

# CLASSIFICATION METHODS



# Nearest-Neighbor classifiers

2

instance-based learning

eager learners

- (1) constructing a classification model from data
- (2) applying the model to test examples

# Nearest-Neighbor classifiers

3

instance-based learning

eager learners

- (1) constructing a classification model from data
- (2) applying the model to test examples

uses specific training instances to make predictions without having to maintain an abstraction (or model) derived from data

Lazy learners

# Nearest-Neighbor classifiers

4

instance-based learning

eager learners

- (1) constructing a classification model from data
- (2) applying the model to test examples

uses specific training instances to make predictions without having to maintain an abstraction (or model) derived from data

Lazy learners

# Nearest-Neighbor classifiers

5

instance-based learning

eager learners

- (1) constructing a classification model from data
- (2) applying the model to test examples

uses specific training instances to make predictions without having to maintain an abstraction (or model) derived from data

Set of Stored Cases

Atr1	.....	AtrN	Class
			A
			B
			B
			C
			A
			C
			B

Unseen Case

Atr1	.....	AtrN

Lazy learners

# Nearest-Neighbor classifiers

6

## Rote-learner

- ❖ Memorizes entire training data and performs classification only if attributes of record match one of the training examples exactly
- ❖ problem

# Nearest-Neighbor classifiers

7

## Rote-learner

- ❖ Memorizes entire training data and performs classification only if attributes of record match one of the training examples exactly
- ❖ problem

## Nearest neighbor

- ❖ Uses  $k$  “closest” points (nearest neighbors) for performing classification

# Nearest Neighbor Classifiers

8

- Basic idea:

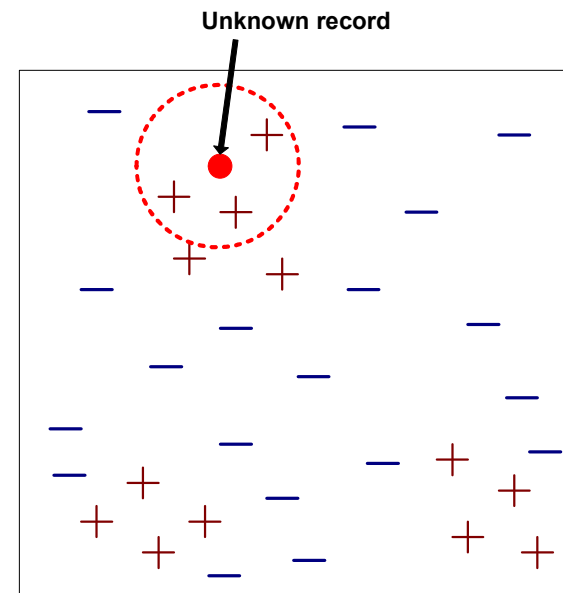
- If it walks like a duck, quacks like a duck, then it's probably a duck



# Nearest-Neighbor Classifiers

9

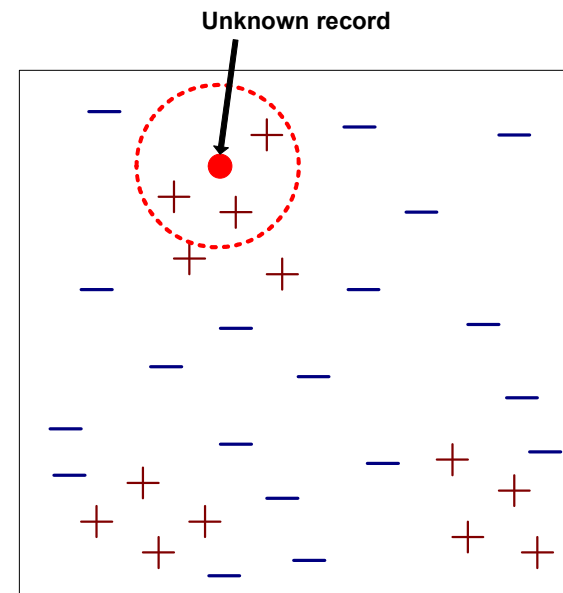
- Requires three things
  - The set of stored records
  - Distance Metric to compute distance between records
  - The value of  $k$ , the number of nearest neighbors to retrieve



# Nearest-Neighbor Classifiers

10

- 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)



---

**Algorithm 5.2** The  $k$ -nearest neighbor classification algorithm.

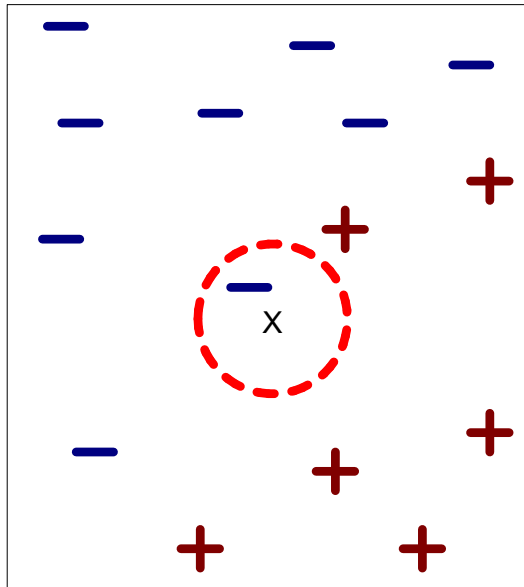
---

- 1: Let  $k$  be the number of nearest neighbors and  $D$  be the set of training examples.
  - 2: **for** each test example  $z = (\mathbf{x}', y')$  **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:    $y' = \operatorname{argmax}_v \sum_{(\mathbf{x}_i, y_i) \in D_z} I(v = y_i)$
  - 6: **end for**
- 

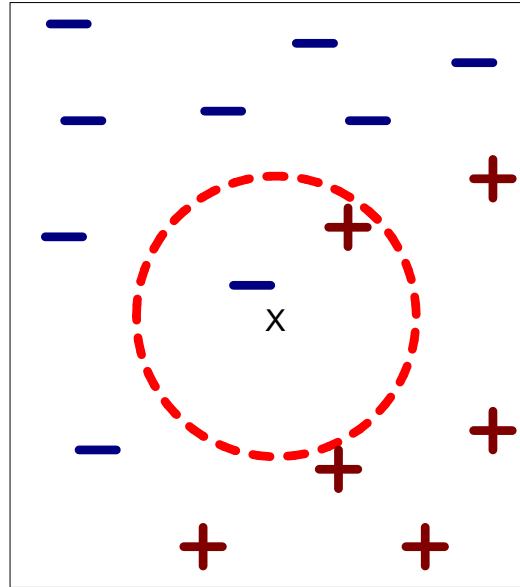
$$\text{Majority Voting: } y' = \operatorname{argmax}_v \sum_{(\mathbf{x}_i, y_i) \in D_z} I(v = y_i)$$

# Definition of Nearest Neighbor

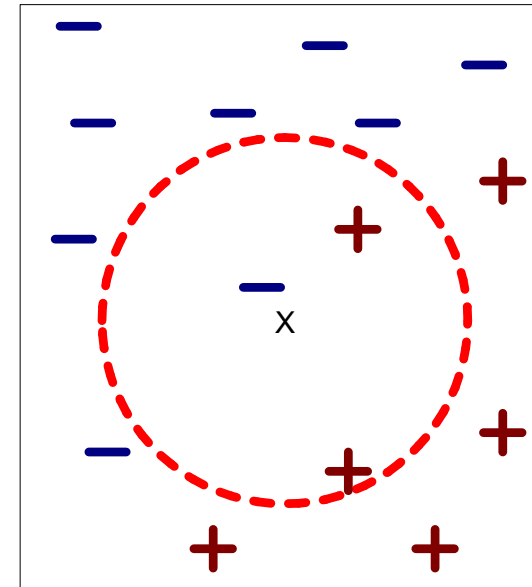
12



(a) 1-nearest neighbor



(b) 2-nearest neighbor



(c) 3-nearest neighbor

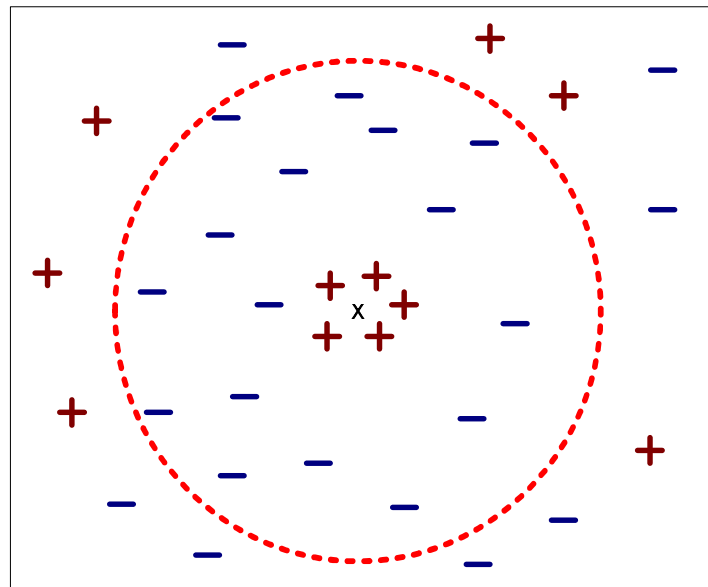
K-nearest neighbors of a record  $x$  are data points that have the  $k$  smallest distance to  $x$

# Nearest-Neighbor Classifiers

13

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



# Nearest-Neighbor Classifiers

14

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


$$w_i = 1/d(\mathbf{x}', \mathbf{x}_i)^2$$

# Nearest-Neighbor Classifiers

15

- ✓ Lazy learners
- ✓ instance-based learning, which uses specific training instances to make predictions without having to maintain an abstraction (or model) derived from data.
- ✓ classifying a test example can be quite expensive
- ✓ eager learners often spend the bulk of their computing resources for model building. Once a model has been built, classifying a test example is extremely fast.
- ✓ Nearest-neighbor classifiers can produce arbitrarily shaped decision boundaries
- ✓ appropriate proximity measure, data preprocessing steps

16

# Bayesian Classifiers



# Bayes Classifier

17

- A probabilistic framework for solving classification problems

- Conditional Probability: 
$$P(C | A) = \frac{P(A, C)}{P(A)}$$

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

- Bayes theorem:

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

# Example of Bayes Theorem

18

## □ 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?

# Example of Bayes Theorem

19

## □ 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?

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

# Bayesian Classifiers

20

- Consider each attribute and class label as random variables
- Given a record with attributes  $(x_1, x_2, \dots, x_n)$ 
  - ▣ Goal is to predict class  $Y$

# Bayesian Classifiers

21

- Consider each attribute and class label as random variables
- Given a record with attributes  $(x_1, x_2, \dots, x_n)$ 
  - ▣ Goal is to predict class  $Y$
  - ▣ Specifically, we want to find the value of  $Y$  that maximizes  $P(Y \mid x_1, x_2, \dots, x_n)$
- Can we estimate  $P(Y \mid x_1, x_2, \dots, x_n)$  directly from data?

# Bayesian Classifiers

22

- Approach:
  - ▣ compute the posterior probability  $P(Y \mid x_1, x_2, \dots, x_n)$  for all values of  $Y$  using the Bayes theorem

$$P(Y \mid x_1 x_2 \dots x_n) = \frac{P(x_1 x_2 \dots x_n \mid Y) P(Y)}{P(x_1 x_2 \dots x_n)}$$

# Bayesian Classifiers

23

## □ Approach:

- compute the posterior probability  $P(Y \mid x_1, x_2, \dots, x_n)$  for all values of  $Y$  using the Bayes theorem

$$P(Y \mid x_1 x_2 \dots x_n) = \frac{P(x_1 x_2 \dots x_n \mid Y) P(Y)}{P(x_1 x_2 \dots x_n)}$$

- Choose value of  $Y$  that maximizes  $P(Y \mid x_1, x_2, \dots, x_n)$

- Equivalent to choosing value of  $Y$  that maximizes  $P(x_1, x_2, \dots, x_n \mid Y) P(Y)$

- How to estimate  $P(x_1, x_2, \dots, x_n \mid Y)$ ?

24

# Naïve Bayes Classifier



# Naïve Bayes Classifier

25

- Assume independence among attributes when class is given:
  - ▣  $P(x_1, x_2, \dots, x_n | Y) = P(x_1 | Y_i) P(x_2 | Y_i) \dots P(x_n | Y_i)$
  - ▣ Can estimate  $P(x_i | Y_i)$  for all  $x_i$  and  $Y_i$ .
  - ▣ New point is classified to  $Y_i$  if  $P(Y_i) \prod P(x_i | Y_i)$  is maximal.

