

Lecture Slides for

INTRODUCTION TO MACHINE LEARNING

3RD EDITION

ETHEM ALPAYDIN
© The MIT Press, 2014

alpaydin@boun.edu.tr
<http://www.cmpe.boun.edu.tr/~ethem/i2ml3e>

CHAPTER 10:

LINEAR DISCRIMINATION

Likelihood- vs. Discriminant-based Classification

3

- **Likelihood-based:** Assume a model for $p(\mathbf{x} | C_i)$, use Bayes' rule to calculate $P(C_i | \mathbf{x})$

$$g_i(\mathbf{x}) = \log P(C_i | \mathbf{x})$$

- **Discriminant-based:** Assume a model for $g_i(\mathbf{x} | \Phi_i)$; no density estimation
- Estimating the boundaries is enough; no need to accurately estimate the densities inside the boundaries

Linear Discriminant

4

- Linear discriminant:

$$g_i(\mathbf{x} | \mathbf{w}_i, w_{i0}) = \mathbf{w}_i^T \mathbf{x} + w_{i0} = \sum_{j=1}^d w_{ij} x_j + w_{i0}$$

- Advantages:
 - ▣ Simple: $O(d)$ space/computation
 - ▣ Knowledge extraction: Weighted sum of attributes; positive/negative weights, magnitudes (credit scoring)
 - ▣ Optimal when $p(\mathbf{x} | C_i)$ are Gaussian with shared cov matrix; useful when classes are (almost) linearly separable

Generalized Linear Model

5

- Quadratic discriminant:

$$g_i(\mathbf{x} | \mathbf{W}_i, \mathbf{w}_i, w_{i0}) = \mathbf{x}^T \mathbf{W}_i \mathbf{x} + \mathbf{w}_i^T \mathbf{x} + w_{i0}$$

- Higher-order (product) terms:

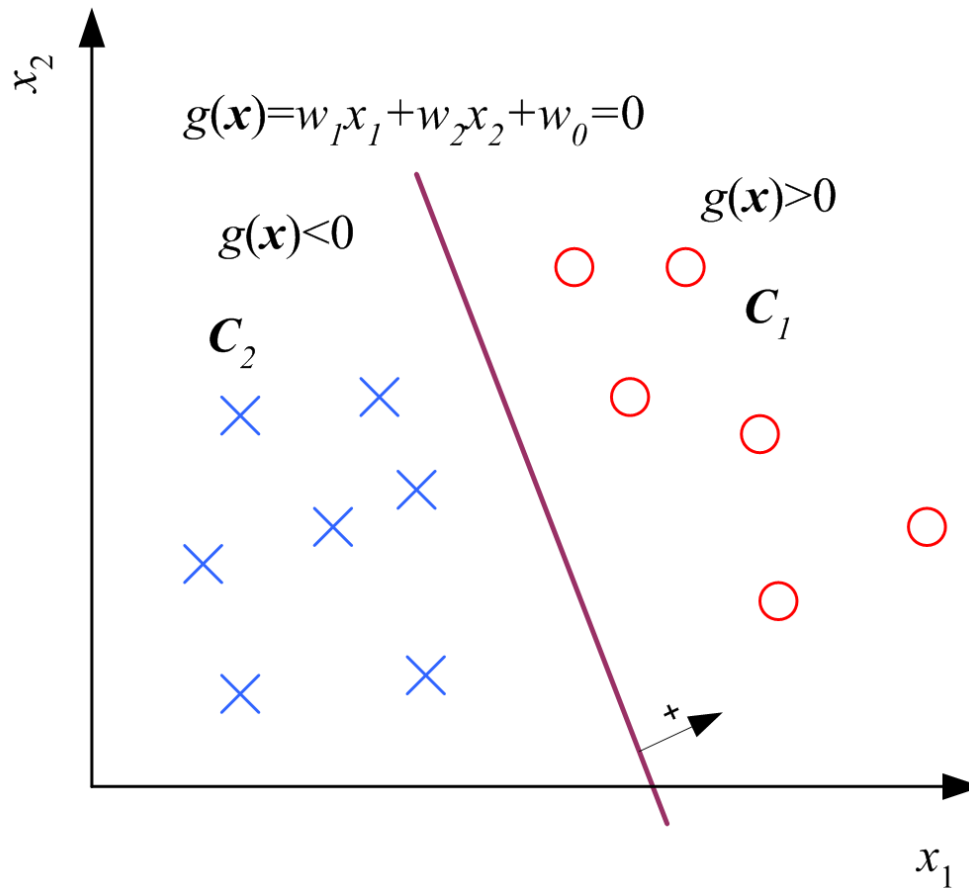
$$z_1 = x_1, z_2 = x_2, z_3 = x_1^2, z_4 = x_2^2, z_5 = x_1 x_2$$

Map from \mathbf{x} to \mathbf{z} using nonlinear basis functions and use a linear discriminant in \mathbf{z} -space

