

# C11: More on Decision Trees, Rule-Based Classifiers, Nearest Neighbor, Bays, SVM Model Evaluation

Ming-Syan Chen

November 20, 2019

## Tentative Class Agenda

- Class 8 – (10/30) Introduction to data
- Class 9 – (11/6) Association (Project announced, HW#4)
- Class 10 – (11/13) Data Exploration, OLAP, Classification
- Class 11 – (11/20) **More on Classification (HW#5)**
- Class 12 – (11/27) Clustering and others; **go over project abs.**
- Class 13 –(12/4) R (HW#6)
- Class 14 – (12/11) Exam (in class, closed book)
- Class 15 – (12/18) Project presentation I
- Class 16 – (12/25) Project presentation II

# More on Decision Tree

## Illustrative Example

RID	age	income	student	credit-rating	Class: buys-computer
1	<=30	high	no	fair	no
2	<=30	high	no	excellent	no
3	31...40	high	no	fair	yes
4	>40	medium	no	fair	yes
5	>40	low	yes	fair	yes
6	>40	low	yes	excellent	no
7	31...40	low	yes	excellent	yes
8	<=30	medium	no	fair	no
9	<=30	low	yes	fair	yes
10	>40	medium	yes	fair	yes
11	<=30	medium	yes	excellent	yes
12	31...40	medium	no	excellent	yes
13	31...40	high	yes	fair	yes
14	>40	medium	no	excellent	no

# Information Gain

- The information gain is :

$$gain(A) = I(y, n) - E(A)$$

- An example :

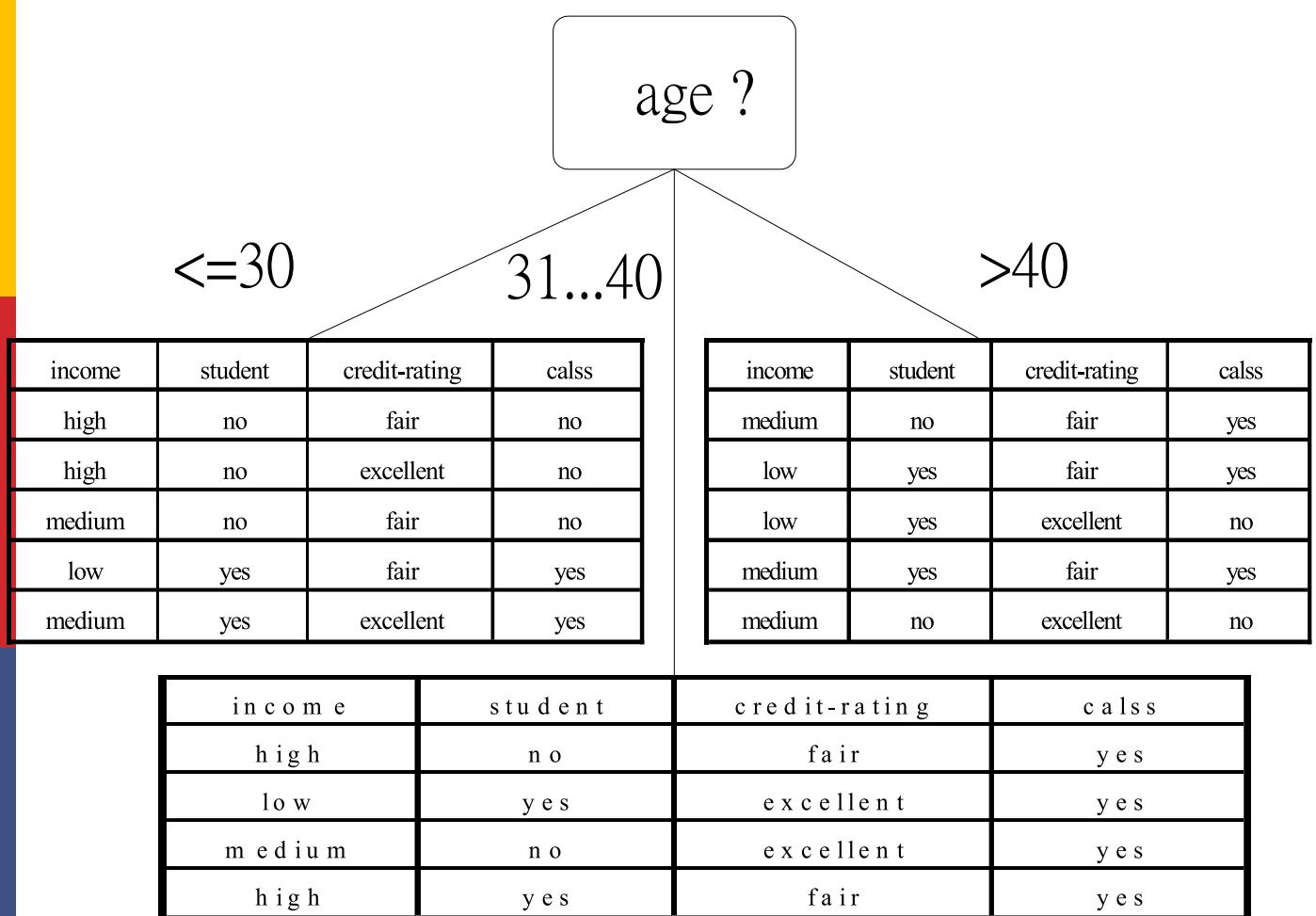
$$I(y, n) = -\frac{9}{14} \log_2 \frac{9}{14} - \frac{5}{14} \log_2 \frac{5}{14} = 0.94$$

$y = 9$	$n = 5$	$I(y, n) = 0.94$	- total	14
$y_1 = 2$	$n_1 = 3$	$I(y_1, n_1) = 0.971$	- age = " $\leq 30$ "	
$y_2 = 4$	$n_2 = 0$	$I(y_2, n_2) = 0$	- age = " $31 - 40$ "	
$y_3 = 3$	$n_3 = 2$	$I(y_3, n_3) = 0.971$	- age = " $> 40$ "	

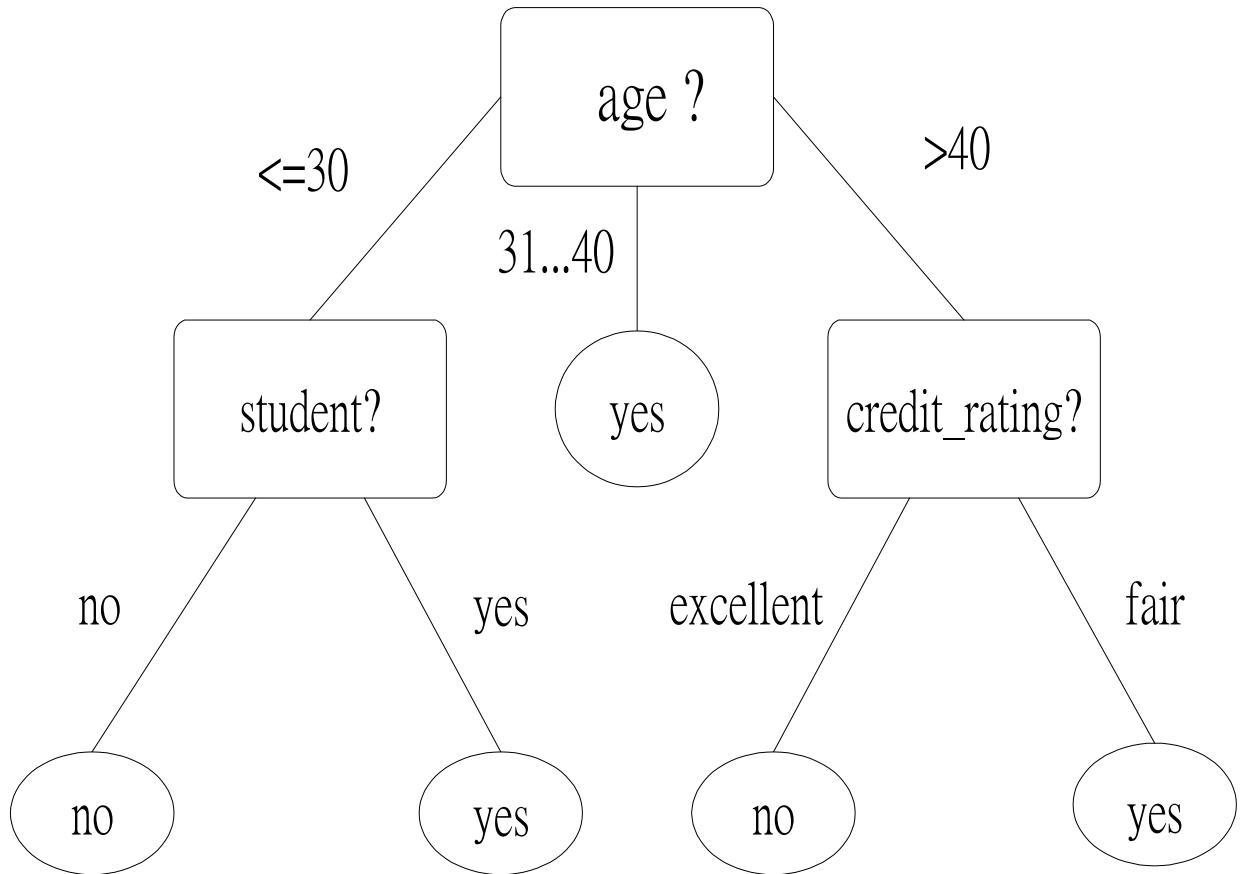
$$E(\text{age}) = \frac{5}{14} I(y_1, n_1) + \frac{4}{14} I(y_2, n_2) + \frac{5}{14} I(y_3, n_3) = 0.694$$

