

FUNDAMENTALS OF MACHINE LEARNING

LINEAR DISCRIMINATION

INTRODUCTION

- Discuss *linear discrimination* methods
- We assume instances of a class are *linearly separable* from instances of other classes
 - ▣ Give example in class of *linearly separable classes*
 - ▣ Give example in class of *linearly non-separable classes*
- A discriminant-based method that **estimates** parameters of the linear discriminant directly from a **given label samples**

Likelihood- vs. Discriminant-based Classification

- In previous chapters, we define a set of discriminant functions $g_j(\mathbf{x}), j = 1, \dots, K$, then we

choose C_i if $g_i(\mathbf{x}) = \max_{j=1,\dots,K} g_j(\mathbf{x})$

- **Likelihood-based method:**

- estimate prior probabilities for each class: $\hat{P}(C_i)$
- estimate the class likelihoods: $\hat{p}(\mathbf{x}|C_i)$
- use Bayes' rule to calculate the posterior densities $P(C_i|\mathbf{x})$

- Then define the discriminant functions as

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

- Used this method in parametric, semi-parametric and nonparametric approaches.

Likelihood- vs. Discriminant-based Classification

- **Discriminant-based:** Assume a model for
$$g_i(x|\Phi_i)$$

no need for density estimation
- Estimating the **boundaries is enough**; no need to accurately estimate the densities inside the boundaries
- Instead of making assumption on the form of the class densities (e.g., Gaussian), it makes assumption on the **form of the boundaries** separating classes
- In discriminant-based classification, it is about **optimization of the model parameters** Φ_i to maximize the **quality of the separation** of classes

Linear Discriminant

- **Linear discriminant for class i :**

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

- **Advantages:**

- **Simple:** $O(d)$ space/computation
- **Knowledge extraction:**
 - Weighted sum of attributes;
 - Magnitude of w_{ij} shows the importance of x_j for the class i discriminant
 - positive (enforcing) or negative (inhibiting) weights (e.g., credit score)
- Linear discriminant is optimal when the likelihood $p(\mathbf{x}|C_i)$ are Gaussian with shared covariance matrix
- Useful when classes are (**almost**) linearly separable

Generalized Linear Model

- **Quadratic discriminant:** $g_i(\mathbf{x}|\mathbf{W}_i, \mathbf{w}_i, w_{i0}) = \mathbf{x}^T \mathbf{W}_i \mathbf{x} + \mathbf{w}_i^T \mathbf{x} + w_{i0}$
 - Complexity is $O(d^2)$, and we have **bias/variance dilemma**
 - Require much larger training sets and may over-fit on small samples

For example, let $d = 2$, $\mathbf{W}_i = \begin{bmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{bmatrix}$, we have

$$\begin{aligned} g_i() &= \mathbf{x}^T \mathbf{W}_i \mathbf{x} + \mathbf{w}_i^T \mathbf{x} + w_{i0} \\ &= [x_1, x_2] \begin{bmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} + \begin{bmatrix} w_1 \\ w_2 \end{bmatrix}^T \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} + w_{i0} \\ &= a_{11}x_1^2 + (a_{12} + a_{21})x_1x_2 + a_{22}x_2^2 + w_1x_1 + w_2x_2 + w_{i0} \end{aligned}$$

- Adding higher-order (product) terms. E.g., two inputs x_1 and x_2

$$Z_1 = x_1, Z_2 = x_2, Z_3 = x_1^2, Z_4 = x_2^2, Z_5 = x_1 x_2$$

take $\mathbf{z} = [Z_1, Z_2, Z_3, Z_4, Z_5]^T$ as the input

Map from \mathbf{x} to \mathbf{z} using **nonlinear basis functions** and use a linear discriminant in \mathbf{z} -space $g_i(\mathbf{x}) = \sum_{j=1} w_j \phi_{ij}(\mathbf{x})$ where $\phi_{ij}(\mathbf{x})$ are *basis functions*

Generalized Linear Model

- Higher-order terms are only one set of possible basis functions
- Other possible basis functions:
 - $\sin(x_1)$
 - $\exp(-(x_1 - m)^2/c)$
 - $\exp(-\|x - \mathbf{m}\|^2/c)$
 - $\log(x_2)$
 - $1(x_1 > c)$
 - $1(ax_1 + bx_2 > c)$

where m, a, b, c are scalars and \mathbf{m} is a d -dimensional vector
and $1(b)$ is an indicating function

Geometry of Linear Discriminant: Two Classes

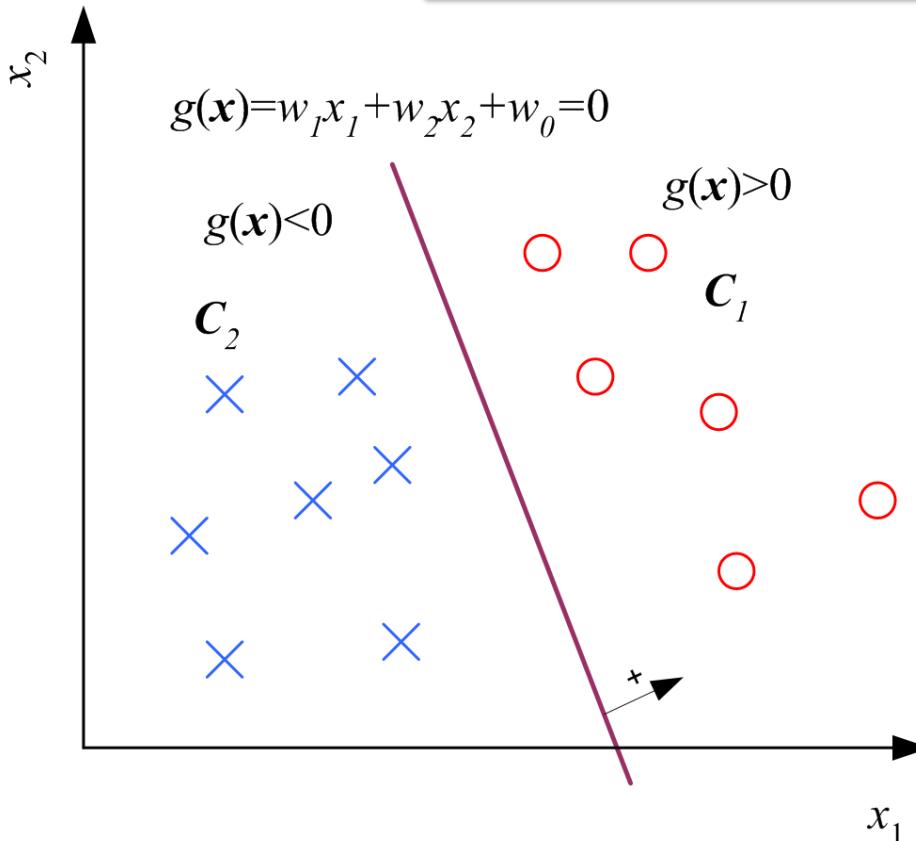
When there are 2 classes, **one** discriminant is enough:

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

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

\mathbf{w} is the *weight vector*

w_0 is the *threshold*.