$$g_i(\mathbf{x}) = \sum_{j=1}^k w_{ij} \phi_j(\mathbf{x})$$

Two Classes

6

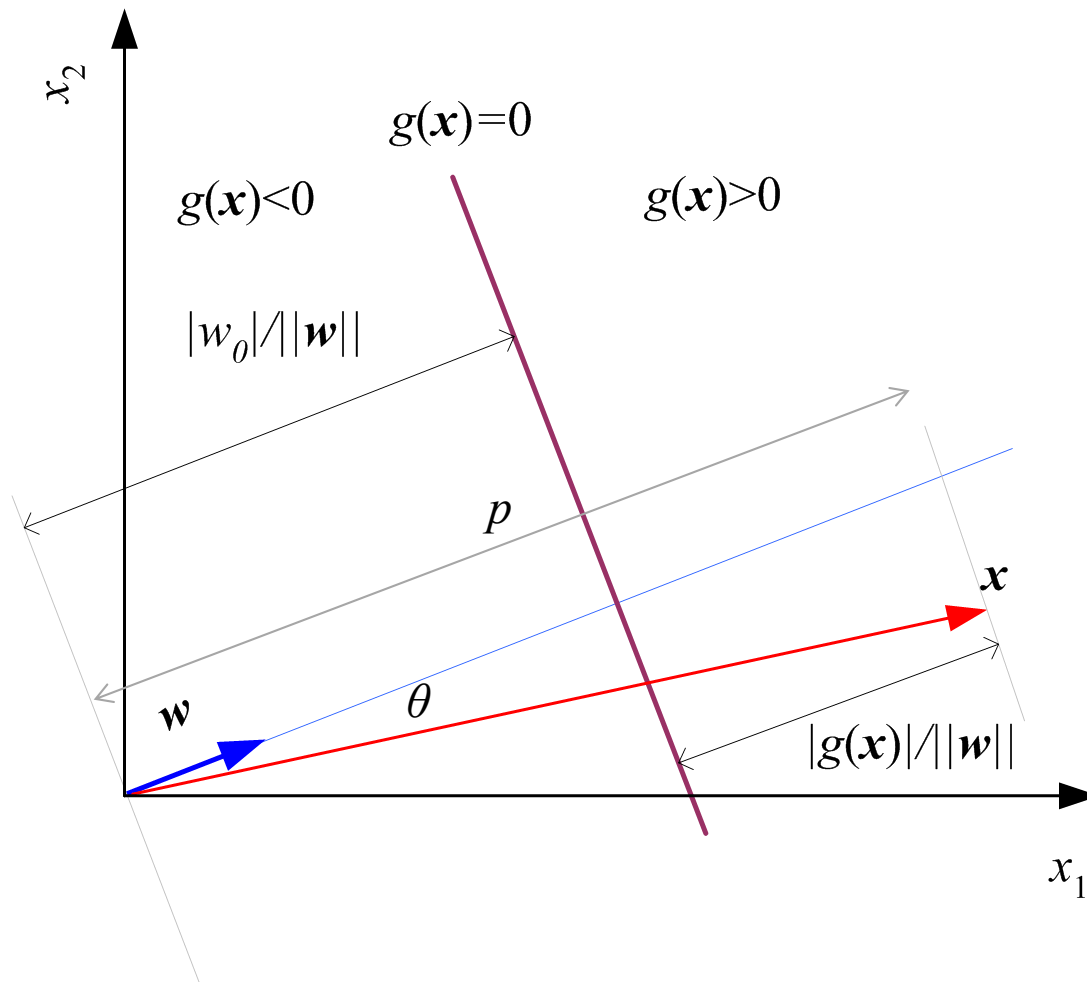


$$\begin{aligned} g(\mathbf{x}) &= g_1(\mathbf{x}) - g_2(\mathbf{x}) \\ &= (\mathbf{w}_1^T \mathbf{x} + w_{10}) - (\mathbf{w}_2^T \mathbf{x} + w_{20}) \\ &= (\mathbf{w}_1 - \mathbf{w}_2)^T \mathbf{x} + (w_{10} - w_{20}) \\ &= \mathbf{w}^T \mathbf{x} + w_0 \end{aligned}$$

choose $\begin{cases} C_1 & \text{if } g(\mathbf{x}) > 0 \\ C_2 & \text{otherwise} \end{cases}$

Geometry

7



$$g(\mathbf{x}) > 0 \text{ iff } \mathbf{w}^T \mathbf{x} + w_0 > 0$$

$$\text{iff } \mathbf{w}^T \mathbf{x} > -w_0$$

$$\text{iff } \|\mathbf{w}\| \|\mathbf{x}\| \cos \theta > -w_0$$

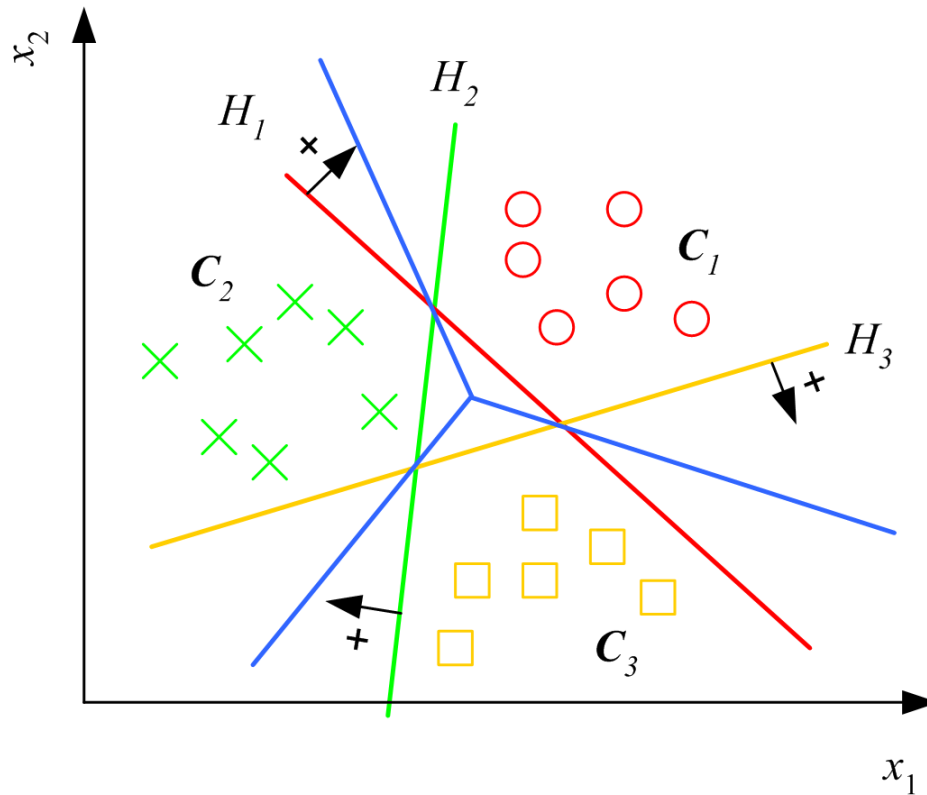
$$\text{iff } p \equiv \|\mathbf{x}\| \cos \theta > -w_0 / \|\mathbf{w}\|$$

