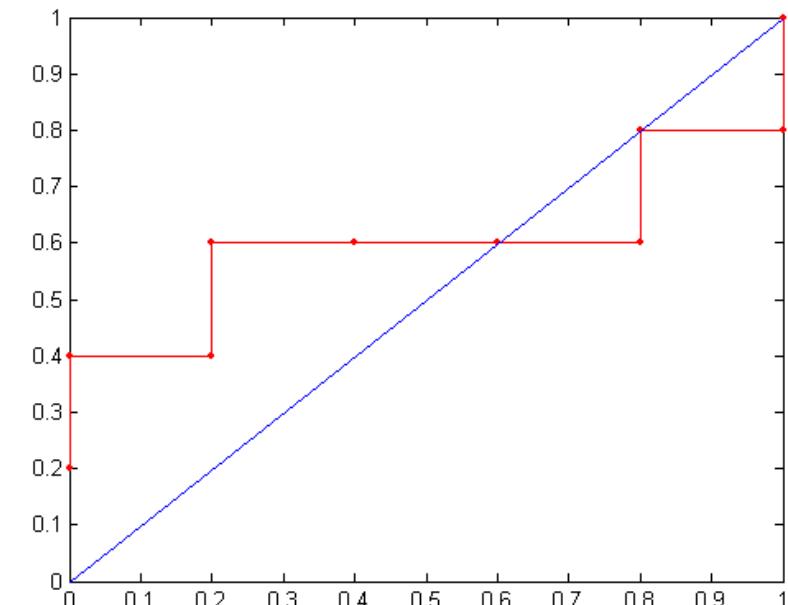
- Use a classifier that produces a continuous-valued score for each instance
  - The more likely it is for the instance to be in the + class, the higher the score
- Sort the instances in decreasing order according to the score
- Apply a threshold at each unique value of the score
- Count the number of TP, FP, TN, FN at each threshold
  - $TPR = TP/(TP+FN)$
  - $FPR = FP/(FP + TN)$



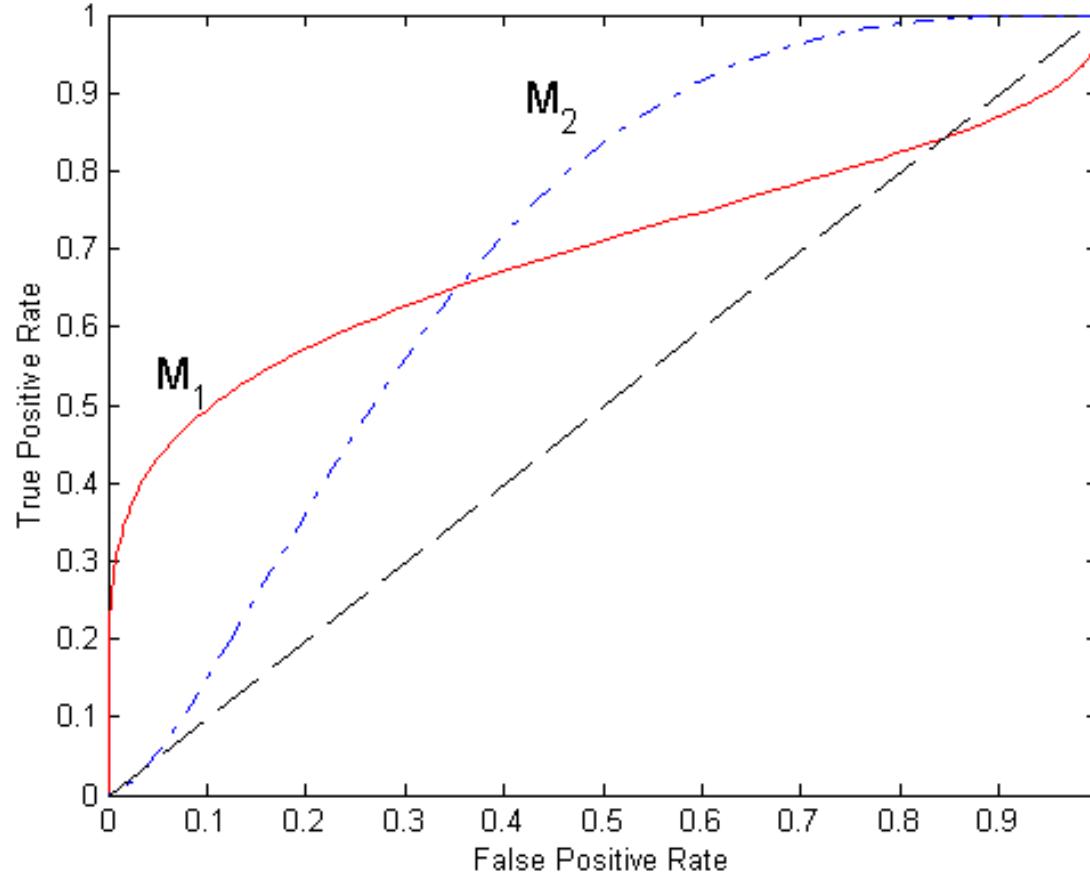
# How to construct an ROC curve

Class	+	-	+	-	-	-	+	-	+	+	
	0.25	0.43	0.53	0.76	0.85	0.85	0.85	0.87	0.93	0.95	1.00
TP	5	4	4	3	3	3	3	2	2	1	0
FP	5	5	4	4	3	2	1	1	0	0	0
TN	0	0	1	1	2	3	4	4	5	5	5
FN	0	1	1	2	2	2	2	3	3	4	5
TPR	1	0.8	0.8	0.6	0.6	0.6	0.6	0.4	0.4	0.2	0
FPR	1	1	0.8	0.8	0.6	0.4	0.2	0.2	0	0	0

**ROC Curve:**



# Using ROC for Model Comparison



- No model consistently outperforms the other
  - M<sub>1</sub> is better for small FPR
  - M<sub>2</sub> is better for large FPR
- Area Under the ROC curve (AUC)
  - Ideal:
    - Area = 1
  - Random guess:
    - Area = 0.5



UNIVERSITÀ  
DEL SALENTO

