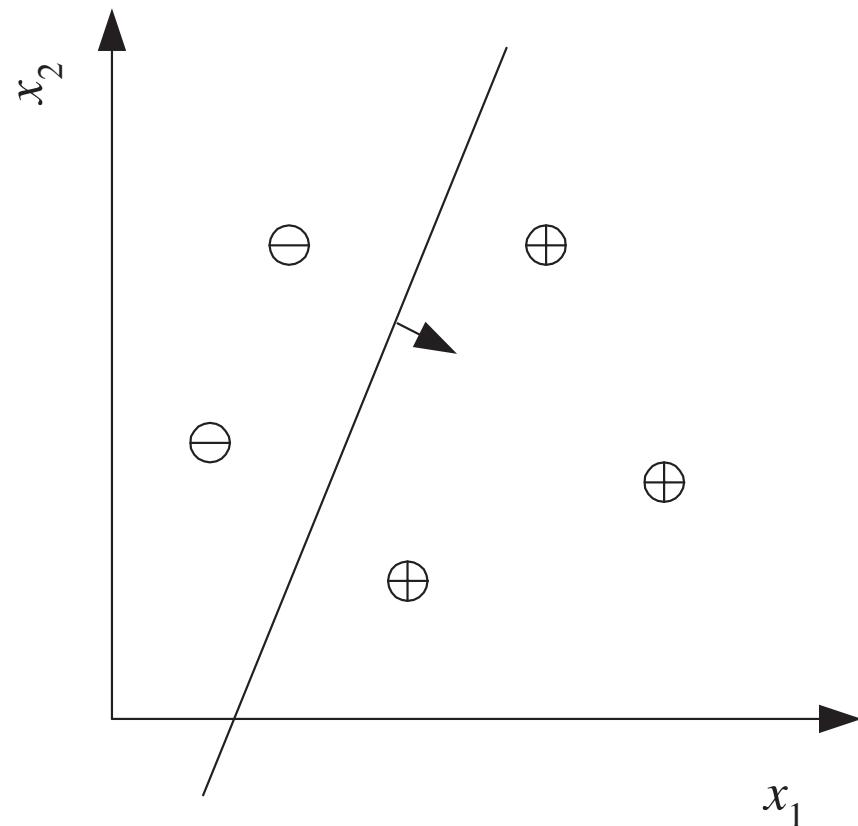


Linear Discrimination

Rui Kuang

Department of Computer Science and Engineering
University of Minnesota

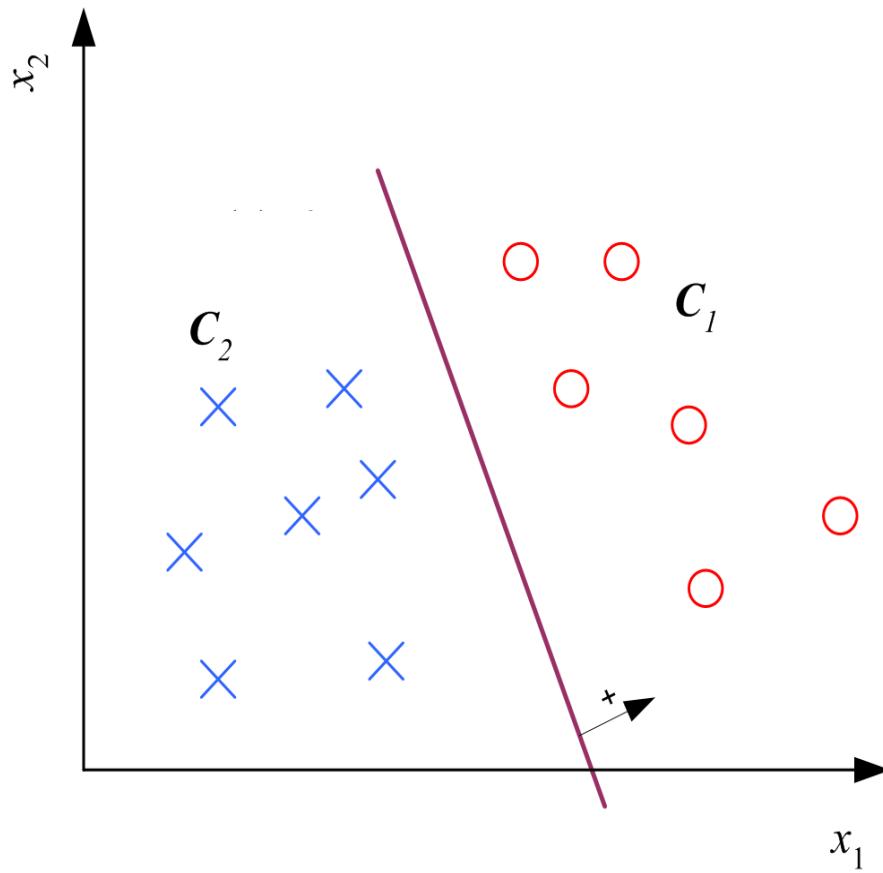
Linear Classifier



$h(x)=\langle w, x \rangle + b$ is a linear classifier

$h(x)>0$ positive
 $h(x)<0$ negative

Likelihood vs Discrimination



Likelihood/Posterior: $g_i(x) = p(x | C_i)P(C_i)$