(edited by BJM)

Multiple Classes

8

$$g_i(\mathbf{x} | \mathbf{w}_i, w_{i0}) = \mathbf{w}_i^T \mathbf{x} + w_{i0}$$

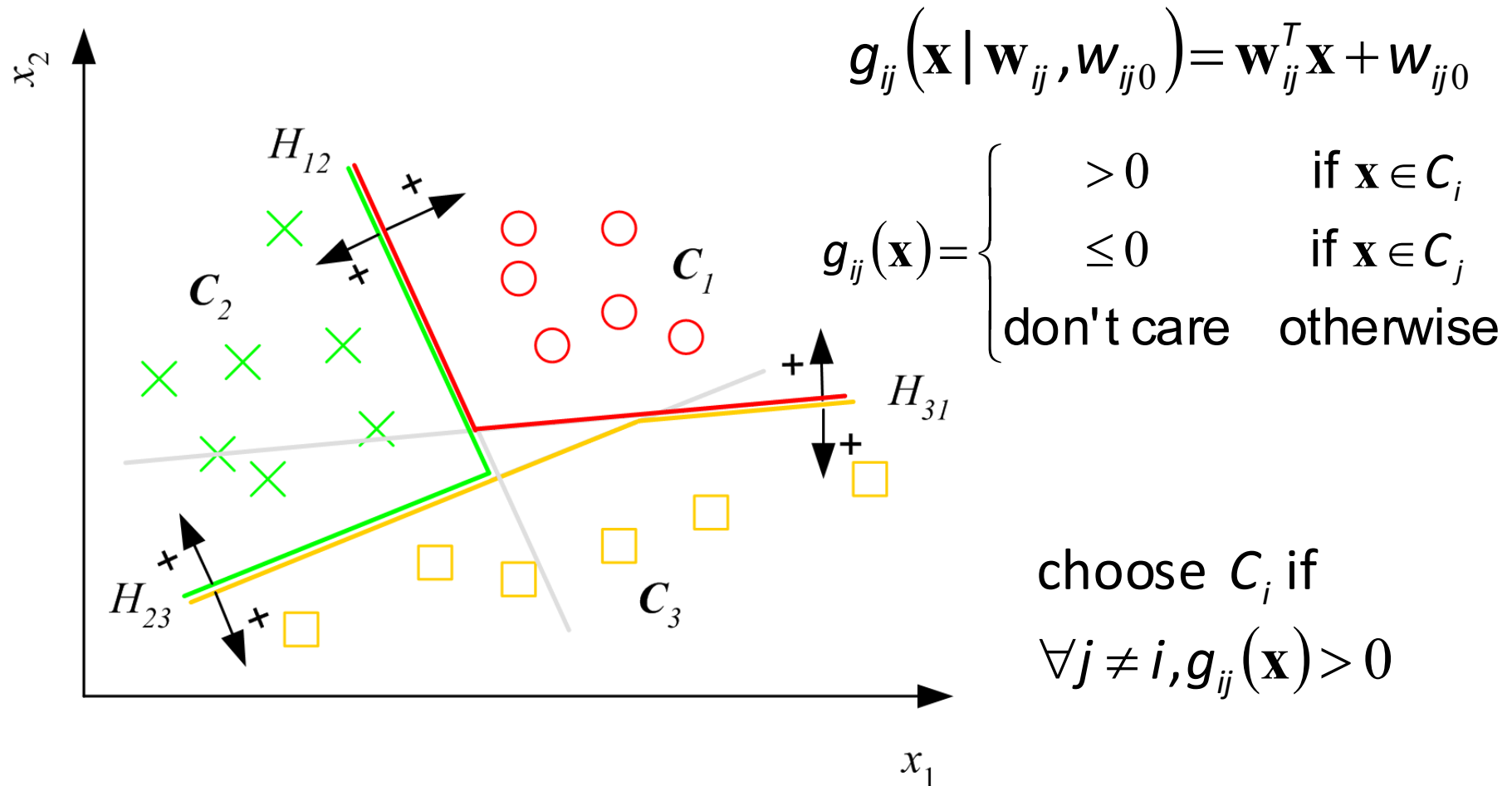


Choose C_i if

$$g_i(\mathbf{x}) = \max_{j=1}^K g_j(\mathbf{x})$$

Classes are
linearly separable

Pairwise Separation



From Discriminants to Posteriors

10

When $p(\mathbf{x} \mid C_i) \sim \mathcal{N}(\boldsymbol{\mu}_i, \Sigma)$

$$g_i(\mathbf{x} \mid \mathbf{w}_i, w_{i0}) = \mathbf{w}_i^T \mathbf{x} + w_{i0}$$

$$\mathbf{w}_i = \Sigma^{-1} \boldsymbol{\mu}_i \quad w_{i0} = -\frac{1}{2} \boldsymbol{\mu}_i^T \Sigma^{-1} \boldsymbol{\mu}_i + \log P(C_i)$$

$$y \equiv P(C_1 \mid \mathbf{x}) \text{ and } P(C_2 \mid \mathbf{x}) = 1 - y$$

$$\text{choose } C_1 \text{ if } \begin{cases} y > 0.5 \\ y/(1-y) > 1 \\ \log[y/(1-y)] > 0 \end{cases} \text{ and } C_2 \text{ otherwise}$$

$$\begin{aligned}
\text{logit}(P(C_1 | \mathbf{x})) &= \log \frac{P(C_1 | \mathbf{x})}{1 - P(C_1 | \mathbf{x})} = \log \frac{P(C_1 | \mathbf{x})}{P(C_2 | \mathbf{x})} \\
&= \log \frac{p(\mathbf{x} | C_1)}{p(\mathbf{x} | C_2)} + \log \frac{P(C_1)}{P(C_2)} \\
&= \log \frac{(2\pi)^{-d/2} |\Sigma|^{-1/2} \exp\left[-(1/2)(\mathbf{x} - \mu_1)^T \Sigma^{-1} (\mathbf{x} - \mu_1)\right]}{(2\pi)^{-d/2} |\Sigma|^{-1/2} \exp\left[-(1/2)(\mathbf{x} - \mu_2)^T \Sigma^{-1} (\mathbf{x} - \mu_2)\right]} + \log \frac{P(C_1)}{P(C_2)} \\
&= \mathbf{w}^T \mathbf{x} + w_0
\end{aligned}$$

$$\text{where } \mathbf{w} = \Sigma^{-1}(\mu_1 - \mu_2) \quad w_0 = -\frac{1}{2}(\mu_1 + \mu_2)^T \Sigma^{-1}(\mu_1 - \mu_2) + \log \frac{P(C_1)}{P(C_2)}$$