# How to Estimate Probabilities from Data?

26

<i>Tid</i>	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

# How to Estimate Probabilities from Data?

27

<i>Tid</i>	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

□ Class:  $P(Y_k) = N_k/N$

□ e.g.,  $P(\text{No}) = 7/10$ ,  
 $P(\text{Yes}) = 3/10$

# How to Estimate Probabilities from Data?

28

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

□ Class:  $P(Y_k) = N_k/N$

□ e.g.,  $P(\text{No}) = 7/10$ ,  
 $P(\text{Yes}) = 3/10$

□ For discrete attributes:

$$P(x_i \mid Y_k) = |x_{ik}| / N_k$$

□ where  $|x_{ik}|$  is number of instances having attribute  $x_i$  and belongs to class  $Y_k$

□ Examples:

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

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

Tid	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

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

# How to Estimate Probabilities from Data?

30

- For continuous attributes:
  - ▣ **Discretize** the range into bins

# How to Estimate Probabilities from Data?

31

- For continuous attributes:
  - ▣ Discretize the range into bins
  - ▣ 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  $P(x_i | Y)$

# How to Estimate Probabilities from Data?

32

<i>Tid</i>	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:

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

□ One for each  $(x_i, Y_i)$  pair



# How to Estimate Probabilities from Data?

33

<i>Tid</i>	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:

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

□ One for each  $(x_i, Y_i)$  pair

□ For (Income, Class=No):

□ If Class=No

■ sample mean = 110

■ sample variance = 2975

# How to Estimate Probabilities from Data?

34

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:

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

□ One for each  $(x_i, Y_i)$  pair

□ For (Income, Class=No):

□ If Class=No

■ sample mean = 110

■ sample variance = 2975

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

# example

35

Tid	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

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

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

**Test data:**  $\mathbf{X} = (\text{Home Owner}=\text{No}, \text{Marital Status} = \text{Married}, \text{Income} = \$120\text{K})$

$$P(\text{No}|\mathbf{X}) \text{ and } P(\text{Yes}|\mathbf{X})$$

New point is classified to  $Y_i$  if  $P(Y_i) \prod P(x_i | Y_i)$  is maximal.

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

For Annual Income:

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

$$P(Y_i) \prod P(x_i | Y_i)$$

**X** = (Home Owner=No, Marital Status = Married, Income = \$120K)

$$P(\text{Yes}) = 0.3 \text{ and } P(\text{No}) = 0.7$$

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

For Annual Income:

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

$P(\text{Yes}) = 0.3$  and  $P(\text{No}) = 0.7$

$$P(Y_i) \prod P(x_i | Y_i)$$

$\mathbf{X} = (\text{Home Owner}=\text{No}, \text{Marital Status} = \text{Married}, \text{Income} = \$120\text{K})$

$$\begin{aligned}
 P(\mathbf{X}|\text{No}) &= P(\text{Home Owner} = \text{No}|\text{No}) \times P(\text{Status} = \text{Married}|\text{No}) \\
 &\quad \times P(\text{Annual Income} = \$120\text{K}|\text{No}) \\
 &= 4/7 \times 4/7 \times 0.0072 = 0.0024.
 \end{aligned}$$

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

For Annual Income:

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

$P(\text{Yes}) = 0.3$  and  $P(\text{No}) = 0.7$

$$P(Y_i) \prod P(x_i | Y_i)$$

$\mathbf{X} = (\text{Home Owner}=\text{No}, \text{Marital Status} = \text{Married}, \text{Income} = \$120\text{K})$

$$\begin{aligned}
 P(\mathbf{X}|\text{No}) &= P(\text{Home Owner} = \text{No}|\text{No}) \times P(\text{Status} = \text{Married}|\text{No}) \\
 &\quad \times P(\text{Annual Income} = \$120\text{K}|\text{No}) \\
 &= 4/7 \times 4/7 \times 0.0072 = 0.0024.
 \end{aligned}$$

$$\begin{aligned}
 P(\mathbf{X}|\text{Yes}) &= P(\text{Home Owner} = \text{No}|\text{Yes}) \times P(\text{Status} = \text{Married}|\text{Yes}) \\
 &\quad \times P(\text{Annual Income} = \$120\text{K}|\text{Yes}) \\
 &= 1 \times 0 \times 1.2 \times 10^{-9} = 0.
 \end{aligned}$$

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

# Bayes Error Rate

39

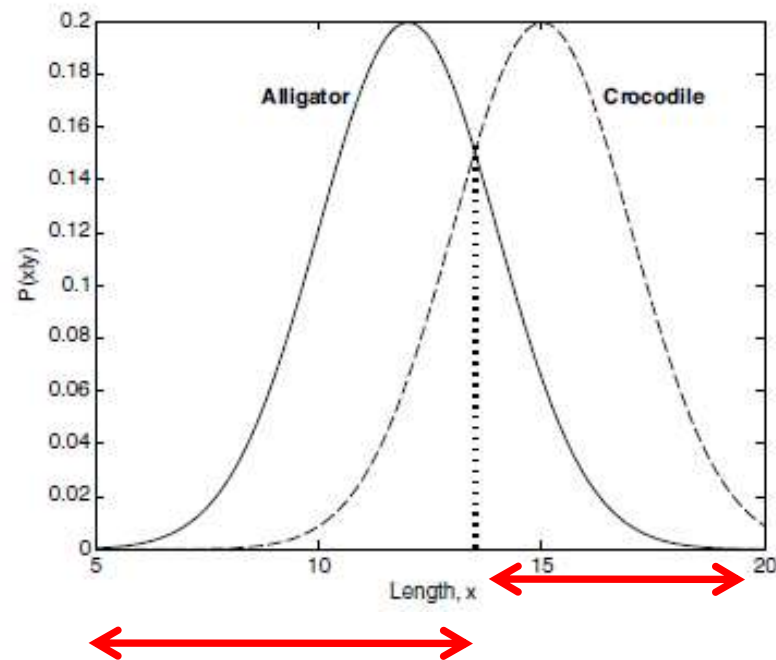
Example:

identifying alligators and crocodiles based on their lengths  
their prior probabilities are the same

$$P(X|\text{Crocodile}) = \frac{1}{\sqrt{2\pi} \cdot 2} \exp \left[ -\frac{1}{2} \left( \frac{X - 15}{2} \right)^2 \right]$$
$$P(X|\text{Alligator}) = \frac{1}{\sqrt{2\pi} \cdot 2} \exp \left[ -\frac{1}{2} \left( \frac{X - 12}{2} \right)^2 \right]$$

# Bayes Error Rate

40





41

## Directed graphical models (Bayesian Belief network)

# Introduction

42

- multiple correlated variables, such as words in a document, pixels in an image, or genes in a microarray
- compactly *represent* the **joint distribution**  $p(\mathbf{x})$  (probabilistic modeling)
- *learn* the parameters of this distribution with a reasonable amount of data (learning)
- *infer* one set of variables given another (inference)

# Introduction

43

- **Chain rule**

$$p(x_{1:V}) = p(x_1)p(x_2|x_1)p(x_3|x_2, x_1)p(x_4|x_1, x_2, x_3) \dots p(x_V|x_{1:V-1})$$

- becomes more and more complicated to represent the conditional distributions

# Introduction

44

## □ Chain rule

$$p(x_{1:V}) = p(x_1)p(x_2|x_1)p(x_3|x_2, x_1)p(x_4|x_1, x_2, x_3) \dots p(x_V|x_{1:V-1})$$

- becomes more and more complicated to represent the conditional distributions
- suppose all the variables have  $K$  states
- $p(x_2|x_1), p(x_3|x_1, x_2)$

# Conditional independence(CI)

45

- The key to efficiently representing large joint distributions is to make some assumptions about conditional independence

$$X \perp Y | Z \iff p(X, Y | Z) = p(X | Z)p(Y | Z)$$

# Conditional independence(CI)

46

- The key to efficiently representing large joint distributions is to make some assumptions about conditional independence

$$X \perp Y | Z \iff p(X, Y | Z) = p(X | Z)p(Y | Z)$$

$$x_{t+1} \perp \mathbf{x}_{1:t-1} | x_t$$

# Conditional independence(CI)

47


- The key to efficiently representing large joint distributions is to make some assumptions about conditional independence

$$X \perp Y | Z \iff p(X, Y | Z) = p(X | Z)p(Y | Z)$$

$$x_{t+1} \perp \mathbf{x}_{1:t-1} | x_t$$

**Markov  
Assumption**

$$p(x_{1:V}) = p(x_1)p(x_2|x_1)p(x_3|x_2, x_1)p(x_4|x_1, x_2, x_3) \dots p(x_V|x_{1:V-1})$$


$$p(\mathbf{x}_{1:V}) = p(x_1) \prod_{t=1}^V p(x_t | x_{t-1})$$

# Graphical Models

48

- first-order Markov assumption is useful for defining distributions on 1d sequences
- A **graphical model (GM)** is a way to represent a joint distribution by making CI assumptions
- nodes in the graph represent random variables
- lack of edges represents CI assumptions



# Graphical Models

49

- first-order Markov assumption is useful for defining distributions on 1d sequences
- A **graphical model (GM)** is a way to represent a joint distribution by making CI assumptions
- nodes in the graph represent random variables
- lack of edges represents CI assumptions

graph  $G = (\mathcal{V}, \mathcal{E})$

$$\mathcal{V} = \{1, \dots, V\}$$

$$\mathcal{E} = \{(s, t) : s, t \in \mathcal{V}\}$$

adjacency matrix

$G(s, t) = 1$  to denote  $(s, t) \in \mathcal{E}$

# Graph terminology

50

- **Parent** For a directed graph, the **parents** of a node is the set of all nodes that feed into it:  
 $\text{pa}(s) \triangleq \{t : G(t, s) = 1\}.$
- **Child** For a directed graph, the **children** of a node is the set of all nodes that feed out of it:  
 $\text{ch}(s) \triangleq \{t : G(s, t) = 1\}.$
- **Root** For a directed graph, a **root** is a node with no parents.
- **Path or trail** A **path** or **trail**  $s \rightsquigarrow t$  is a series of directed edges leading from  $s$  to  $t$ .

# Graph terminology

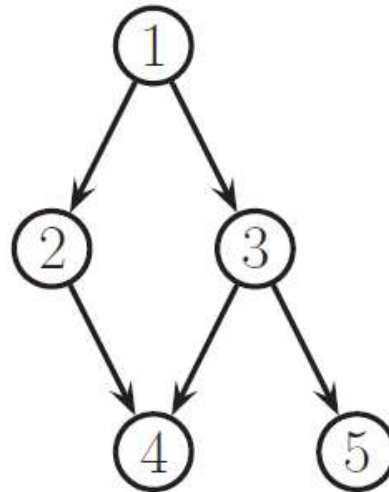
51

- **Ancestors** For a directed graph, the **ancestors** are the parents, grand-parents, etc of a node. That is, the ancestors of  $t$  is the set of nodes that connect to  $t$  via a trail:  $\text{anc}(t) \triangleq \{s : s \rightsquigarrow t\}$ .
- **Descendants** For a directed graph, the **descendants** are the children, grand-children, etc of a node. That is, the descendants of  $s$  is the set of nodes that can be reached via trails from  $s$ :  $\text{desc}(s) \triangleq \{t : s \rightsquigarrow t\}$ .
- **Cycle or loop** For any graph, we define a **cycle** or **loop** to be a series of nodes such that we can get back to where we started by following edges

# Graph terminology

52

- **DAG** A **directed acyclic graph** or **DAG** is a directed graph with no directed cycles.



- **Topological ordering** For a DAG, a **topological ordering** or **total ordering** is a numbering of the nodes such that parents have lower numbers than their children.

# Directed graphical models

53

- A **directed graphical model** or **DGM** is a GM whose graph is a DAG
- Bayesian Network, Causal Network
- topological ordering
- **ordered Markov property**

# Directed graphical models

54

- A **directed graphical model** or **DGM** is a GM whose graph is a DAG
- Bayesian Network, Causal Network
- topological ordering
- **ordered Markov property**

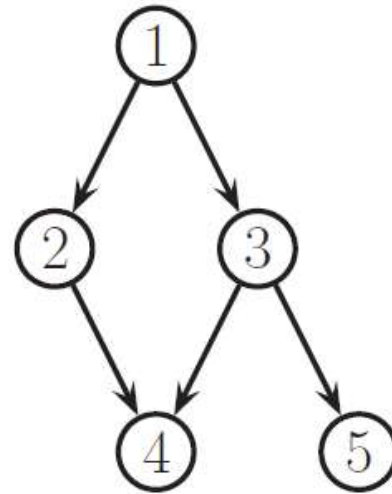
**Property 1 (Conditional Independence).** *A node in a Bayesian network is conditionally independent of its non-descendants, if its parents are known.*



natural generalization of the  
first-order Markov property

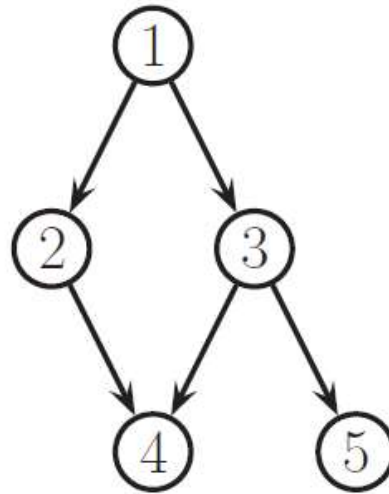
# Directed graphical models

55



# Directed graphical models

56

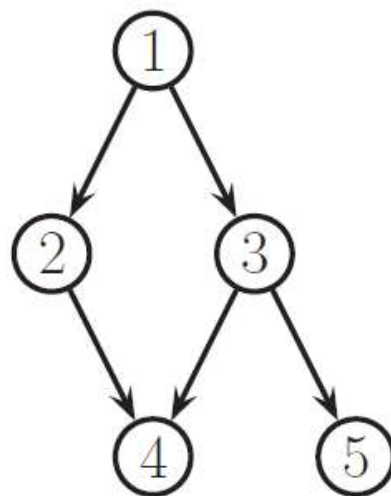


$$p(\mathbf{x}_{1:5}) = p(x_1)p(x_2|x_1)p(x_3|x_1, \cancel{x_2})p(x_4|\cancel{x_1}, x_2, x_3)p(x_5|\cancel{x_1}, \cancel{x_2}, x_3, \cancel{x_4})$$



