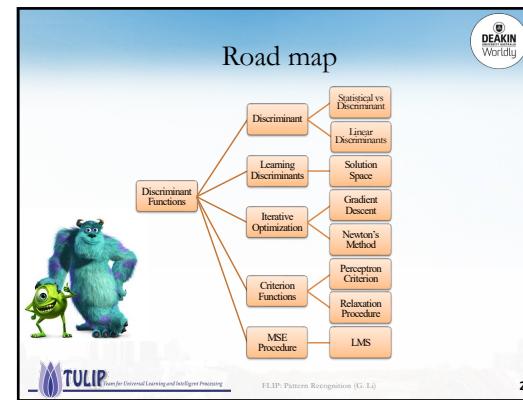
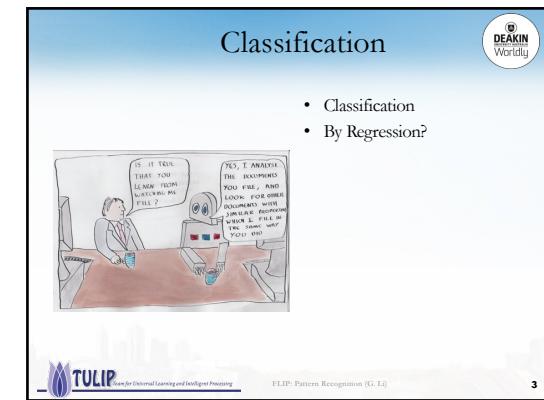




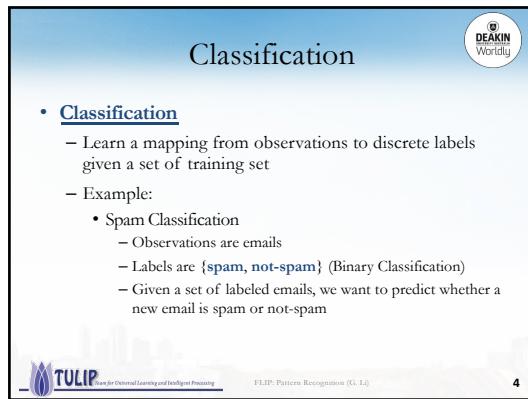
1



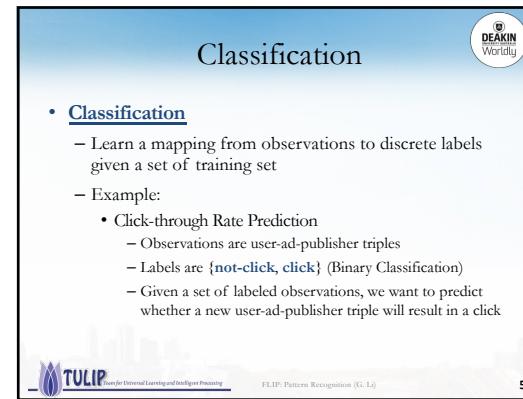
2



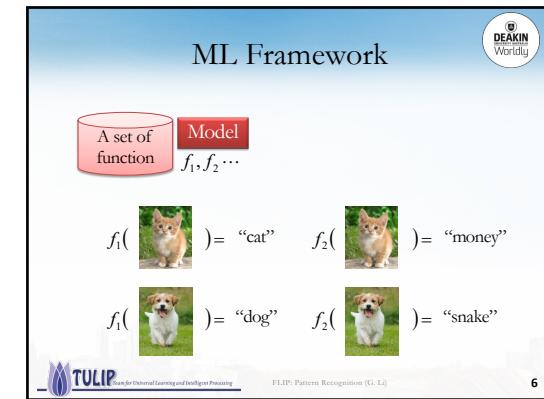
3



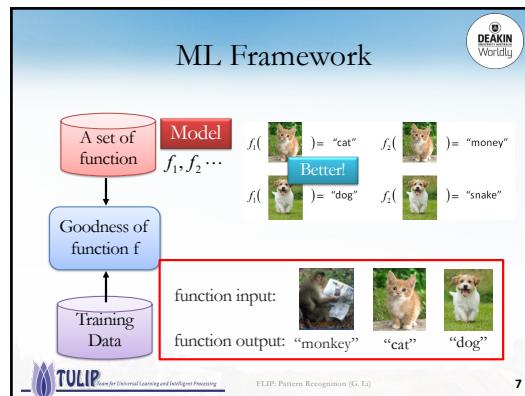
4



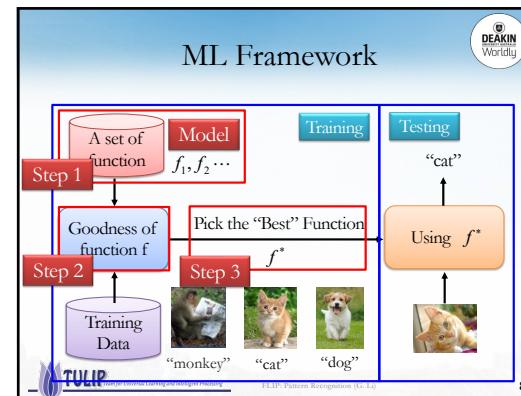
5



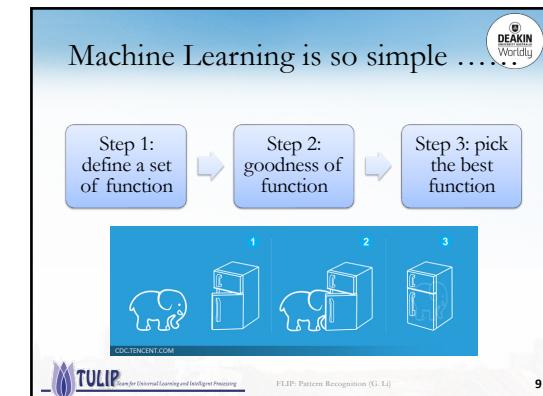
6



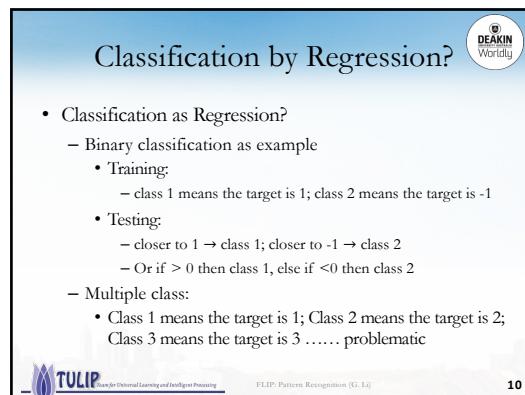
7



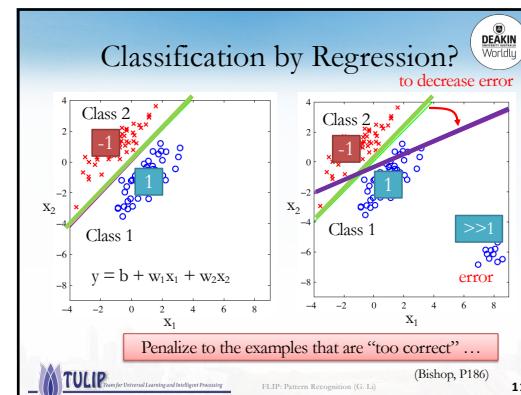
8



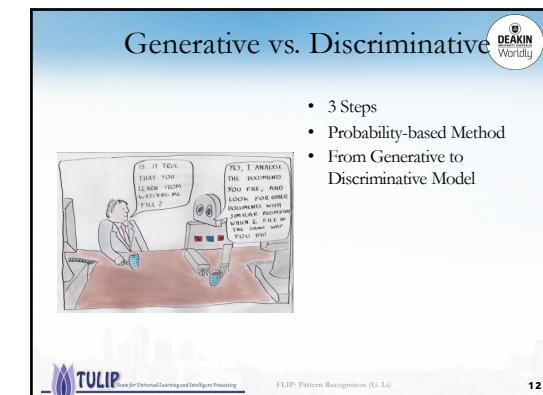
9



10



11



12

## Three Steps

- Function Set (Generative Model):

$$f(x) = \begin{cases} \text{class 1} & \text{if } P(C_1|x) > 0.5 \\ \text{class 2} & \text{else if } P(C_1|x) \leq 0.5 \end{cases}$$

$P(C_i|x) = \frac{P(x|C_i)P(C_i)}{P(x|C_1)P(C_1) + P(x|C_2)P(C_2)}$

- Goodness of a function:

– The parameters of the distribution that maximizing the likelihood (the probability of generating data)

- Find the best function:
  - MLE
  - Bayesian Estimation



FLIP: Pattern Recognition (G. Li)

13

## Generative v.s. Discriminative

- Example

Training Data  
Class 1: 1 X 4  
Class 2: 1 X 4  
X 4  
X 4

Testing Data  
1 Class 1?  
1 Class 2?