Geometrical Interpretation for 2 Classes

Take two points \mathbf{x}_1 and \mathbf{x}_2 both on the decision surface

$$g(\mathbf{x}_1) = g(\mathbf{x}_2) = 0$$

$$\mathbf{w}^T \mathbf{x}_1 + w_0 = \mathbf{w}^T \mathbf{x}_2 + w_0$$

$$\mathbf{w}^T (\mathbf{x}_1 - \mathbf{x}_2) = 0$$

\mathbf{w} is normal to any vector

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

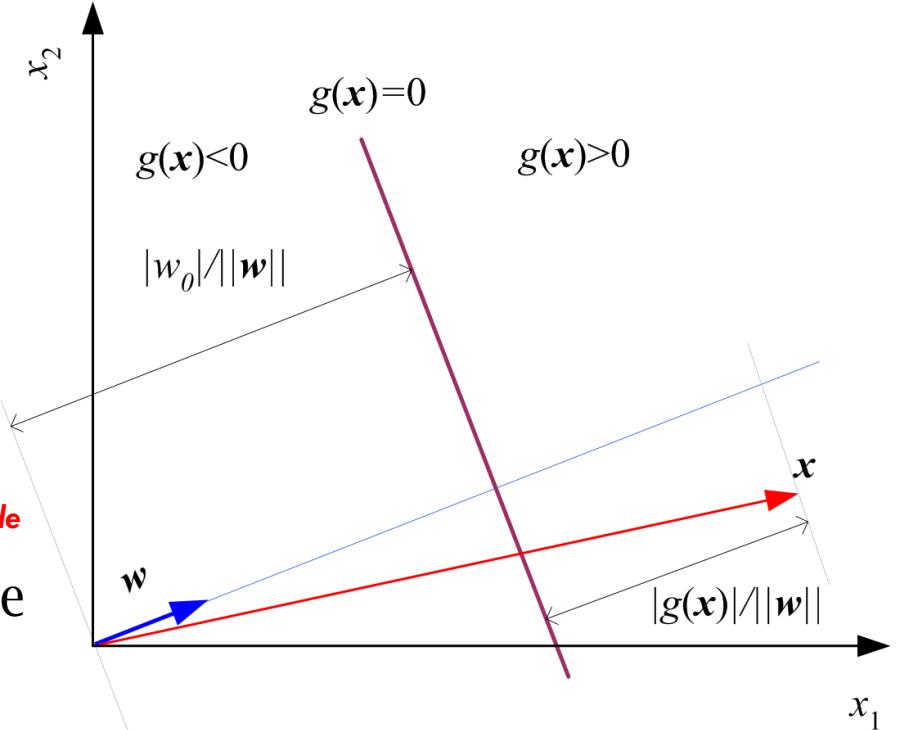
\mathbf{x}_p is the normal projection of \mathbf{x}

onto the hyperplane

r gives the distance from \mathbf{x} to the hyperplane
since $g(\mathbf{x}_p) = 0$.

See *derivation* in the next page:

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



Apply it to the origin $\mathbf{x} = 0$ and
following similar derivation:

$$r_0 = \frac{w_0}{\|\mathbf{w}\|}$$

Geometrical Interpretation for 2 Classes (derivation)

We have $g(\mathbf{x}_p) = \mathbf{w}^T \mathbf{x}_p + w_0 = 0$.

Since $\mathbf{x} = \mathbf{x}_p + \frac{r}{\|\mathbf{w}\|} \mathbf{w}$, or $\mathbf{x}_p = \mathbf{x} - \frac{r}{\|\mathbf{w}\|} \mathbf{w}$.

Apply \mathbf{x}_p onto the $g()$, we have:

$$\mathbf{w}^T \left(\mathbf{x} - \frac{r}{\|\mathbf{w}\|} \mathbf{w} \right) + w_0 = 0$$

$$\mathbf{w}^T \mathbf{x} - \frac{r}{\|\mathbf{w}\|} \mathbf{w}^T \mathbf{w} + w_0 = 0$$

$$\mathbf{w}^T \mathbf{x} - \frac{r}{\|\mathbf{w}\|} \|\mathbf{w}\|^2 + w_0 = 0$$

$$\mathbf{w}^T \mathbf{x} + w_0 = r \|\mathbf{w}\|$$

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

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

Multiple Classes ($K > 2$)

there are K discriminant functions

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

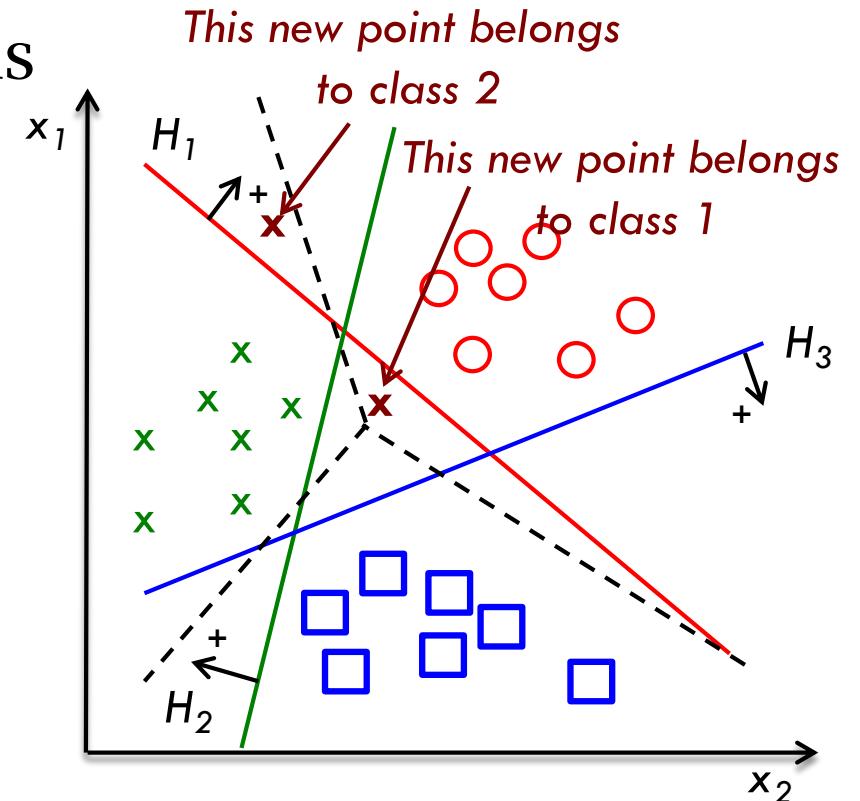
where

$$g_i(\mathbf{x} | \mathbf{w}_i, w_{i0}) = \begin{cases} > 0 & \text{if } \mathbf{x} \in C_i \\ \leq 0 & \text{otherwise} \end{cases}$$

Classes are
linearly separable

To relax “linear separability”

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



$|g_i(\mathbf{x})| / \|\mathbf{w}_i\|$ is the distance from the input point to the hyperplane i . Assume \mathbf{w}_i has the same length. It assigns the point to the class with the largest distance.