$$gain(\text{age}) = I(p, n) - E(\text{age}) = 0.246$$

5



6



7

## Rule Based Classifier

# Rule-Based Classifier

- Classify records by using a collection of “if...then...” rules
- Rule:  $(Condition) \rightarrow y$ 
  - where
    - *Condition* is a conjunctions of attributes
    - $y$  is the class label
  - *LHS*: rule antecedent or condition
  - *RHS*: rule consequent
  - Examples of classification rules:
    - $(\text{Blood Type}=\text{Warm}) \wedge (\text{Lay Eggs}=\text{Yes}) \rightarrow \text{Birds}$
    - $(\text{Taxable Income} < 50K) \wedge (\text{Refund}=\text{Yes}) \rightarrow \text{Evade}=\text{No}$

9

## Rule-based Classifier (Example)

Name	Blood Type	Give Birth	Can Fly	Live in Water	Class
human	warm	yes	no	no	mammals
python	cold	no	no	no	reptiles
salmon	cold	no	no	yes	fishes
whale	warm	yes	no	yes	mammals
frog	cold	no	no	sometimes	amphibians
komodo	cold	no	no	no	reptiles
bat	warm	yes	yes	no	mammals
pigeon	warm	no	yes	no	birds
cat	warm	yes	no	no	mammals
leopard shark	cold	yes	no	yes	fishes
turtle	cold	no	no	sometimes	reptiles
penguin	warm	no	no	sometimes	birds
porcupine	warm	yes	no	no	mammals
eel	cold	no	no	yes	fishes
salamander	cold	no	no	sometimes	amphibians
gila monster	cold	no	no	no	reptiles
platypus	warm	no	no	no	mammals
owl	warm	no	yes	no	birds
dolphin	warm	yes	no	yes	mammals
eagle	warm	no	yes	no	birds

R1:  $(\text{Give Birth} = \text{no}) \wedge (\text{Can Fly} = \text{yes}) \rightarrow \text{Birds}$

R2:  $(\text{Give Birth} = \text{no}) \wedge (\text{Live in Water} = \text{yes}) \rightarrow \text{Fishes}$

R3:  $(\text{Give Birth} = \text{yes}) \wedge (\text{Blood Type} = \text{warm}) \rightarrow \text{Mammals}$

R4:  $(\text{Give Birth} = \text{no}) \wedge (\text{Can Fly} = \text{no}) \rightarrow \text{Reptiles}$

R5:  $(\text{Live in Water} = \text{sometimes}) \rightarrow \text{Amphibians}$

# Applying Rule-Based Classifier

- A rule  $r$  **covers** an instance  $x$  if the attributes of the instance satisfy the condition of the rule

R1: (Give Birth = no)  $\wedge$  (Can Fly = yes)  $\rightarrow$  Birds

R2: (Give Birth = no)  $\wedge$  (Live in Water = yes)  $\rightarrow$  Fishes

R3: (Give Birth = yes)  $\wedge$  (Blood Type = warm)  $\rightarrow$  Mammals

R4: (Give Birth = no)  $\wedge$  (Can Fly = no)  $\rightarrow$  Reptiles

R5: (Live in Water = sometimes)  $\rightarrow$  Amphibians

Name	Blood Type	Give Birth	Can Fly	Live in Water	Class
hawk	warm	no	yes	no	?
grizzly bear	warm	yes	no	no	?

The rule R1 covers a hawk => Bird

The rule R3 covers the grizzly bear => Mammal

11

## Rule Coverage and Accuracy

- Coverage of a rule:
  - Fraction of records that satisfy the antecedent of a rule
- Accuracy of a rule:
  - Fraction of records that satisfy both the antecedent and consequent of a rule

Tid	Refund	Marital Status	Taxable Income	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

(Status=Single)  $\rightarrow$  No

Coverage = 40%, Accuracy = 50%

12

# How Rule-based Classifier Work?

R1: (Give Birth = no)  $\wedge$  (Can Fly = yes)  $\rightarrow$  Birds

R2: (Give Birth = no)  $\wedge$  (Live in Water = yes)  $\rightarrow$  Fishes

R3: (Give Birth = yes)  $\wedge$  (Blood Type = warm)  $\rightarrow$  Mammals

R4: (Give Birth = no)  $\wedge$  (Can Fly = no)  $\rightarrow$  Reptiles

R5: (Live in Water = sometimes)  $\rightarrow$  Amphibians

Name	Blood Type	Give Birth	Can Fly	Live in Water	Class
lemur	warm	yes	no	no	?
turtle	cold	no	no	sometimes	?
dogfish shark	cold	yes	no	yes	?

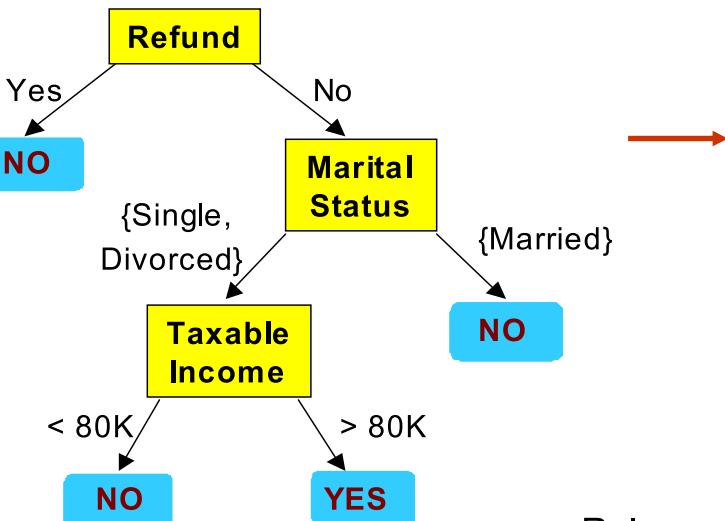
A lemur triggers rule R3, so it is classified as a mammal

A turtle triggers both R4 and R5

A dogfish shark triggers none of the rules

13

## From Decision Trees To Rules



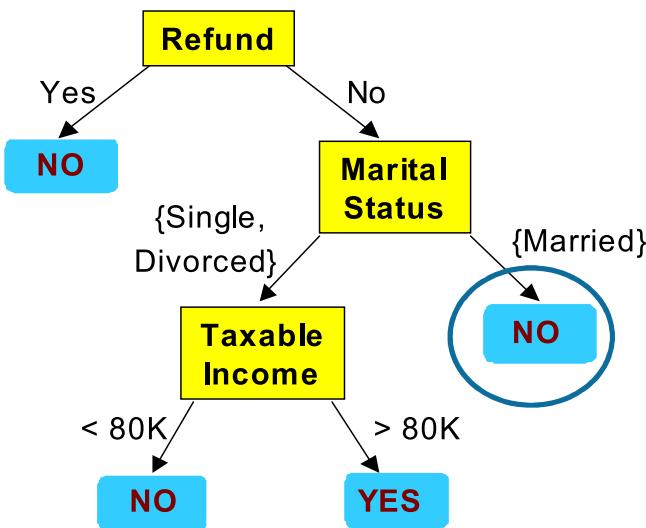
Classification Rules
(Refund=Yes) ==> No
(Refund=No, Marital Status={Single,Divorced}, Taxable Income<80K) ==> No
(Refund=No, Marital Status={Single,Divorced}, Taxable Income>80K) ==> Yes
(Refund=No, Marital Status={Married}) ==> No