How about Naïve Bayes?  
 $P(x|C_i) = P(x_1|C_i)P(x_2|C_i)$



FLIP: Pattern Recognition (G. Li)

14

## Generative v.s. Discriminative

- Example

Training Data  
Class 1: 1 X 4  
Class 2: 1 X 4  
X 4  
X 4

$$\begin{aligned} P(C_1) &= \frac{1}{13} & P(x_1 = 1|C_1) &= 1 & P(x_2 = 1|C_1) &= 1 \\ P(C_2) &= \frac{12}{13} & P(x_1 = 1|C_2) &= \frac{1}{3} & P(x_2 = 1|C_2) &= \frac{1}{3} \end{aligned}$$

FLIP: Pattern Recognition (G. Li)

15

## Generative v.s. Discriminative

- Example

Training Data  
Class 1: 1 X 4  
Class 2: 1 X 4  
X 4  
X 4

Testing Data  
1  
1

$P(C_1|x) = \frac{P(x|C_1)P(C_1)}{P(x|C_1)P(C_1) + P(x|C_2)P(C_2)}$

$$= \frac{1}{1 + \frac{P(x|C_2)P(C_2)}{P(x|C_1)P(C_1)}} = \frac{1}{1 + \exp(-z)} = \sigma(z)$$

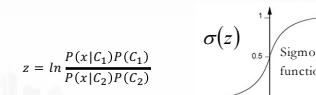
$$z = \ln \frac{P(x|C_1)P(C_1)}{P(x|C_2)P(C_2)}$$


FLIP: Pattern Recognition (G. Li)

16

## Posterior Probability

$$\begin{aligned} P(C_1|x) &= \frac{P(x|C_1)P(C_1)}{P(x|C_1)P(C_1) + P(x|C_2)P(C_2)} \\ &= \frac{1}{1 + \frac{P(x|C_2)P(C_2)}{P(x|C_1)P(C_1)}} = \frac{1}{1 + \exp(-z)} = \sigma(z) \end{aligned}$$



FLIP: Pattern Recognition (G. Li)

17

## Posterior Probability

$$P(C_1|x) = \sigma(z) \quad \text{sigmoid} \quad z = \ln \frac{P(x|C_1)P(C_1)}{P(x|C_2)P(C_2)}$$

$$z = \ln \frac{P(x|C_1)}{P(x|C_2)} + \ln \frac{P(C_1)}{P(C_2)} \rightarrow \frac{\frac{N_1}{N_1+N_2}}{\frac{N_2}{N_1+N_2}} = \frac{N_1}{N_2}$$

$$P(x|C_1) = \frac{1}{(2\pi)^{D/2} |\Sigma^1|^{1/2}} \exp \left\{ -\frac{1}{2}(x - \mu^1)^T (\Sigma^1)^{-1} (x - \mu^1) \right\}$$

$$P(x|C_2) = \frac{1}{(2\pi)^{D/2} |\Sigma^2|^{1/2}} \exp \left\{ -\frac{1}{2}(x - \mu^2)^T (\Sigma^2)^{-1} (x - \mu^2) \right\}$$

FLIP: Pattern Recognition (G. Li)

18

16

17

18

### Generative vs Discriminative Model

$$\begin{aligned} z &= \ln \frac{P(x|C_1)}{P(x|C_2)} + \ln \frac{P(C_1)}{P(C_2)} \\ &\quad \cancel{\ln \frac{1}{(2\pi)^{N/2} |\Sigma^1|^{1/2}} \exp \left\{ -\frac{1}{2} (x - \mu^1)^T (\Sigma^1)^{-1} (x - \mu^1) \right\}} \\ &\quad \cancel{\ln \frac{1}{(2\pi)^{N/2} |\Sigma^2|^{1/2}} \exp \left\{ -\frac{1}{2} (x - \mu^2)^T (\Sigma^2)^{-1} (x - \mu^2) \right\}} \\ &= \ln \frac{|\Sigma^2|^{1/2}}{|\Sigma^1|^{1/2}} - \frac{1}{2} [(x - \mu^1)^T (\Sigma^1)^{-1} (x - \mu^1) - (x - \mu^2)^T (\Sigma^2)^{-1} (x - \mu^2)] \end{aligned}$$



FLIP: Pattern Recognition (G. Li)

19

### Generative vs Discriminative Model

$$\begin{aligned} z &= \ln \frac{P(x|C_1)}{P(x|C_2)} + \ln \frac{P(C_1)}{P(C_2)} \\ &\quad (x - \mu^1)^T (\Sigma^1)^{-1} (x - \mu^1) \\ &= x^T (\Sigma^1)^{-1} x - x^T (\Sigma^1)^{-1} \mu^1 - (\mu^1)^T (\Sigma^1)^{-1} x + (\mu^1)^T (\Sigma^1)^{-1} \mu^1 \\ &= x^T (\Sigma^1)^{-1} x - 2(\mu^1)^T (\Sigma^1)^{-1} x + (\mu^1)^T (\Sigma^1)^{-1} \mu^1 \\ &\quad (x - \mu^2)^T (\Sigma^2)^{-1} (x - \mu^2) \\ &= x^T (\Sigma^2)^{-1} x - 2(\mu^2)^T (\Sigma^2)^{-1} x + (\mu^2)^T (\Sigma^2)^{-1} \mu^2 \end{aligned}$$

$$\begin{aligned} z &= \ln \frac{|\Sigma^2|^{1/2}}{|\Sigma^1|^{1/2}} - \frac{1}{2} x^T (\Sigma^1)^{-1} x + (\mu^1)^T (\Sigma^1)^{-1} x - \frac{1}{2} (\mu^1)^T (\Sigma^1)^{-1} \mu^1 \\ &\quad + \frac{1}{2} x^T (\Sigma^2)^{-1} x - (\mu^2)^T (\Sigma^2)^{-1} x + \frac{1}{2} (\mu^2)^T (\Sigma^2)^{-1} \mu^2 + \ln \frac{N_1}{N_2} \end{aligned}$$



FLIP: Pattern Recognition (G. Li)

20

### Generative vs Discriminative Model

$$\begin{aligned} z &= \ln \frac{|\Sigma^2|^{1/2}}{|\Sigma^1|^{1/2}} - \frac{1}{2} x^T (\Sigma^1)^{-1} x + (\mu^1)^T (\Sigma^1)^{-1} x - \frac{1}{2} (\mu^1)^T (\Sigma^1)^{-1} \mu^1 \\ &\quad + \frac{1}{2} x^T (\Sigma^2)^{-1} x - (\mu^2)^T (\Sigma^2)^{-1} x + \frac{1}{2} (\mu^2)^T (\Sigma^2)^{-1} \mu^2 + \ln \frac{N_1}{N_2} \\ \Sigma_1 &= \Sigma_2 = \Sigma \\ z &= (\mu^1 - \mu^2)^T \Sigma^{-1} x - \frac{1}{2} (\mu^1)^T \Sigma^{-1} \mu^1 + \frac{1}{2} (\mu^2)^T \Sigma^{-1} \mu^2 + \ln \frac{N_1}{N_2} \end{aligned}$$

$$\mathbf{w}^T \quad b$$

$$P(C_1|x) = \sigma(w \cdot x + b)$$

How about directly find  $w$  and  $b$ ?In generative model, we estimate  $N_1, N_2, \mu^1, \mu^2, \Sigma$ Then we have  $w$  and  $b$ 

FLIP: Pattern Recognition (G. Li)

21

### Generative vs. Discriminative

$$P(C_1|x) = \sigma(w \cdot x + b)$$

directly find  $w$  and  $b$

Will we obtain the same set of  $w$  and  $b$ ?