Pairwise Separation

- If the classes are not linearly separable, we divide it into a set of linear problems, e.g., *pairwise separation*

- Let K be # of classes. We use $K(K - 1)/2$ linear discriminants, $g_{ij}(\mathbf{x})$, one for every pair of distinct classes, where

$$g_{ij}(\mathbf{x} | \mathbf{w}_{ij}, w_{ij0}) = \mathbf{w}_{ij}^T \mathbf{x} + w_{ij0}$$

- We compute the parameters, \mathbf{w}_{ij} , during the training

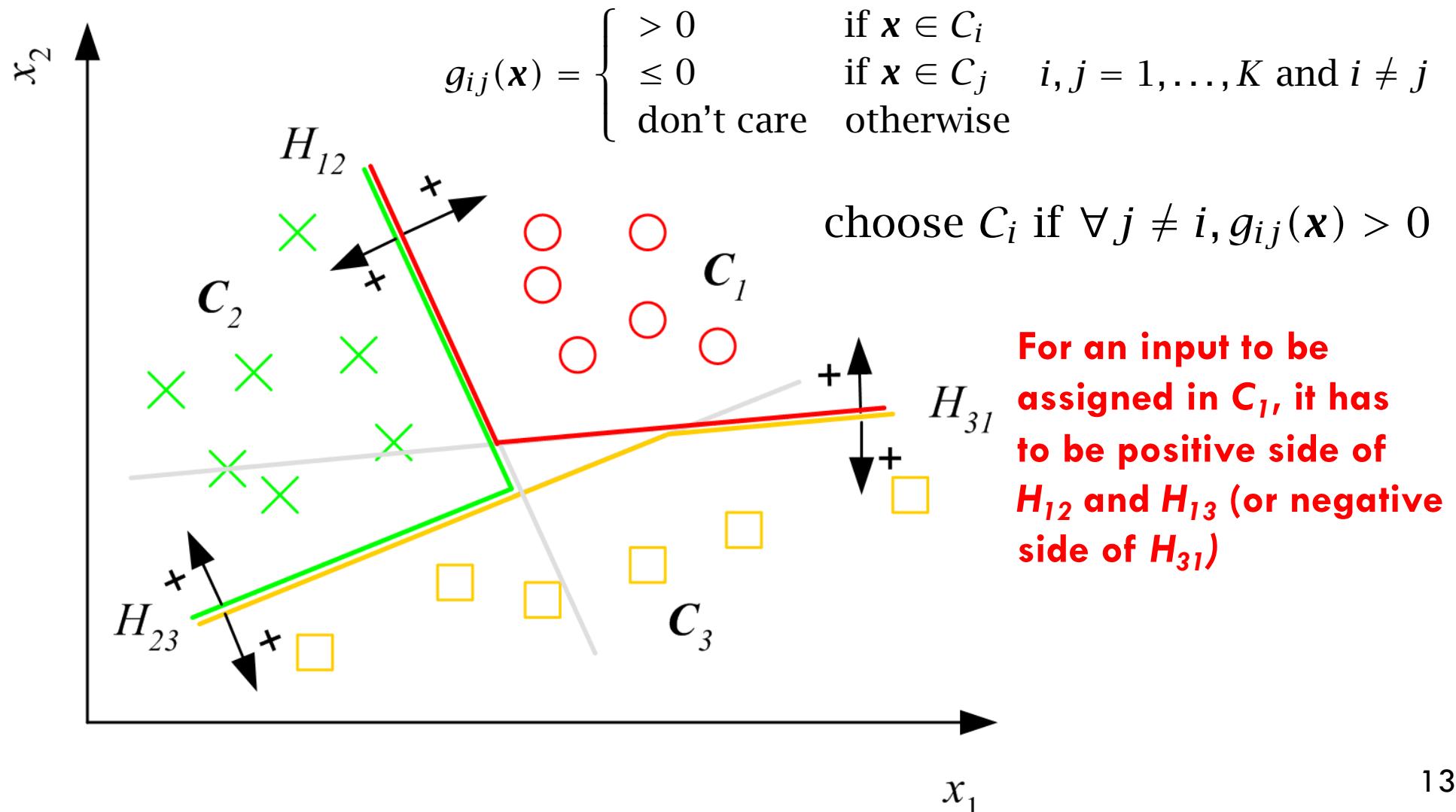
$$g_{ij}(\mathbf{x}) = \begin{cases} > 0 & \text{if } \mathbf{x} \in C_i \\ \leq 0 & \text{if } \mathbf{x} \in C_j \quad i, j = 1, \dots, K \text{ and } i \neq j \\ \text{don't care} & \text{otherwise} \end{cases}$$

If $\mathbf{x}^t \in C_k$ and $k \neq i, k \neq j$, then \mathbf{x}^t is not used in training of $g_{ij}(\mathbf{x})$

- During test, we choose C_i if $\forall j \neq i, g_{ij}(\mathbf{x}) > 0$

Pairwise Separation

$$g_{ij}(\mathbf{x} | \mathbf{w}_{ij}, w_{ij0}) = \mathbf{w}_{ij}^T \mathbf{x} + w_{ij0}$$



Parametric Discrimination Revisited

In parametric ML, if the class densities, $p(\mathbf{x}|C_i)$, are Gaussian and **share a common covariance matrix** or $p(\mathbf{x}|C_i) \sim N(\boldsymbol{\mu}_i, \Sigma)$, we have:

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

$$\mathbf{w}_i = \Sigma^{-1} \boldsymbol{\mu}_i$$

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

Given the data, we “estimate” $\boldsymbol{\mu}_i$ and Σ via \mathbf{m}_i and \mathbf{S} , then use the last two equations to get \mathbf{w}_i and w_{i0}

Parametric Discrimination Revisited

Consider a special case of $K=2$ classes. We define

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

In classification, we have

$$\text{choose } C_1 \text{ if } \begin{cases} y > 0.5 \\ \frac{y}{1-y} > 1 \\ \log \frac{y}{1-y} > 0 \end{cases} \text{ and } C_2 \text{ otherwise}$$

$\log y / (1 - y)$ is known as the *logit* transformation or *log odds* of y

Note that $\log \frac{y}{1-y}$ goes from $-\infty$ to $+\infty$.

Assume two normal classes share a common covariance matrix

$$\begin{aligned}
 \text{logit}(P(C_1|\mathbf{x})) &= \log \frac{P(C_1|\mathbf{x})}{1 - P(C_1|\mathbf{x})} = \log \frac{P(C_1|\mathbf{x})}{P(C_2|\mathbf{x})} \\
 &= \log \frac{p(\mathbf{x}|C_1)}{p(\mathbf{x}|C_2)} + \log \frac{P(C_1)}{P(C_2)} \\
 &= \log \frac{(2\pi)^{-d/2} |\Sigma|^{-1/2} \exp[-(1/2)(\mathbf{x} - \boldsymbol{\mu}_1)^T \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu}_1)]}{(2\pi)^{-d/2} |\Sigma|^{-1/2} \exp[-(1/2)(\mathbf{x} - \boldsymbol{\mu}_2)^T \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu}_2)]} + \log \frac{P(C_1)}{P(C_2)} \\
 &= \mathbf{w}^T \mathbf{x} + w_0
 \end{aligned}$$

where

$$\mathbf{w} = \Sigma^{-1}(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)$$