Rely on estimation of the probability densities.

Discriminant: $g(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + w_0$

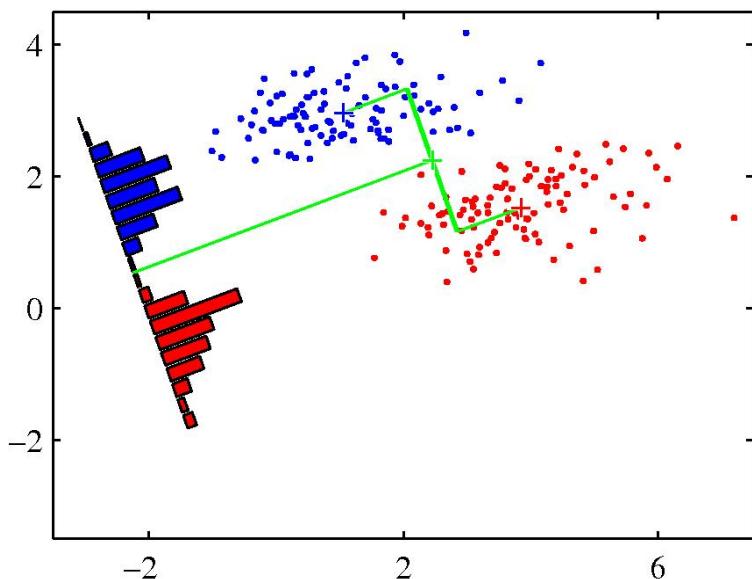
Directly estimate the discrimination function.

Claim: Estimating class densities is harder than estimating the discriminant.

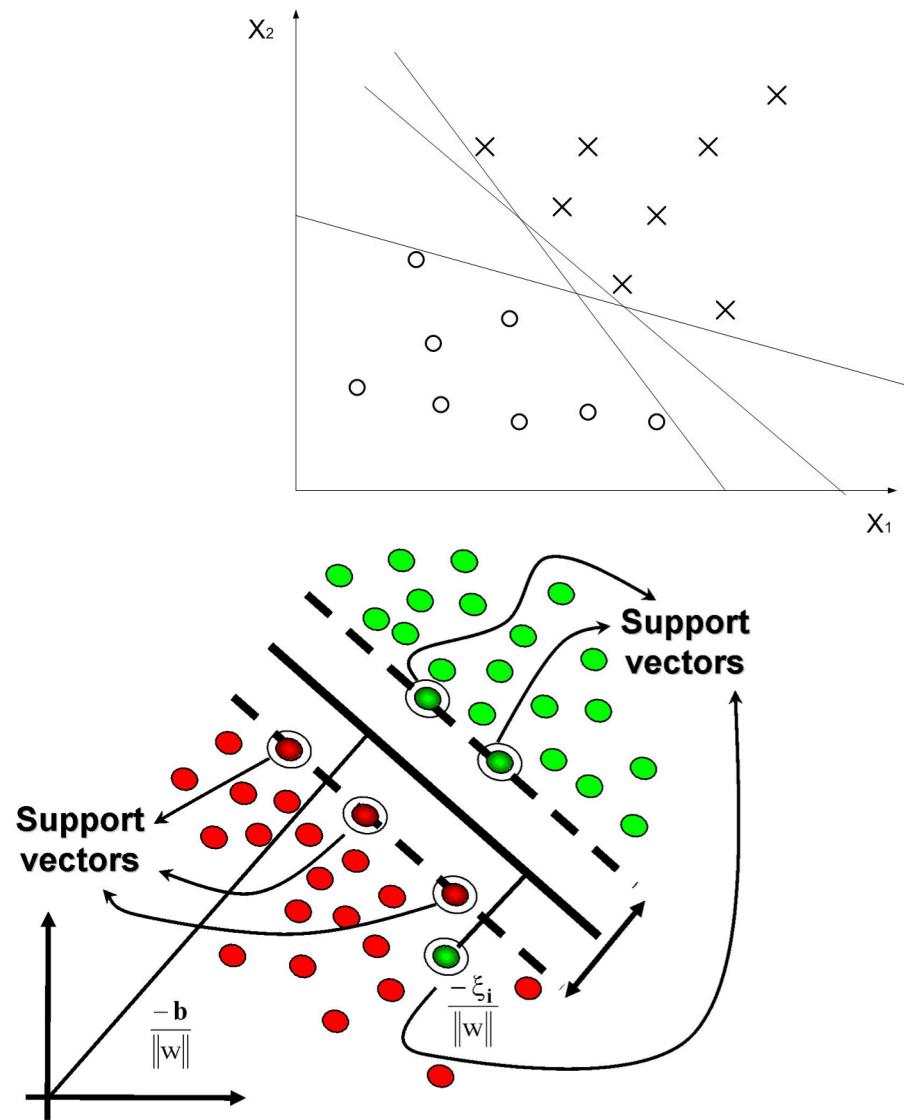
Examples:

Logistic discriminant, LDA,
perceptron, SVM...

Linear Discrimination Examples



LDA, perceptron, SVM...



Likelihood vs Discrimination

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

From Discriminants to Posteriors

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

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

From Discriminants to Posteriors

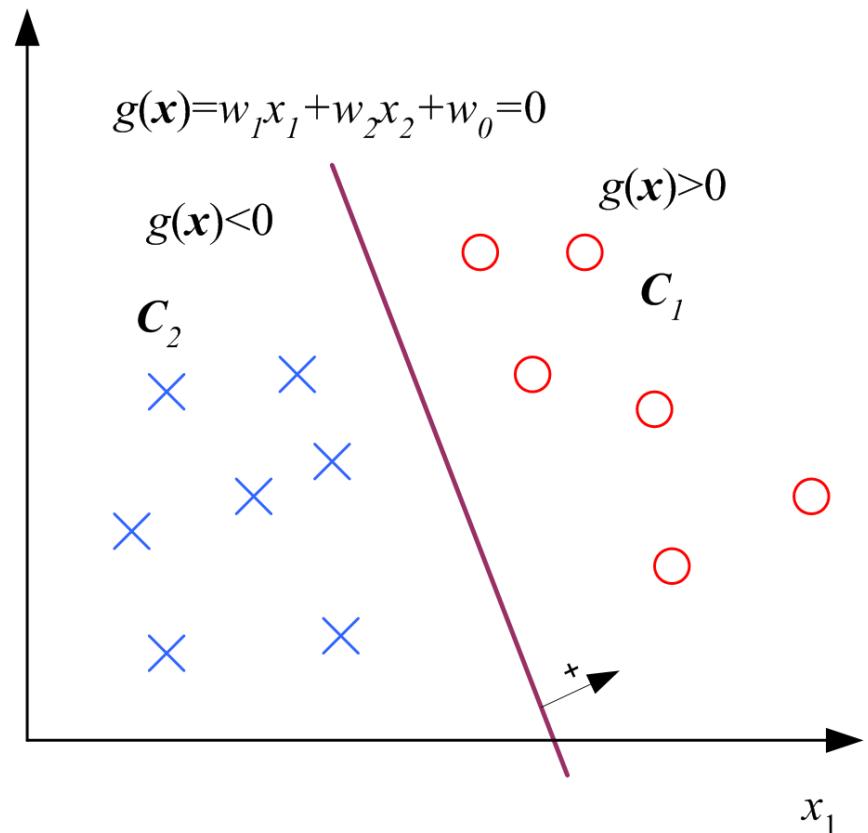
When $p(\mathbf{x} | C_i) \sim N(\mu_i, \Sigma)$, linear discrimination is a log likelihood ratio.