# Directed graphical models

57



$$\begin{aligned} p(\mathbf{x}_{1:5}) &= p(x_1)p(x_2|x_1)p(x_3|x_1, \cancel{x_2})p(x_4|\cancel{x_1}, x_2, x_3)p(x_5|\cancel{x_1}, \cancel{x_2}, x_3, \cancel{x_4}) \\ &= p(x_1)p(x_2|x_1)p(x_3|x_1)p(x_4|x_2, x_3)p(x_5|x_3) \end{aligned}$$

$$p(\mathbf{x}_{1:V}|G) = \prod_{t=1}^V p(x_t | \mathbf{x}_{\text{pa}(t)})$$

**What are the random variables in Bayesian learning that we want to jointly model them?**

**What are the random variables in Bayesian learning that we want to jointly model them?**

$$P(Y \mid x_1 x_2 \dots x_n) = \frac{P(x_1 x_2 \dots x_n \mid Y) P(Y)}{P(x_1 x_2 \dots x_n)}$$

# Example: Naive Bayes classifiers

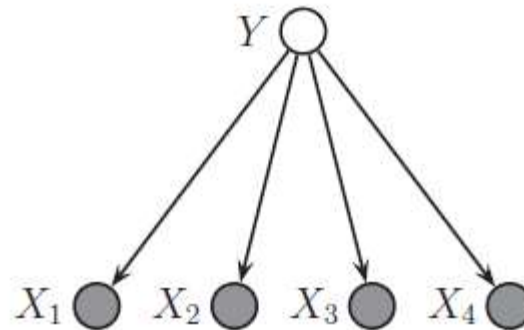
60

features are conditionally independent given the class label

# Example: Naive Bayes classifiers

61

features are conditionally independent given the class label

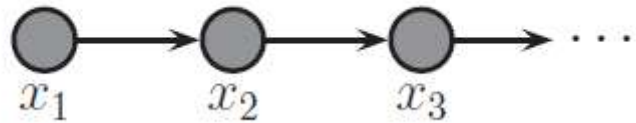


$$p(y, \mathbf{x}) = p(y) \prod_{j=1}^D p(x_j | y)$$

# Example: Markov Chain

62

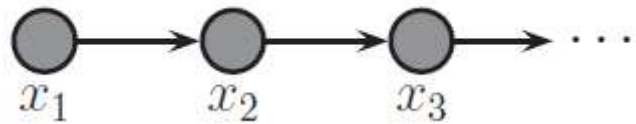
first-order Markov chain



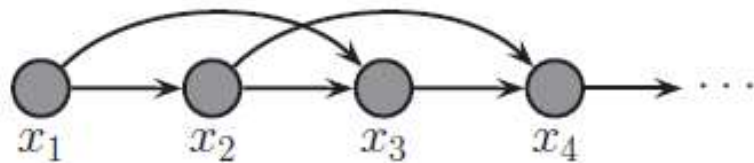
# Example: Markov Chain

63

first-order Markov chain



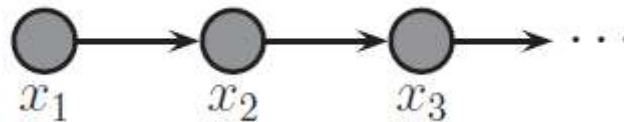
second order Markov chain



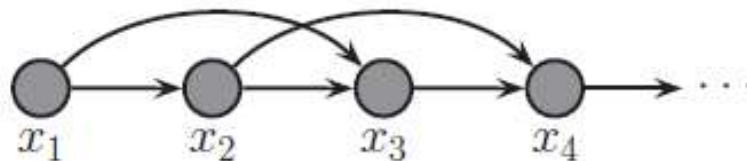
# Example: Markov Chain

64

first-order Markov chain



second order Markov chain

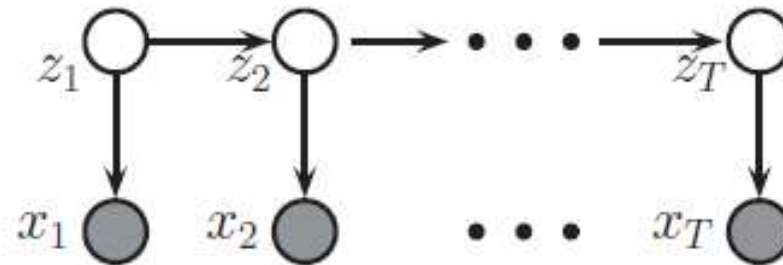


$$p(\mathbf{x}_{1:T}) = p(x_1, x_2)p(x_3|x_1, x_2)p(x_4|x_2, x_3) \dots = p(x_1, x_2) \prod_{t=3}^T p(x_t|x_{t-1}, x_{t-2})$$



# Example: HMM

65



- ❖ hidden variables often represent quantities of interest
- ❖ estimate the hidden state given the data

# Inference

66

- ✓ graphical models : compact way to define joint probability distributions
- ✓ The main use for such a joint distribution is to perform **probabilistic inference**

task of estimating unknown quantities from known quantities

set of correlated random variables with joint distribution

$$p(\mathbf{x}_{1:V} | \theta)$$

visible variables and hidden variables

# Inference

67

$$p(\mathbf{x}_h | \mathbf{x}_v, \boldsymbol{\theta}) = \frac{p(\mathbf{x}_h, \mathbf{x}_v | \boldsymbol{\theta})}{p(\mathbf{x}_v | \boldsymbol{\theta})} = \frac{p(\mathbf{x}_h, \mathbf{x}_v | \boldsymbol{\theta})}{\sum_{\mathbf{x}'_h} p(\mathbf{x}'_h, \mathbf{x}_v | \boldsymbol{\theta})}$$

In Classification : class label :Y

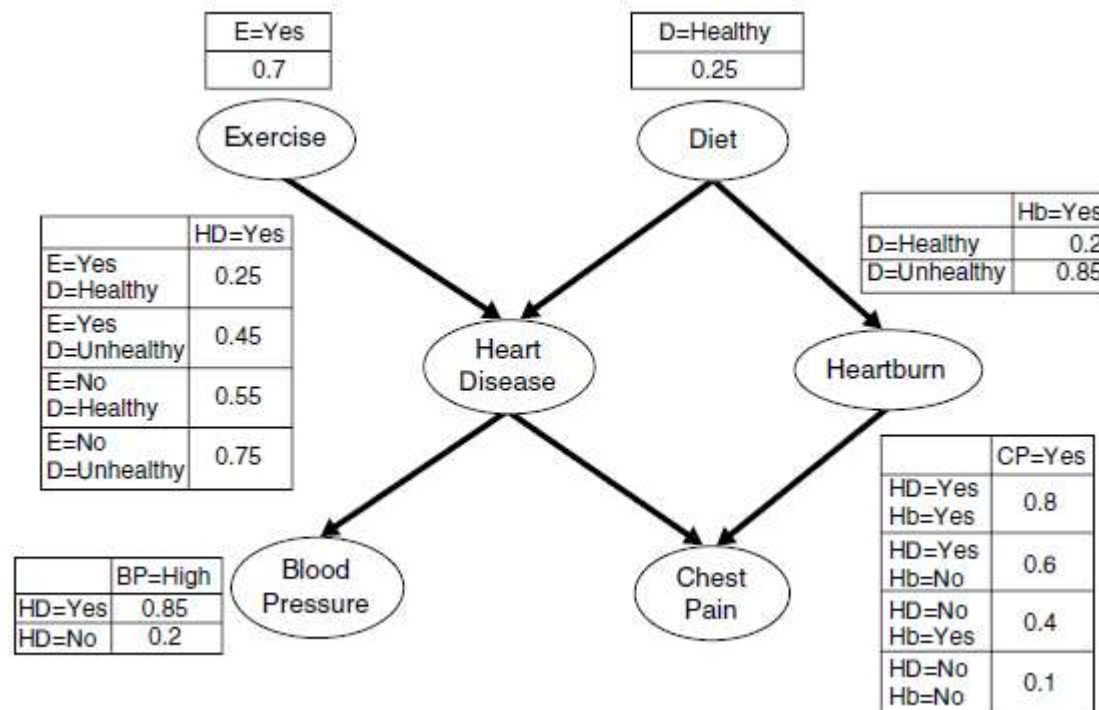
Model Paramters

# Example of Inferencing Using

68

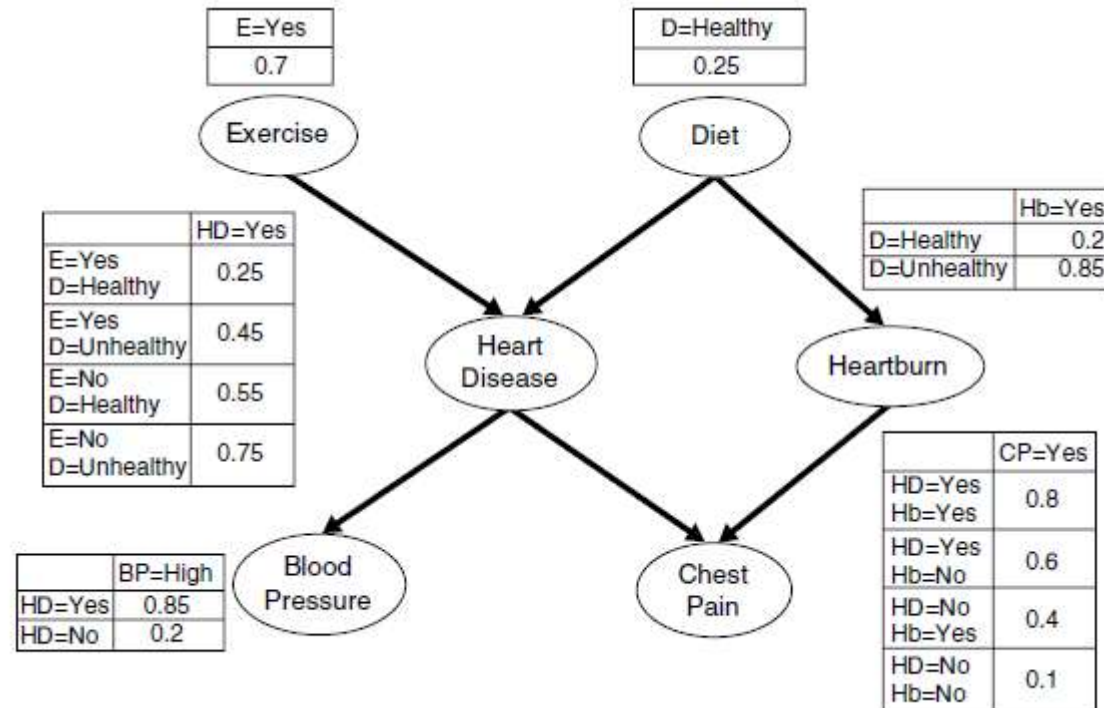
## Case 1: No Prior Information

Without any prior information, we can determine whether the person is likely to have heart disease



# Example of Inferencing Using

69



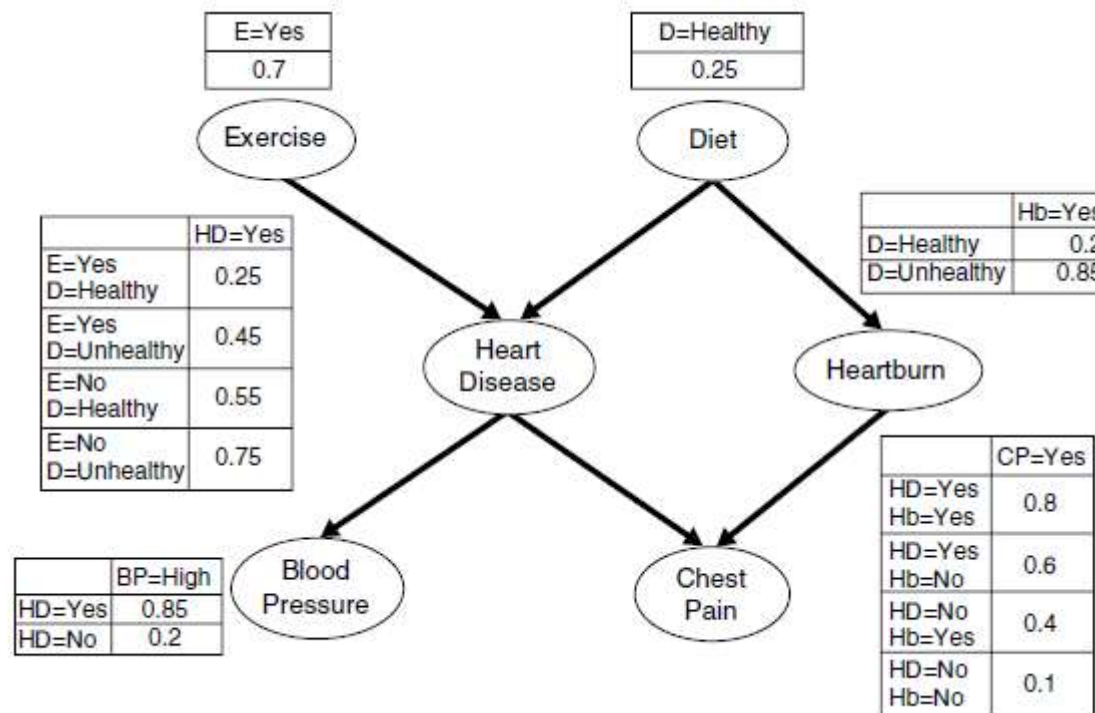
$$\begin{aligned}
 P(\text{HD} = \text{Yes}) &= \sum_{\alpha} \sum_{\beta} P(\text{HD} = \text{Yes} | E = \alpha, D = \beta) P(E = \alpha, D = \beta) \\
 &= \sum_{\alpha} \sum_{\beta} P(\text{HD} = \text{Yes} | E = \alpha, D = \beta) P(E = \alpha) P(D = \beta) \\
 &= 0.25 \times 0.7 \times 0.25 + 0.45 \times 0.7 \times 0.75 + 0.55 \times 0.3 \times 0.25 \\
 &\quad + 0.75 \times 0.3 \times 0.75 \\
 &= 0.49.
 \end{aligned}$$