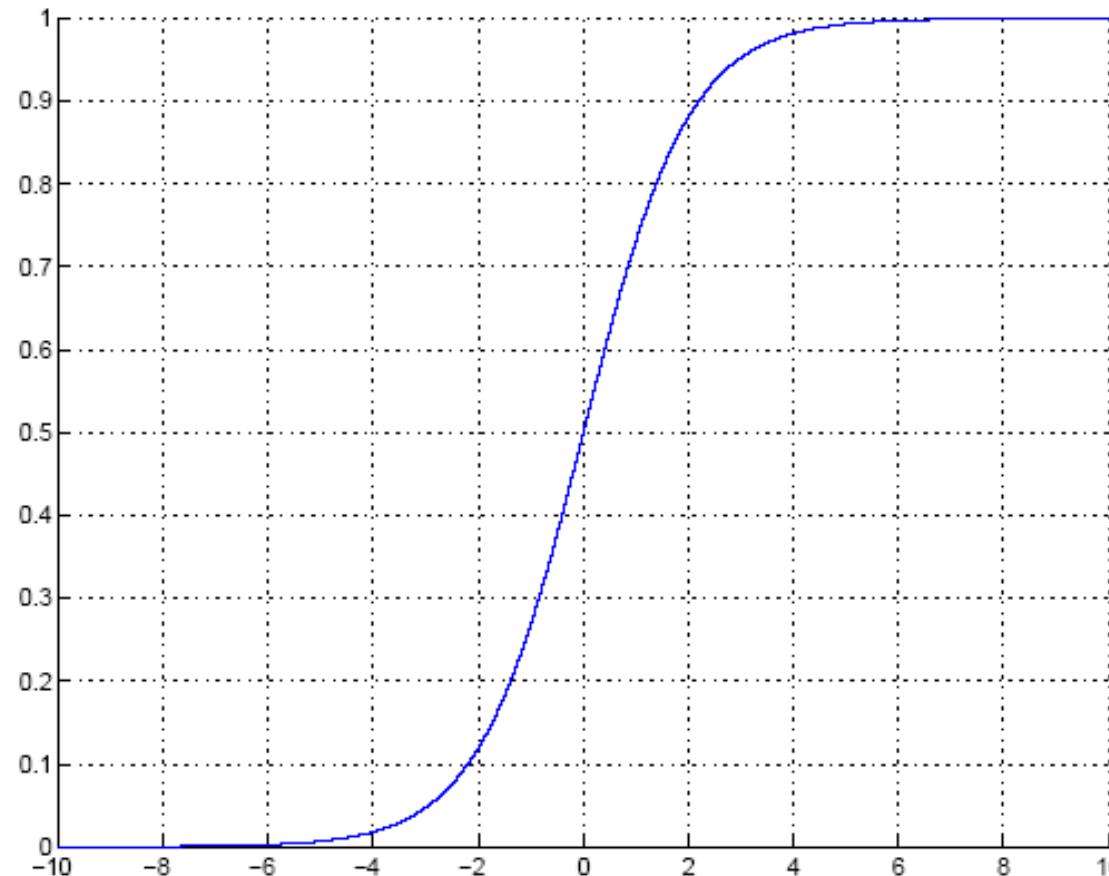
$$w_0 = -\frac{1}{2}(\boldsymbol{\mu}_1 + \boldsymbol{\mu}_2)^T \Sigma^{-1}(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2) + \log \frac{P(C_1)}{P(C_2)}$$

1. estimate $\boldsymbol{\mu}_1, \boldsymbol{\mu}_2, \Sigma$
2. calculate $g(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + w_0$ and choose C_1 if $g(\mathbf{x}) > 0$

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

Show in page 21 how we derive this from logit !!!!

Sigmoid (Logistic) Function



1. calculate $g(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + w_0$ and choose C_1 if $g(\mathbf{x}) > 0$, or
2. calculate $y = \text{sigmoid}(\mathbf{w}^T \mathbf{x} + w_0)$ and choose C_1 if $y > 0.5$,

Will show later how to handle $K > 2$

Gradient-Descent

- In discriminant-based approach, we **optimized** the parameters of the discriminants to reduce classification error
- Define $E(\mathbf{w}|X)$ as error with parameters \mathbf{w} on sample X

$$\mathbf{w}^* = \operatorname{argmin}_{\mathbf{w}} E(\mathbf{w}|X)$$

- No analytical solution, resort to *iterative method*
- **Gradient descent**

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

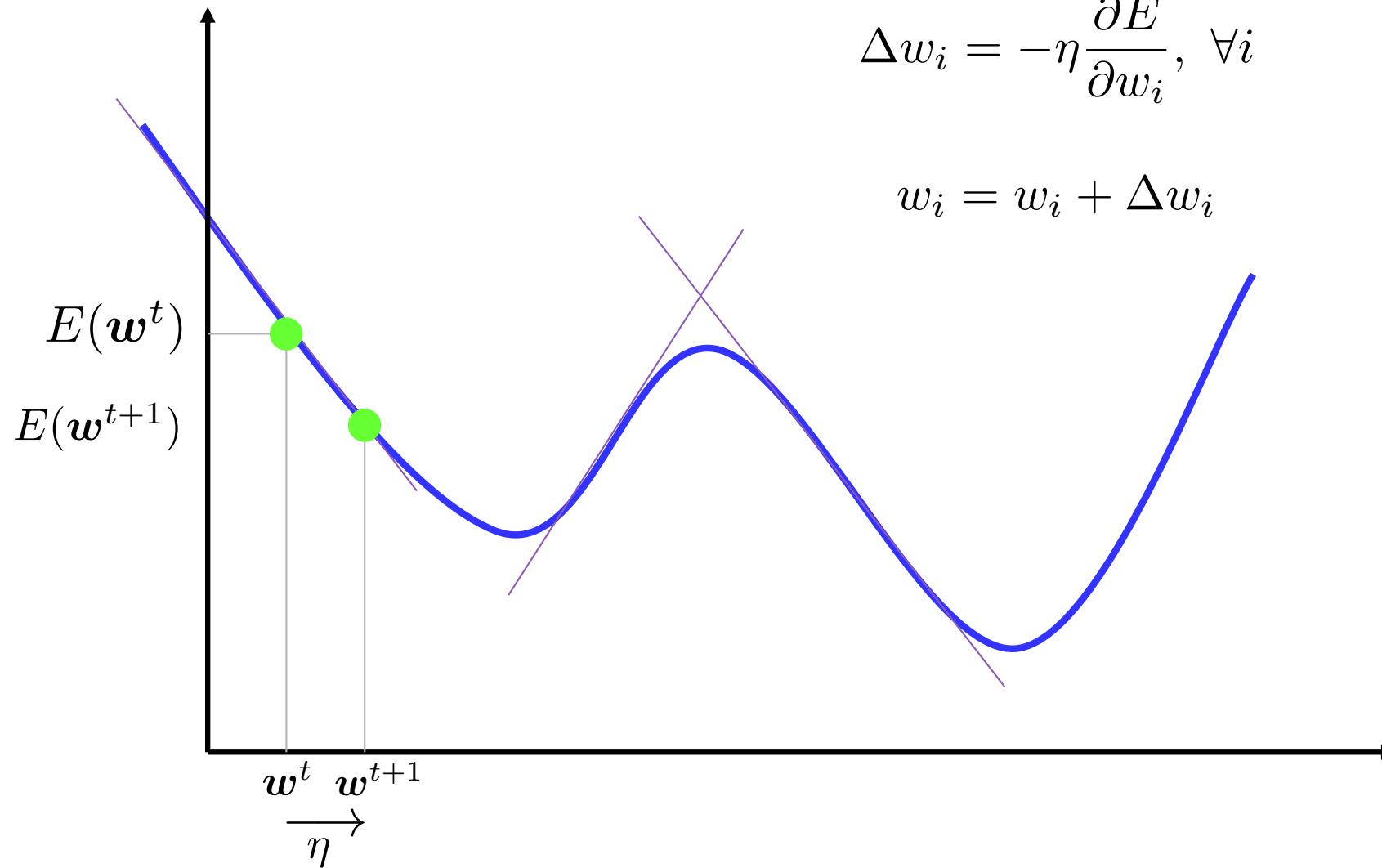
- **Gradient-descent:** Starts from random \mathbf{w} , and updates \mathbf{w} iteratively in the **negative direction** of gradient

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

$$w_i = w_i + \Delta w_i$$

η is step-size, or learning factor

Gradient-Descent



Logistic Discrimination

- In logistic discrimination, we don't model class-conditional densities $p(\mathbf{x}|C_i)$, but instead their **ratio**.
- Assume 2 classes and assume linear log likelihood ratio

$$\log \frac{p(\mathbf{x}|C_1)}{p(\mathbf{x}|C_2)} = \mathbf{w}^T \mathbf{x} + w_0^o \quad (\text{true when class-conditional densities are normal})$$

- Use Bayes' rule, we have

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

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

- Rearranging terms, our estimator of $P(C_1|\mathbf{x})$

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

See next page for derivation

Derivation

We have:

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

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

How to learn w and w_0 for $K = 2$ Classes ?