Rules are mutually exclusive and exhaustive

Rule set contains as much information as the tree

14

# Rules Can Be Simplified



Tid	Refund	Marital Status	Taxable Income	Cheat
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

Initial Rule:  $(\text{Refund}=\text{No}) \wedge (\text{Status}=\text{Married}) \rightarrow \text{No}$

Simplified Rule:  $(\text{Status}=\text{Married}) \rightarrow \text{No}$

15

## Advantages of Rule-Based Classifiers

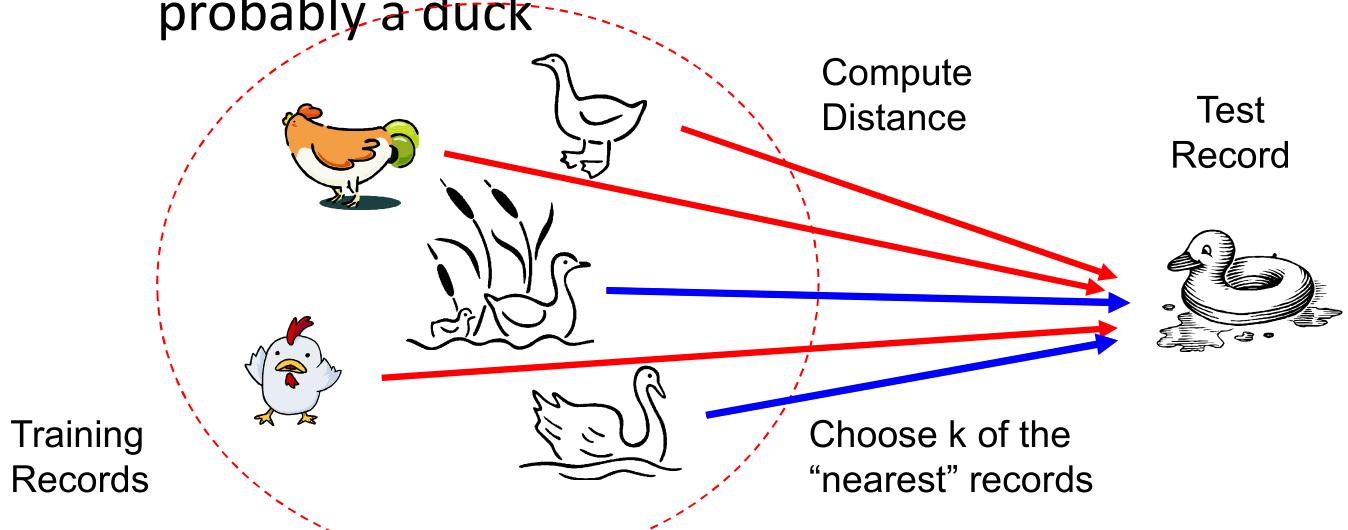
- As highly expressive as decision trees
- Easy to interpret
- Easy to generate
- Can classify new instances rapidly
- Performance comparable to decision trees

16

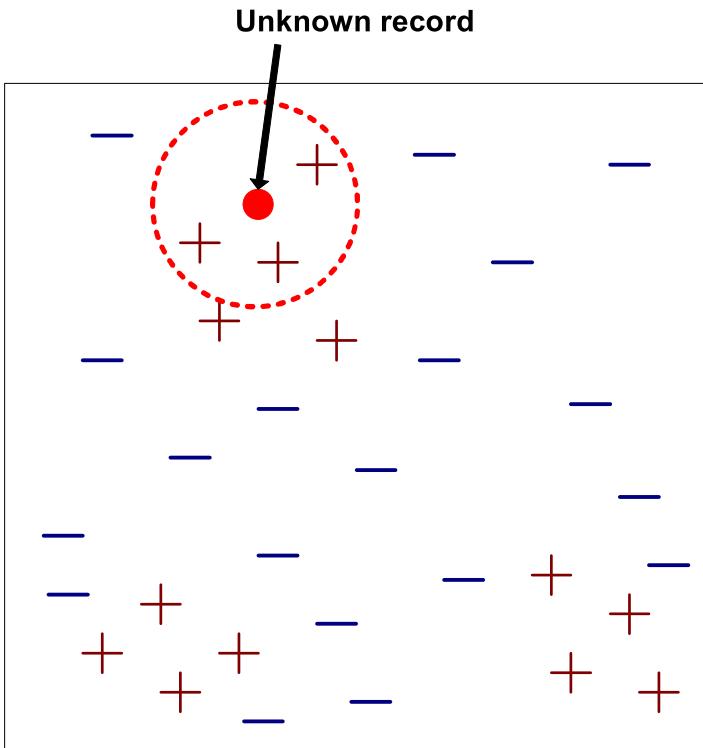
# Nearest Neighbor Classifier

# Nearest Neighbor Classifiers

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



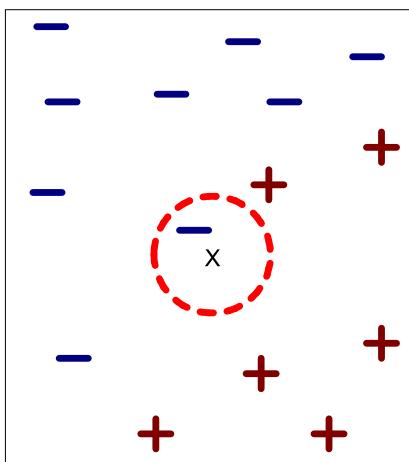
# Nearest-Neighbor Classifiers



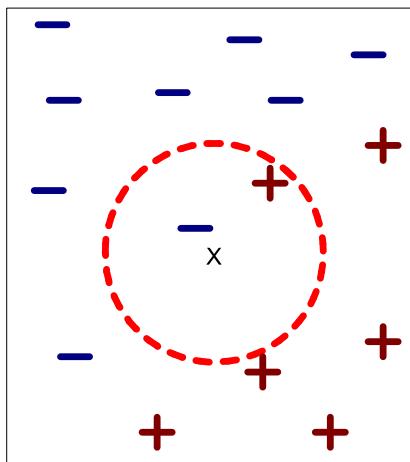
- Requires three things
  - The set of stored records
  - Distance Metric to compute distance between records
  - The value of  $k$ , the number of nearest neighbors to retrieve
- To classify an unknown record:
  - Compute distance to other training records
  - Identify  $k$  nearest neighbors
  - Use class labels of nearest neighbors to determine the class label of unknown record (e.g., by taking majority vote)

19

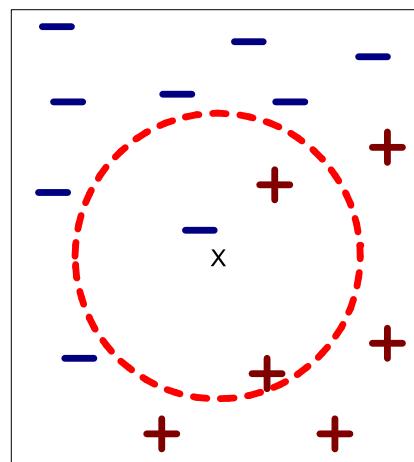
## Definition of Nearest Neighbor



(a) 1-nearest neighbor



(b) 2-nearest neighbor



(c) 3-nearest neighbor

K-nearest neighbors of a record  $x$  are data points that have the  $k$  smallest distance to  $x$

20

# Nearest Neighbor Classification

- Compute distance between two points:
  - Euclidean distance

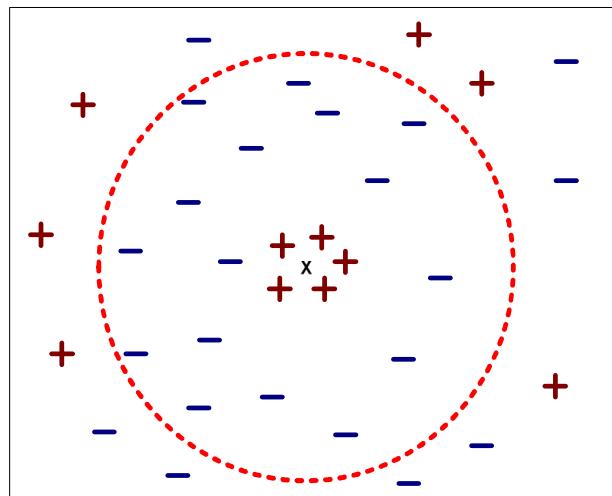
$$d(p, q) = \sqrt{\sum_i (p_i - q_i)^2}$$