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

choose C_1 if $\log [y / (1 - y)] > 0$

C_2 otherwise

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



Sigmoid (Logistic) Function

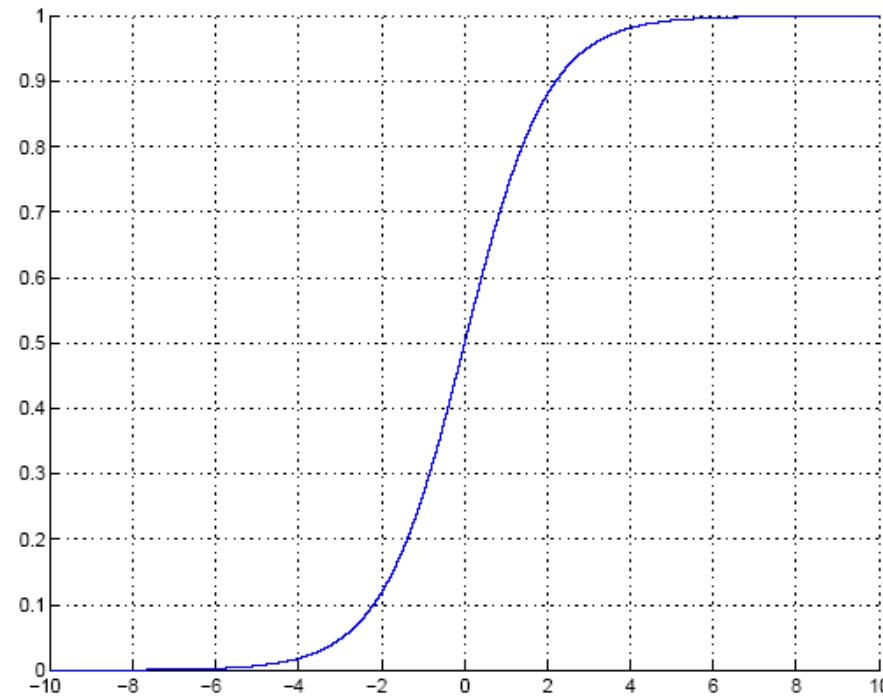
The logit($P(C_1|\mathbf{x})$) is actual a log odds

The inverse of logit is logistic/sigmoid function

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

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

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 $y > 0.5$

Training: Two Classes

- Directly minimize negative log Posterior function

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

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

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

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

Gradient-Descent

- $E(\mathbf{w}|\mathbf{X})$ is error with parameters \mathbf{w} on sample \mathbf{X}

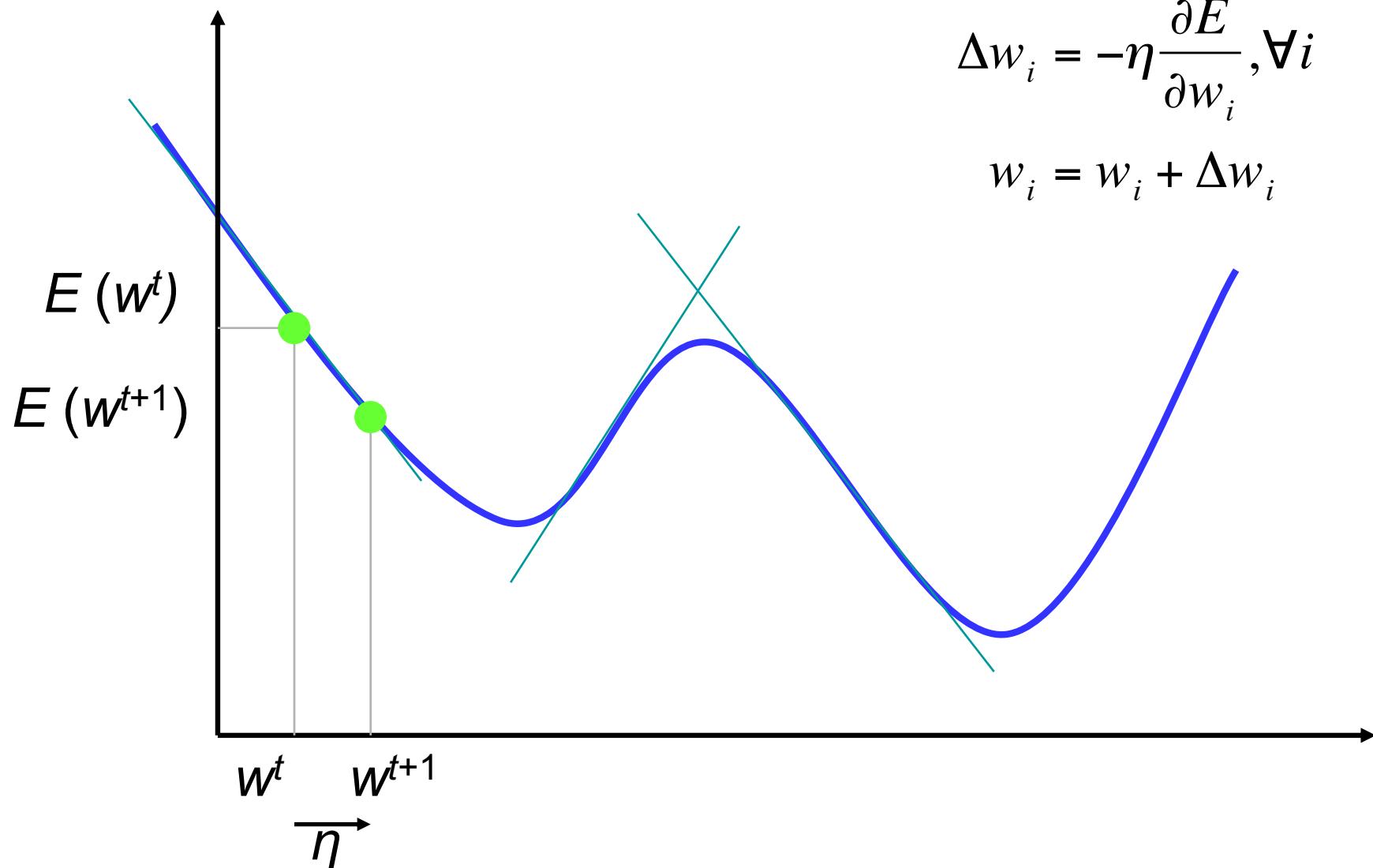
$$\mathbf{w}^* = \arg \min_{\mathbf{w}} E(\mathbf{w} | \mathbf{X})$$