The inverse of logit

$$\log \frac{P(C_1 | \mathbf{x})}{1 - P(C_1 | \mathbf{x})} = \mathbf{w}^T \mathbf{x} + w_0$$

$$P(C_1 | \mathbf{x}) = \text{sigmoid}(\mathbf{w}^T \mathbf{x} + w_0) = \frac{1}{1 + \exp\left[-(\mathbf{w}^T \mathbf{x} + w_0)\right]}$$

$$\log \frac{P(C_1|x)}{1-P(C_1|x)} \equiv \log \frac{P}{1-P} = w^T x + w_0$$

$$\therefore \frac{P}{1-P} = \exp(w^T x + w_0)$$

$$P = (1 - P) \exp(w^T x + w_0)$$

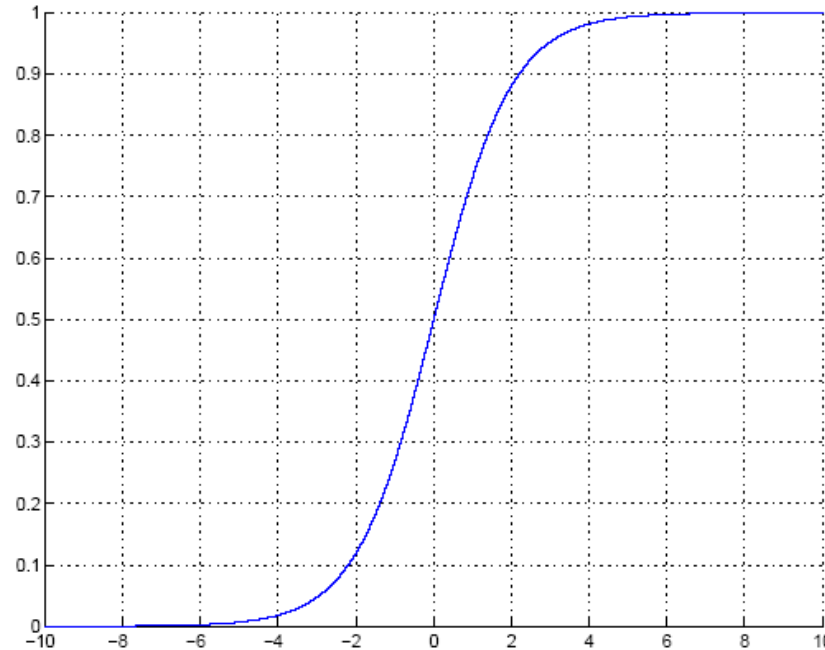
$$= \exp(w^T x + w_0) - P \exp(w^T x + w_0)$$

$$P[1 + \exp(w^T x + w_0)] = \exp(w^T x + w_0)$$

$$P = \frac{\exp(w^T x + w_0)}{1 + \exp(w^T x + w_0)} = \frac{1}{1 + \exp[-(w^T x + w_0)]}$$

Sigmoid (Logistic) Function

13



Calculate $g(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + w_0$ and choose C_1 if $g(\mathbf{x}) > 0$, or

Calculate $y = \text{sigmoid}(\mathbf{w}^T \mathbf{x} + w_0)$ and choose C_1 if $y > 0.5$

Gradient-Descent

14

- $E(\mathbf{w} \mid \mathcal{X})$ is error with parameters \mathbf{w} on sample \mathcal{X}
 $\mathbf{w}^* = \arg \min_{\mathbf{w}} E(\mathbf{w} \mid \mathcal{X})$