Find  $\mu^1, \mu^2, \Sigma^{-1}$ 

$$w^T = (\mu^1 - \mu^2)^T \Sigma^{-1}$$

$$\begin{aligned} b &= -\frac{1}{2} (\mu^1)^T (\Sigma^1)^{-1} \mu^1 \\ &\quad + \frac{1}{2} (\mu^2)^T (\Sigma^2)^{-1} \mu^2 + \ln \frac{N_1}{N_2} \end{aligned}$$

The same model (function set), but different function may be selected by the same training data.



FLIP: Pattern Recognition (G. Li)

22

### Generative vs. Discriminative

- Usually, people believe discriminative model is better
- Benefit of generative model
  - With the assumption of probability distribution
    - less training data is needed
    - more robust to the noise
  - Priors and class-dependent probabilities can be estimated from different sources.



FLIP: Pattern Recognition (G. Li)

23

### Discriminant Approach

- Discriminant Approach
- Discriminant Function
- Linear Discriminant Function



FLIP: Pattern Recognition (G. Li)

24

## Statistical vs Discriminant Approach

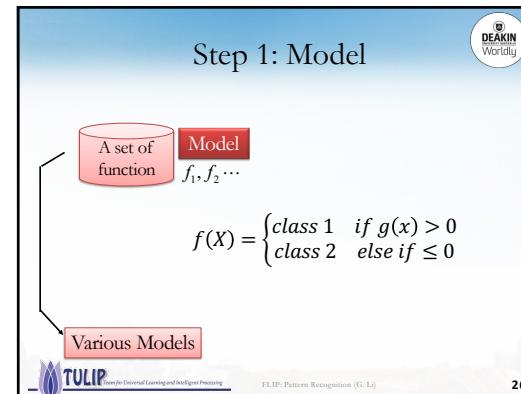
- Statistical** approaches find the decision boundary by first estimating the probability distribution of the patterns belonging to each class.
- In the **discriminant** approach, the decision boundary is constructed explicitly without assuming a probability distribution.



FLIP: Pattern Recognition (G. Li)

25

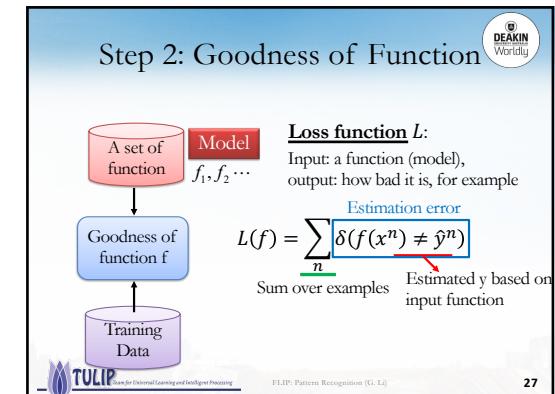
## Step 1: Model



FLIP: Pattern Recognition (G. Li)

26

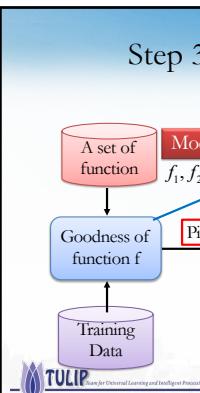
## Step 2: Goodness of Function



FLIP: Pattern Recognition (G. Li)

27

## Step 3: Best Function



FLIP: Pattern Recognition (G. Li)

28

## Discriminant Approach

- Specify **parametric form** of the decision boundary (e.g., linear or quadratic).
- Find the **best** decision boundary of the specified form using a set of training examples  $x_k$ .
  - This is performed by **minimizing** a criterion function (e.g., “training error” or “sample risk”):

$$J(w) = \frac{1}{n} \sum_{k=1}^n [z_k - g(x_k, w)]^2$$

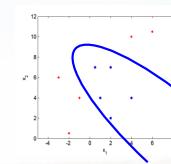
correct class                            predicted class

FLIP: Pattern Recognition (G. Li)

29

## Discriminant Functions: two-categories case

- Decide  $w_1$  if  $g(\mathbf{x}) > 0$  and  $w_2$  if  $g(\mathbf{x}) < 0$ 
  - If  $g(\mathbf{x})=0$ , then  $\mathbf{x}$  lies on the decision boundary and can be assigned to either class.
- Classification is viewed as finding a decision boundary that separates the data belonging to different classes.



FLIP: Pattern Recognition (G. Li)

30

## Linear Discriminant Functions: two-categories case

- A linear discriminant function has the following form:
$$g(\mathbf{x}) = \mathbf{w}'\mathbf{x} + w_0 = \sum_{i=1}^d w_i x_i + w_0$$
- The decision boundary, given by  $g(\mathbf{x})=0$ , is a **hyperplane** where the orientation of the hyperplane is determined by  $\mathbf{w}$  and its location by  $w_0$ .
  - $\mathbf{w}$  is the normal to the hyperplane
  - If  $w_0=0$ , the hyperplane passes through the origin

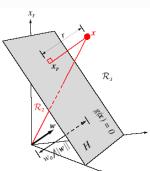


FLIP: Pattern Recognition (G. Li)

31

## Geometric Interpretation of $g(\mathbf{x})$

- $g(\mathbf{x})$  provides **an algebraic measure of the distance** of  $\mathbf{x}$  from the hyperplane.



$\mathbf{x}$  can be expressed as follows:

$$\mathbf{x} = \mathbf{x}_p + r \frac{\mathbf{w}}{\|\mathbf{w}\|}$$

• direction of  $r$

$$g(\mathbf{x}) = \mathbf{w}'\mathbf{x} + w_0 = \mathbf{w}'(\mathbf{x}_p + r \frac{\mathbf{w}}{\|\mathbf{w}\|}) + w_0$$

$$= \mathbf{w}'\mathbf{x}_p + r \frac{\mathbf{w}'\mathbf{w}}{\|\mathbf{w}\|} + w_0 = r \|\mathbf{w}\|$$



FLIP: Pattern Recognition (G. Li)

32

## Geometric Interpretation of $g(\mathbf{x})$

- Therefore, the distance of  $\mathbf{x}$  from the hyperplane is given by:

$$r = \frac{g(\mathbf{x})}{\|\mathbf{w}\|}$$

- $w_0$  determines the distance of the hyperplane from the origin:

• setting  $\mathbf{x}=0$ :

$$r = \frac{w_0}{\|\mathbf{w}\|}$$



FLIP: Pattern Recognition (G. Li)

33

## Linear Discriminant Functions: multi-category case

- There are several ways to devise multi-category classifiers using linear discriminant functions:
  - One against the rest** (i.e.,  $c-1$  two-class problems)

problem  
• ambiguous regions



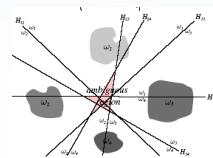
FLIP: Pattern Recognition (G. Li)

34

## Linear Discriminant Functions: multi-category case

- There are several ways to devise multi-category classifiers using linear discriminant functions:
  - One against another** (i.e.,  $c(c-1)/2$  pairs of classes)

problem  
• ambiguous regions

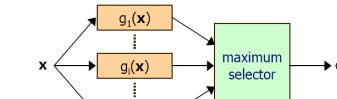


FLIP: Pattern Recognition (G. Li)

35

## Linear Discriminant Functions: multi-category case

- To avoid the problem of ambiguous regions:
  - Define  $c$  linear discriminant functions
  - Assign  $\mathbf{x}$  to  $\omega_i$  if  $g_i(\mathbf{x}) > g_j(\mathbf{x})$  for all  $j \neq i$ .
- The resulting classifier is called a **linear machine**



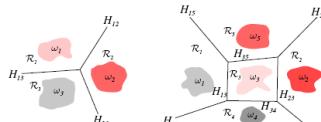
FLIP: Pattern Recognition (G. Li)

36

## Linear Discriminant Functions: multi-category case

- A linear machine divides the feature space in  $c$  convex decisions regions.