- 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

Gradient-Descent



Training: Gradient-Descent

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

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

$$a = \mathbf{w}^T \mathbf{x} + w_0 \quad \frac{da}{dw_j} = x_j$$

$$\frac{\partial E}{\partial w_j} = \frac{\partial E}{\partial y} \frac{\partial y}{\partial a} \frac{\partial a}{\partial w_j} = -\sum_t \left(\frac{r^t}{y^t} - \frac{1 - r^t}{1 - y^t} \right) y^t (1 - y^t) x_j^t$$

$$\frac{\partial E}{\partial w_0} = \frac{\partial E}{\partial y} \frac{\partial y}{\partial a} \frac{\partial a}{\partial w_0} = -\sum_t \left(\frac{r^t}{y^t} - \frac{1 - r^t}{1 - y^t} \right) y^t (1 - y^t)$$

Training: Gradient-Descent

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

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

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

```

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

```

Matrix Version of Gradeint Descent

$$w = \begin{bmatrix} w_0 \\ w_1 \\ \vdots \\ w_d \end{bmatrix}, r = \begin{bmatrix} r^1 \\ r^2 \\ \vdots \\ r^N \end{bmatrix}, y = \begin{bmatrix} y^1 \\ y^2 \\ \vdots \\ y^N \end{bmatrix}, X = \begin{bmatrix} 1 & x_1^1 & x_2^1 & \cdots & x_d^1 \\ 1 & x_1^2 & x_2^2 & \cdots & x_d^2 \\ \vdots & & & & \\ 1 & x_1^N & x_2^N & \cdots & x_d^N \end{bmatrix}$$

$$\Delta w = \eta X^T (r - y)$$

$$w^{new} = w^{old} + \eta X^T (r - y)$$

Perceptron with Error Function

- Perceptron algorithm

for each example (\mathbf{x}^t, r^t) :

$$\text{if } (\langle \mathbf{w}, \mathbf{x}^t \rangle * r^t \leq 0)$$

$$\mathbf{w} = \mathbf{w} + \eta r^t \mathbf{x}^t$$

- One-sample logistic regression

$$y^t = \text{sigmoid}(\mathbf{w}^T \mathbf{x}^t)$$

$$E^t(\mathbf{w} | \mathbf{x}^t, \mathbf{r}^t) = -r^t \log y^t - (1 - r^t) \log (1 - y^t)$$