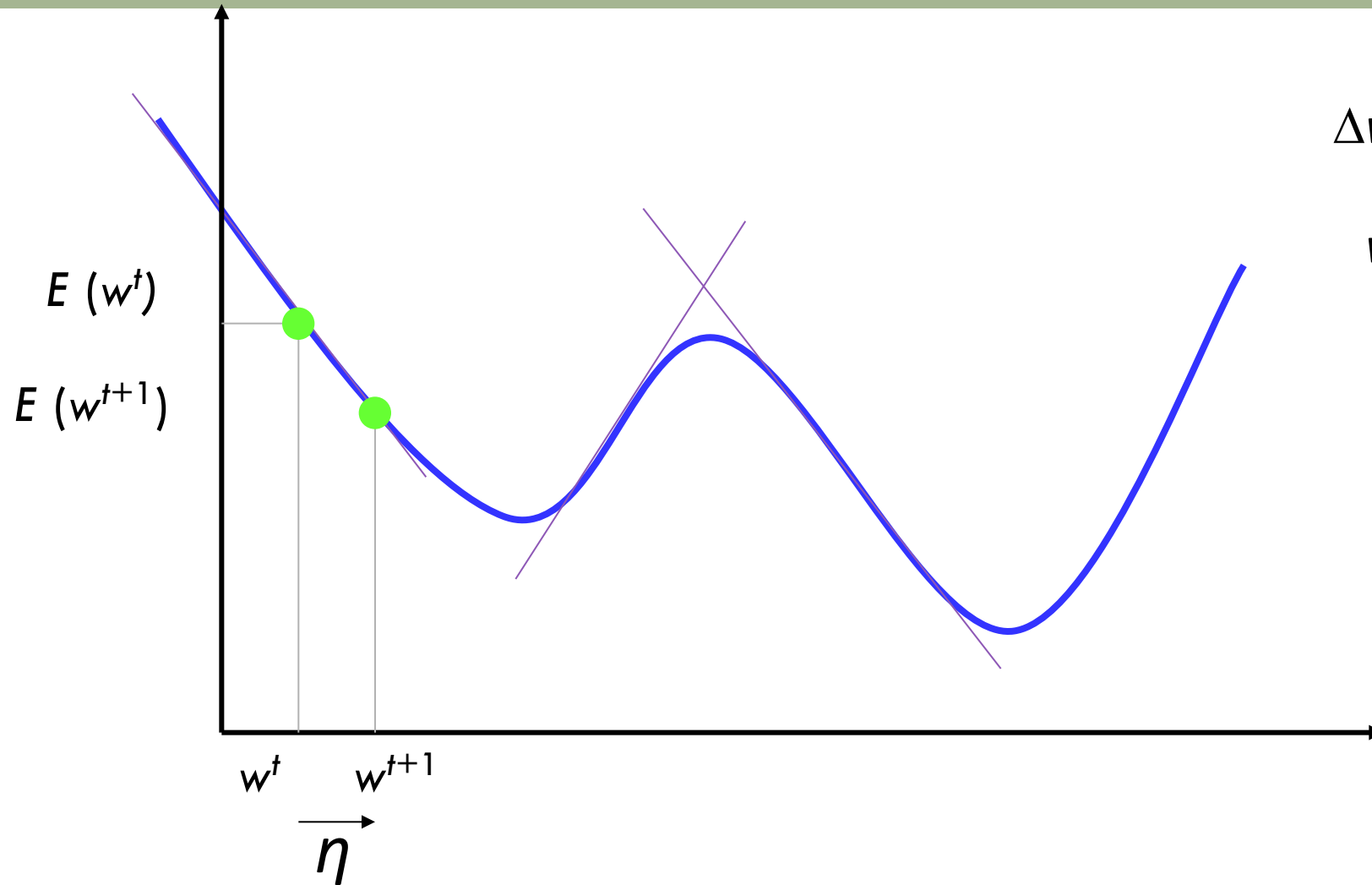
- Gradient
$$\nabla_{\mathbf{w}} E = \left[\frac{\partial E}{\partial w_1}, \frac{\partial E}{\partial w_2}, \dots, \frac{\partial E}{\partial w_d} \right]^T$$

- Gradient-descent:

Starts from random \mathbf{w} and updates \mathbf{w} iteratively in the negative direction of gradient (steepest descent)

Gradient-Descent

15



$$\Delta w_i = -\eta \frac{\partial E}{\partial w_i}, \forall i$$

$$w_i = w_i + \Delta w_i$$

Logistic Discrimination

16

Two classes: Assume log likelihood ratio is linear

$$\log \frac{p(\mathbf{x} | C_1)}{p(\mathbf{x} | C_2)} = \mathbf{w}^T \mathbf{x} + w_0^o$$

$$\begin{aligned} \text{logit}(P(C_1 | \mathbf{x})) &= \log \frac{P(C_1 | \mathbf{x})}{1 - P(C_1 | \mathbf{x})} = \log \frac{p(\mathbf{x} | C_1)}{p(\mathbf{x} | C_2)} + \log \frac{P(C_1)}{P(C_2)} \\ &= \mathbf{w}^T \mathbf{x} + w_0 \end{aligned}$$

$$\text{where } w_0 = w_0^o + \log \frac{P(C_1)}{P(C_2)}$$

$$y = \hat{P}(C_1 | \mathbf{x}) = \frac{1}{1 + \exp[-(\mathbf{w}^T \mathbf{x} + w_0)]}$$

Training: Two Classes

17

$$\mathcal{X} = \{\mathbf{x}^t, r^t\}_t \quad r^t \mid \mathbf{x}^t \sim \text{Bernoulli}(y^t)$$

$$y = P(C_1 \mid \mathbf{x}) = \frac{1}{1 + \exp\left[-(\mathbf{w}^T \mathbf{x} + w_0)\right]}$$

$$l(\mathbf{w}, w_0 \mid \mathcal{X}) = \prod_t (y^t)^{(r^t)} (1 - y^t)^{(1-r^t)}$$

$$E = -\log l$$

$$E(\mathbf{w}, w_0 \mid \mathcal{X}) = -\sum_t r^t \log y^t + (1 - r^t) \log (1 - y^t)$$

CROSS ENTROPY

- The optimal code for signal s_k with probability q_k has length $-\lg q_k$ bits
- The average length is the entropy $H(q) = E_q[-\lg q_k] = -\sum_k q_k \lg q_k$ bits
- But suppose the actual probabilities are p_k
- Then the average length is given by the cross-entropy:
 $H(p \parallel q) = E_p[-\lg q_k] = -\sum_k p_k \lg q_k$ bits
- The cross-entropy is minimized when the probability distributions are equal:
 $H(p \parallel q) = H(p) = H(q)$

Training: Gradient-Descent

19

$$E(\mathbf{w}, w_0 \mid \mathcal{X}) = -\sum_t r^t \log y^t + (1 - r^t) \log (1 - y^t)$$

$$\text{If } y = \text{sigmoid}(a) \quad \frac{dy}{da} = y(1 - y)$$

$$\begin{aligned} \Delta w_j &= -\eta \frac{\partial E}{\partial w_j} = \eta \sum_t \left(\frac{r^t}{y^t} - \frac{1 - r^t}{1 - y^t} \right) y^t (1 - y^t) x_j^t \\ &= \eta \sum_t (r^t - y^t) x_j^t, j = 1, \dots, d \end{aligned}$$

$$\Delta w_0 = -\eta \frac{\partial E}{\partial w_0} = \eta \sum_t (r^t - y^t)$$

$$\begin{aligned} \frac{\partial}{\partial w_j} r^t \log y^t &= \frac{r^t}{y^t} \frac{\partial y^t}{\partial w_j} \\ &= \frac{r^t}{y^t} \frac{\partial}{\partial w_j} \text{sigmoid}(a) \\ &= \frac{r^t}{y^t} y^t (1 - y^t) \frac{\partial a}{\partial w_j} \\ &= \frac{r^t}{y^t} y^t (1 - y^t) \frac{\partial (\mathbf{w}^T \mathbf{x}^t + w_0)}{\partial w_j} \\ &= \frac{r^t}{y^t} y^t (1 - y^t) \mathbf{x}_j^t \end{aligned}$$

