

# LINEAR MODELS FOR CLASSIFICATION

①

For  $K$  classes the target can be represented as

$$t = (0, 0, \dots, \underset{j}{\downarrow}, 1, 0, \dots, 0) \in \mathbb{R}^K$$

indicating that the corresponding  $x \in C_j$ . We will consider models of the form

$$y(x) = f(w^T x + w_0)$$

where  $f$  is a nonlinear function. For  $x \in \mathbb{R}^D$ , the space is divided into regions where  $w^T x + w_0 = \text{const}$ . These  $D-1$  hyperplanes are the decision boundaries.

## Discriminant Functions

A discriminant function is a rule that assigns  $x$  to only one class  $C_k$ , for  $k = 1, 2, \dots, K$ . Consider

$$y(x) = w^T x + w_0$$

For a 2-class problem we define

$$\begin{cases} x \in C_1 & \text{if } y(x) > 0 \\ x \in C_2 & \text{if } y(x) \leq 0 \end{cases}$$

Thus  $y=0$  is the decision boundary. If  $x_A$  and  $x_B$  are any two points on the boundary,  $w^T(x_A - x_B) = 0$ , thus  $w$  is normal to this  $D-1$  dimensional hyperplane.

For any  $x$  we can write

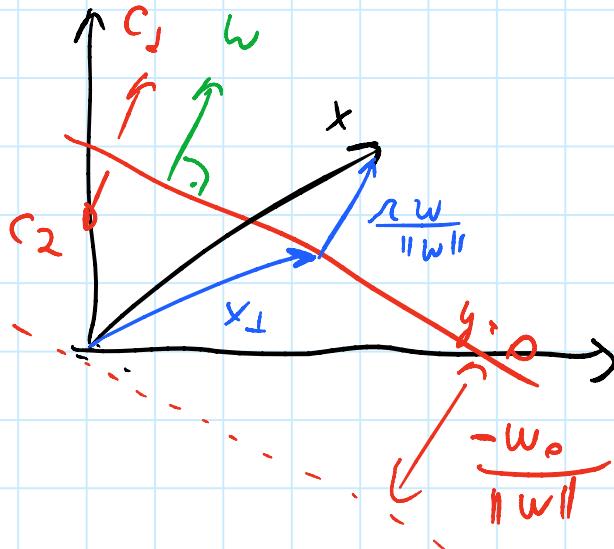
$$x = x_1 + r \frac{w}{\|w\|}$$

where  $y(x_1) = 0$ .

]

$$\text{Now } y(x) = \mathbf{w}^T \mathbf{x} + w_0 = \underbrace{\mathbf{w}^T \mathbf{x}_\perp}_{y(\mathbf{x}_\perp) = 0} + w_0 + \frac{r}{\|\mathbf{w}\|} \|\mathbf{w}\|$$
(2)

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



For a K-class problem, we can avoid ambiguities in defining the decision regions by considering K discriminant functions:

$$y_k(x) = \mathbf{w}_k^T \mathbf{x} + w_{k0} \quad (k=1, \dots, K)$$

Then,

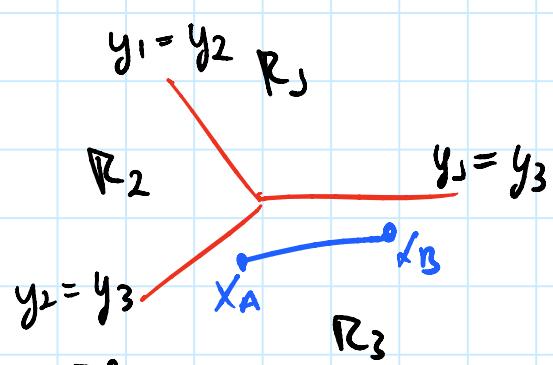
$$x \in C_k \text{ if } y_k(x) > y_j(x), \forall j \neq k.$$

The decision boundary between classes  $C_i$  and  $C_j$  is given by  $y_k(x) = y_j(x)$  which is a  $D-1$ -dimensional hyperplane determined by

$$(\mathbf{w}_k - \mathbf{w}_j)^T \mathbf{x} + (w_{k0} - w_{j0}) = 0$$

Moreover, any region  $R_K$  is convex. To see this, let  $x_A$  and  $x_B$  be in  $R_K$ , and consider the line segment  $\hat{x} = \lambda x_A + (1-\lambda)x_B$ ,  $\lambda \in [0, 1]$ . By linearity

$$\begin{aligned} y_K(\hat{x}) &= \lambda y_K(x_A) + (1-\lambda)y_K(x_B) \\ &> \lambda y_j(x_A) + (1-\lambda)y_j(x_B) = y_j(\hat{x}) \quad \therefore \hat{x} \in R_K \end{aligned}$$



③

### LEAST SQUARES

Now we consider the least-squares solution for  $k=1, \dots, K$

Let  $w_k \rightarrow \begin{pmatrix} w_{k0} \\ w_k \end{pmatrix} \equiv w_k$  and  $x \rightarrow \begin{pmatrix} 1 \\ x \end{pmatrix} \equiv x$ , such that

$$y_k(x) = w_k^T x \rightarrow y(x) = w^T x$$

where

$$y(x) = \begin{pmatrix} y_1(x) \\ \vdots \\ y_K(x) \end{pmatrix}_{K \times 1} \quad w = \begin{pmatrix} w_1 & \cdots & w_K \end{pmatrix}_{(D+1) \times K}$$

We have training data  $\{(x_m, t_m)\}_{m=1}^N$ , and we want to minimize

$$E_D(w) = \frac{1}{2} \sum_{m=1}^N \|t_m - w^T x_m\|^2$$

Define  $T = \begin{pmatrix} t_1^T \\ \vdots \\ t_N^T \end{pmatrix}$  and  $X = \begin{pmatrix} x_1^T \\ \vdots \\ x_N^T \end{pmatrix}$ .