If  $x$  is in region  $R_i$ , the  $g_i(x)$  is the largest.



$c(c-1)/2$  pairs of regions but typically less decision boundaries



FLIP: Pattern Recognition (G. Li)

37

## Linear Discriminant Functions: multi-category case

- The boundary between two regions  $R_i$  and  $R_j$  is a portion of the hyperplane  $H_{ij}$  given by:

$$g_i(\mathbf{x}) = g_j(\mathbf{x}) \text{ or } g_i(\mathbf{x}) - g_j(\mathbf{x}) = 0$$

$$\text{or } (\mathbf{w}_i - \mathbf{w}_j)^T \mathbf{x} + (w_{i0} - w_{j0}) = 0$$

- $(\mathbf{w}_i - \mathbf{w}_j)$  is normal to  $H_{ij}$  and the signed distance from  $\mathbf{x}$  to  $H_{ij}$  is

$$r = \frac{g_i(\mathbf{x}) - g_j(\mathbf{x})}{\|\mathbf{w}_i - \mathbf{w}_j\|}$$



FLIP: Pattern Recognition (G. Li)

38

## Higher Order Discriminant Functions

- Can produce more complicated decision boundaries than linear discriminant functions.

$$\text{Linear discriminant: } g(\mathbf{x}) = \mathbf{w}' \mathbf{x} + w_0 = \sum_{i=1}^d w_i x_i + w_0$$

Quadratic discriminant: obtained by adding terms corresponding to products of pairs of components of  $\mathbf{x}$

$$g(\mathbf{x}) = w_0 + \sum_{i=1}^d w_i x_i + \sum_{i=1}^d \sum_{j=i+1}^d x_i x_j w_{ij}$$

Polynomial discriminant: obtained by adding terms such as  $x_i x_j x_k w_{ijk}$ .



FLIP: Pattern Recognition (G. Li)

39

## Generalized discriminants

- Generalized linear discriminant functions

$$g(\mathbf{x}) = \sum_{i=1}^{\hat{d}} a_i y_i(\mathbf{x}) \text{ or } g(\mathbf{x}) = \mathbf{a}' \mathbf{y}$$

$\mathbf{a}$  is a  $\hat{d}$ -dimensional weight vector

$y_i(\mathbf{x})$  functions map points from the  $d$ -dimensional  $\mathbf{x}$ -space to the  $\hat{d}$ -dimensional  $\mathbf{y}$ -space

usually,  $\hat{d} >> d$

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_d \end{bmatrix} \Rightarrow \begin{bmatrix} y_1(\mathbf{x}) \\ y_2(\mathbf{x}) \\ \vdots \\ y_{\hat{d}}(\mathbf{x}) \end{bmatrix}$$



FLIP: Pattern Recognition (G. Li)

40

## Generalized discriminants

- Generalized linear discriminant functions

$$g(\mathbf{x}) = \sum_{i=1}^{\hat{d}} a_i y_i(\mathbf{x}) \text{ or } g(\mathbf{x}) = \mathbf{a}' \mathbf{y}$$

The resulting discriminant function is not linear in  $\mathbf{x}$ , but it is linear in  $\mathbf{y}$ .

The generalized discriminant separates points in the transformed space by a hyperplane passing through the origin.



FLIP: Pattern Recognition (G. Li)

41

## Generalized discriminants

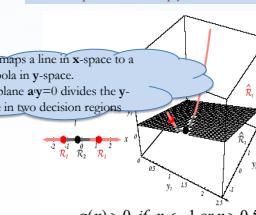
- Example

$$g(\mathbf{x}) = -1 + x + 2x^2$$

$$\mathbf{a} = \begin{bmatrix} -1 \\ 1 \\ 2 \end{bmatrix}$$

$$\mathbf{y} = \begin{bmatrix} y_1(\mathbf{x}) \\ y_2(\mathbf{x}) \\ y_3(\mathbf{x}) \end{bmatrix} = \begin{bmatrix} 1 \\ x \\ x^2 \end{bmatrix}$$

The corresponding decision regions  $R_1, R_2$  in the  $\mathbf{x}$ -space are not simply connected!



FLIP: Pattern Recognition (G. Li)

42

**Learning Discriminants**

- Linear Separable
- Solution Space

TULIP: Transfer for Universal Learning and Intelligent Processing

43

**Learning: two-category, linearly separable case**

- Task: Given a linear discriminant function, **learn** the weights  $\mathbf{w}$  using a set of  $n$  labeled samples  $\mathbf{x}_i$ 
  - i.e., each  $\mathbf{x}_i$  has a class label  $\omega_1$  or  $\omega_2$

$$g(\mathbf{x}) = \mathbf{w}' \mathbf{x} + w_0 = \sum_{i=1}^d w_i x_i + w_0$$

TULIP: Transfer for Universal Learning and Intelligent Processing

44

**Augmented feature/weight vectors**

- Increasing dimensionality simplifies things:
  - Find a single weight vector  $\mathbf{a}$  in  $\mathbf{y}$ -space.
  - Decision hyperplane, given by  $g(\mathbf{y})=0$ , in  $\mathbf{y}$ -space passes through the origin (but could be at any position in  $\mathbf{x}$ -space).

$$g(\mathbf{x}) = \mathbf{w}' \mathbf{x} + w_0 = \sum_{i=1}^d w_i x_i + x_0 w_0 = \sum_{i=0}^d w_i x_i = \mathbf{a}' \mathbf{y}$$

$$\mathbf{w} = \begin{bmatrix} w_1 \\ w_2 \\ \dots \\ w_d \end{bmatrix} \Rightarrow \mathbf{a} = \begin{bmatrix} w_0 \\ w_1 \\ \dots \\ w_d \end{bmatrix} \quad \mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ \dots \\ x_d \end{bmatrix} \Rightarrow \mathbf{y} = \begin{bmatrix} 1 \\ x_1 \\ x_2 \\ \dots \\ x_d \end{bmatrix} \quad (\text{i.e., } x_0 = 1)$$

• dimensionality:  $d \rightarrow (d+1)$

TULIP: Transfer for Universal Learning and Intelligent Processing

45

**Classification in augmented space**

- Given a discriminant function  $g(\mathbf{x})=\mathbf{a}'\mathbf{y}$ , the classification rule is:
  - If  $\mathbf{a}'\mathbf{y} > 0$  assign  $\mathbf{y}_i$  to  $\omega_1$
  - else if  $\mathbf{a}'\mathbf{y} < 0$  assign  $\mathbf{y}_i$  to  $\omega_2$

TULIP: Transfer for Universal Learning and Intelligent Processing

46

**Learning: two-category, linearly separable case**

- Given a discriminant function  $g(\mathbf{x})=\mathbf{a}'\mathbf{y}$ , the goal is to **learn** the weights  $\mathbf{a}$  using a set of  $n$  labeled samples  $\mathbf{y}_i$ 
  - Every training sample  $\mathbf{y}_i$  places a constraint on the weight vector  $\mathbf{a}$ .
    - $\mathbf{a}'\mathbf{y}_i = 0$  defines a hyperplane in solution space having  $\mathbf{y}_i$  as a normal vector.
  - Given  $n$  examples, the solution  $\mathbf{a}$  must lie on the intersection of  $n$  half-spaces

TULIP: Transfer for Universal Learning and Intelligent Processing

47

**Effect of training examples on solution**

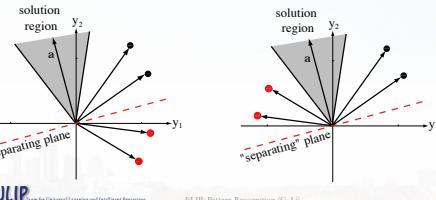
- Solution can be visualized either in the solution space or the feature space
  - solution space ( $a_1, a_2$ )
  - feature space ( $y_1, y_2$ )

TULIP: Transfer for Universal Learning and Intelligent Processing

48

## Normalized Version

- If  $y_i$  in  $\omega_2$ , replace  $y_i$  by  $-y_i$
- Find  $a$  such that:  $a'y_i > 0$

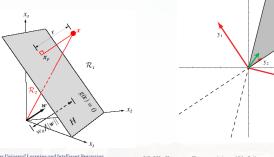


49

## Solution Uniqueness/Constraints

- Solution vector is usually not unique; we can impose certain constraints, e.g.:

**Method 1:** find unit-length weight vector that maximizes the minimum distance from the samples to the separating plane.