(edited by BJM)

```
For  $j = 0, \dots, d$   
     $w_j \leftarrow \text{rand}(-0.01, 0.01)$   
Repeat  
    For  $j = 0, \dots, d$   
         $\Delta w_j \leftarrow 0$   
    For  $t = 1, \dots, N$   
         $o \leftarrow 0$   
        For  $j = 0, \dots, d$   
             $o \leftarrow o + w_j x_j^t$   
         $y \leftarrow \text{sigmoid}(o)$   
         $\Delta w_j \leftarrow \Delta w_j + (r^t - y)x_j^t$   
    For  $j = 0, \dots, d$   
         $w_j \leftarrow w_j + \eta \Delta w_j$   
Until convergence
```



$K > 2$ Classes

22

$$\mathcal{X} = \{\mathbf{x}^t, \mathbf{r}^t\}_t \quad r^t | \mathbf{x}^t \sim \text{Mult}_K(1, \mathbf{y}^t)$$

$$\log \frac{p(\mathbf{x} | C_i)}{p(\mathbf{x} | C_K)} = \mathbf{w}_i^T \mathbf{x} + w_{i0}$$

$$y = \hat{P}(C_i | \mathbf{x}) = \frac{\exp[\mathbf{w}_i^T \mathbf{x} + w_{i0}]}{\sum_{j=1}^K \exp[\mathbf{w}_j^T \mathbf{x} + w_{j0}]}, i = 1, \dots, K \quad \text{softmax}$$

$$l(\{\mathbf{w}_i, w_{i0}\}_i | \mathcal{X}) = \prod_t \prod_i (y_i^t)^{(r_i^t)}$$

$$E(\{\mathbf{w}_i, w_{i0}\}_i | \mathcal{X}) = - \sum_t r_i^t \log y_i^t \quad \text{cross-entropy}$$

$$\Delta \mathbf{w}_j = \eta \sum_t (r_j^t - y_j^t) \mathbf{x}^t \quad \Delta w_{j0} = \eta \sum_t (r_j^t - y_j^t)$$

```

For  $i = 1, \dots, K$ , For  $j = 0, \dots, d$ ,  $w_{ij} \leftarrow \text{rand}(-0.01, 0.01)$ 
Repeat
  For  $i = 1, \dots, K$ , For  $j = 0, \dots, d$ ,  $\Delta w_{ij} \leftarrow 0$ 
  For  $t = 1, \dots, N$ 
    For  $i = 1, \dots, K$ 
       $o_i \leftarrow 0$ 
      For  $j = 0, \dots, d$ 
         $o_i \leftarrow o_i + w_{ij} x_j^t$ 
      For  $i = 1, \dots, K$ 
         $y_i \leftarrow \exp(o_i) / \sum_k \exp(o_k)$ 
      For  $i = 1, \dots, K$ 
        For  $j = 0, \dots, d$ 
           $\Delta w_{ij} \leftarrow \Delta w_{ij} + (r_i^t - y_i) x_j^t$ 
    For  $i = 1, \dots, K$ 
      For  $j = 0, \dots, d$ 
         $w_{ij} \leftarrow w_{ij} + \eta \Delta w_{ij}$ 
  Until convergence

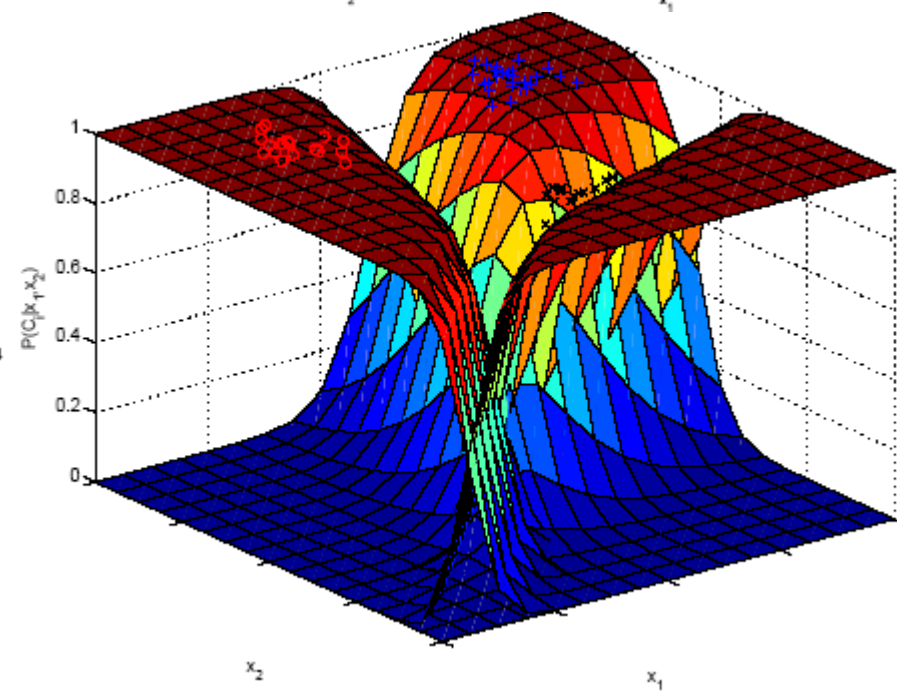
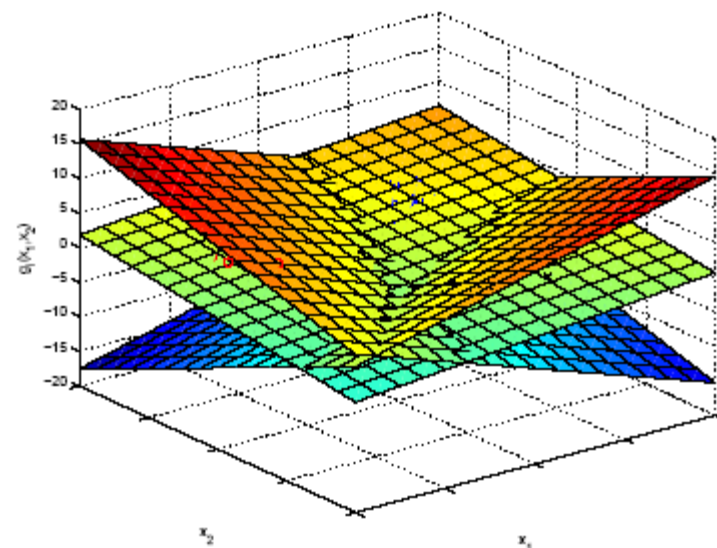
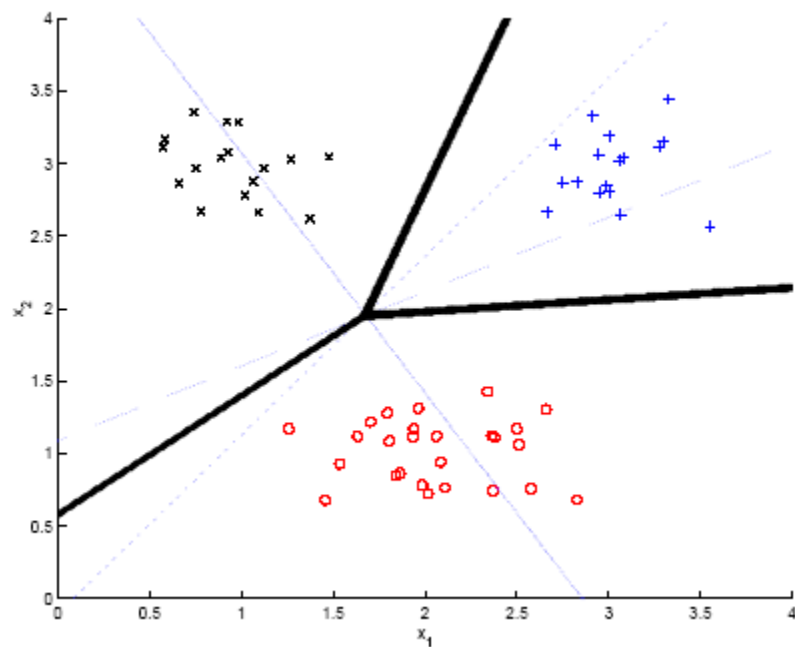
```