# Example

70

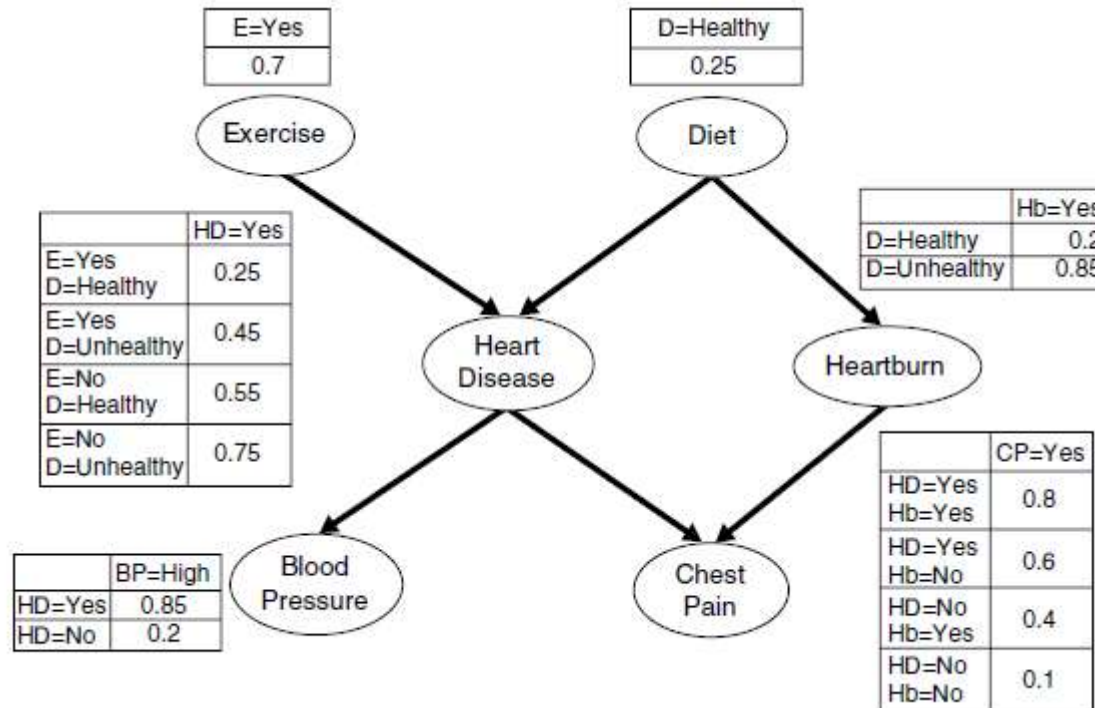
## Case 2: High Blood Pressure

If the person has high blood pressure, we can make a diagnosis about heart disease



# Example

71

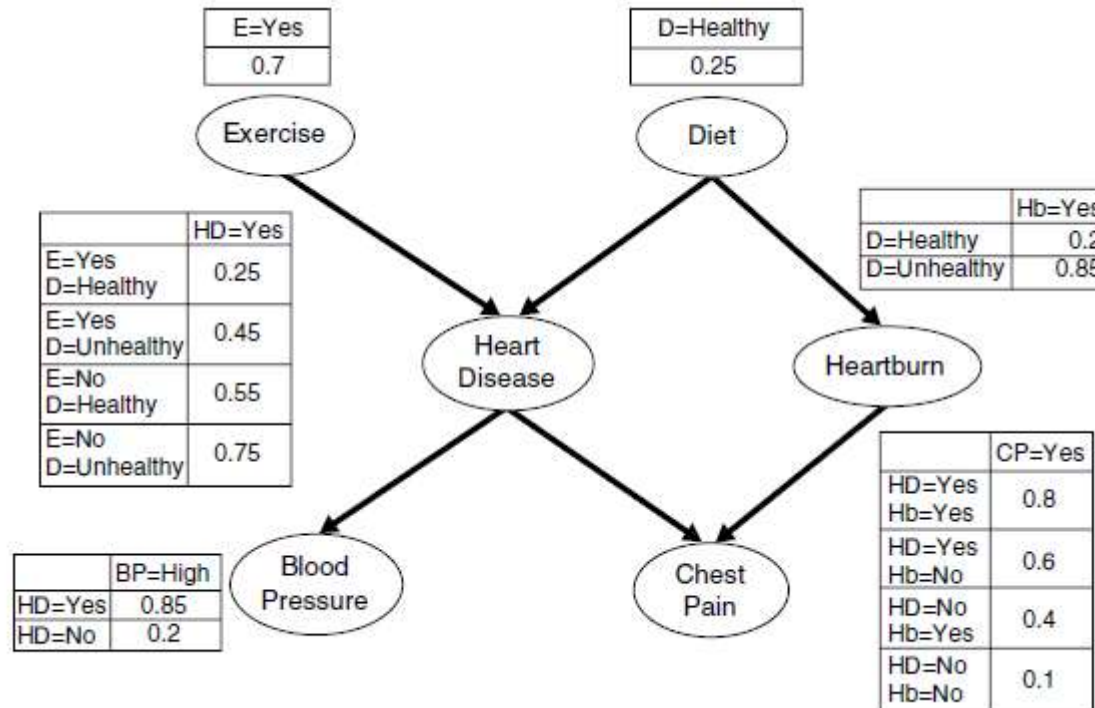


$$P(\text{HD} = \text{Yes} | \text{BP} = \text{High})$$

$$\begin{aligned}
 &= \frac{P(\text{BP} = \text{High} | \text{HD} = \text{Yes}) P(\text{HD} = \text{Yes})}{P(\text{BP} = \text{High})} \\
 &= \frac{0.85 \times 0.49}{0.5185} = 0.8033.
 \end{aligned}$$

# Example

72



$$P(\text{HD} = \text{Yes} | \text{BP} = \text{High})$$

$$\begin{aligned} P(\text{BP} = \text{High}) &= \sum_{\gamma} P(\text{BP} = \text{High} | \text{HD} = \gamma) P(\text{HD} = \gamma) \\ &= 0.85 \times 0.49 + 0.2 \times 0.51 = 0.5185. \end{aligned}$$

$$\begin{aligned} &= \frac{P(\text{BP} = \text{High} | \text{HD} = \text{Yes}) P(\text{HD} = \text{Yes})}{P(\text{BP} = \text{High})} \\ &= \frac{0.85 \times 0.49}{0.5185} = 0.8033. \end{aligned}$$

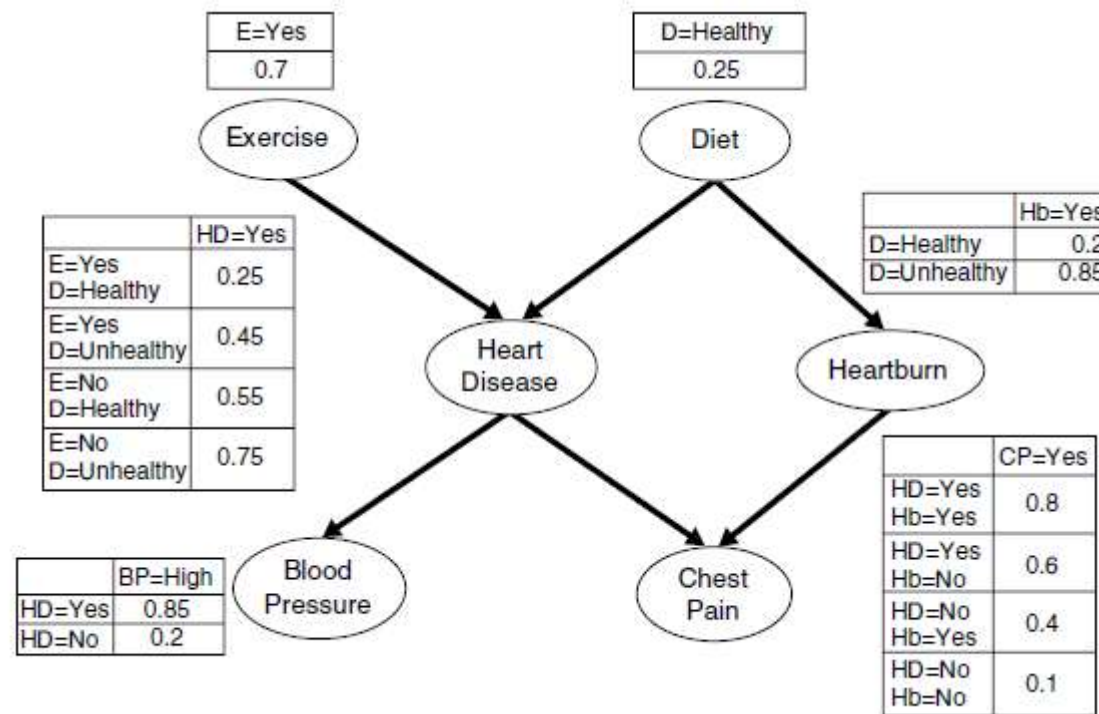


# Example

73

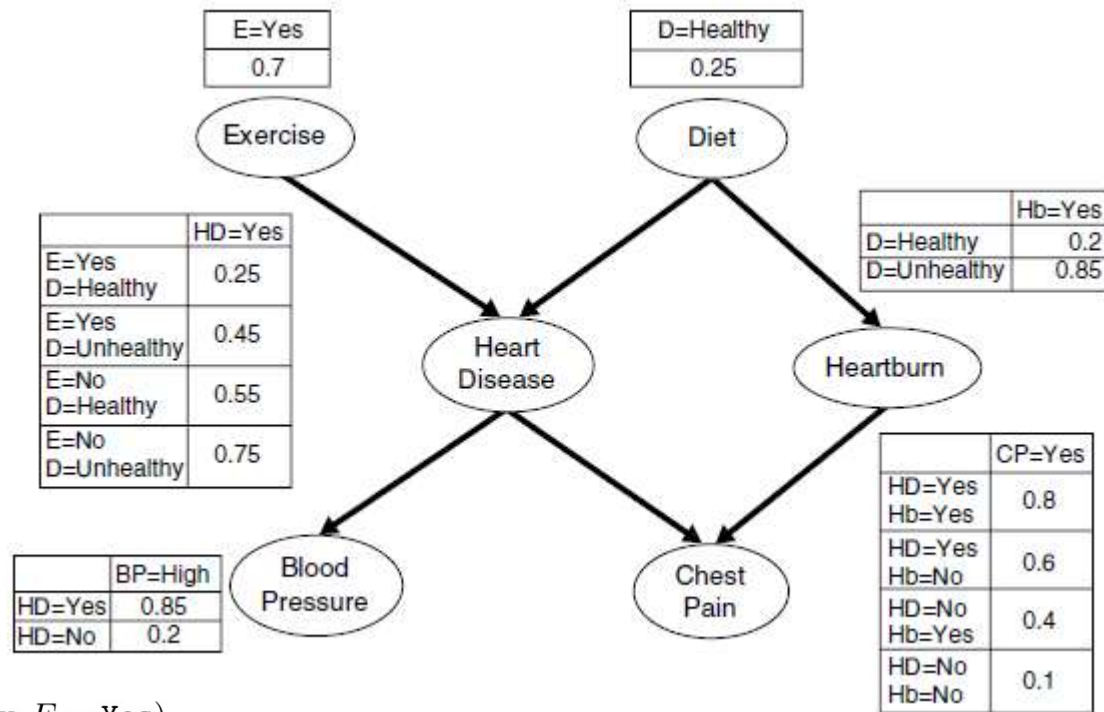
## Case 3: High Blood Pressure, Healthy Diet, and Regular Exercise

Suppose we are told that the person exercises regularly and eats a healthy diet. How does the new information affect our diagnosis?