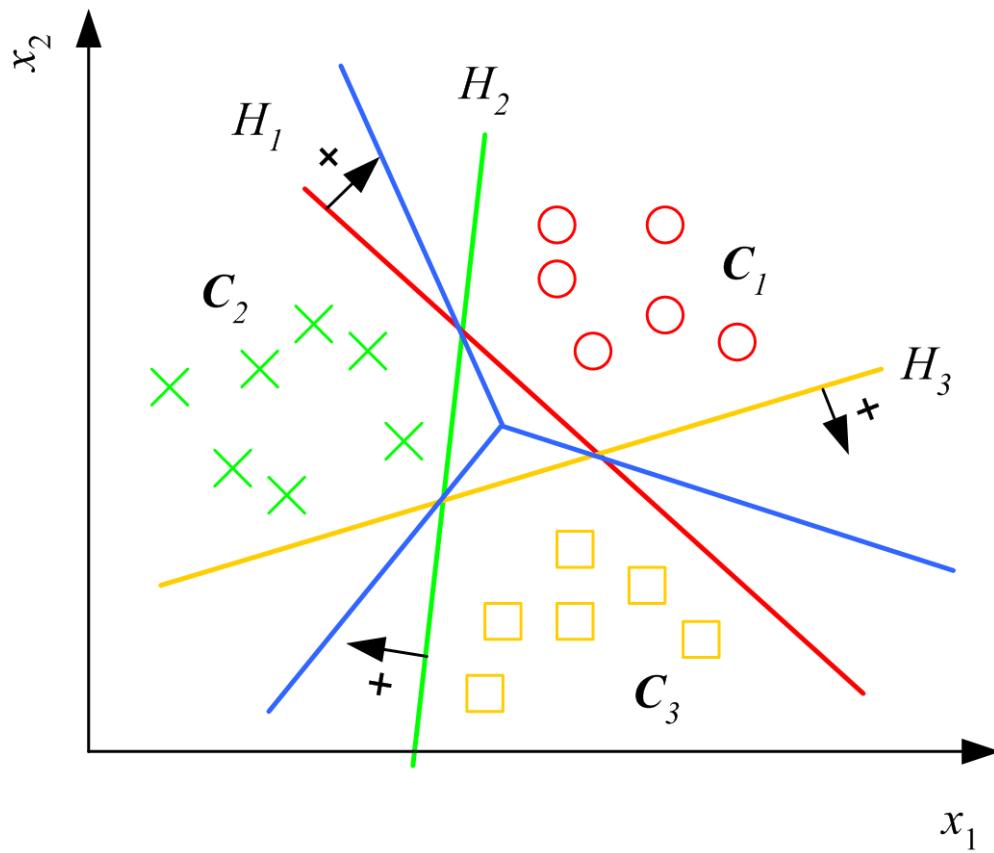
$$\mathbf{w} = \mathbf{w} + \eta (r^t - y^t) \mathbf{x}^t$$

Notes on Logistic Regression

- Optimal when $p(\mathbf{x}|C_i)$ are Gaussian with shared cov matrix but no need to worry about cov estimation.
- Useful when classes are (almost) linearly separable
- The error function is concave with a unique solution with no closed-form solution.
- More efficient algorithms with higher order gradient such as Newton-Raphson (see Bishop's).