- Input: sample of two classes:

$$X = \{ \mathbf{x}^t, r^t \}, \text{ where } r^t = 1 \text{ if } \mathbf{x}^t \in C_1 \text{ and } r^t = 0 \text{ if } \mathbf{x}^t \in C_2$$

- Assume r^t , given \mathbf{x}^t , is Bernoulli, or $r^t | \mathbf{x}^t \sim \text{Bernoulli}(y^t)$

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

- The sample likelihood is Use sequence of coin toss as example

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

- Maximize likelihood can be turned into minimize error function as $E = -\log l$

- In this case, we have cross-entropy:

$$E(\mathbf{w}, w_0 | \mathcal{X}) = - \left[\sum_t r^t \log(y^t) + (1 - r^t) \log(1 - y^t) \right]$$

Training: Gradient-Descent

- We use **gradient descent** to minimize cross-entropy (or maximize likelihood or log likelihood)
- Note that if

$$y = \text{sigmoid}(a) = 1/(1 + \exp(-a))$$

$$\frac{dy}{da} = y(1 - y)$$

and we get the following update equations:

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

See derivation

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

Derivation

$$y^t = \frac{1}{1 + e^{-(\mathbf{w}^T \mathbf{x} + w_0)}} \quad \text{and} \quad (1 - y^t) = \frac{e^{-(\mathbf{w}^T \mathbf{x} + w_0)}}{1 + e^{-(\mathbf{w}^T \mathbf{x} + w_0)}}$$

$$E = -\log l = -\sum_{t=1}^N [r^t \log y^t + (1 - r^t) \log(1 - y^t)]$$

Remember the following rules:

$\frac{d \log(x)}{dx} = \frac{1}{x} dx$	$\frac{de^{ax}}{dx} = ae^{ax} dx$
---	-----------------------------------

$$\begin{aligned} \frac{\partial \log(l)}{\partial w_j} &= \sum_{t=1}^N \left[\left(\frac{r^t}{y^t} \right) \left(\frac{x_j e^{-(\mathbf{w}^t \mathbf{x} + w_0)}}{(1 + e^{-(\mathbf{w}^t \mathbf{x} + w_0)})^2} \right) - \left(\frac{1 - r^t}{1 - y^t} \right) \left(\frac{x_j e^{-(\mathbf{w}^t \mathbf{x} + w_0)}}{(1 + e^{-(\mathbf{w}^t \mathbf{x} + w_0)})^2} \right) \right] \\ &= \sum_{t=1}^N \left[\left(\frac{r^t}{y^t} \right) x_j y^t (1 - y^t) - \left(\frac{1 - r^t}{1 - y^t} \right) x_j y^t (1 - y^t) \right] \\ &= \sum_{t=1}^N \left[\left(\frac{r^t}{y^t} - \frac{1 - r^t}{1 - y^t} \right) x_j y^t (1 - y^t) \right] \quad \text{What about “-”?} \\ &= \sum_{t=1}^N \left[\frac{(r^t - y^t)}{y^t(1 - y^t)} y^t (1 - y^t) x_j \right] = \sum_{i=1}^N [(r^t - y^t) x_j] \quad \text{Because } E \text{ is defined with “-”, so the two “negatives” cancel out.} \end{aligned}$$

```

For  $j = 0, \dots, d$ 
     $w_j \leftarrow \text{rand}(-0.01, 0.01)$ 

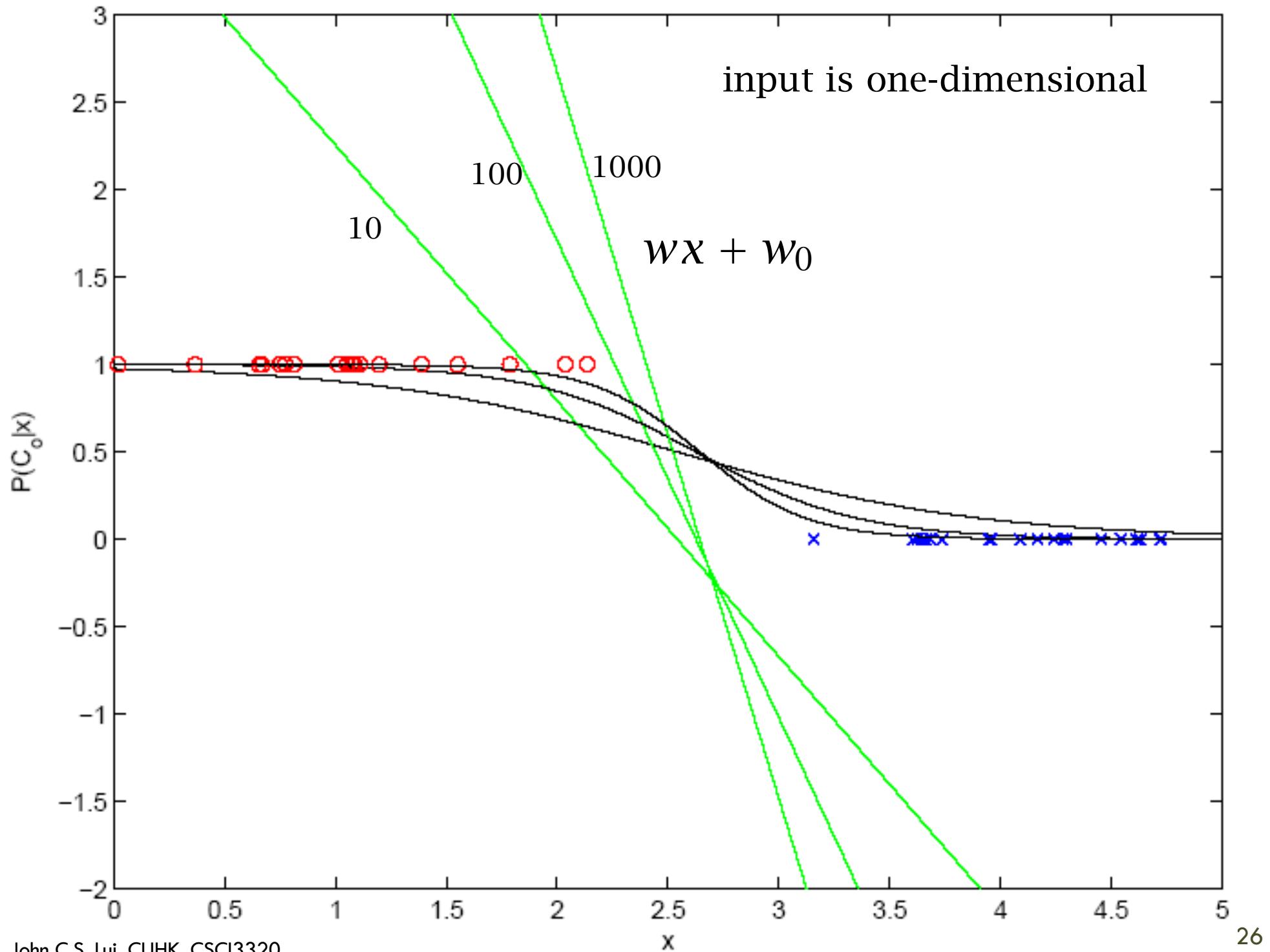
Repeat
    For  $j = 0, \dots, d$ 
         $\Delta w_j \leftarrow 0$ 
    For  $t = 1, \dots, N$ 
         $o \leftarrow 0$ 
        For  $j = 0, \dots, d$ 
             $o \leftarrow o + w_j x_j^t$ 
         $y \leftarrow \text{sigmoid}(o)$ 
        For  $j = 0, \dots, d$ 
             $\Delta w_j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

```

Why?