- Determine the class from nearest neighbor list
  - take the majority vote of class labels among the k-nearest neighbors
  - Weigh the vote according to distance
    - weight factor,  $w = 1/d^2$

21

## Nearest Neighbor Classification...

- Choosing the value of k:
  - If k is too small, sensitive to noise points
  - If k is too large, neighborhood may include points from other classes



22

# Nearest Neighbor Classification...

- Scaling issues
  - 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

23

# Nearest neighbor Classification...

- k-NN classifiers are lazy learners
  - It does not build models explicitly
  - Unlike eager learners such as decision tree induction and rule-based systems
  - Classifying unknown records are relatively expensive

24

# Bays Classifiers

## Bayes Classifier

- A probabilistic framework for solving classification problems

$$P(C | A) = \frac{P(A, C)}{P(A)}$$

- Conditional Probability:

$$P(A | C) = \frac{P(A, C)}{P(C)}$$

- Bayes theorem:

$$P(C | A) = \frac{P(A | C)P(C)}{P(A)}$$

# Example of Bayes Theorem

- Given:
  - A doctor knows that meningitis causes stiff neck 50% of the time
  - Prior probability of any patient having meningitis is 1/50,000
  - Prior probability of any patient having stiff neck is 1/20
- If a patient has stiff neck, what's the probability he/she has meningitis?

$$P(M | S) = \frac{P(S | M)P(M)}{P(S)} = \frac{0.5 \times 1/50000}{1/20} = 0.0002$$

27

## Another Example

- Suppose we have 3 round apples (2 red, one green), 2 bananas, 3 litchis (red and round), and 4 mangoes
- X: a red and round object  
H: being an apple
- $P(X)=5/12$ ,  $P(H)=3/12$ ,  $P(X|H)=2/3$
- Then,  $P(H|X)=P(X|H)P(H)/P(X)=2/5$

28

# Bayesian Classifiers

- Consider each attribute and class label as random variables
- Given a record with attributes  $(A_1, A_2, \dots, A_n)$ 
  - Goal is to predict class C
  - Specifically, we want to find the value of C that maximizes  $P(C | A_1, A_2, \dots, A_n)$
- Can we estimate  $P(C | A_1, A_2, \dots, A_n)$  directly from data?

29

# Bayesian Classifiers

- Approach:
  - compute the posterior probability  $P(C | A_1, A_2, \dots, A_n)$  for all values of C using the Bayes theorem
  - Choose value of C that maximizes  $P(C | A_1, A_2, \dots, A_n)$
  - Equivalent to choosing value of C that maximizes  $P(A_1, A_2, \dots, A_n | C) P(C)$
- How to estimate  $P(A_1, A_2, \dots, A_n | C)$ ?

30

# Naïve Bayes Classifier

- Assume independence among attributes  $A_i$  when class is given:
  - $P(A_1, A_2, \dots, A_n | C_j) = P(A_1 | C_j) P(A_2 | C_j) \dots P(A_n | C_j)$
  - Can estimate  $P(A_i | C_j)$  for all  $A_i$  and  $C_j$ .
  - New point is classified to  $C_j$  if  $P(C_j) \prod P(A_i | C_j)$  is maximal.

31

## To Estimate Probabilities from Data

- Class:  $P(C) = N_c/N$ 
  - e.g.,  $P(\text{No}) = 7/10$ ,  $P(\text{Yes}) = 3/10$
- For discrete attributes:
 
$$P(A_i | C_k) = |A_{ik}| / N_c$$
  - where  $|A_{ik}|$  is number of instances having attribute  $A_i$  and belongs to class  $C_k$
  - 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

32

# To Estimate Probabilities from Data

- For continuous attributes:
  - Discretize the range into bins
    - one ordinal attribute per bin
    - violates independence assumption
  - Two-way split:  $(A < v)$  or  $(A > v)$ 
    - choose only one of the two splits as new attribute
  - 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, can use it to estimate the conditional probability  $P(A_i | c)$

33

# To Estimate Probabilities 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

- Normal distribution:
- $$P(A_i | c_j) = \frac{1}{\sqrt{2\pi\sigma_{ij}^2}} e^{-\frac{(A_i - \mu_{ij})^2}{2\sigma_{ij}^2}}$$
- One for each  $(A_i, c_j)$  pair
- For (Income, Class=No):
    - If Class=No
      - sample mean = 110
      - sample variance = 2975

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

34

# Example of Naïve Bayes Classifier

Given a Test Record:

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

naive Bayes Classifier:

$P(\text{Refund}=\text{Yes}|\text{No}) = 3/7$   
 $P(\text{Refund}=\text{No}|\text{No}) = 4/7$   
 $P(\text{Refund}=\text{Yes}|\text{Yes}) = 0$   
 $P(\text{Refund}=\text{No}|\text{Yes}) = 1$   
 $P(\text{Marital Status}=\text{Single}|\text{No}) = 2/7$   
 $P(\text{Marital Status}=\text{Divorced}|\text{No}) = 1/7$   
 $P(\text{Marital Status}=\text{Married}|\text{No}) = 4/7$   
 $P(\text{Marital Status}=\text{Single}|\text{Yes}) = 2/7$   
 $P(\text{Marital Status}=\text{Divorced}|\text{Yes}) = 1/7$   
 $P(\text{Marital Status}=\text{Married}|\text{Yes}) = 0$

For taxable income:

If class=No: sample mean=110  
sample variance=2975  
If class=Yes: sample mean=90  
sample variance=25

- $P(X|\text{Class}=\text{No}) = P(\text{Refund}=\text{No}|\text{Class}=\text{No}) \times P(\text{Married}|\text{Class}=\text{No}) \times P(\text{Income}=120\text{K}|\text{Class}=\text{No}) = 4/7 \times 4/7 \times 0.0072 = 0.0024$
- $P(X|\text{Class}=\text{Yes}) = P(\text{Refund}=\text{No}|\text{Class}=\text{Yes}) \times P(\text{Married}|\text{Class}=\text{Yes}) \times P(\text{Income}=120\text{K}|\text{Class}=\text{Yes}) = 1 \times 0 \times 1.2 \times 10^{-9} = 0$

Since  $P(X|\text{No})P(\text{No}) > P(X|\text{Yes})P(\text{Yes})$

Therefore  $P(\text{No}|X) > P(\text{Yes}|X)$   
=> Class = No

35

## Naïve Bayes Classifier

- If one of the conditional probability is zero, then the entire expression becomes zero
- Probability estimation:

$$\text{Original: } P(A_i | C) = \frac{N_{ic}}{N_c}$$

c: number of classes

$$\text{Laplace: } P(A_i | C) = \frac{N_{ic} + 1}{N_c + c}$$

p: prior probability

$$\text{m - estimate: } P(A_i | C) = \frac{N_{ic} + mp}{N_c + m}$$

m: parameter

36

# Example of Naïve Bayes Classifier

Name	Give Birth	Can Fly	Live in Water	Have Legs	Class
human	yes	no	no	yes	mammals
python	no	no	no	no	non-mammals
salmon	no	no	yes	no	non-mammals
whale	yes	no	yes	no	mammals
frog	no	no	sometimes	yes	non-mammals
komodo	no	no	no	yes	non-mammals
bat	yes	yes	no	yes	mammals
pigeon	no	yes	no	yes	non-mammals
cat	yes	no	no	yes	mammals
leopard shark	yes	no	yes	no	non-mammals
turtle	no	no	sometimes	yes	non-mammals
penguin	no	no	sometimes	yes	non-mammals
porcupine	yes	no	no	yes	mammals
eel	no	no	yes	no	non-mammals
salamander	no	no	sometimes	yes	non-mammals
gila monster	no	no	no	yes	non-mammals
platypus	no	no	no	yes	mammals
owl	no	yes	no	yes	non-mammals
dolphin	yes	no	yes	no	mammals
eagle	no	yes	no	yes	non-mammals

A: attributes

M: mammals

N: non-mammals

$$P(A | M) = \frac{6}{7} \times \frac{6}{7} \times \frac{2}{7} \times \frac{2}{7} = 0.06$$

$$P(A | N) = \frac{1}{13} \times \frac{10}{13} \times \frac{3}{13} \times \frac{4}{13} = 0.0042$$

$$P(A | M)P(M) = 0.06 \times \frac{7}{20} = 0.021$$

$$P(A | N)P(N) = 0.004 \times \frac{13}{20} = 0.0027$$

$$\begin{aligned} P(A|M)P(M) &> P(A|N)P(N) \\ \Rightarrow \text{Mammals} \end{aligned}$$

Give Birth	Can Fly	Live in Water	Have Legs	Class
yes	no	yes	no	?