FLIP: Pattern Recognition (G. Li)

50

## Solution Uniqueness/Constraints

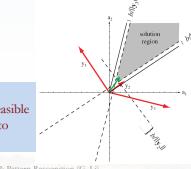
- Solution vector is usually not unique; we can impose certain constraints, e.g.:

**Method 2:** find the minimum-length weight vector satisfying the constraint below where  $b$  is a positive constant called **margin**.

$$g(\mathbf{x}) = \mathbf{a}'\mathbf{y}_i \geq b > 0$$

### Main objective:

- move solution to the center of the feasible region as this solution is more likely to classify new test samples correctly.



FLIP: Pattern Recognition (G. Li)

51

## Iterative Optimization

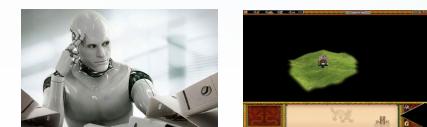
- Gradient Descent
- Why it works?
- Newton's Method



52

## Gradient Descent

- Gradient Descent is the “learning” of machines in deep learning .....
- Even AlphaGo using this approach.



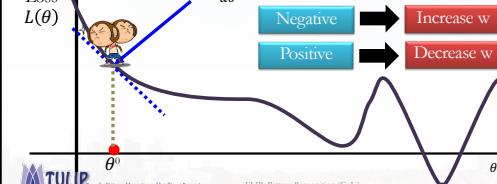
FLIP: Pattern Recognition (G. Li)

53

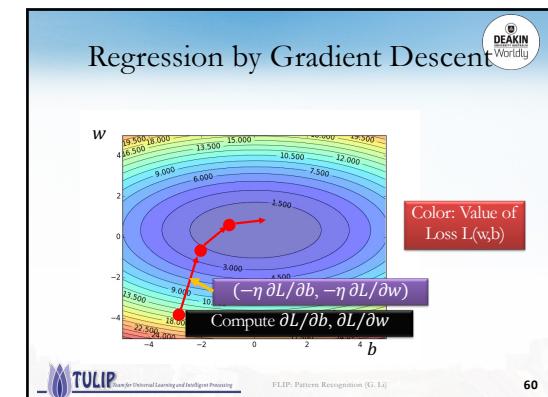
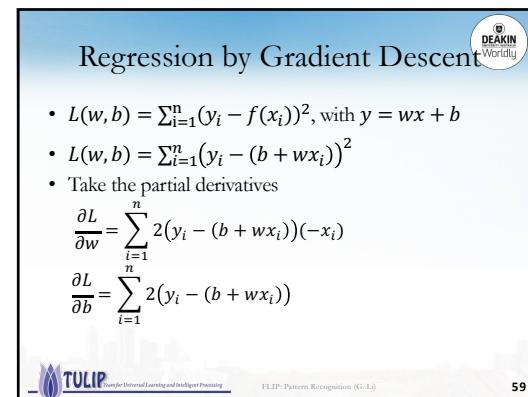
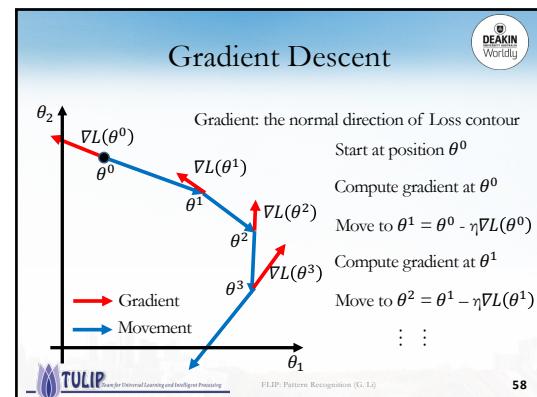
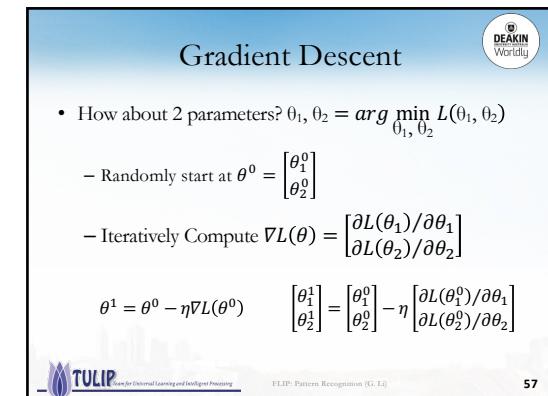
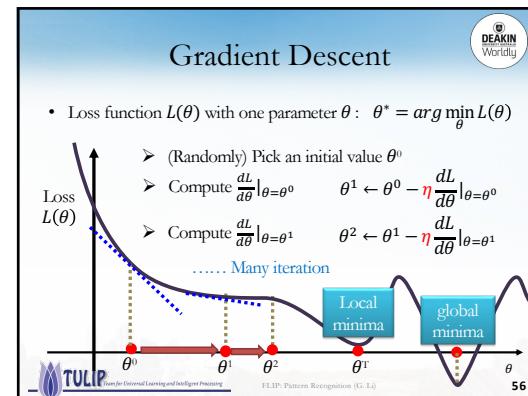
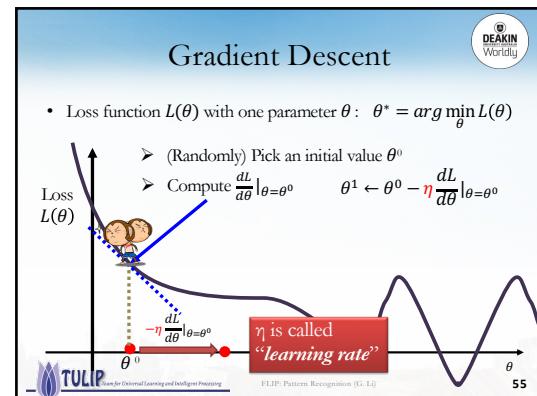
## Gradient Descent

- Loss function  $L(\theta)$  with one parameter  $\theta$ :  $\theta^* = \arg \min_{\theta} L(\theta)$

➤ (Randomly) Pick an initial value  $\theta^0$   
➤ Compute  $\frac{dL}{d\theta}|_{\theta=\theta^0}$



54



## Gradient Descent

- When solving:
$$\theta^* = \arg \min_{\theta} L(\theta) \text{ by gradient descent}$$
  - Each time we update the parameters, we obtain  $\theta$  that makes  $L(\theta)$  smaller.
- $L(\theta^0) > L(\theta^1) > L(\theta^2) > \dots$

Is this statement correct?