# Example

74



$$\begin{aligned}
 & P(\text{HD} = \text{Yes} | \text{BP} = \text{High}, D = \text{Healthy}, E = \text{Yes}) \\
 = & \left[ \frac{P(\text{BP} = \text{High} | \text{HD} = \text{Yes}, D = \text{Healthy}, E = \text{Yes})}{P(\text{BP} = \text{High} | D = \text{Healthy}, E = \text{Yes})} \right] \\
 & \times P(\text{HD} = \text{Yes} | D = \text{Healthy}, E = \text{Yes}) \\
 = & \frac{P(\text{BP} = \text{High} | \text{HD} = \text{Yes}) P(\text{HD} = \text{Yes} | D = \text{Healthy}, E = \text{Yes})}{\sum_{\gamma} P(\text{BP} = \text{High} | \text{HD} = \gamma) P(\text{HD} = \gamma | D = \text{Healthy}, E = \text{Yes})} \\
 = & \frac{0.85 \times 0.25}{0.85 \times 0.25 + 0.2 \times 0.75} \\
 = & 0.5862,
 \end{aligned}$$

75

# Support Vector Machine (SVM)

# Hyperplanes

76

A *hyperplane* is a set of the form

$$\{x \mid a^T x = b\},$$

$a \in \mathbf{R}^n$ ,  $a \neq 0$ , and  $b \in \mathbf{R}$ .

$a$  is the normal vector

# Hyperplanes

77

A *hyperplane* is a set of the form

$$\{x \mid a^T x = b\},$$

$a \in \mathbf{R}^n$ ,  $a \neq 0$ , and  $b \in \mathbf{R}$ .

$a$  is the normal vector

Geometrical interpretation

$x_0$  is any point in the hyperplane

$$a^T x_0 = b$$

$$\{x \mid a^T (x - x_0) = 0\},$$

# Hyperplanes

78

A *hyperplane* is a set of the form

$$\{x \mid a^T x = b\},$$

$a \in \mathbf{R}^n$ ,  $a \neq 0$ , and  $b \in \mathbf{R}$ .