Multiple Classes

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

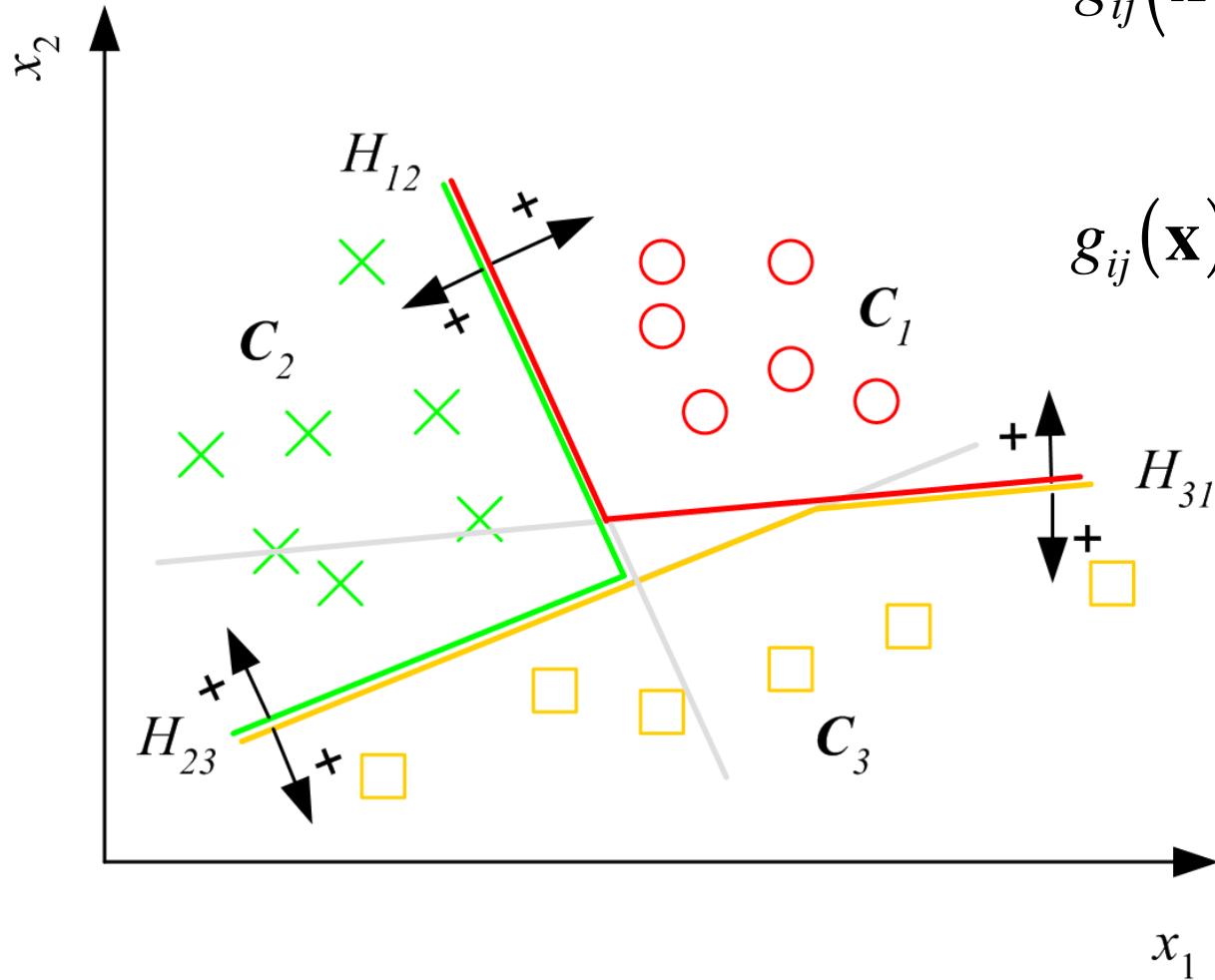


Choose C_i if

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

Classes are
linearly separable

Pairwise Separation



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

K>2 Classes

- Learn a w for each class;
- Choose a reference class C_K

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

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

$$P(C_K | \mathbf{x}) = \frac{1}{1 + \sum_{j=1}^{K-1} \exp[\mathbf{w}_j^T \mathbf{x} + w_{j0}]}$$

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

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

Softmax: taking the max on a differentiable function.