37

## Naïve Bayes (Summary)

- Robust to isolated noise points
- Handle missing values by ignoring the instance during probability estimate calculations
- Robust to irrelevant attributes
- Independence assumption may not hold for some attributes

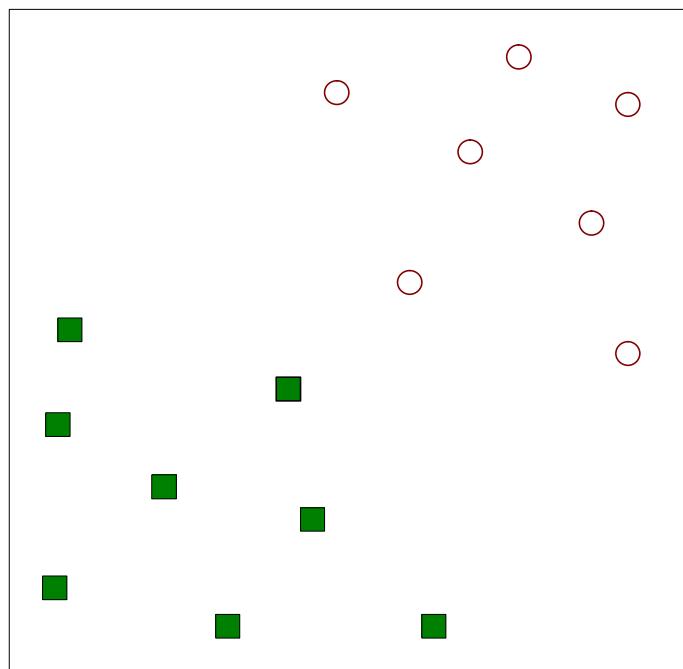
38

# SVM

## Support Vector Machine

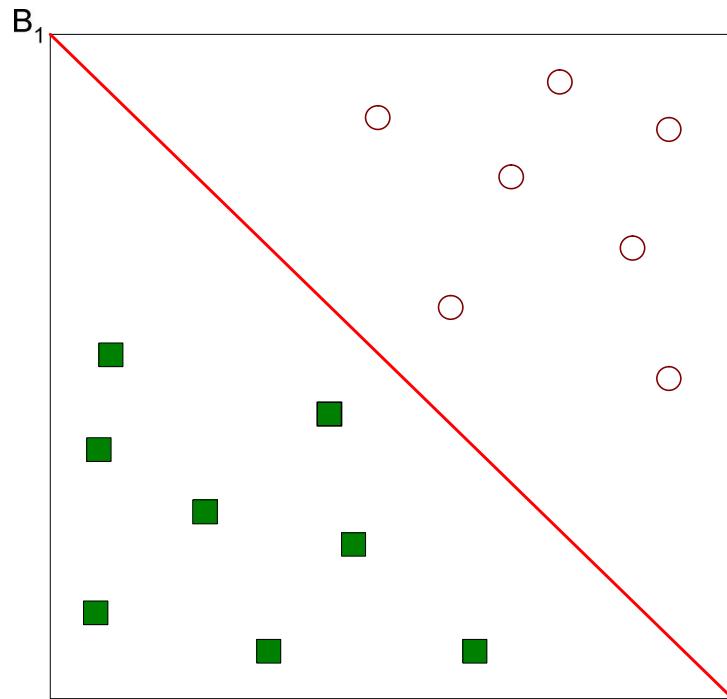
39

## Support Vector Machines



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

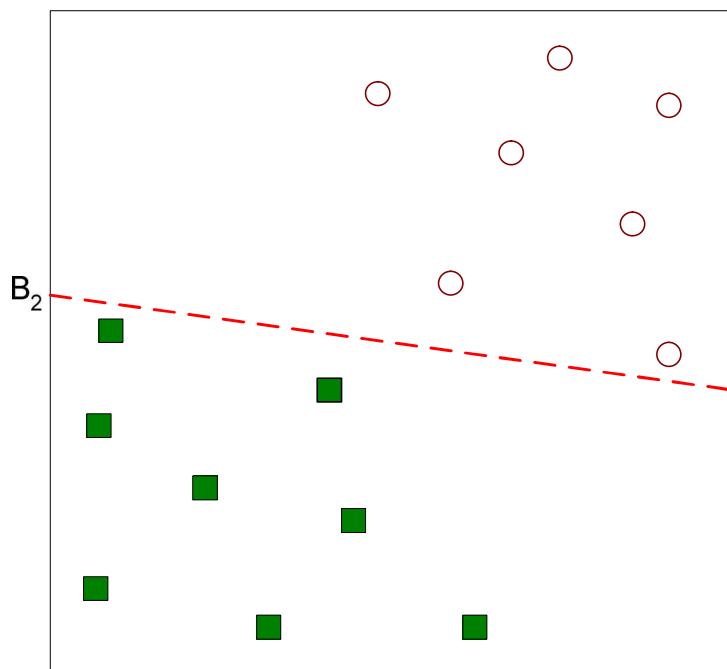
# Support Vector Machines



- One Possible Solution

41

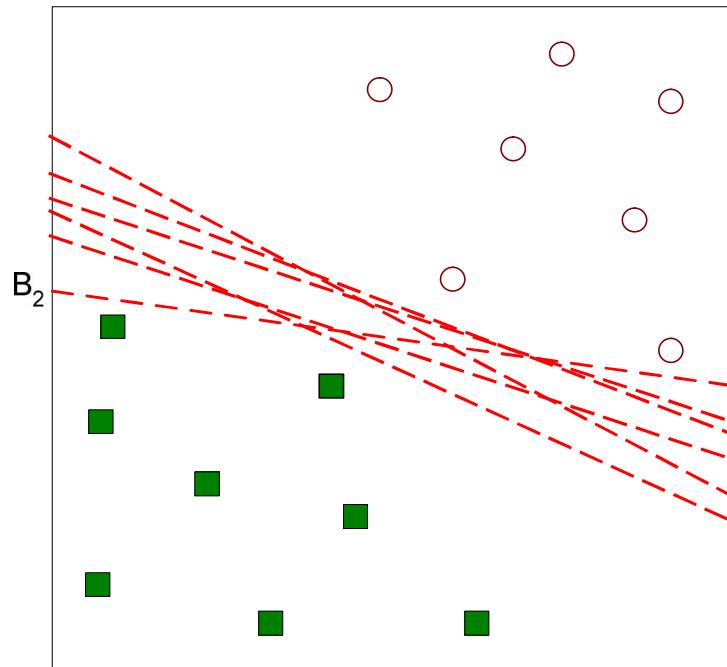
# Support Vector Machines



- Another possible solution

42

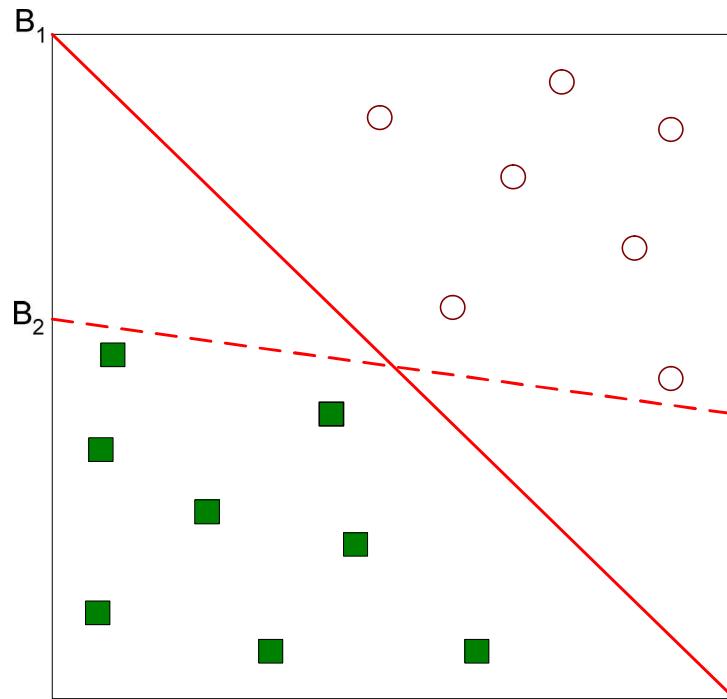
# Support Vector Machines



- Other possible solutions

43

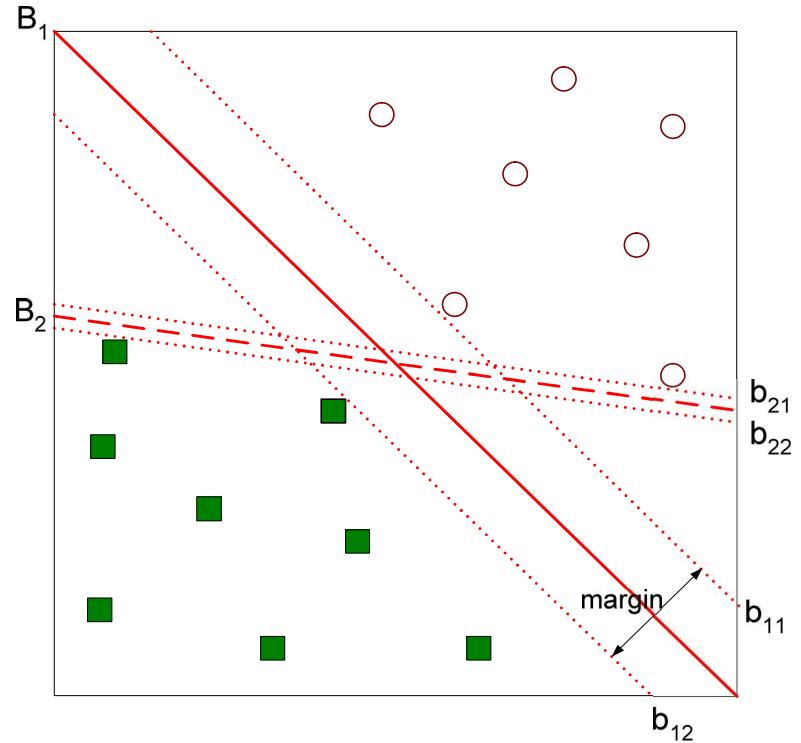
# Support Vector Machines



- Which one is better?  $B_1$  or  $B_2$ ?
- How do you define better?

44

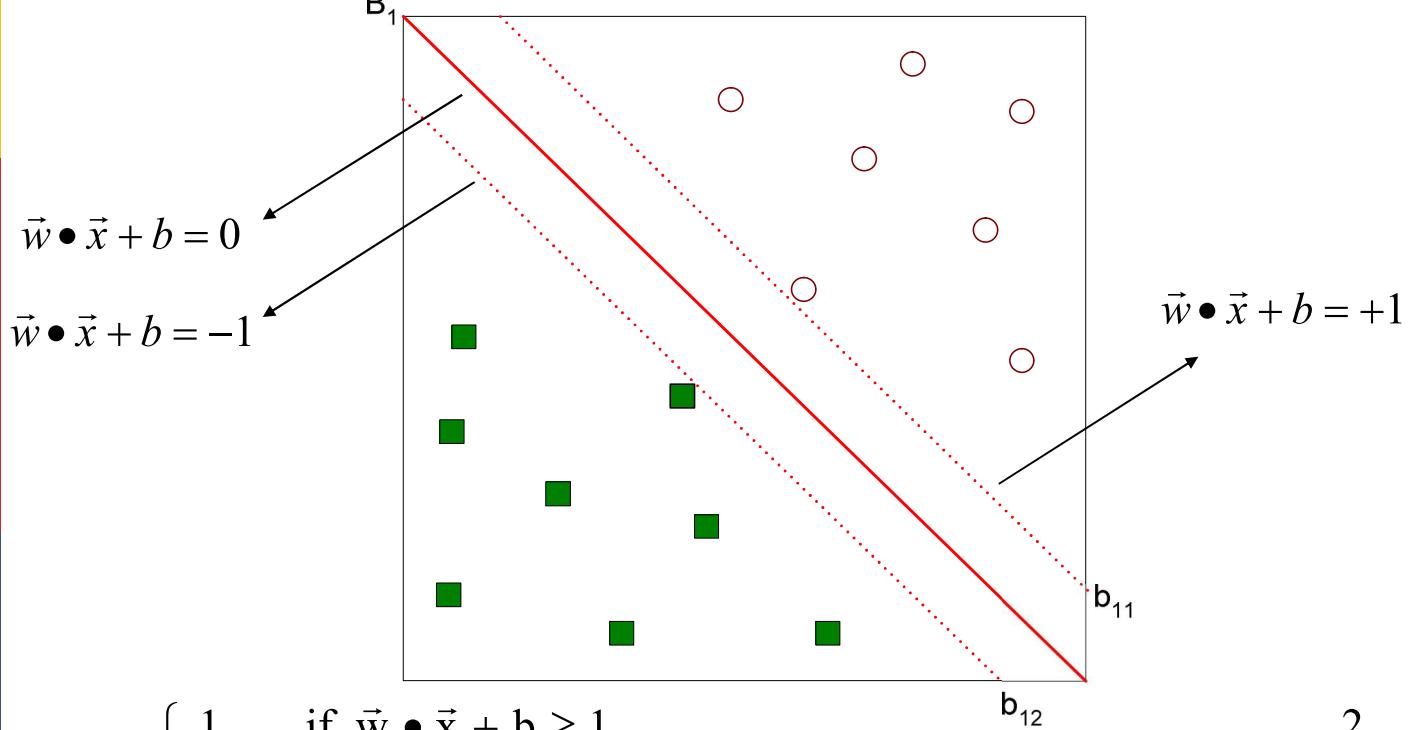
# Support Vector Machines



- Find hyperplane **maximizes** the margin => B1 is better than B2

45

# Support Vector Machines



$$f(\vec{x}) = \begin{cases} 1 & \text{if } \vec{w} \bullet \vec{x} + b \geq 1 \\ -1 & \text{if } \vec{w} \bullet \vec{x} + b \leq -1 \end{cases}$$