Note that we are looping over N training points

input is one-dimensional



$K > 2$ Classes

Input: $X = \{\mathbf{x}^t, \mathbf{r}^t\}_t$ $\mathbf{r}^t \mid \mathbf{x}^t \sim \text{Multi}_k(1, \mathbf{y}^t)$, where $y_i^t = P(C_i \mid \mathbf{x}^t)$

Take class C_K as the reference class and assume

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

Then we have

$$\frac{P(C_i \mid \mathbf{x})}{P(C_K \mid \mathbf{x})} = \exp[\mathbf{w}_i^T \mathbf{x} + w_{i0}]$$

with $w_{i0} = w_{i0}^o + \log P(C_i) / P(C_K)$

$$\begin{aligned} \sum_{i=1}^{K-1} \frac{P(C_i \mid \mathbf{x})}{P(C_K \mid \mathbf{x})} &= \frac{1 - P(C_K \mid \mathbf{x})}{P(C_K \mid \mathbf{x})} = \sum_{i=1}^{K-1} \exp[\mathbf{w}_i^T \mathbf{x} + w_{i0}] \\ \Rightarrow P(C_K \mid \mathbf{x}) &= \frac{1}{1 + \sum_{i=1}^{K-1} \exp[\mathbf{w}_i^T \mathbf{x} + w_{i0}]} \end{aligned}$$

$K > 2$ classes

- Since $\frac{P(C_i|x)}{P(C_K|x)} = \exp[\mathbf{w}_i^T x + w_{i0}]$, and so

$$P(C_i|x) = \frac{\exp[\mathbf{w}_i^T x + w_{i0}]}{1 + \sum_{j=1}^{K-1} \exp[\mathbf{w}_j^T x + w_{j0}]}, \quad i = 1, 2, \dots, K-1$$

- Treat all classes uniformly, we have the “**softmax**”

$$y_i = \hat{P}(C_i|x) = \frac{\exp[\mathbf{w}_i^T x + w_{i0}]}{\sum_{j=1}^K \exp[\mathbf{w}_j^T x + w_{j0}]}, \quad i = 1, 2, \dots, K$$

$K > 2$ classes: Softmax to Sigmoid

- **Softmax** can be viewed as an extension of the Sigmoid function for $K=2$ classes

$$y_i = \hat{P}(C_i | \mathbf{x}) = \frac{\exp[\mathbf{w}_i^T \mathbf{x} + w_{i0}]}{\sum_{j=1}^K \exp[\mathbf{w}_j^T \mathbf{x} + w_{j0}]}, \quad i = 1, 2, \dots, K$$

- Consider $K=2$, define $a_i = \mathbf{w}_i^T \mathbf{x} + w_{i0}$, $i = 1, 2$.

$$\square y_1 = \frac{e^{a_1}}{e^{a_1} + e^{a_2}} = \frac{1}{1 + e^{a_2 - a_1}} = \frac{1}{1 + e^{-(a_1 - a_2)}} = \frac{1}{1 + e^{-a}}$$

$$\square y_2 = \frac{e^{a_2}}{e^{a_1} + e^{a_2}} = \frac{e^{a_2 - a_1}}{1 + e^{a_2 - a_1}} = \frac{e^{-(a_1 - a_2)}}{1 + e^{-(a_1 - a_2)}} = \frac{e^{-a}}{1 + e^{-a}}$$

- Where $a = \mathbf{w}^T \mathbf{x} + w_0$

- This is exactly the result we have for $K=2$

- How to learn the parameters via the sample likelihood?

- Learn the parameters via the sample likelihood:

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

$K > 2$ Classes: how to learn the parameters

The sample likelihood is

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

the error function is again cross-entropy:

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

use gradient descent. If $y_i = \exp(a_i) / \sum_j \exp(a_j)$

$$\frac{\partial y_i}{\partial a_j} = y_i (\delta_{ij} - y_j) \quad \begin{aligned} \delta_{ij} \text{ is the Kronecker delta} \\ \delta_{ij} = 1 \text{ if } i = j \text{ and 0 otherwise} \end{aligned}$$