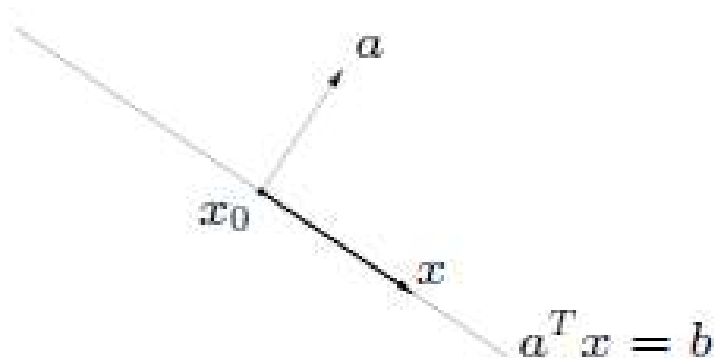
$a$  is the normal vector

Geometrical interpretation

$x_0$  is any point in the hyperplane

$$a^T x_0 = b$$

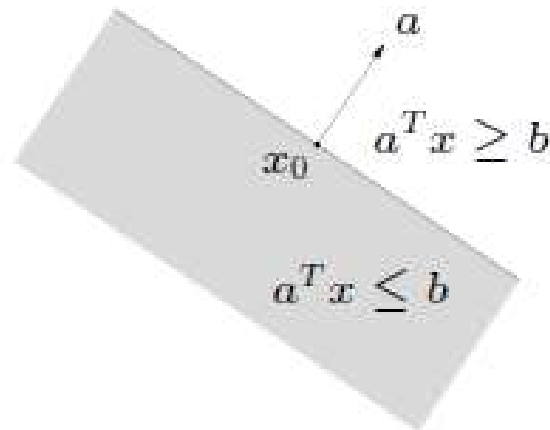
$$\{x \mid a^T (x - x_0) = 0\},$$



# halfspaces

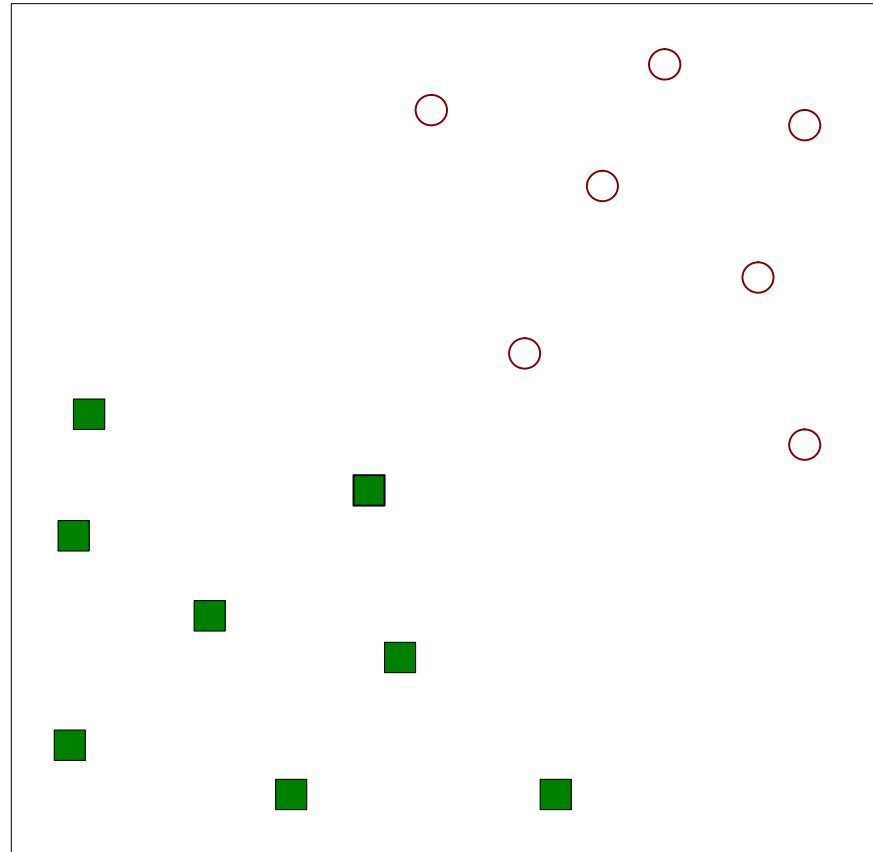
79

**halfspace:** set of the form  $\{x \mid a^T x \leq b\}$  ( $a \neq 0$ )



# Maximum Margin Hyperplanes

80

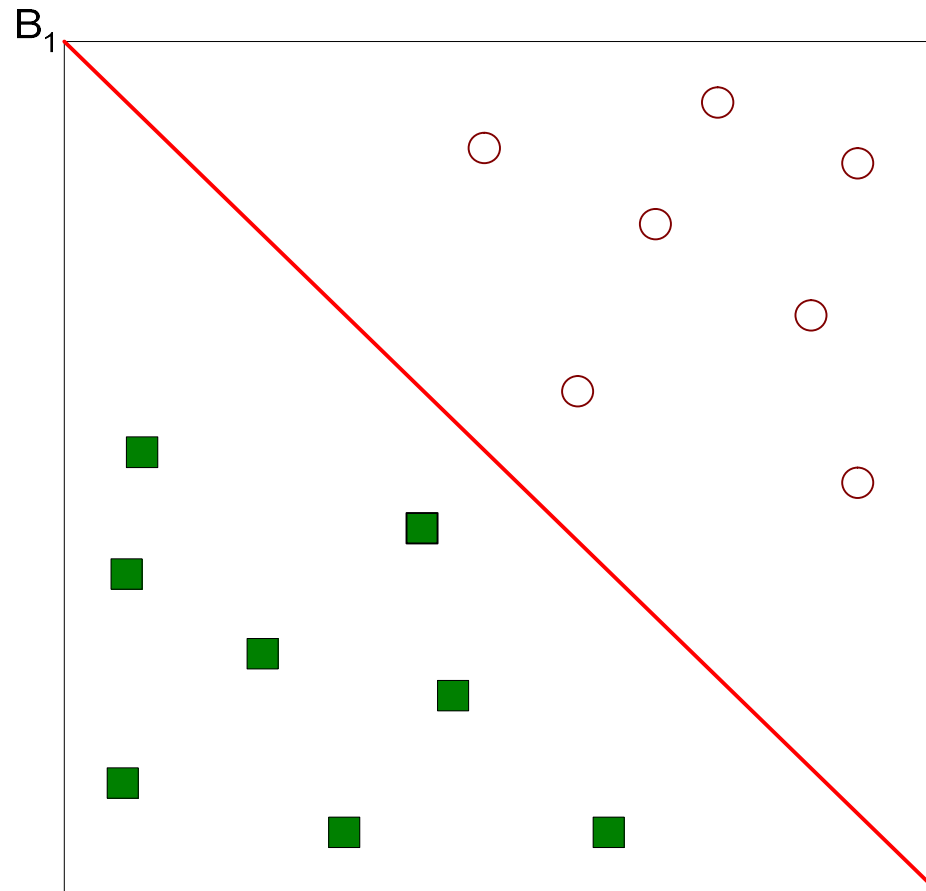


- Find a linear hyperplane (decision boundary) that will separate the data



# Maximum Margin Hyperplanes

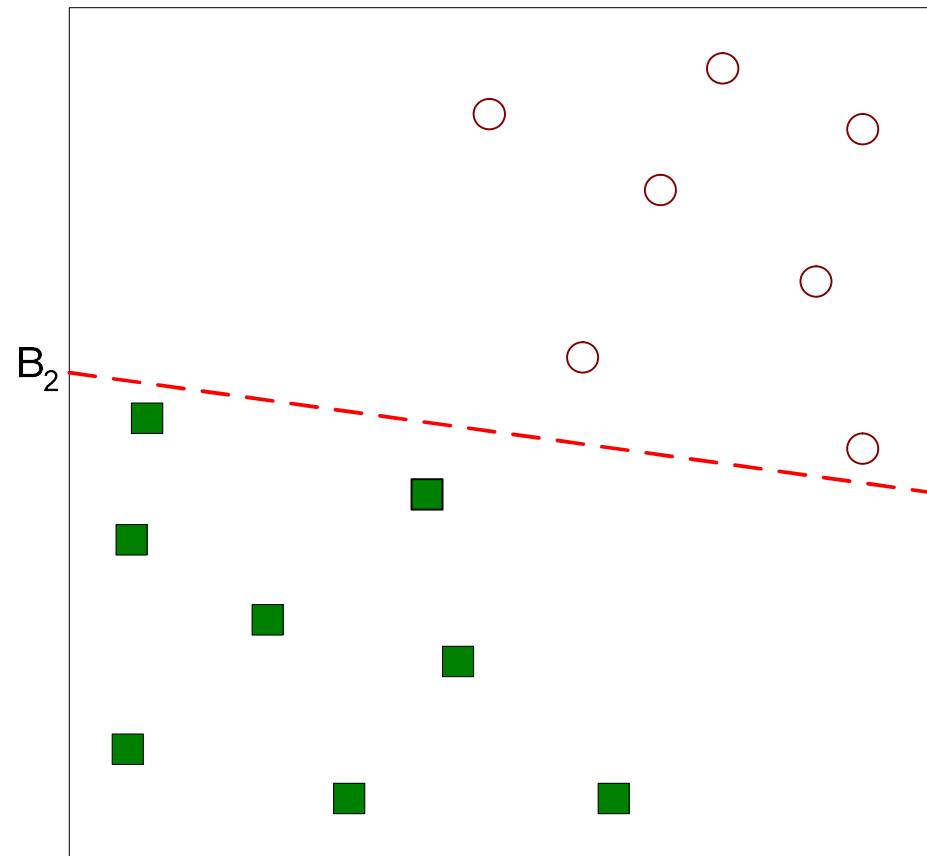
81



□ One Possible Solution

# Maximum Margin Hyperplanes

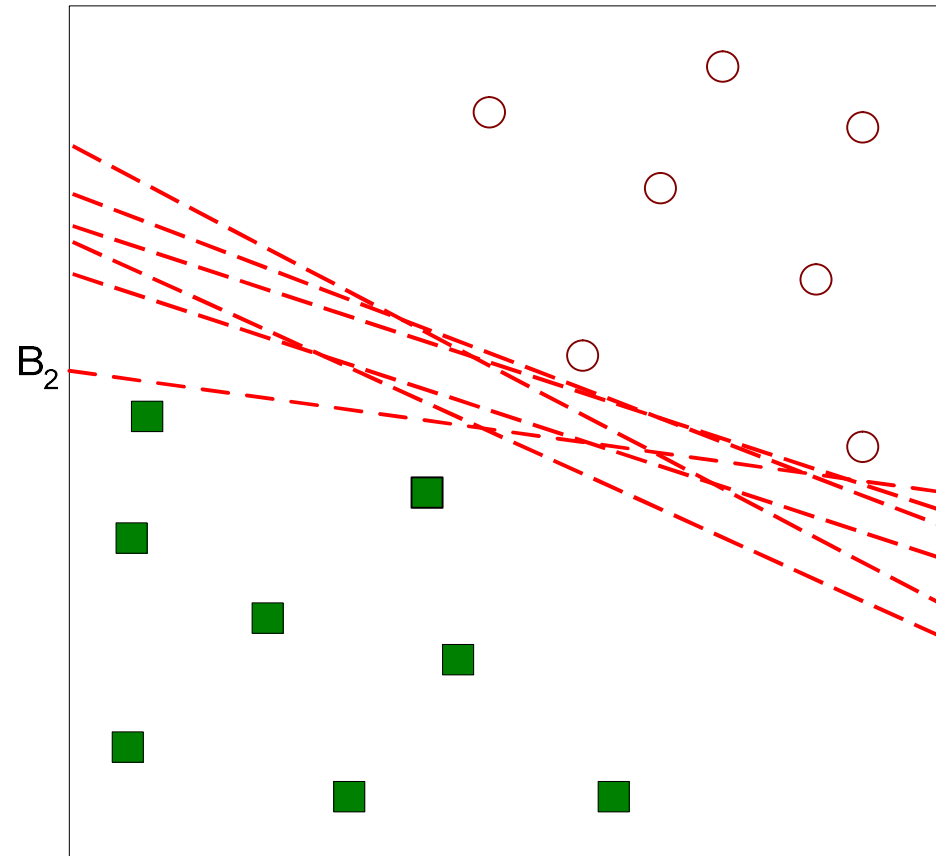
82



□ Another possible solution

# Maximum Margin Hyperplanes

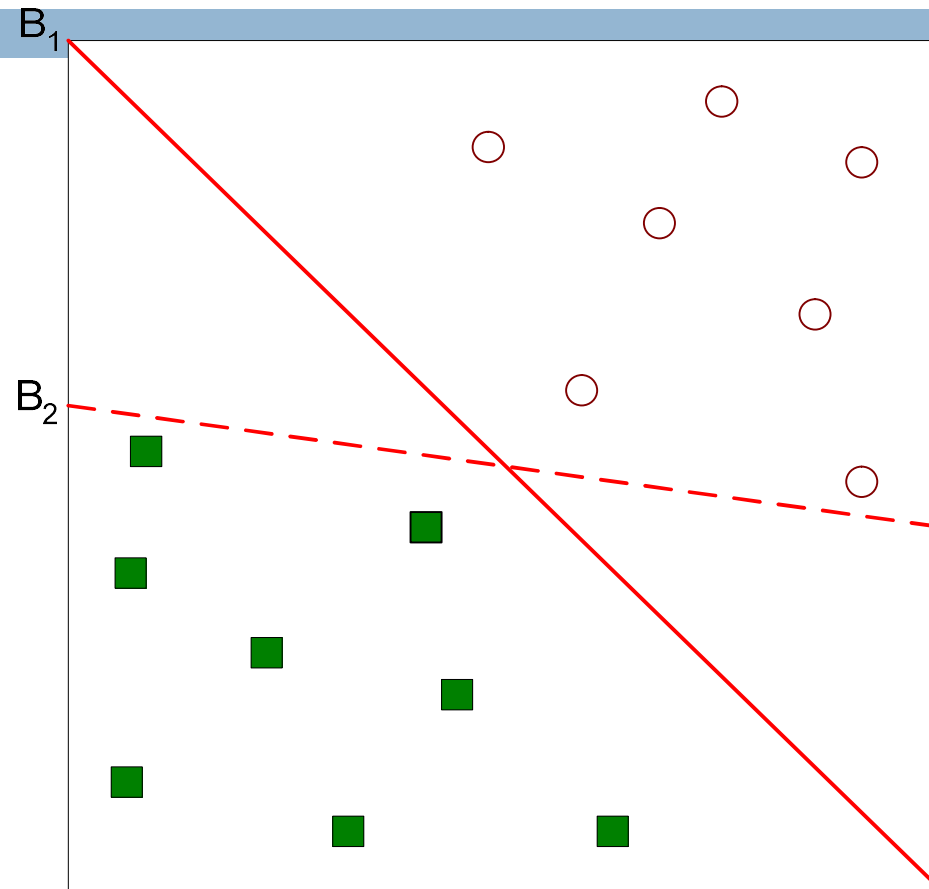
83



□ Other possible solutions

# Maximum Margin Hyperplanes

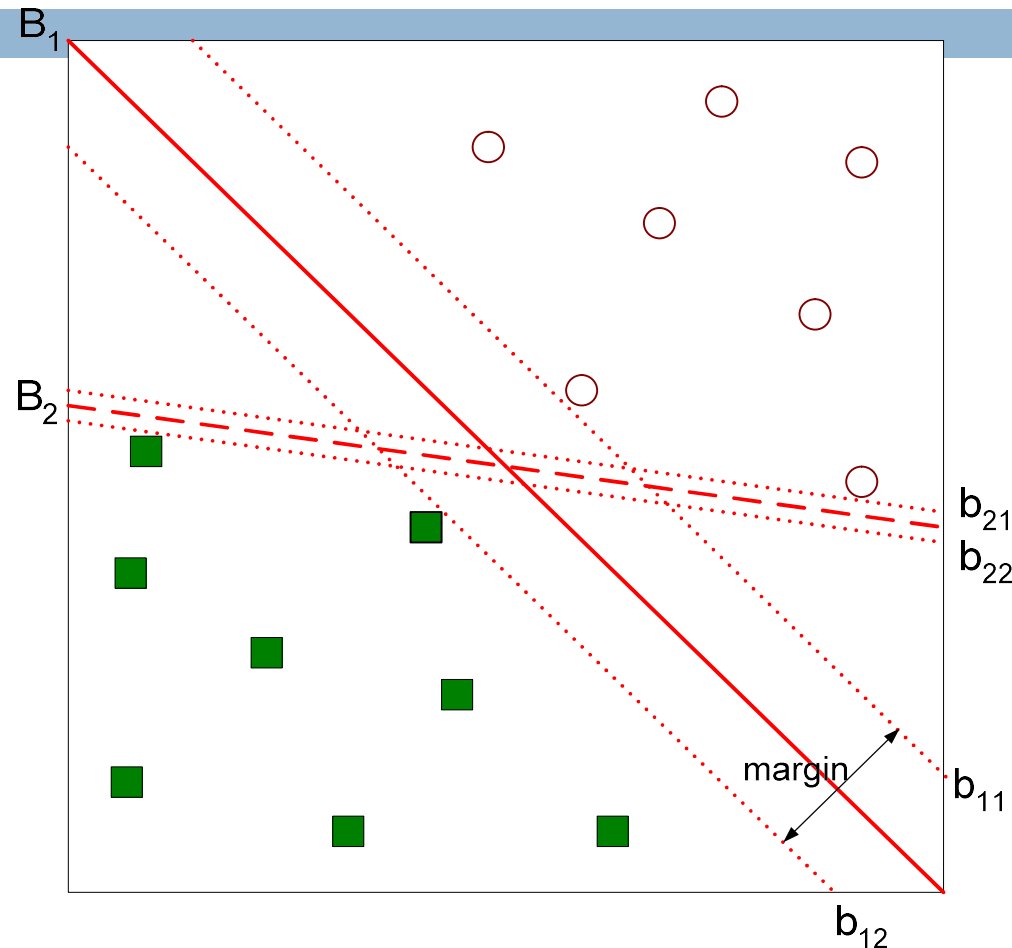
84



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

# Maximum Margin Hyperplanes

85

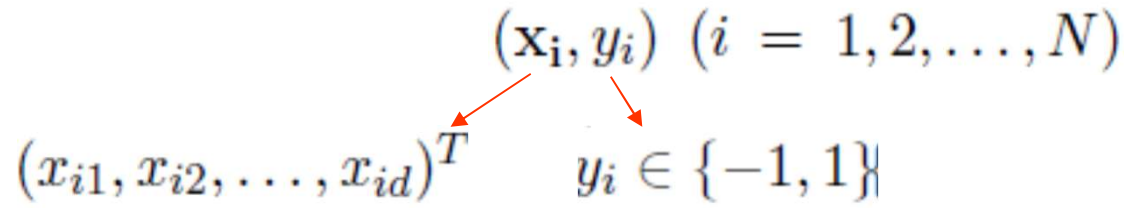


- Find hyperplane **maximizes** the margin  $\Rightarrow$  B1 is better than B2
- Generalization error

# Linear SVM: Separable Case

86

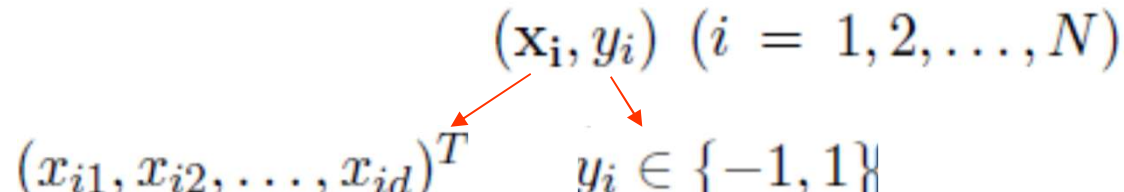
A linear SVM is a classifier that searches for a hyperplane with the largest margin

$$(x_i, y_i) \quad (i = 1, 2, \dots, N)$$

$$(x_{i1}, x_{i2}, \dots, x_{id})^T \quad y_i \in \{-1, 1\}$$

# Linear SVM: Separable Case

87

A linear SVM is a classifier that searches for a hyperplane with the largest margin

$$(x_i, y_i) \quad (i = 1, 2, \dots, N)$$

$$(x_{i1}, x_{i2}, \dots, x_{id})^T \quad y_i \in \{-1, 1\}$$

decision boundary of a linear classifier

$$\mathbf{w} \cdot \mathbf{x} + b = 0,$$

$\mathbf{w}$  and  $b$  are parameters of the model

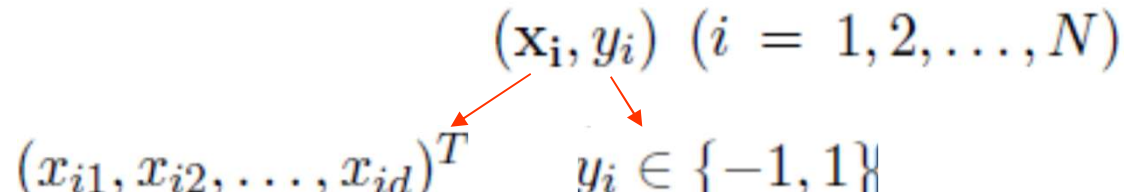
$$\mathbf{w} \cdot \mathbf{x}_s + b = k, \quad k > 0.$$

$$\mathbf{w} \cdot \mathbf{x}_c + b = k', \quad k' < 0.$$

# Linear SVM: Separable Case

88

A linear SVM is a classifier that searches for a hyperplane with the largest margin

$$(x_i, y_i) \quad (i = 1, 2, \dots, N)$$

$$(x_{i1}, x_{i2}, \dots, x_{id})^T \quad y_i \in \{-1, 1\}$$

decision boundary of a linear classifier

$$\mathbf{w} \cdot \mathbf{x} + b = 0,$$

$\mathbf{w}$  and  $b$  are parameters of the model

$$\mathbf{w} \cdot \mathbf{x}_s + b = k, \quad k > 0.$$

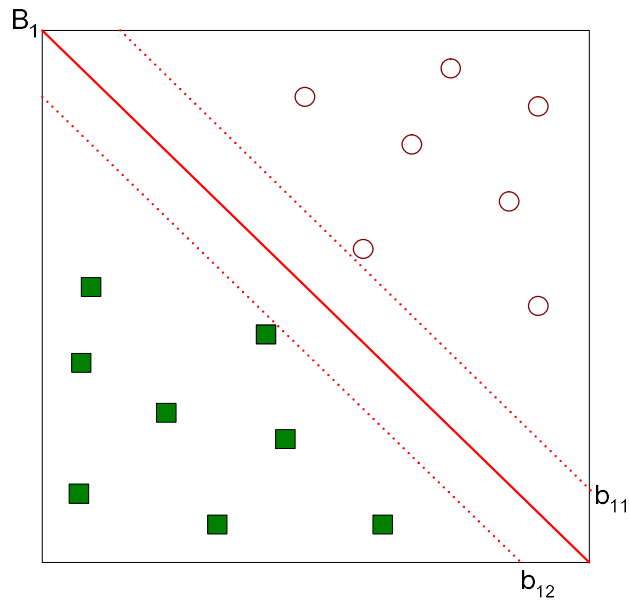
$$\mathbf{w} \cdot \mathbf{x}_c + b = k', \quad k' < 0.$$

$$y = \begin{cases} 1, & \text{if } \mathbf{w} \cdot \mathbf{z} + b > 0; \\ -1, & \text{if } \mathbf{w} \cdot \mathbf{z} + b < 0. \end{cases}$$