$$\text{Margin} = \frac{2}{\|\vec{w}\|^2}$$

46

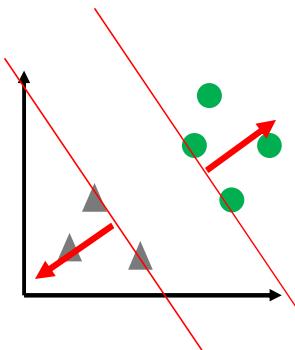
# Support Vector Machines

- We want to maximize: Margin =  $\frac{2}{\|\vec{w}\|^2}$ 
  - Which is equivalent to minimizing:  $L(w) = \frac{\|\vec{w}\|^2}{2}$
  - But subject to the following constraints:
$$f(\vec{x}_i) = \begin{cases} 1 & \text{if } \vec{w} \bullet \vec{x}_i + b \geq 1 \\ -1 & \text{if } \vec{w} \bullet \vec{x}_i + b \leq -1 \end{cases}$$
- This is a constrained optimization problem
  - Numerical approaches to solve it (e.g., quadratic programming)

47

$$\begin{aligned} w^T x_i - b &\leq -1 & \forall y_i = -1 \\ w^T x_i - b &\geq +1 & \forall y_i = +1 \end{aligned}$$

$$\Rightarrow y_i(w^T x_i - b) \geq 1$$



$\Rightarrow$  objective

$$\Rightarrow \text{maximize } \frac{2}{\|w\|^2}$$

$$\text{subject to } y_i(w^T x_i - b) - 1 \geq 0 \quad \forall x_i$$

$\Rightarrow$  minimize (use "Larange Multiplier Method")

$$L(w, b, \alpha) = \frac{1}{2} \|w\|^2 - \sum_{i=1}^N \alpha_i [y_i(w^T x_i - b) - 1]$$

where  $\alpha_i$  is "Largrange Multiplier"

48

$$L(w, b, \alpha) = \frac{1}{2} \|w\|^2 - \sum_{i=1}^N \alpha_i [y_i(w^T x_i - b) - 1]$$