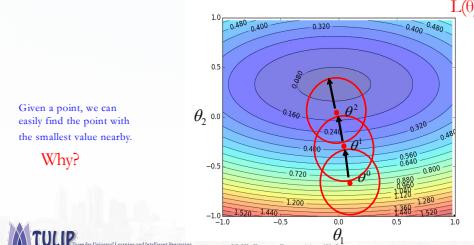
FLIP: Pattern Recognition (G. Li)

61



## Formal Derivation

- Suppose that  $\theta$  has two variables  $\{\theta_1, \theta_2\}$

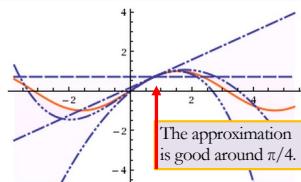


62

61

## Taylor Series

$$\sin(x) = \frac{1}{\sqrt{2}} + \frac{x - \frac{\pi}{4}}{\sqrt{2}} - \frac{(x - \frac{\pi}{4})^2}{2\sqrt{2}} + \frac{(x - \frac{\pi}{4})^3}{6\sqrt{2}} - \frac{(x - \frac{\pi}{4})^4}{24\sqrt{2}} + \frac{(x - \frac{\pi}{4})^5}{120\sqrt{2}} - \frac{(x - \frac{\pi}{4})^6}{720\sqrt{2}} - \dots$$



FLIP: Pattern Recognition (G. Li)

64



## Multivariable Taylor Series

$$h(x, y) = h(x_0, y_0) + \frac{\partial h(x_0, y_0)}{\partial x}(x - x_0) + \frac{\partial h(x_0, y_0)}{\partial y}(y - y_0) + \text{something related to } (x-x_0)^2 \text{ and } (y-y_0)^2 + \dots$$

When  $x$  and  $y$  are close to  $x_0$  and  $y_0$

$$h(x, y) \approx h(x_0, y_0) + \frac{\partial h(x_0, y_0)}{\partial x}(x - x_0) + \frac{\partial h(x_0, y_0)}{\partial y}(y - y_0)$$



FLIP: Pattern Recognition (G. Li)

65

65

## Taylor Series

- Taylor series:** Let  $h(x)$  be any function infinitely differentiable around  $x = x_0$ .

$$h(x) = \sum_{k=0}^{\infty} \frac{h^{(k)}(x_0)}{k!} (x - x_0)^k = h(x_0) + h'(x_0)(x - x_0) + \frac{h''(x_0)}{2!} (x - x_0)^2 + \dots$$

When  $x$  is close to  $x_0$   $\Rightarrow h(x) \approx h(x_0) + h'(x_0)(x - x_0)$



FLIP: Pattern Recognition (G. Li)

63

63

## Back to Formal Derivation

Based on Taylor Series:

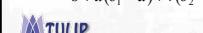
- If the red circle is small enough, in the red circle

$$L(\theta) \approx L(a, b) + \frac{\partial L(a, b)}{\partial \theta_1} (\theta_1 - a) + \frac{\partial L(a, b)}{\partial \theta_2} (\theta_2 - b)$$

$$s = L(a, b)$$

$$u = \frac{\partial L(a, b)}{\partial \theta_1}, v = \frac{\partial L(a, b)}{\partial \theta_2}$$

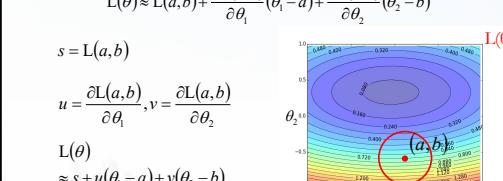
$$L(\theta) \approx s + u(\theta_1 - a) + v(\theta_2 - b)$$



FLIP: Pattern Recognition (G. Li)

66

66



64

## Back to Formal Derivation

Based on Taylor Series:

- If the red circle is small enough, in the red circle

$$L(\theta) \approx s + u(\theta_1 - a) + v(\theta_2 - b)$$

$$s = L(a, b)$$

$$u = \frac{\partial L(a, b)}{\partial \theta_1}, v = \frac{\partial L(a, b)}{\partial \theta_2}$$

Find  $\theta_1$  and  $\theta_2$  in the red circle minimizing  $L(\theta)$

$$(\theta_1 - a)^2 + (\theta_2 - b)^2 \leq d^2$$

**Simple, right?**

FLIP: Pattern Recognition (G, L)

67

## Gradient descent – two variables

Red Circle: (If the radius is small)

$$L(\theta) \approx s + u(\theta_1 - a) + v(\theta_2 - b)$$

$$\Delta \theta_1 \quad \Delta \theta_2$$

Find  $\theta_1$  and  $\theta_2$  in the red circle minimizing  $L(\theta)$

$$(\theta_1 - a)^2 + (\theta_2 - b)^2 \leq d^2$$

To minimize  $L(\theta)$

$$\begin{bmatrix} \Delta \theta_1 \\ \Delta \theta_2 \end{bmatrix} = -\eta \begin{bmatrix} u \\ v \end{bmatrix} \rightarrow \begin{bmatrix} \theta_1 \\ \theta_2 \end{bmatrix} = \begin{bmatrix} a \\ b \end{bmatrix} - \eta \begin{bmatrix} u \\ v \end{bmatrix}$$

FLIP: Pattern Recognition (G, L)

68

## Back to Formal Derivation

Based on Taylor Series:

- If the red circle is small enough, in the red circle

$$s = L(a, b)$$

$$u = \frac{\partial L(a, b)}{\partial \theta_1}, v = \frac{\partial L(a, b)}{\partial \theta_2}$$

Find  $\theta_1$  and  $\theta_2$  yielding the smallest value of  $L(\theta)$  in the circle

$$\begin{bmatrix} \theta_1 \\ \theta_2 \end{bmatrix} = \begin{bmatrix} a \\ b \end{bmatrix} - \eta \begin{bmatrix} u \\ v \end{bmatrix} = \begin{bmatrix} a \\ b \end{bmatrix} - \eta \begin{bmatrix} \frac{\partial L(a, b)}{\partial \theta_1} \\ \frac{\partial L(a, b)}{\partial \theta_2} \end{bmatrix}$$

This is gradient descent.

Not satisfied if the red circle (learning rate) is not small enough  
You can consider the second order term, e.g. Newton's method.

FLIP: Pattern Recognition (G, L)

69

## Newton's method

- Using second order approximation
  - Make sure  $L(\theta)$  reaches the lowest point along the decreasing direction

$$L(\theta) \approx L(\theta_0) + L'(\theta_0)(\theta - \theta_0) + \frac{L''(\theta_0)}{2}(\theta - \theta_0)^2$$

*Descent (1<sup>st</sup> derivatives):  $\nabla L$*

*Hessian (2<sup>nd</sup> derivatives):  $H$*

*optimum learning rate*

If  $L(\theta)$  is quadratic, then  $H$  is constant which implies that the learning rate is constant!

FLIP: Pattern Recognition (G, L)

70

## Newton's method

- We will try to update  $\theta$  to minimize  $L(\theta)$

$$L(\theta) \approx L(\theta_0) + L'(\theta_0)(\theta - \theta_0) + \frac{L''(\theta_0)}{2}(\theta - \theta_0)^2$$

- Take the first derivative:

$$\frac{\partial L}{\partial \theta} = L'(\theta_0) + L''(\theta_0)(\theta - \theta_0) = 0$$

- Then we have:  $(\theta - \theta_0) = -\frac{L'(\theta_0)}{L''(\theta_0)}$  hence the update rule:

$$\theta = \theta_0 - \frac{L'(\theta_0)}{L''(\theta_0)}$$

*optimum learning rate*

- If  $L(\theta)$  is quadratic, then  $H$  is constant which implies that the learning rate is constant!