# Linear SVM: Separable Case

89

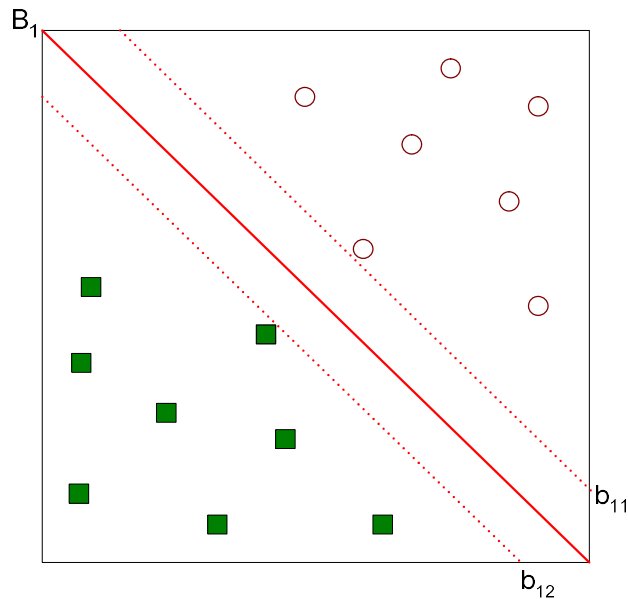


$$b_{i1} : \mathbf{w} \cdot \mathbf{x} + b = 1,$$

$$b_{i2} : \mathbf{w} \cdot \mathbf{x} + b = -1.$$

# Linear SVM: Separable Case

90



$$b_{i1} : \mathbf{w} \cdot \mathbf{x} + b = 1,$$

$$b_{i2} : \mathbf{w} \cdot \mathbf{x} + b = -1.$$

margin of the decision boundary  
is given by the distance between  
these two hyperplanes

$$d = \frac{2}{\|\mathbf{w}\|}$$

# SVM

91

$$\begin{aligned} \mathbf{w} \cdot \mathbf{x}_i + b &\geq 1 \text{ if } y_i = 1, \\ \mathbf{w} \cdot \mathbf{x}_i + b &\leq -1 \text{ if } y_i = -1. \end{aligned}$$

$$y_i(\mathbf{w} \cdot \mathbf{x}_i + b) \geq 1, \quad i = 1, 2, \dots, N.$$

# SVM

92

$$\begin{aligned} \mathbf{w} \cdot \mathbf{x}_i + b &\geq 1 \text{ if } y_i = 1, \\ \mathbf{w} \cdot \mathbf{x}_i + b &\leq -1 \text{ if } y_i = -1. \end{aligned} \quad y_i(\mathbf{w} \cdot \mathbf{x}_i + b) \geq 1, \quad i = 1, 2, \dots, N.$$

**Definition 5.1 (Linear SVM: Separable Case).** The learning task in SVM can be formalized as the following constrained optimization problem:

$$\begin{aligned} \min_{\mathbf{w}} \quad & \frac{\|\mathbf{w}\|^2}{2} \\ \text{subject to} \quad & y_i(\mathbf{w} \cdot \mathbf{x}_i + b) \geq 1, \quad i = 1, 2, \dots, N. \end{aligned}$$

This is a constrained optimization problem

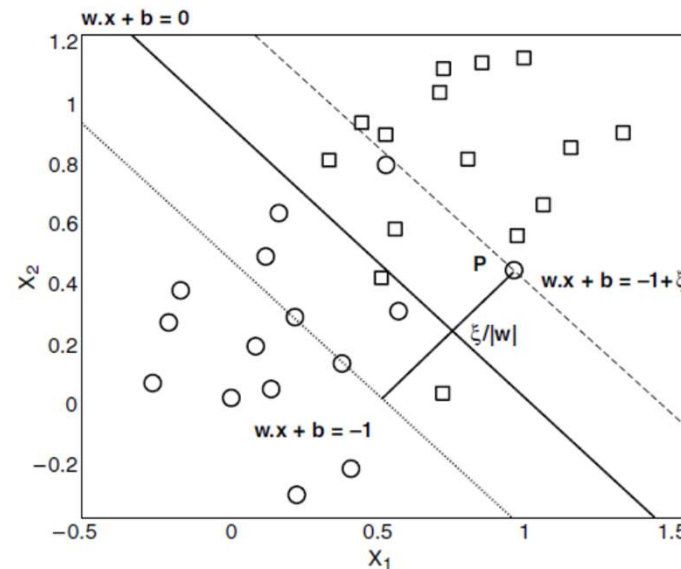
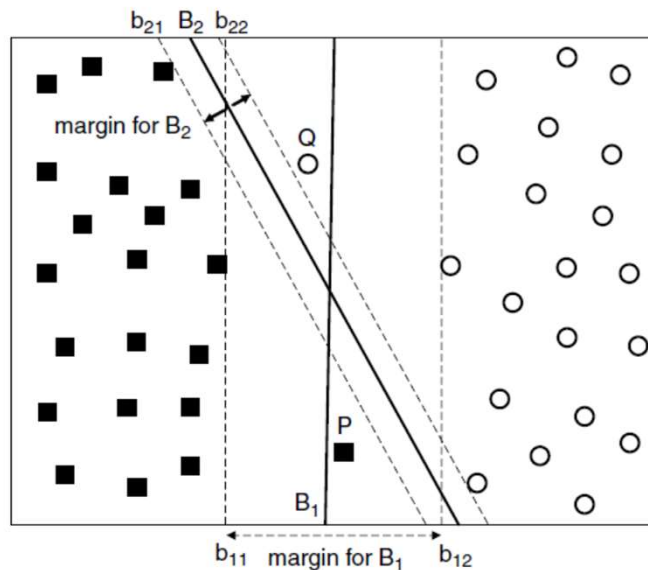
Numerical approaches to solve it (e.g., quadratic programming)

$$L_D = \sum_{i=1}^N \lambda_i - \frac{1}{2} \sum_{i,j} \lambda_i \lambda_j y_i y_j \mathbf{x}_i \cdot \mathbf{x}_j.$$

$$\left( \sum_{i=1}^N \lambda_i y_i \mathbf{x}_i \cdot \mathbf{x} \right) + b = 0.$$

# Linear SVM: Nonseparable Case

94

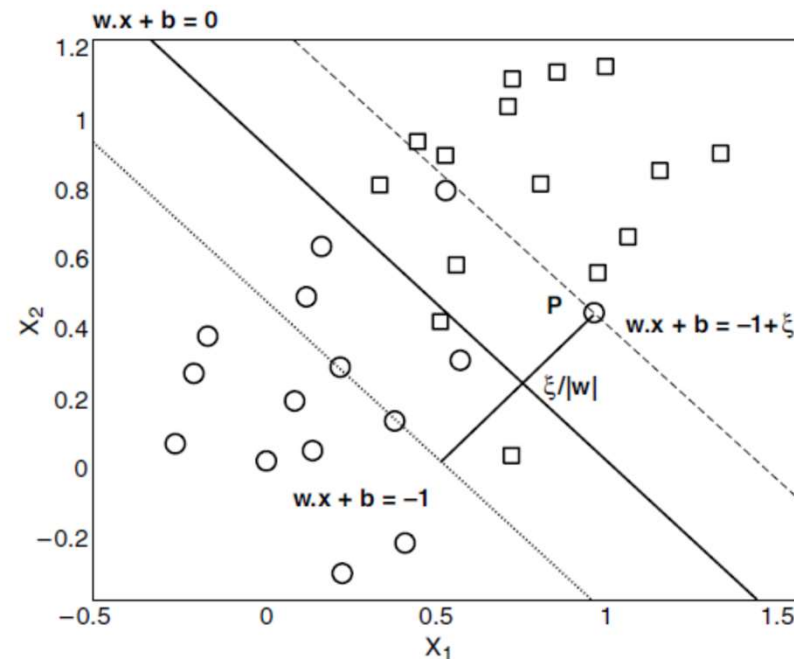


while  $B_2$  classifies them correctly, this does not mean that  $B_2$  is a better decision boundary than  $B_1$

# Linear SVM: Nonseparable Case

95

- ✓ learn a decision boundary that is tolerable to small training errors
- ✓ construct a linear decision boundary even in situations where the classes are not linearly separable
- ✓ inequality constraints must therefore be relaxed to accommodate the nonlinearly separable data

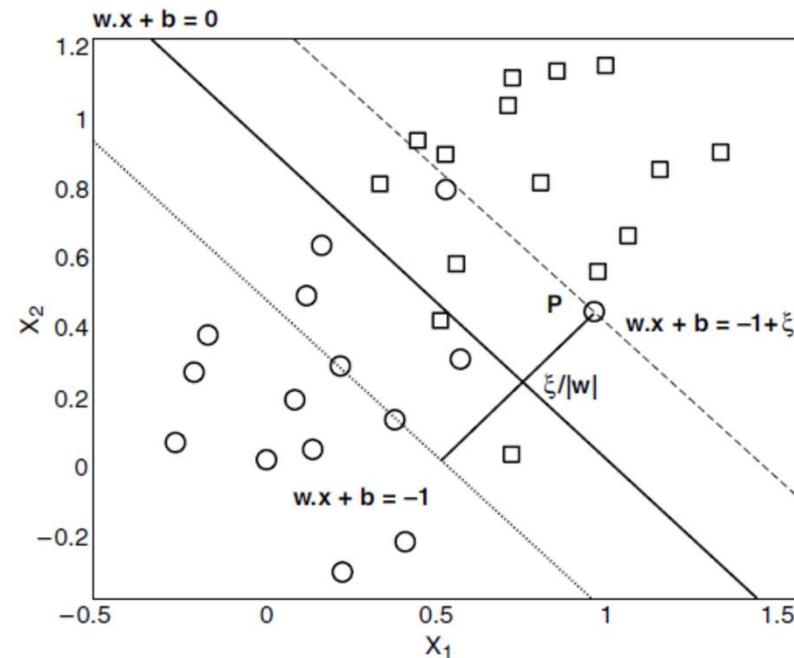


# Linear SVM: Nonseparable Case

96

- ✓ learn a decision boundary that is tolerable to small training errors
- ✓ construct a linear decision boundary even in situations where the classes are not linearly separable
- ✓ inequality constraints must therefore be relaxed to accommodate the nonlinearly separable data

$$w \cdot x_i + b \geq 1 - \xi_i \text{ if } y_i = 1,$$
$$w \cdot x_i + b \leq -1 + \xi_i \text{ if } y_i = -1,$$



$\xi$  provides an estimate of the error of the decision boundary



# Linear SVM: Nonseparable Case

97

$$\min_w \frac{\|w\|^2}{2}$$

$$w \cdot x_i + b \geq 1 - \xi_i \text{ if } y_i = 1,$$

$$w \cdot x_i + b \leq -1 + \xi_i \text{ if } y_i = -1,$$



problem

# Linear SVM: Nonseparable Case

98

$$\min_{\mathbf{w}} \frac{\|\mathbf{w}\|^2}{2}$$

$$\mathbf{w} \cdot \mathbf{x}_i + b \geq 1 - \xi_i \quad \text{if } y_i = 1,$$

$$\mathbf{w} \cdot \mathbf{x}_i + b \leq -1 + \xi_i \quad \text{if } y_i = -1,$$



problem

$$f(\mathbf{w}) = \frac{\|\mathbf{w}\|^2}{2} + C \left( \sum_{i=1}^N \xi_i \right)^k,$$

$C$  and  $k$  are user-specified parameters

# Nonlinear SVM

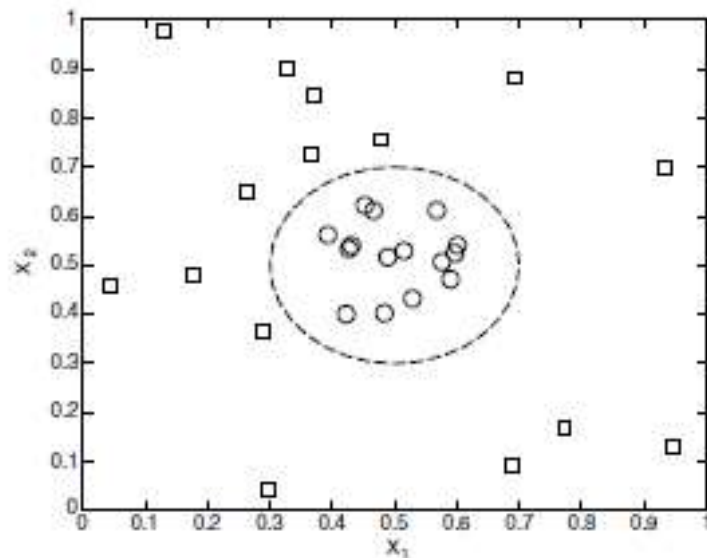
99

- ✓ applying SVM to data sets that have nonlinear decision boundaries
- ✓ transform the data  $\mathbf{x}$  into a new space  $\Phi(\mathbf{x})$  so that a linear decision boundary can be used to separate the instances in the transformed space

# Nonlinear SVM

100

- ✓ applying SVM to data sets that have nonlinear decision boundaries
- ✓ transform the data  $\mathbf{x}$  into a new space  $\Phi(\mathbf{x})$  so that a linear decision boundary can be used to separate the instances in the transformed space



(a) Decision boundary in the original two-dimensional space.