Note when  $y_i(w^T x_i - b) - 1 = 0 \quad \alpha_i \geq 0$

when  $y_i(w^T x_i - b) - 1 > 0 \quad \alpha_i = 0$

$\because$ 求極小

對  $w$  微分  $\Rightarrow w - \sum_{i=1}^N \alpha_i y_i x_i = 0$

對  $b$  微分  $\Rightarrow \sum_{i=1}^N \alpha_i y_i = 0$

整體之 constraints 稱  
KKT(Karush-Kuhn-Tucker)  
conditions

49

參考 “Data Mining: Concepts and Techniques, 2e” Eq.6.39 in p.341  
“..., 3e” Eq.9.19 in p.412

$$d(X^T) = \sum_{i=1}^l \alpha_i y_i X_i X^T + b_0$$

$$\begin{aligned} \Rightarrow f(x) &= w^{*T} x - b^* \\ &= \sum_{i=1}^N \alpha_i y_i x_i x - b^* \end{aligned}$$

$\uparrow$   
*test vector*

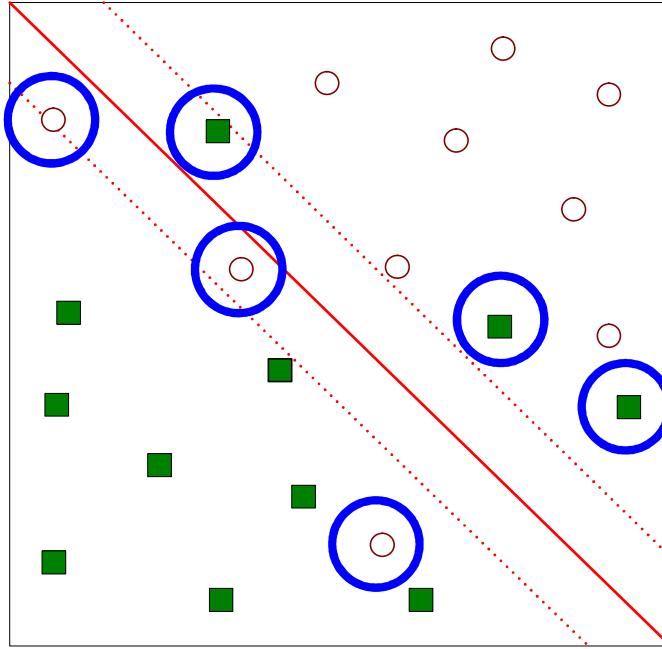
Read Vipin's book: Chapter 5.5  
Try a tiny numerical example!

*i: support vectors*

50

# Support Vector Machines

- What if the problem is not linearly separable?



51

# Support Vector Machines

- What if the problem is not linearly separable?

– Introduce slack variables

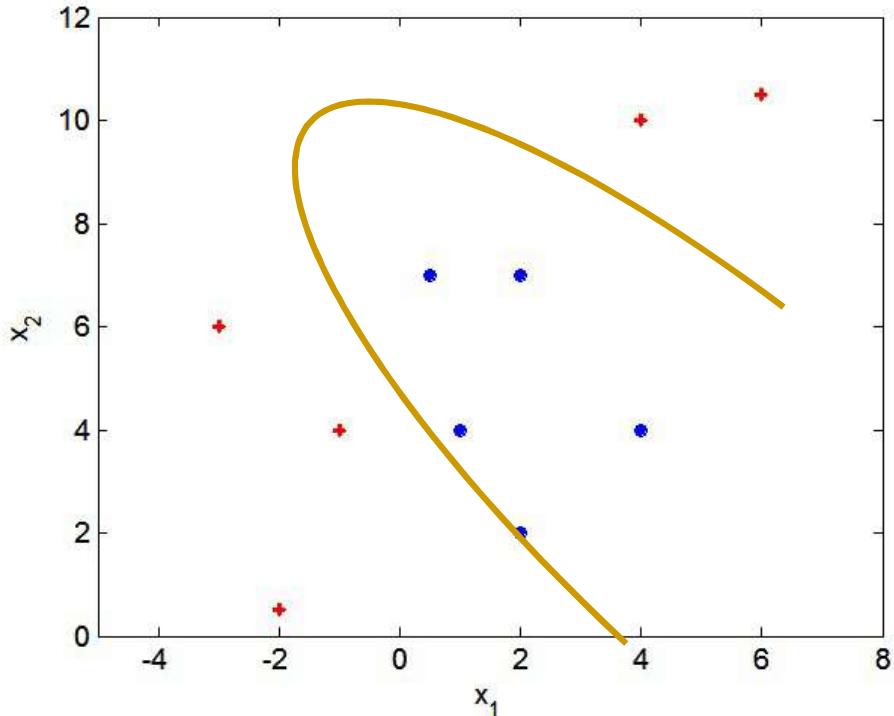
- Need to minimize: 
$$L(w) = \frac{\|\vec{w}\|^2}{2} + C \left( \sum_{i=1}^N \xi_i^k \right)$$

- Subject to:  
$$f(\vec{x}_i) = \begin{cases} 1 & \text{if } \vec{w} \bullet \vec{x}_i + b \geq 1 - \xi_i \\ -1 & \text{if } \vec{w} \bullet \vec{x}_i + b \leq -1 + \xi_i \end{cases}$$

52

# Nonlinear SVM

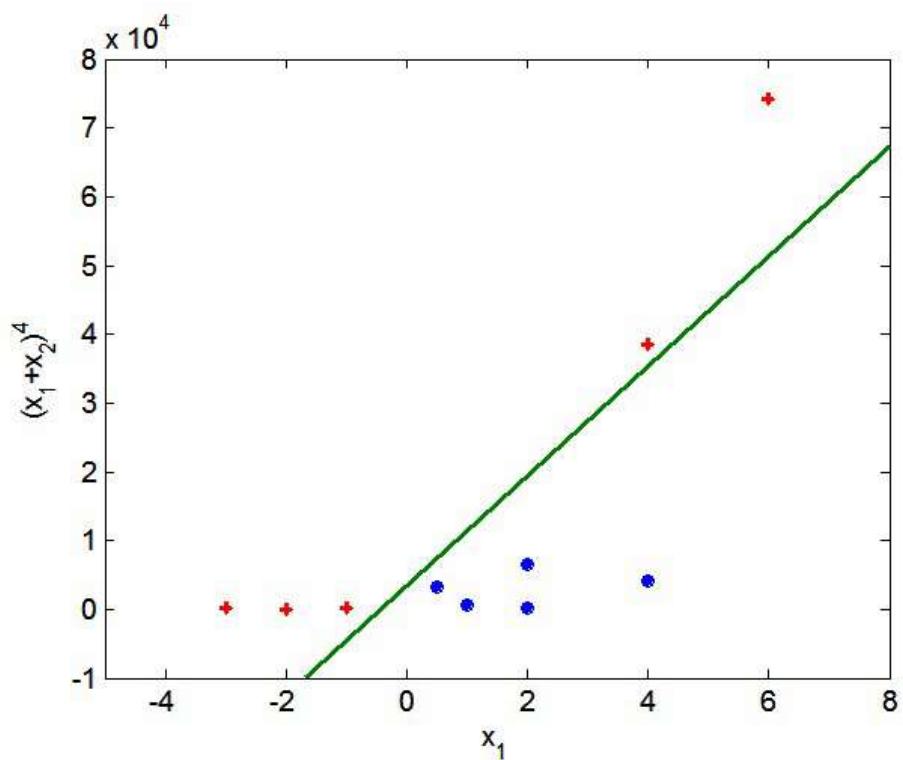
- What if decision boundary is not linear?



53

# Nonlinear SVM

- Transform data into higher dimensional space



54

# Nonlinear SVM: Attribute Transform

- To make the data linearly separable we could:
  - Project the data from the input space to a new space called “feature (or transformed) space”
  - This feature space having more dimensions than the input space we could separate the data in the feature space

55

## One possible transform

- Let's use an example projection:

$$P: I \Rightarrow F$$

$$\begin{pmatrix} x_1 \\ x_2 \end{pmatrix} \rightarrow \begin{pmatrix} x_1^2 \\ x_1 x_2 \\ x_2^2 \end{pmatrix}$$

- The inner product of two vectors  $x$  and  $y$  projected in the space  $F$  becomes:

$$\begin{pmatrix} x_1^2 \\ x_1 x_2 \\ x_2^2 \end{pmatrix} \cdot \begin{pmatrix} y_1^2 \\ y_1 y_2 \\ y_2^2 \end{pmatrix} = x_1^2 y_1^2 + x_1 x_2 y_1 y_2 + x_2^2 y_2^2$$

# Observation

- Decision boundary decided in the higher space represented by inner products after the transformation
  - Computationally costly
  - Curse of dimensionality
- This inner product in the transformed space, which is computed in the original space is known as Kernel function

57

# Kernel Function

- Kernel Function is the function that represents the inner product of some space (feature) in another (original) space
- Support Vector Machine finds an hyperplace separating the training set in a feature space induced by a kernel function used as the inner product in the algorithm.