FLIP: Pattern Recognition (G, L)

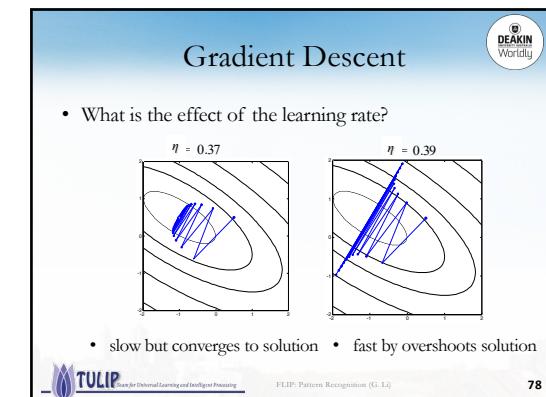
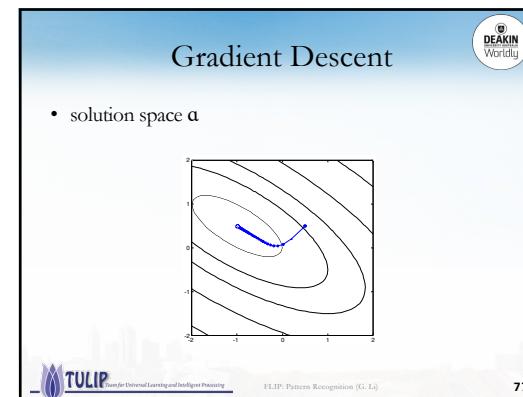
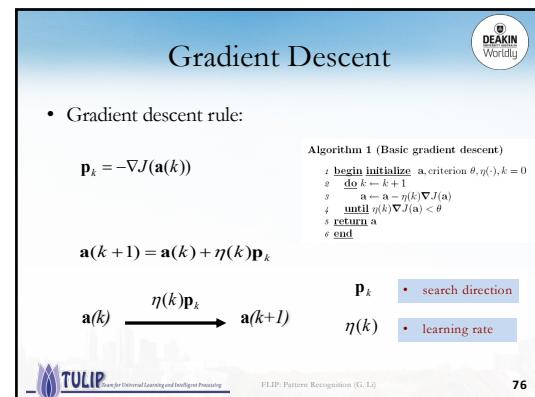
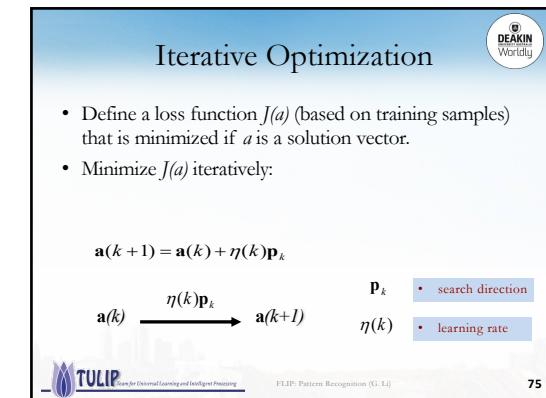
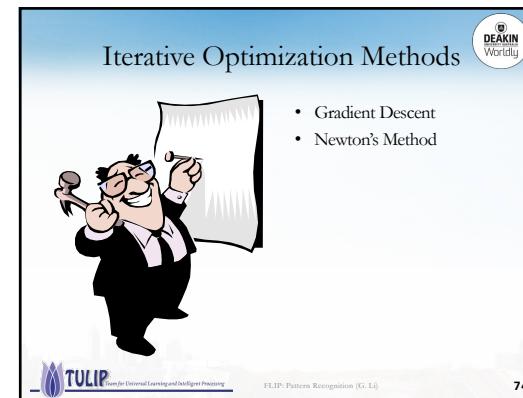
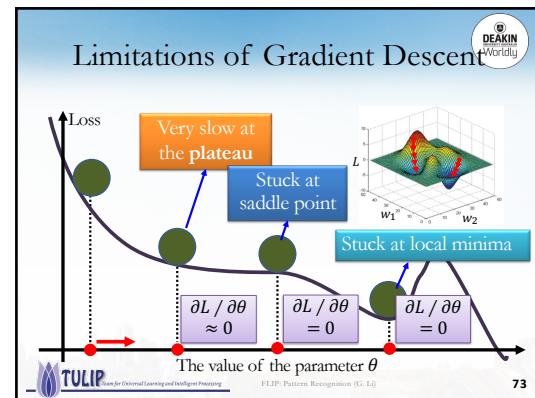
71

## Newton's method

- If  $L(\theta)$  is quadratic, Newton's method converges in **one step!**

FLIP: Pattern Recognition (G, L)

72



## Gradient Descent

- How to choose the learning rate  $\eta(k)$ ?
  - Make sure  $J(a(k+1))$  reaches the lowest point along the decreasing direction

*Taylor series approximation*

$$J(\mathbf{a}) \simeq J(\mathbf{a}(k)) + \nabla J^t(\mathbf{a} - \mathbf{a}(k)) + \frac{1}{2}(\mathbf{a} - \mathbf{a}(k))^t \mathbf{H} (\mathbf{a} - \mathbf{a}(k)),$$

*Hessian (2nd derivatives)*

*optimum learning rate*

- If  $J(\mathbf{a})$  is quadratic, then  $\mathbf{H}$  is constant which implies that the learning rate is constant!

**TULIP** Team for Universal Learning and Intelligent Processing FLIP: Pattern Recognition (G. Li) 79

## Newton's method

- Newton's descent rule:

$$\mathbf{p}_k = -\mathbf{H}^{-1}\nabla J(\mathbf{a}(k))$$

requires inverting  $\mathbf{H}$

$$\mathbf{a}(k+1) = \mathbf{a}(k) + \eta(k)\mathbf{p}_k$$

$\mathbf{p}_k$  • search direction  
 $\eta(k)$  • learning rate

**Algorithm 2 (Newton descent)**

```

1 begin initialize a, criterion θ
2   do
3     a ← a - H-1∇J(a)
4     until H-1∇J(a) < θ
5   return a
6 end

```

**TULIP** Team for Universal Learning and Intelligent Processing FLIP: Pattern Recognition (G. Li) 80

## Newton's method

- If  $J(\mathbf{a})$  is quadratic, Newton's method converges in **one step**!

**TULIP** Team for Universal Learning and Intelligent Processing FLIP: Pattern Recognition (G. Li) 81

79

80

81

## Criterion Function

- Perceptron Criterion Function
- Relaxation Procedure

**TULIP** Team for Universal Learning and Intelligent Processing FLIP: Pattern Recognition (G. Li) 82

## Criterion Function

- Constructing a **criterion function** for solving the linear inequalities  $\mathbf{a}^t \mathbf{y} \geq 0$ .
  - The simplest choice is let  $J(\mathbf{a})$  be the number of samples misclassified by  $\mathbf{a}$ .
  - This leads to a piecewise constant, not suitable for gradient search

**TULIP** Team for Universal Learning and Intelligent Processing FLIP: Pattern Recognition (G. Li) 83

82

83

## Perceptron Criterion Function

- Constructing a **criterion function** for solving the linear inequalities  $\mathbf{a}^t \mathbf{y} \geq 0$ .
  - Perceptron Criterion Function:**

$$J_p(\mathbf{a}) = \sum_{\mathbf{y} \in Y(\mathbf{a})} (-\mathbf{a}^t \mathbf{y})$$

- where  $Y(\mathbf{a})$  is the set of samples misclassified by  $\mathbf{a}$ .
- If  $Y(\mathbf{a})$  is empty,  $J_p(\mathbf{a})=0$ ; otherwise,  $J_p(\mathbf{a})>0$  (normalized version)

**TULIP** Team for Universal Learning and Intelligent Processing FLIP: Pattern Recognition (G. Li) 84

84

## Perceptron rule (Batch Correction)

- The gradient of  $J_p(\mathbf{a})$  is:

$$\nabla J_p = \sum_{\mathbf{y} \in Y(\mathbf{a})} (-\mathbf{y})$$

- The perceptron update rule is obtained using gradient descent:

$$\begin{aligned} \mathbf{a}(k+1) &= \mathbf{a}(k) - \eta(k) \nabla J(\mathbf{a}(k)) \\ \mathbf{a}(k+1) &= \mathbf{a}(k) + \eta(k) \sum_{\mathbf{y} \in Y(\mathbf{a})} \mathbf{y} \end{aligned}$$

Algorithm 3 (Batch Perceptron)

```

1 begin initialize a,  $\eta()$ , criterion  $\theta$ ,  $k = 0$ 
2   do  $k \leftarrow k + 1$ 
3     a  $\leftarrow a + \eta(k) \sum_{\mathbf{y} \in Y_k} \mathbf{y}$ 
4     until  $\eta(k) \sum_{\mathbf{y} \in Y_k} \mathbf{y} < \theta$ 
5   return a
6 end
• consider all examples misclassified

```



FLIP: Pattern Recognition (G. Li)

85

## Perceptron rule (Single-Sample Correction)

- Fix the learning rate  $\eta(k) = 1$
- Modify the weight vector upon each misclassified single sample
- The misclassified sample will be input infinitely often
- The perceptron update rule is obtained using gradient descent:

$$\mathbf{a}(k+1) = \mathbf{a}(k) + \eta(k) \mathbf{y}$$



FLIP: Pattern Recognition (G. Li)

86

## Perceptron rule (Single-Sample Correction)

### • Perceptron Convergence Theorem:

- If training samples are linearly separable, then the sequence of weight vectors will terminate at a solution vector in a finite number of steps.