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

$$\Delta w_{j0} = \eta \sum_t (r_j^t - y_j^t)$$

```
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

Example

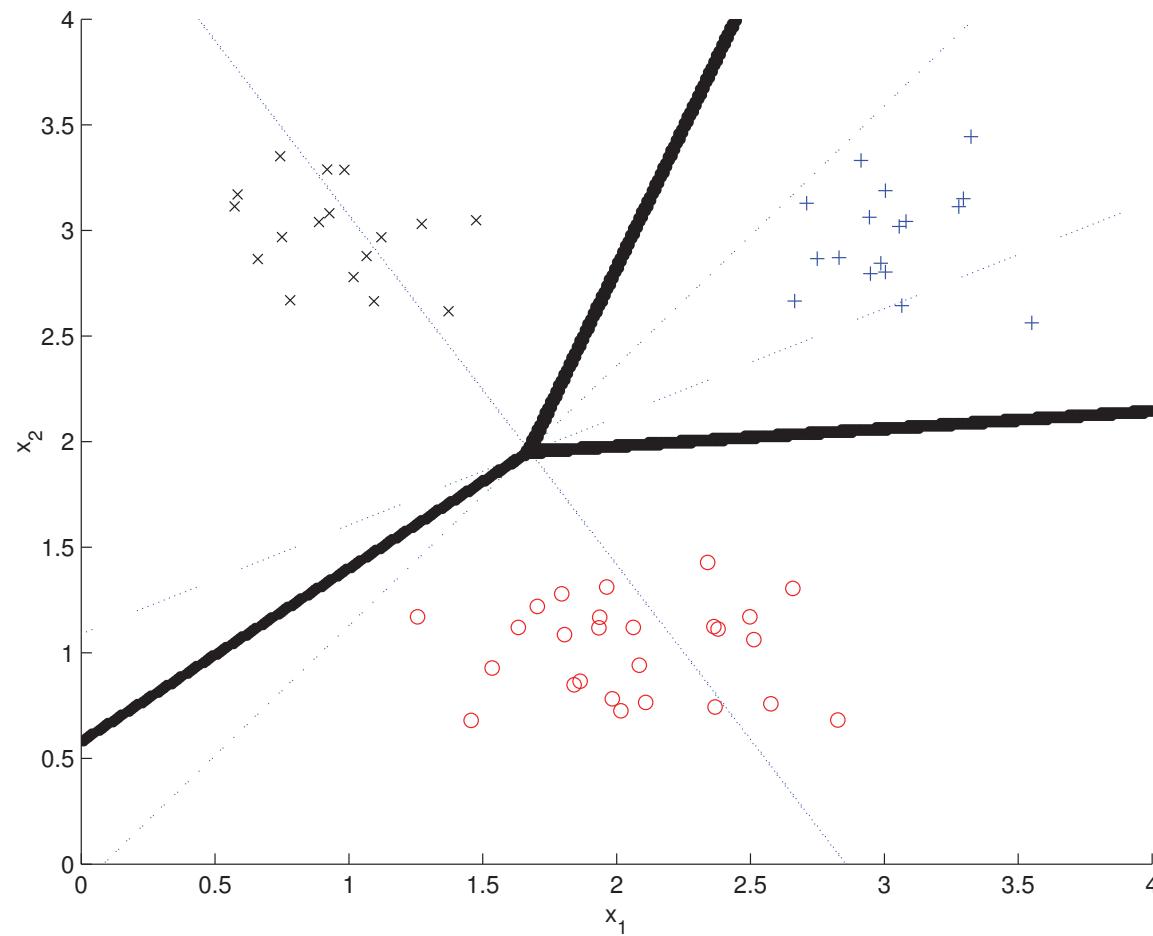
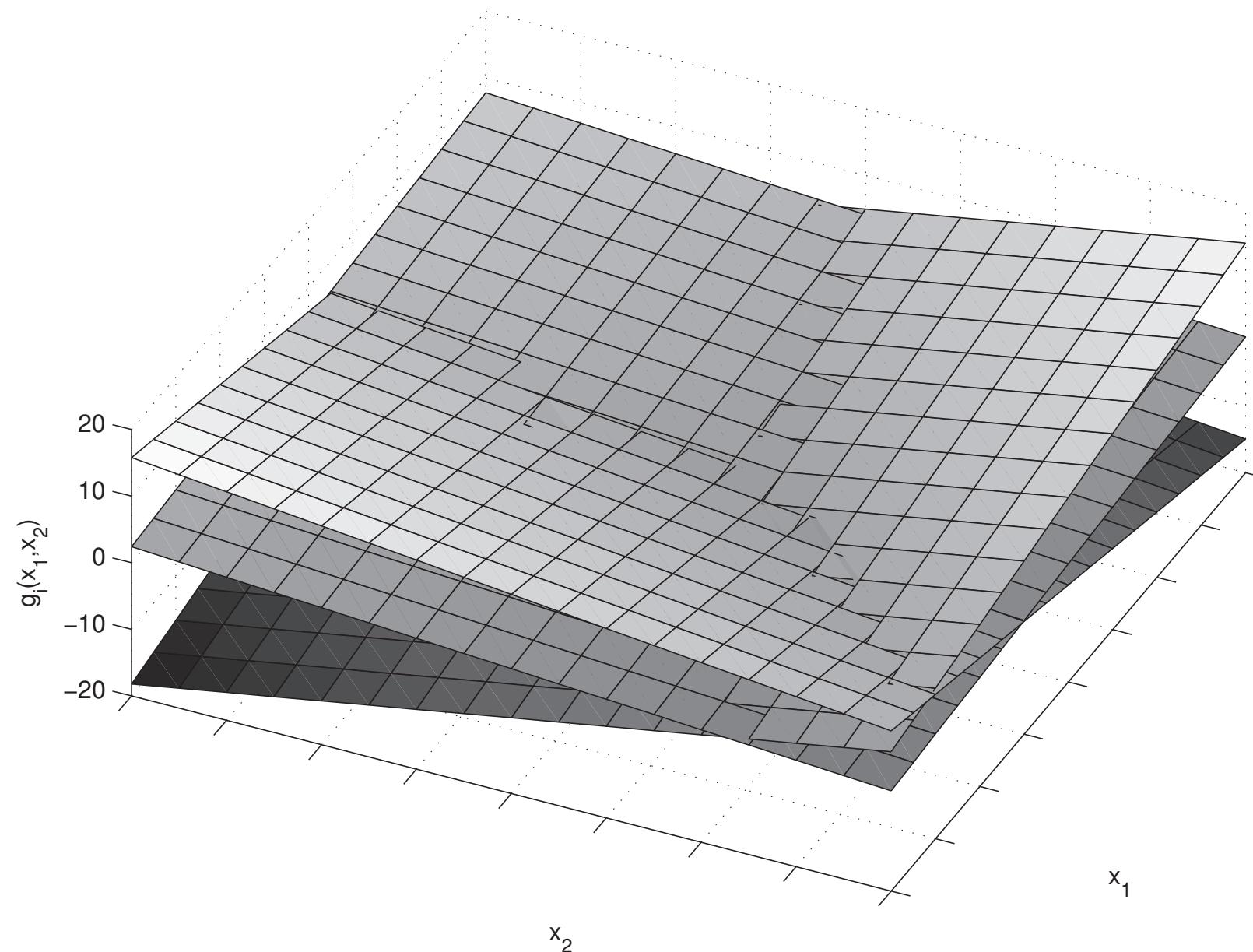


Figure 10.9 For a two-dimensional problem with three classes, the solution found by logistic discrimination. Thin lines are where $g_i(\mathbf{x}) = 0$, and the thick line is the boundary induced by the linear classifier choosing the maximum.

Example



Example

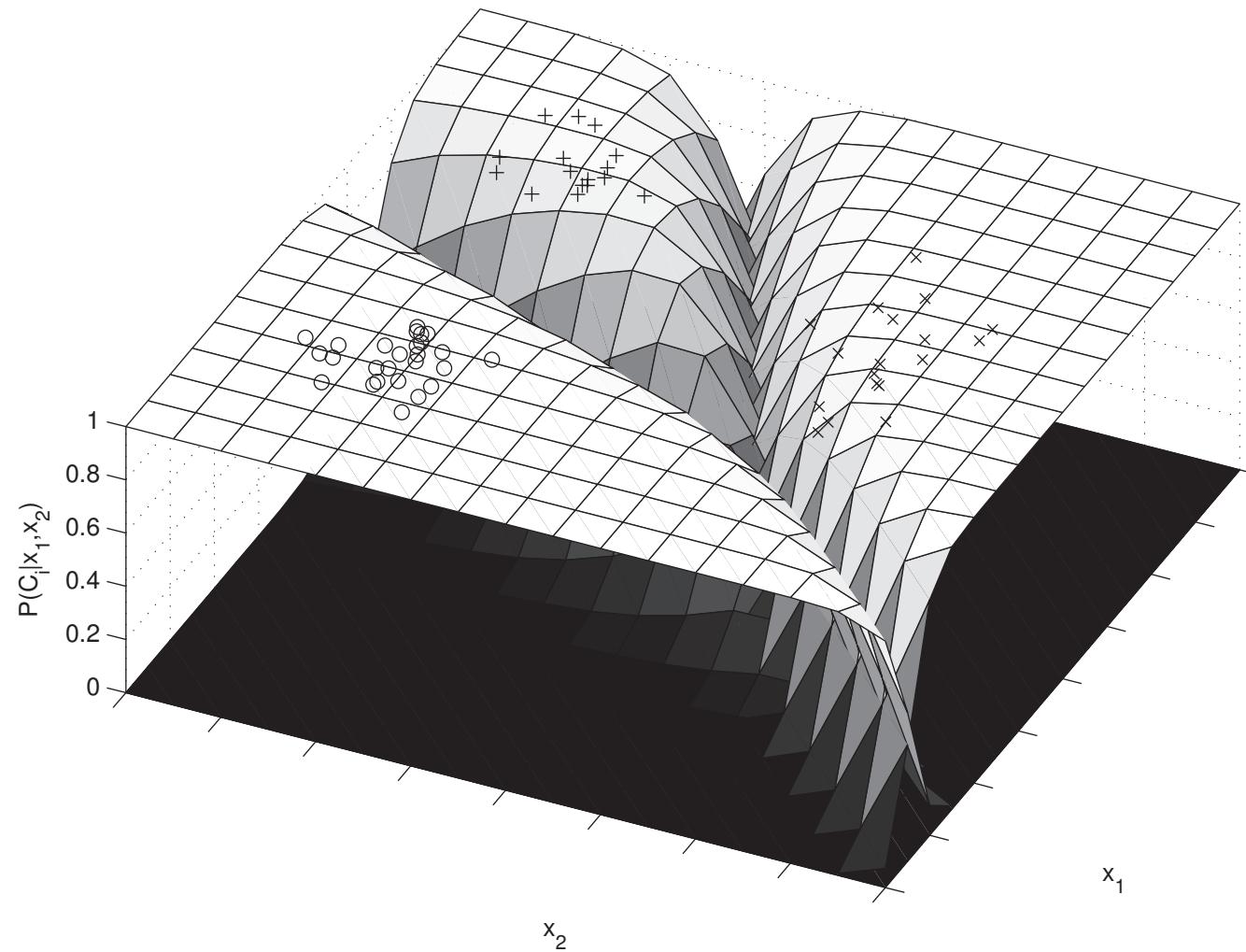


Figure 10.10 For the same example in figure 10.9, the linear discriminants (top), and the posterior probabilities after the softmax (bottom).

Generalizing the Linear Model

- 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:

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

where $\boldsymbol{\phi}(\cdot)$ are basis functions. E.g., *sigmoid function*.

Discrimination by Regression

- Classes are NOT mutually exclusive and exhaustive

- In regression, the probabilistic model is

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

- If $r^t \in \{0, 1\}$, then y^t can be constrained to lie in the range of a sigmoid function

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

- The sample *likelihood* in regression, assume $r|\mathbf{x} \sim \mathcal{N}(y, \sigma^2)$