$\mathbf{W} \leftarrow \text{rand}_{K \times d}(-0.01, +0.01)$
 Repeat
 $\Delta \mathbf{W} \leftarrow \text{zeroes}(K, d)$
 For $t = 1, \dots, N$
 $\mathbf{o} \leftarrow \mathbf{W} \mathbf{x}^t$
 $\mathbf{z} \leftarrow \exp(\mathbf{o})$ (component-wise)
 $\mathbf{y} \leftarrow \mathbf{z} / \text{sum}(\mathbf{z})$
 $\Delta \mathbf{W} \leftarrow (\mathbf{r}^t - \mathbf{y})(\mathbf{x}^t)^T$
 $\mathbf{W} \leftarrow \mathbf{W} + \Delta \mathbf{W}$
 until convergence

Given \mathbf{X} ($N \times d$) and \mathbf{R} ($N \times K$)
 $\mathbf{W} \leftarrow \text{rand}_{K \times d}(-0.01, +0.01)$
 Repeat
 $\mathbf{O} \leftarrow \mathbf{X} \mathbf{W}^T$ ($N \times k$)
 $\mathbf{Z} \leftarrow \exp(\mathbf{O})$ (component-wise)
 $\mathbf{S} \leftarrow \mathbf{Z} \mathbf{1}_{N \times k}$ ($N \times k$)
 $\mathbf{Y} \leftarrow \mathbf{Z} / \mathbf{S}$ (component-wise)
 $\Delta \mathbf{W} \leftarrow (\mathbf{R} - \mathbf{Y})^T \mathbf{X}$ ($K \times d$)
 $\mathbf{W} \leftarrow \mathbf{W} + \Delta \mathbf{W}$
 until convergence

Example

25



Generalizing the Linear Model

26

- Quadratic:

$$\log \frac{p(\mathbf{x} | C_i)}{p(\mathbf{x} | C_K)} = \mathbf{x}^T \mathbf{W}_i \mathbf{x} + \mathbf{w}_i^T \mathbf{x} + w_{i0}$$

- Sum of basis functions:

$$\log \frac{p(\mathbf{x} | C_i)}{p(\mathbf{x} | C_K)} = \mathbf{w}_i^T \phi(\mathbf{x}) + w_{i0}$$

where $\phi(\mathbf{x})$ are basis functions. Examples:

- ▣ Hidden units in neural networks (Chapters 11 and 12)
- ▣ Kernels in SVM (Chapter 13)

Discrimination by Regression

27

- Classes are NOT mutually exclusive and exhaustive

$$r^t = y^t + \varepsilon \text{ where } \varepsilon \sim \mathcal{N}(0, \sigma^2)$$

$$y^t = \text{sigmoid}(\mathbf{w}^T \mathbf{x}^t + w_0) = \frac{1}{1 + \exp[-(\mathbf{w}^T \mathbf{x}^t + w_0)]}$$

$$l(\mathbf{w}, w_0 | \mathcal{X}) = \prod_t \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left[-\frac{(r^t - y^t)^2}{2\sigma^2}\right]$$

$$E(\mathbf{w}, w_0 | \mathcal{X}) = \frac{1}{2} \sum_t (r^t - y^t)^2$$

$$\Delta \mathbf{w} = \eta \sum_t (r^t - y^t) y^t (1 - y^t) \mathbf{x}^t$$

Learning to Rank

28

- Ranking: A different problem than classification or regression
- Let us say \mathbf{x}^u and \mathbf{x}^v are two instances, e.g., two movies