K>2 Classes

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

$$\frac{\partial y_i}{\partial a_j} = y_i(\delta_{ij} - y_j)$$

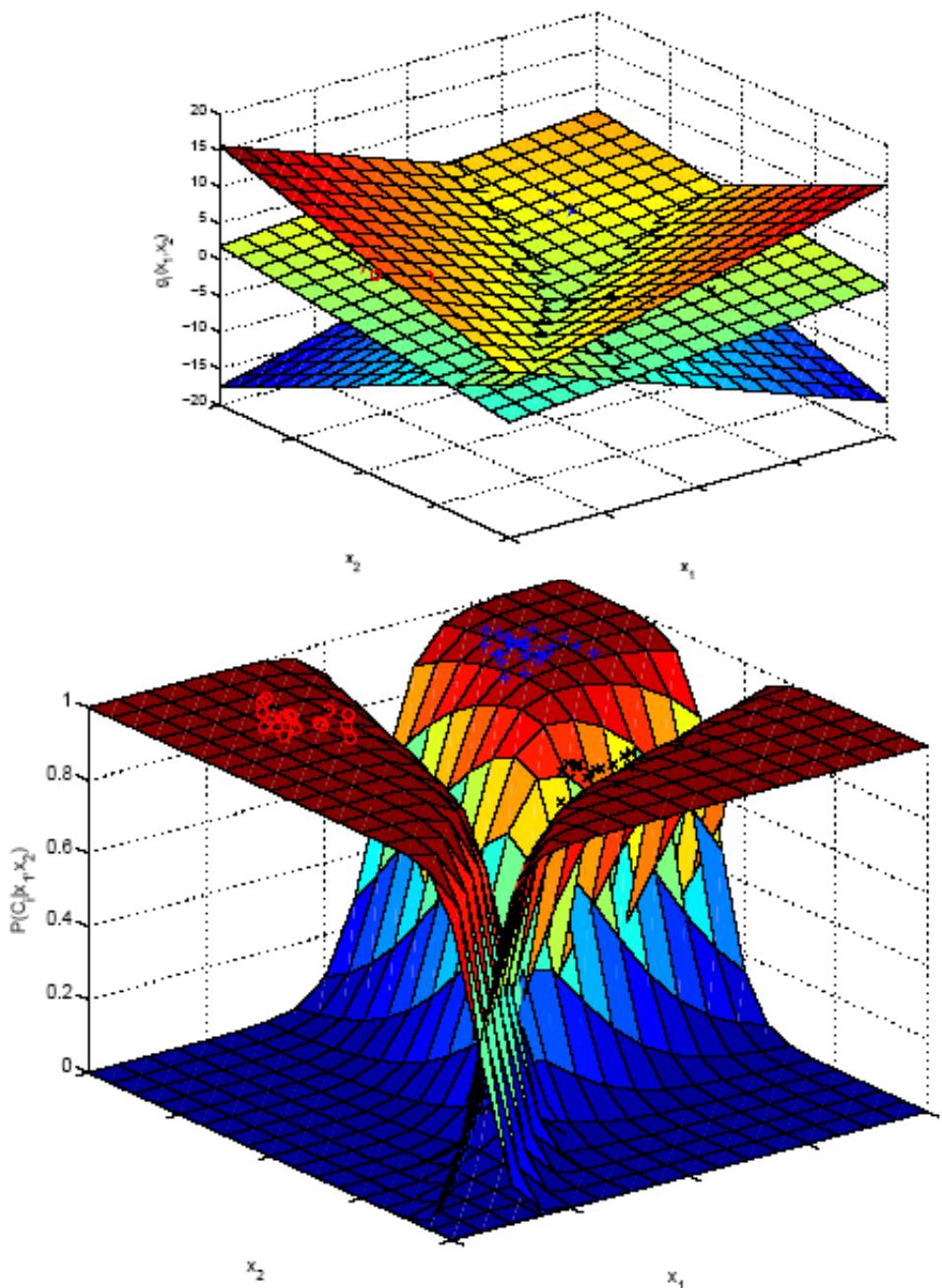
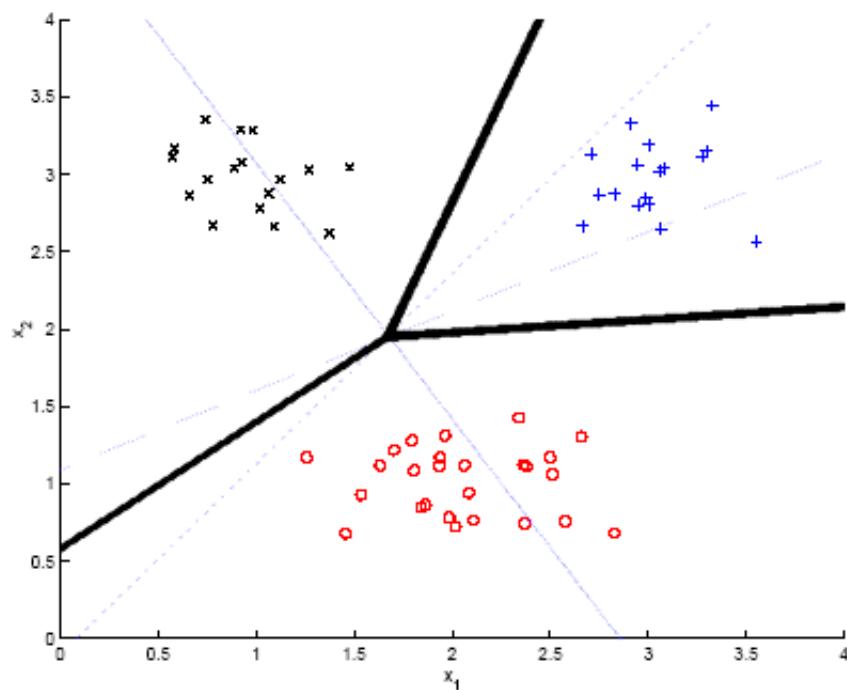
$$l(\{\mathbf{w}_i, w_{i0}\}_i | \mathcal{X}) = \prod_t \prod_i (y_i^t)^{(r_i^t)}$$

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

$$\Delta \mathbf{w}_j = \eta \sum_t \sum_i \frac{r_i^t}{y_i^t} y_i^t (\delta_{ij} - y_j^t) \mathbf{x}^t = \eta \sum_t (r_j^t - y_j^t) \mathbf{x}^t$$

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

Example



Summary of Logistic Regression

■ The computational trick:

$$\text{Chain rule : } \frac{\partial E}{\partial w} = \frac{\partial E}{\partial y} \frac{\partial y}{\partial a} \frac{\partial a}{\partial w}$$

Decided by the error function

$$y = \text{sigmoid}(a)$$

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

$$a = w^T x + w_0$$
$$\frac{da}{dw} = x$$

$$\text{Iterative update : } w^{new} = w^{old} + \Delta w = w^{old} - \eta \frac{\partial E}{\partial w}$$

Variations

■ Other link functions:

$$\text{Chain rule : } \frac{\partial E}{\partial w} = \frac{\partial E}{\partial y} \frac{\partial y}{\partial a} \frac{\partial a}{\partial w}$$

Binomial distribution

$$y = a \quad \frac{dy}{da} = 1$$

$$y = \log(a) \quad \frac{dy}{da} = \frac{1}{a}$$

$$a = w^T x + w_0$$
$$\frac{da}{dw} = x$$

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

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