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

$$x_1^2 - x_1 + x_2^2 - x_2 = -0.46$$

$$\Phi : (x_1, x_2) \longrightarrow (x_1^2, x_2^2, \sqrt{2}x_1, \sqrt{2}x_2, 1).$$

$$\mathbf{W} = (w_0, w_1, \dots, w_4)$$

$$w_4 x_1^2 + w_3 x_2^2 + w_2 \sqrt{2}x_1 + w_1 \sqrt{2}x_2 + w_0 = 0.$$

# Nonlinear SVM

101

**Definition 5.2 (Nonlinear SVM).** The learning task for a nonlinear SVM can be formalized as the following optimization problem:

$$\begin{aligned} & \min_{\mathbf{w}} \frac{\|\mathbf{w}\|^2}{2} \\ \text{subject to} \quad & y_i(\mathbf{w} \cdot \Phi(\mathbf{x}_i) + b) \geq 1, \quad i = 1, 2, \dots, N. \end{aligned}$$

# Nonlinear SVM

102

**Definition 5.2 (Nonlinear SVM).** The learning task for a nonlinear SVM can be formalized as the following optimization problem:

$$\begin{aligned} \min_{\mathbf{w}} \quad & \frac{\|\mathbf{w}\|^2}{2} \\ \text{subject to} \quad & y_i(\mathbf{w} \cdot \Phi(\mathbf{x}_i) + b) \geq 1, \quad i = 1, 2, \dots, N. \end{aligned}$$

$$L_D = \sum_{i=1}^N \lambda_i - \frac{1}{2} \sum_{i,j} \lambda_i \lambda_j y_i y_j \mathbf{x}_i \cdot \mathbf{x}_j.$$

$$L_D = \sum_{i=1}^n \lambda_i - \frac{1}{2} \sum_{i,j} \lambda_i \lambda_j y_i y_j \Phi(\mathbf{x}_i) \cdot \Phi(\mathbf{x}_j)$$

# Nonlinear SVM

103

**Definition 5.2 (Nonlinear SVM).** The learning task for a nonlinear SVM can be formalized as the following optimization problem:

$$\begin{aligned} \min_{\mathbf{w}} \quad & \frac{\|\mathbf{w}\|^2}{2} \\ \text{subject to} \quad & y_i(\mathbf{w} \cdot \Phi(\mathbf{x}_i) + b) \geq 1, \quad i = 1, 2, \dots, N. \end{aligned}$$

$$\begin{aligned} \Phi(\mathbf{u}) \cdot \Phi(\mathbf{v}) &= (u_1^2, u_2^2, \sqrt{2}u_1, \sqrt{2}u_2, 1) \cdot (v_1^2, v_2^2, \sqrt{2}v_1, \sqrt{2}v_2, 1) \\ &= u_1^2 v_1^2 + u_2^2 v_2^2 + 2u_1 v_1 + 2u_2 v_2 + 1 \\ &= (\mathbf{u} \cdot \mathbf{v} + 1)^2. \end{aligned}$$

$$K(\mathbf{u}, \mathbf{v}) = \Phi(\mathbf{u}) \cdot \Phi(\mathbf{v}) = (\mathbf{u} \cdot \mathbf{v} + 1)^2$$

# Nonlinear SVM

104

The main requirement for the kernel function used in nonlinear SVM is that there must exist a corresponding transformation such that the kernel function computed for a pair of vectors is equivalent to the dot product between the vectors in the transformed space.

$$K(\mathbf{x}, \mathbf{y}) = (\mathbf{x} \cdot \mathbf{y} + 1)^p$$

$$K(\mathbf{x}, \mathbf{y}) = e^{-\|\mathbf{x} - \mathbf{y}\|^2 / (2\sigma^2)}$$

$$K(\mathbf{x}, \mathbf{y}) = \tanh(k\mathbf{x} \cdot \mathbf{y} - \delta)$$



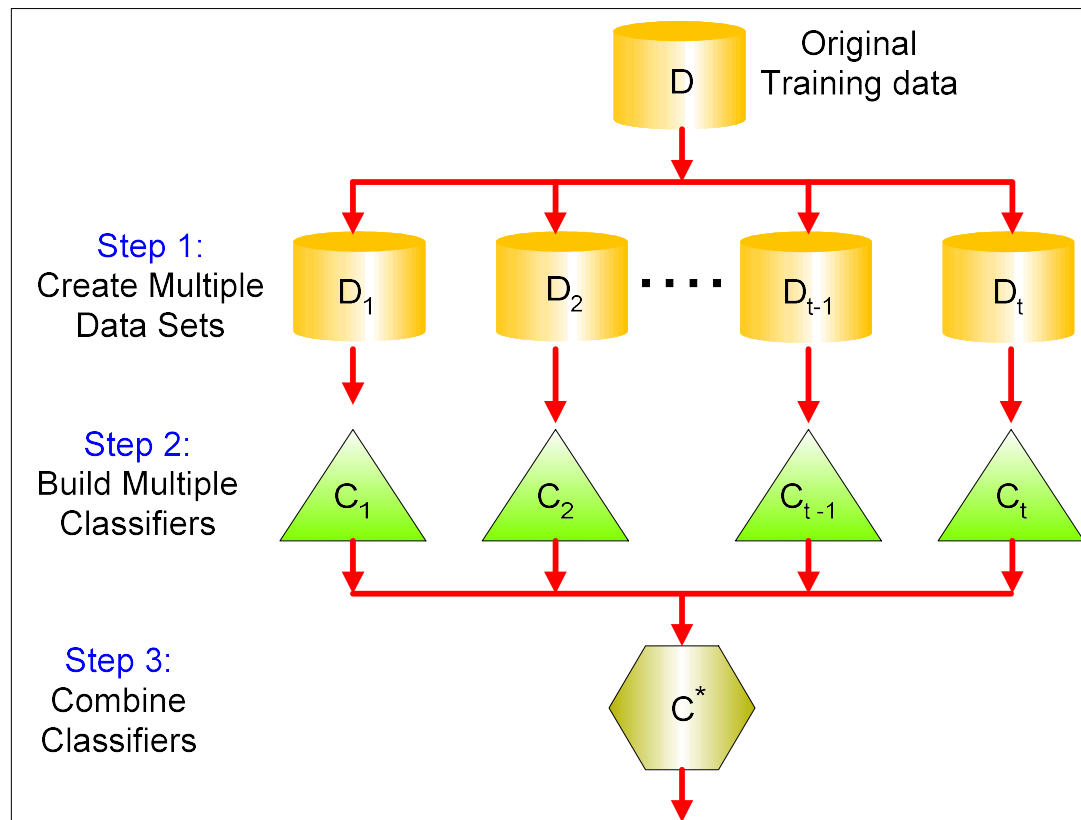
105

# Ensemble Methods

# Ensemble Methods

106

- ✓ improving classification accuracy by aggregating the predictions of multiple classifiers
- ✓ Construct a set of base classifiers from the training data
- ✓ Predict class label of previously unseen records by aggregating predictions made by multiple classifiers



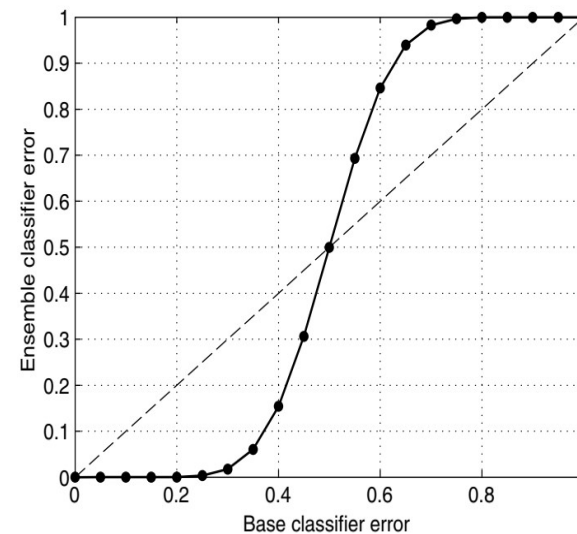
# Ensemble Methods

107

- Suppose there are 25 base classifiers
  - ▣ Each classifier has error rate,  $\varepsilon = 0.35$
  - ▣ Assume classifiers are independent
  - ▣ Probability that the ensemble classifier makes a wrong prediction:

$$\sum_{i=13}^{25} \binom{25}{i} \varepsilon^i (1 - \varepsilon)^{25-i} = 0.06$$

base classifiers should be independent of each other  
base classifiers should do better than a classifier that performs random guessing.



How to generate an ensemble of classifiers?

- Bagging: bootstrap aggregating
- Boosting

# Bagging

109

- Sampling with replacement
- Build classifier on each bootstrap sample

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

---

## Algorithm 5.6 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:  $C^*(x) = \operatorname{argmax}_y \sum_i \delta(C_i(x) = y)$ .  
     $\{\delta(\cdot) = 1 \text{ if its argument is true and } 0 \text{ otherwise}\}.$
-

# Boosting

110

- An iterative procedure to adaptively change distribution of training data by focusing more on previously misclassified records
  - ▣ Initially, all  $N$  records are assigned equal weights
  - ▣ Unlike bagging, weights may change at the end of boosting round
- 1. how the weights of the training examples are updated at the end of each boosting round
- 2. how the predictions made by each classifier are combined.

# Boosting: AdaBoost

111

Let  $\{(\mathbf{x}_j, y_j) \mid j = 1, 2, \dots, N\}$  denote a set of  $N$  training examples.  
Unlike bagging, importance of a base classifier  $C_i$  depends on its error rate

$$\epsilon_i = \frac{1}{N} \left[ \sum_{j=1}^N w_j I(C_i(\mathbf{x}_j) \neq y_j) \right], \quad \alpha_i = \frac{1}{2} \ln \left( \frac{1 - \epsilon_i}{\epsilon_i} \right)$$

weight assigned to example  $(\mathbf{x}_i, y_i)$   
during the  $j$ th boosting round

$$w_i^{(j+1)} = \frac{w_i^{(j)}}{Z_j} \times \begin{cases} \exp^{-\alpha_j} & \text{if } C_j(\mathbf{x}_i) = y_i \\ \exp^{\alpha_j} & \text{if } C_j(\mathbf{x}_i) \neq y_i \end{cases}$$

Normalization factor

intermediate rounds produce an error rate higher than 50%, the weights  $w_i = 1/N$

# AdaBoost

112

---

**Algorithm 5.7** 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:    $\epsilon_i = \frac{1}{N} [\sum_j w_j \delta(C_i(x_j) \neq y_j)]$    {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 5.69.
14: end for
15:  $C^*(\mathbf{x}) = \operatorname{argmax}_y \sum_{j=1}^T \alpha_j \delta(C_j(\mathbf{x}) = y)$ .
```

---



# Example

113

$x$	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

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

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

# Random Forest

114

- ✓ ensemble methods specifically designed for decision tree classifiers
- ✓ multiple decision trees where each tree is generated based on the values of an independent set of random vectors

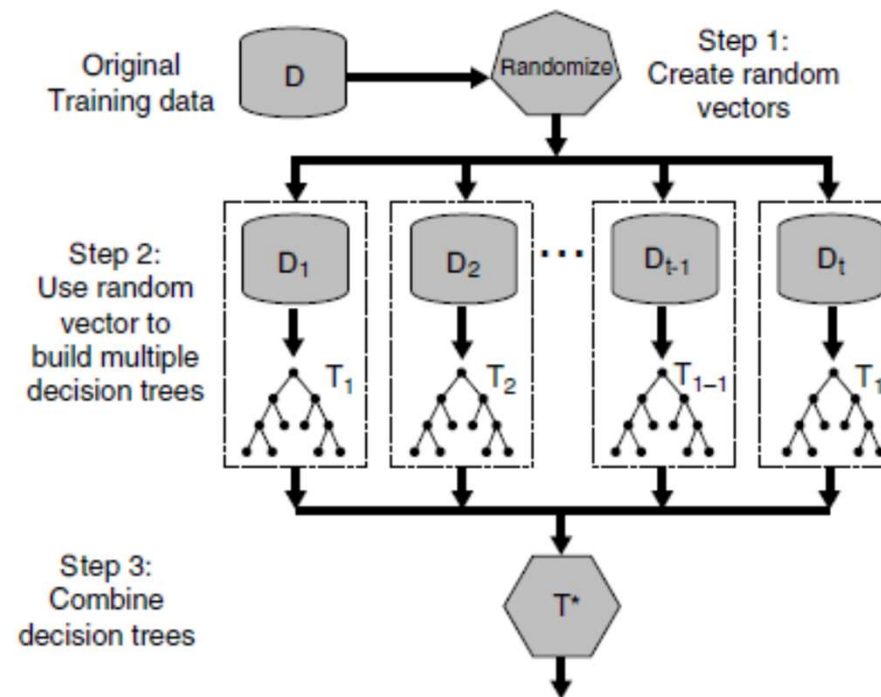


Figure 5.40. Random forests.

# Random Forest

115

- ✓ Bagging using decision trees is a special case of random forests
- ✓ randomly select  $F$  input features to split at each node of the decision tree
- ✓ majority voting scheme
- ✓ To increase randomness, bagging can also be used to generate bootstrap samples for Forest-R1

116

## Metrics for class imbalance problem

# Imbalance

117

- ✓ Data sets with imbalanced class distributions
- ✓ in credit card fraud detection, fraudulent transactions are outnumbered by legitimate transactions
- ✓ accuracy measure, used extensively for classifiers, may not be well suited for evaluating models derived from imbalanced data sets

example : 1% of the credit card transactions fraudulent,  
a model that predicts every transaction as legitimate  
accuracy 99%

it fails to detect any of the fraudulent activities.

binary classification, the rare class is often denoted as the positive class against negative class

		Predicted Class	
		+	−
Actual Class	+	$f_{++}$ (TP)	$f_{+-}$ (FN)
	−	$f_{-+}$ (FP)	$f_{--}$ (TN)

confusion matrix

# Imbalance

118

Precision : fraction of records that actually turns out to be positive in the group the classifier has declared as a positive class

$$\text{Precision, } p = \frac{TP}{TP + FP}$$

Recall measures the fraction of positive examples correctly predicted by the classifier

$$\text{Recall, } r = \frac{TP}{TP + FN}$$

maximizes both precision and recall

# Imbalance

119

Precision and recall can be summarized into another metric known as the F1 measure

$$F_1 = \frac{2}{\frac{1}{r} + \frac{1}{p}}.$$

tends to be closer to the smaller of the two numbers  
a high value of  $F_1$ -measure ensures that both precision and recall are reasonably high

$$\text{Weighted accuracy} = \frac{w_1 TP + w_4 TN}{w_1 TP + w_2 FP + w_3 FN + w_4 TN}.$$