Algorithm 4 (Fixed-increment single-sample Perceptron)

```

1 begin initialize a,  $k = 0$ 
2   do  $k \leftarrow (k + 1) \bmod n$ 
3     if  $\mathbf{y}_k$  is misclassified by  $a$  then  $a \leftarrow a + \mathbf{y}_k$ 
4     until all patterns properly classified
5   return a
6 end
• consider one example at a time

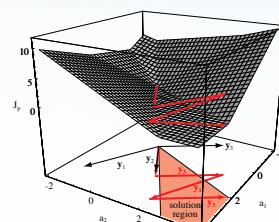
```



FLIP: Pattern Recognition (G. Li)

87

## Perceptron rule



FLIP: Pattern Recognition (G. Li)

88

## Perceptron rule (Margins)

- Generalizations:
  - Variable learning rate and a margin

Algorithm 5 (Variable increment Perceptron with margin)

```

1 begin initialize a, criterion  $\theta$ , margin b,  $\eta()$ ,  $k = 0$ 
2   do  $k \leftarrow k + 1$ 
3     if  $\mathbf{a}^T \mathbf{y}_k + b < 0$  then  $a \leftarrow a + \eta(k) \mathbf{y}_k$ 
4     until  $\mathbf{a}^T \mathbf{y}_k + b \leq 0$  for all  $k$ 
5   return a
6 end

```



FLIP: Pattern Recognition (G. Li)

89

## Relaxation Procedures

- Note that different criterion functions exist
- One possible choice is:

$$J_q(a) = \sum_{\mathbf{y} \in Y} (a^T \mathbf{y})^2$$

- Where  $Y$  is again the set of the training samples that are misclassified by  $a$

- However, there are two problems with this criterion

- The function is too smooth and can converge to  $a=0$
- $J_q$  is dominated by training samples with large magnitude



FLIP: Pattern Recognition (G. Li)

90

## Relaxation Procedures

- A modified version that avoids two problems is

$$J_r(a) = \frac{1}{2} \sum_{y \in Y} \frac{(a^T y - b)^2}{\|y\|^2}$$

- Here  $Y$  is the set of samples for which its gradient is given by:

$$\nabla J_r = \sum_{y \in Y} \frac{(a^T y - b)}{\|y\|^2} y$$



FLIP: Pattern Recognition (G. Li)

91



91

## Relaxation Procedures

### Algorithm 8 (Batch relaxation with margin)

```

1 begin initialize a,  $\eta(\cdot), k = 0$ 
2   do  $k \leftarrow k + 1$ 
3      $\mathcal{Y}_k = \{\}$ 
4      $j = 0$ 
5     do  $j \leftarrow j + 1$ 
6       if  $y_j$  is misclassified then Append  $y_j$  to  $\mathcal{Y}_k$ 
7       until  $j = n$ 
8      $a \leftarrow a + \eta(k) \sum_{y \in \mathcal{Y}_k} \frac{b - a^T y}{\|y\|^2} y$ 
9     until  $\mathcal{Y}_k = \{\}$ 
10    return a
11  end
```



FLIP: Pattern Recognition (G. Li)

92

92



## Relaxation Procedures

### Algorithm 9 (Single-sample relaxation with margin)

```

1 begin initialize a,  $\eta(\cdot), k = 0$ 
2   do  $k \leftarrow k + 1$ 
3     if  $y_k$  is misclassified then  $a \leftarrow a + \eta(k) \frac{b - a^T y_k}{\|y_k\|^2} y_k$ 
4     until all patterns properly classified
5   return a
6  end
```



FLIP: Pattern Recognition (G. Li)

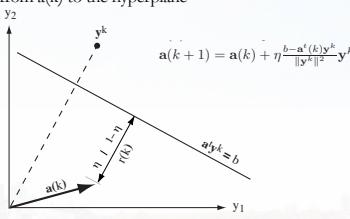
93

93



## Relaxation Procedures

- Since  $y/\|y\|$  is the unit normal vector for the hyperplane. The correction calls for  $a(k)$  to be moved a certain fraction of  $\eta$  of the distance from  $a(k)$  to the hyperplane



FLIP: Pattern Recognition (G. Li)

94

94



## Nonseparable Behavior

- Perceptron and relaxation procedures are all error-correction procedure.
  - Their success on separable problems is largely due to its relentless search for an error-free solution.
  - For non-separable problems, correction may never converge.



FLIP: Pattern Recognition (G. Li)

95

95



## MSE Procedure

- MSE and Discrimants
- Pseudoinverse



FLIP: Pattern Recognition (G. Li)

96

16

## Minimum Squared Error Procedures

- Define a criterion function for all samples, rather than the misclassified samples
- Try to find the vector such that  $\mathbf{a}^T \mathbf{y}_i > b_i$ 
  - Where  $b_i$  is any positive value

$$\begin{pmatrix} Y_{10} & Y_{11} & \cdots & Y_{1d} \\ Y_{20} & Y_{21} & \cdots & Y_{2d} \\ \vdots & \vdots & & \vdots \\ \vdots & \vdots & & \vdots \\ Y_{n0} & Y_{n1} & \cdots & Y_{nd} \end{pmatrix} \begin{pmatrix} a_0 \\ a_1 \\ \vdots \\ a_d \end{pmatrix} = \begin{pmatrix} b_1 \\ b_2 \\ \vdots \\ b_n \end{pmatrix}$$



FLIP: Pattern Recognition (G. Li)

97

97

## Minimum Squared Error Procedures

- Minimum squared error and pseudoinverse
  - The problem is to find a weight vector  $\mathbf{a}$  satisfying  $\mathbf{Y}\mathbf{a} = \mathbf{b}$
  - If we have more equations than unknowns,  $\mathbf{a}$  is over-determined.
  - We want to choose the one that minimizes the sum-of-squared-error criterion function

$$J_s(\mathbf{a}) = \|\mathbf{Y}\mathbf{a} - \mathbf{b}\|^2 = \sum_{i=1}^n (\mathbf{a}^T \mathbf{y}_i - b_i)^2.$$



FLIP: Pattern Recognition (G. Li)

98

98

## Minimum Squared Error Procedures

- Pseudoinverse

$$\begin{aligned} \mathbf{a} &= (\mathbf{Y}^T \mathbf{Y})^{-1} \mathbf{Y}^T \mathbf{b} \\ &= \mathbf{Y}^{\dagger} \mathbf{b}, \end{aligned}$$

### Algorithm 10 (LMS)

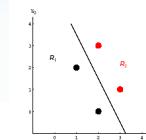
```

1 begin initialize  $\mathbf{a}, \mathbf{b}$ , criterion  $\theta, \eta(\cdot), k = 0$ 
2   do  $k \leftarrow k + 1$ 
3      $\mathbf{a} \leftarrow \mathbf{a} + \eta(k)(b_k - \mathbf{a}^T \mathbf{y}^k) \mathbf{y}^k$ 
4   until  $\eta(k)(b_k - \mathbf{a}^T \mathbf{y}^k) \mathbf{y}^k < \theta$ 
5   return  $\mathbf{a}$ 
6 end

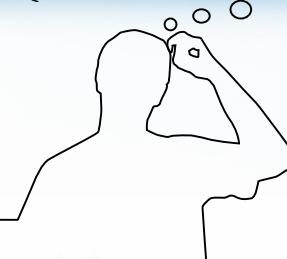
```

FLIP: Pattern Recognition (G. Li)

99



Questions?



FLIP: Pattern Recognition (G. Li)

101

101