58

# Another transform

- Let's use an example projection:

$$P: I \Rightarrow F$$

$$\begin{pmatrix} x_1 \\ x_2 \end{pmatrix} \rightarrow \begin{pmatrix} x_1^2 \\ x_1 x_2 \\ x_2^2 \end{pmatrix}$$

- The inner product of two vectors  $x$  and  $y$  projected in the space  $F$  becomes

$$\begin{pmatrix} x_1^2 \\ x_1 x_2 \\ x_2^2 \end{pmatrix} \cdot \begin{pmatrix} y_1^2 \\ y_1 y_2 \\ y_2^2 \end{pmatrix} = x_1^2 y_1^2 + x_1 x_2 y_1 y_2 + x_2^2 y_2^2$$

## SVM—Kernel functions

- Instead of computing the dot product on the transformed data tuples, it is mathematically equivalent to instead applying a kernel function  $K(X_i, X_j)$  to the original data, i.e.,  $K(X_i, X_j) = \Phi(X_i) \cdot \Phi(X_j)$
- Typical Kernel Functions
- (no need to have the form of  $X_i \cdot X_j$ )

Polynomial kernel of degree  $h$  :  $K(X_i, X_j) = (X_i \cdot X_j + 1)^h$

Gaussian radial basis function kernel :  $K(X_i, X_j) = e^{-\|X_i - X_j\|^2 / 2\sigma^2}$

Sigmoid kernel :  $K(X_i, X_j) = \tanh(\kappa X_i \cdot X_j - \delta)$

- SVM can also be used for classifying multiple ( $> 2$ ) classes and for regression analysis (with additional user parameters)

# Model Evaluation

# Model Evaluation

- Metrics for Performance Evaluation
  - How to evaluate the performance of a model?
- Methods for Performance Evaluation
  - How to obtain reliable estimates?

# Model Evaluation

- Metrics for Performance Evaluation
  - How to evaluate the performance of a model?
- Methods for Performance Evaluation
  - How to obtain reliable estimates?

63

## Metrics for Performance

		PREDICTED CLASS	
ACTUAL CLASS		Class=Yes	Class>No
	Class=Yes	a (TP)	b (FN)
	Class>No	c (FP)	d (TN)

- Most widely-used metric:

$$\text{Accuracy} = \frac{a + d}{a + b + c + d} = \frac{TP + TN}{TP + TN + FP + FN}$$

64

# Limitation of Accuracy

- Consider a 2-class problem
  - Number of Class 0 examples = 9990
  - Number of Class 1 examples = 10
- If model predicts everything to be class 0, accuracy is  $9990/10000 = 99.9\%$ 
  - Accuracy is misleading because model does not detect any class 1 example

65

## Cost Matrix

		PREDICTED CLASS	
		C(i j)	Class=Yes
ACTUAL CLASS	Class=Yes	C(Yes Yes)	C(No Yes)
	Class=No	C(Yes No)	C(No No)

$C(i|j)$ : Cost of misclassifying class j example as class i

66

# Computing Cost of Classification

Cost Matrix		PREDICTED CLASS		
		C(i j)	+	-
ACTUAL CLASS	+	-1	100	
	-	1	0	

Model M <sub>1</sub>	PREDICTED CLASS		
	+	-	
ACTUAL CLASS	+	150	40
	-	60	250

Accuracy = 80%

Cost = 3910

Model M <sub>2</sub>	PREDICTED CLASS		
	+	-	
ACTUAL CLASS	+	250	45
	-	5	200

Accuracy = 90%

Cost = 4255

67

## Cost vs Accuracy

Count	PREDICTED CLASS		
		Class=Yes	Class>No
ACTUAL CLASS	Class=Yes	a	b
	Class>No	c	d

Accuracy is proportional to cost if  
 1.  $C(\text{Yes}|\text{No})=C(\text{No}|\text{Yes}) = q$   
 2.  $C(\text{Yes}|\text{Yes})=C(\text{No}|\text{No}) = p$

$$N = a + b + c + d$$

$$\text{Accuracy} = (a + d)/N$$

Cost	PREDICTED CLASS		
		Class=Yes	Class>No
ACTUAL CLASS	Class=Yes	p	q
	Class>No	q	p

$$\begin{aligned}
 \text{Cost} &= p(a + d) + q(b + c) \\
 &= p(a + d) + q(N - a - d) \\
 &= qN - (q - p)(a + d) \\
 &= N[q - (q-p) \times \text{Accuracy}]
 \end{aligned}$$

# Cost-Sensitive Measures

$$\text{Precision (p)} = \frac{a}{a+c}$$

$$\text{Recall (r)} = \frac{a}{a+b}$$

$$\text{F - measure (F)} = \frac{2rp}{r+p} = \frac{2a}{2a+b+c}$$

- Precision is biased towards C(Yes|Yes) & C(Yes|No)
- Recall is biased towards C(Yes|Yes) & C(No|Yes)
- F-measure is biased towards all except C(No|No)

$$\text{Weighted Accuracy} = \frac{w_1 a + w_4 d}{w_1 a + w_2 b + w_3 c + w_4 d}$$

69

# Model Evaluation

- Metrics for Performance Evaluation
  - How to evaluate the performance of a model?
- **Methods for Performance Evaluation**
  - How to obtain reliable estimates?

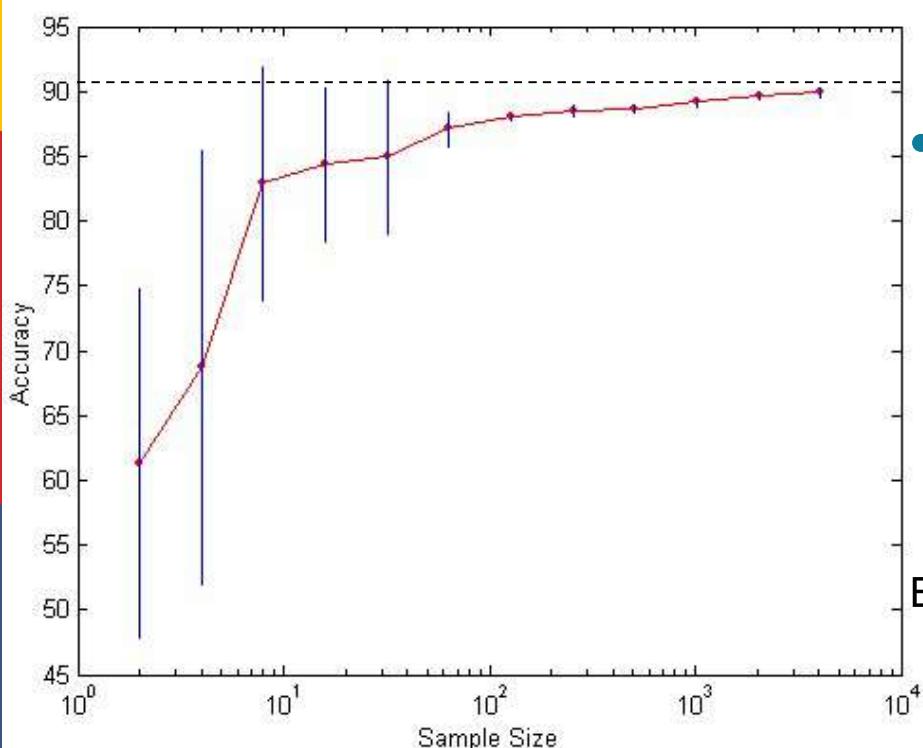
70

# Methods for Performance Evaluation

- How to obtain a reliable estimate of performance?
- Performance of a model may depend on other factors besides the learning algorithm:
  - Class distribution
  - Cost of misclassification
  - Size of training and test sets

71

## Learning Curve



- Learning curve shows how accuracy changes with varying sample size
  - Requires a sampling schedule for creating learning curve:
    - Arithmetic sampling (Langley, et al)
    - Geometric sampling (Provost et al)
- Effect of small sample size:
- Bias in the estimate
  - Variance of estimate

72

# Methods of Estimation

- Holdout
  - Reserve 2/3 for training and 1/3 for testing
- Random subsampling
  - Repeated holdout
- Cross validation
  - Partition data into  $k$  disjoint subsets
  - $k$ -fold: train on  $k-1$  partitions, test on the remaining one
  - Leave-one-out:  $k=n$
- Bootstrap
  - Sampling with replacement