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

- Maximize *log likelihood* is to minimize sum of square errors $E(\mathbf{w}, w_0 | \mathcal{X}) = \frac{1}{2} \sum_t (r^t - y^t)^2$. Using gradient decent:

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

$$\Delta w_0 = \eta \sum_t (r^t - y^t) y^t (1 - y^t)$$

Learning to Rank

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

We prefer u to v implies that $g(x^u) > g(x^v)$ where $g(x)$ is a score function, here linear:

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

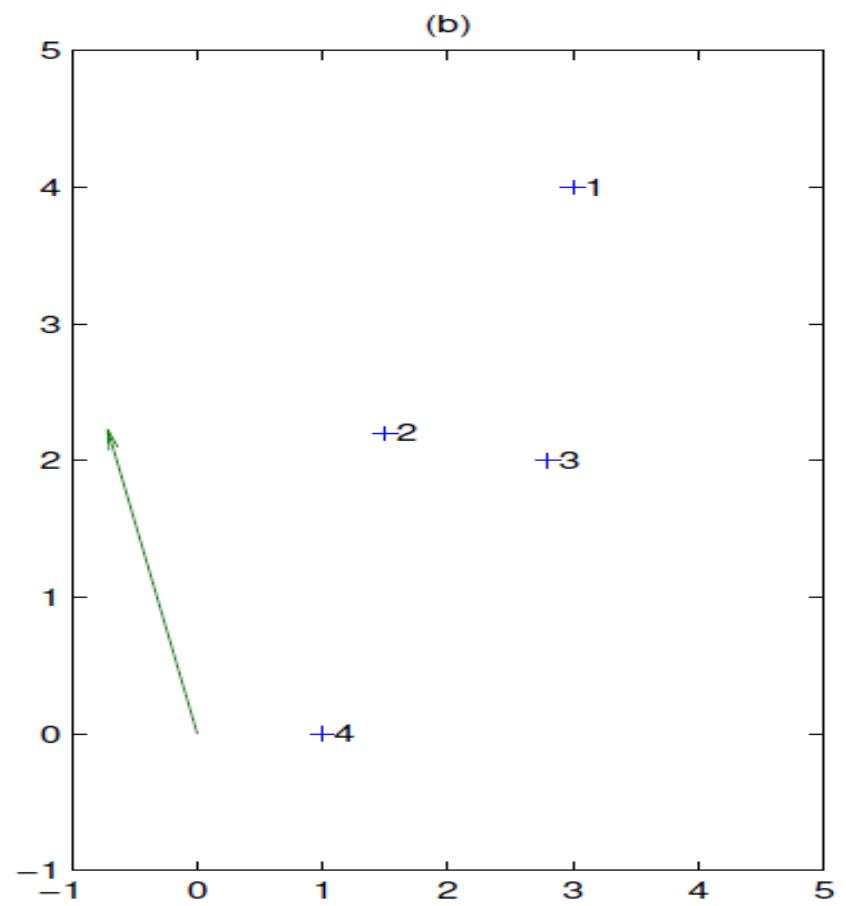
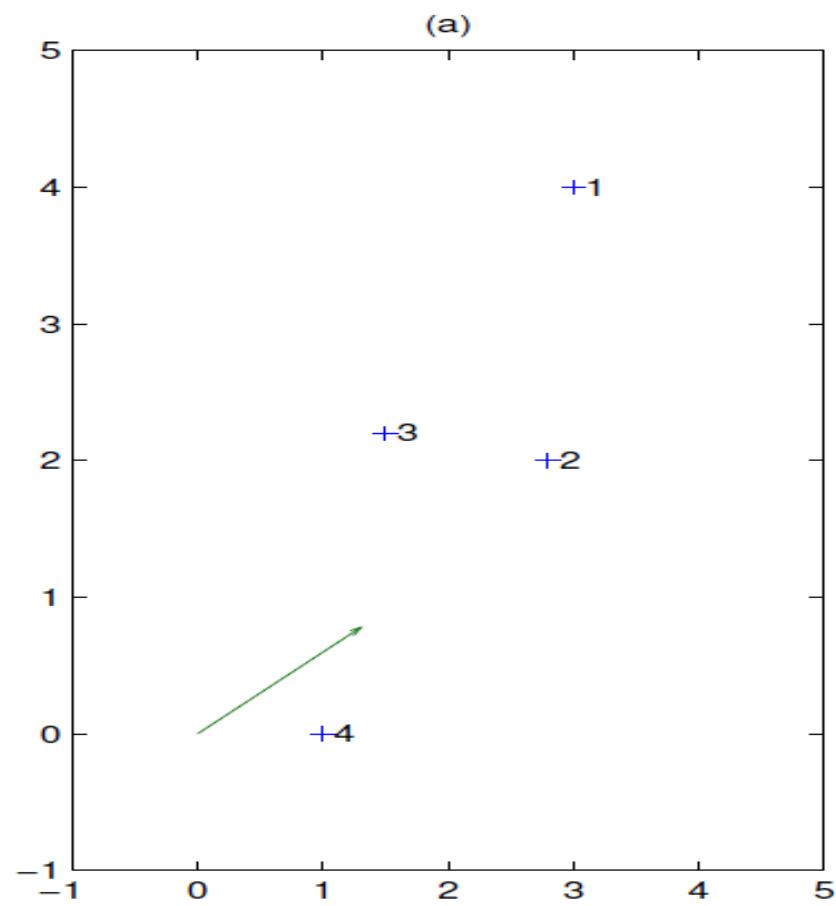
- Find a direction \mathbf{w} such that we get the desired ranks when instances **are projected** along \mathbf{w}

Ranking Error

- 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 $g(\mathbf{x}^u) < g(\mathbf{x}^v)$

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

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



Conclusion

- Geometric interpretation
- Sigmoid and softmax
- Likelihood and error functions
- Gradient decent method
- Depending on our error function, we have different update rules for w
- We can change variables, to consider quadratic (or higher order) discriminants
- We consider logistic regression.