We prefer u to v implies that $g(\mathbf{x}^u) > g(\mathbf{x}^v)$

where $g(\mathbf{x})$ is a score function, here linear:

$$g(\mathbf{x}) = \mathbf{w}^\top \mathbf{x}$$

- Find a direction \mathbf{w} such that we get the desired ranks when instances are projected along \mathbf{w}

Ranking Error

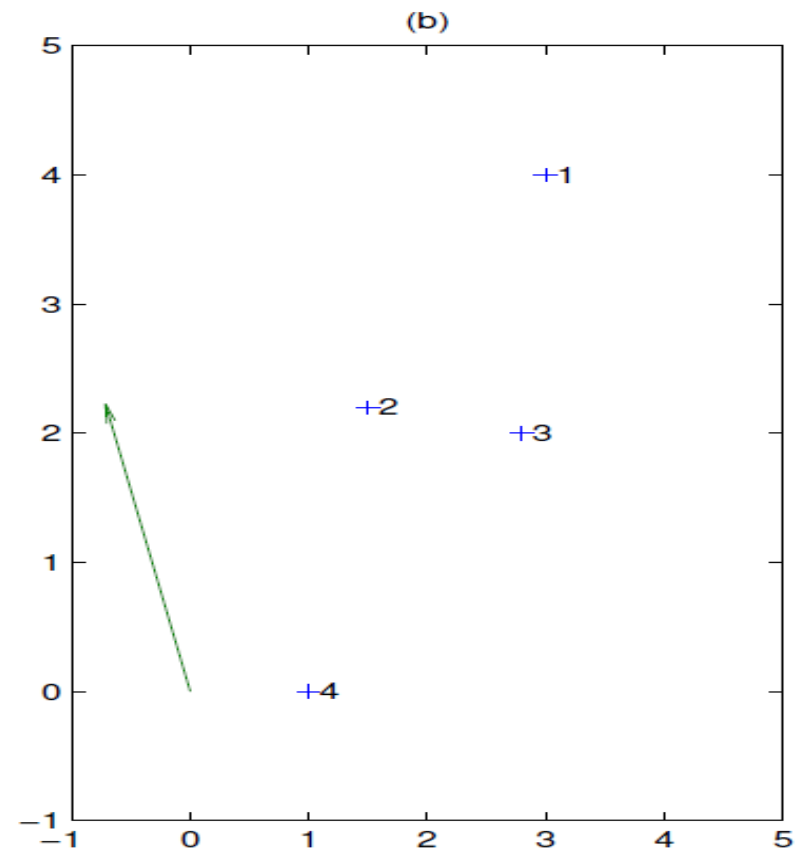
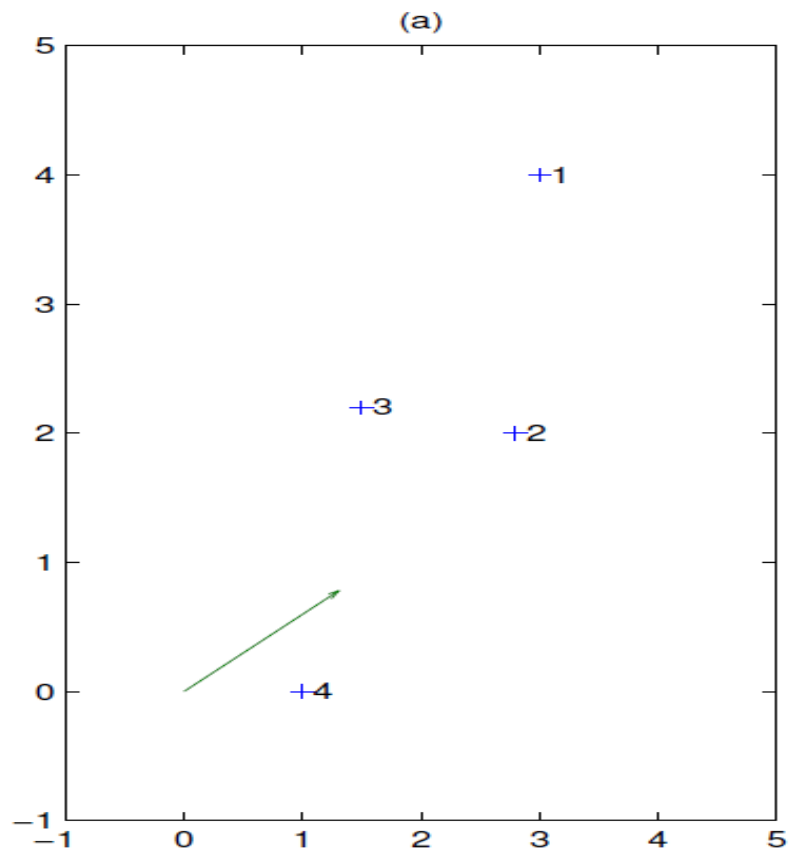
29

- We prefer u to v implies that $g(\mathbf{x}^u) > g(\mathbf{x}^v)$, so error is $g(\mathbf{x}^v) - g(\mathbf{x}^u)$, if $g(\mathbf{x}^u) < g(\mathbf{x}^v)$

$$E(\mathbf{w}|\{r^u, r^v\}) = \sum_{r^u < r^v} [g(\mathbf{x}^v|\theta) - g(\mathbf{x}^u|\theta)]_+$$

where a_+ is equal to a if $a \geq 0$ and 0 otherwise.

- Linear case: $E(\mathbf{w}|\{r^u, r^v\}) = \sum_{r^u < r^v} [\mathbf{w}^T(\mathbf{x}^v - \mathbf{x}^u)]_+$
- Gradient descent update when ranked wrong:
 $\Delta w_j = -\eta(x_j^v - x_j^u), j = 1, \dots, d$



READ ALPAYDIN CH. 11