Notice that  $w^T x_m = (x_m^T w)^T = (X^T w)_m$  and  $t_m = T_m^T$ , so the summand is  $\|(T - Xw)_m^T\|^2$ . For any set of vectors  $\{v_m\}$  we have  $\sum_m \|v_m\|^2 = \sum_m v_m^T v_m = \text{Tr}(V^T V)$  where

$$V = \begin{pmatrix} v_1^T \\ \vdots \\ v_N^T \end{pmatrix}.$$

—

Thus the least-squares error can be written as

$$E_{\text{LS}}(w) = \frac{1}{2} \text{Tr} \left\{ (T - Xw)^T (T - Xw) \right\} \quad (4)$$

Then we have

$$\begin{aligned} \frac{\partial E_{\text{LS}}(w)}{\partial w} &= \frac{\partial}{\partial w} \frac{1}{2} \text{Tr} \left\{ (T - Xw)^T (T - Xw) \right\} \\ &= -X^T T + X^T X w = 0 \end{aligned}$$

where we used  $\frac{\partial}{\partial x} \text{Tr}(XA) = A^T$  and  $\frac{\partial}{\partial x} \text{Tr}(X^T BX) = (B + B^T)X$ .

The solution is thus

$$w = (X^T X)^{-1} X^T T = X^+ T //$$

where  $X^+ \equiv (X^T X)^{-1} X^T$  is the left pseudo-inverse,  $X^+ X = I$ .

The discriminant function is

$$y(x) = (X^+ T)^T x$$

Least squares suffers from being too sensitive to outliers, and also, sometimes, it cannot find the correct regions even when the data is linearly separable.

## FISHER LINEAR DISCRIMINANT (LDA)

Suppose we have  $N_1$  points from class  $C_1$  and  $N_2$  points from class  $C_2$ . We want to project the data into 1D:  $y = w^T x$ . If we build a classifier such that  $y > -w_0 \Rightarrow x \in C_1$ , and  $x \in C_2$  otherwise, we obtain the previous linear classifier.

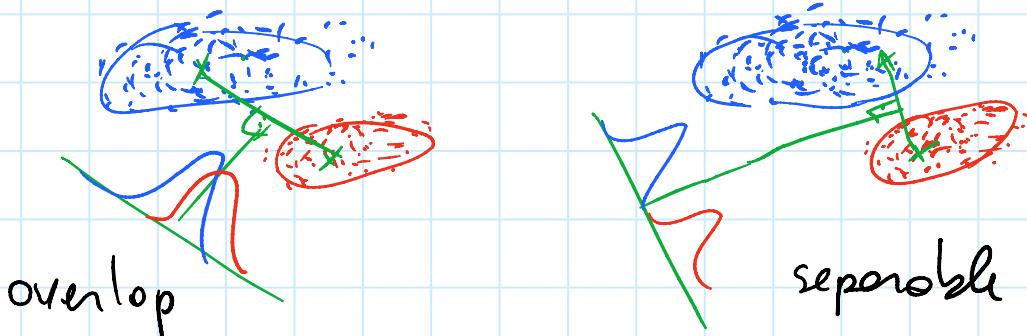
]

One idea is to maximize the projected means: (5) ↗

$$w^T(m_1 - m_2)$$

$$\text{where } m_1 = \frac{1}{N_1} \sum_{x_i \in C_1} x_i, \quad m_2 = \frac{1}{N_2} \sum_{x_i \in C_2} x_i.$$

It can be shown that  $w \propto m_1 - m_2$  is the solution to this problem. This represents a serious drawback if the data has high "off-diagonal" covariance:



Even when the data is linearly separable, the projection on  $m_1 - m_2$  may be highly overlapping.

We need to find a direction such that  $\|m_1 - m_2\|$  is maximized, while keeping the variance of the within-class projected data as small as possible.

Let  $\mu_k = w^T m_k$  be the projected mean and  $y_n = w^T x_n$  the projected data point. The within-class variance is

$$s_k^2 = \sum_{y_n \in C_k} (y_n - \mu_k)^2$$

For a two-class problem the Fisher criterion is

$$J(w) = \frac{(\mu_2 - \mu_1)^2}{S_1^2 + S_2^2}.$$

(6)

Re-writing this:

$$(\mu_2 - \mu_1)^2 = w^T(m_2 - m_1)(m_2 - m_1)^T w \equiv w^T S_B w$$

$$\begin{aligned} S_1^2 + S_2^2 &= \sum_{x_m \in C_1} w^T(x_m - m_1)(x_m - m_1)^T w \\ &\quad + \sum_{x_m \in C_2} w^T(x_m - m_2)(x_m - m_2)^T w \\ &= w^T S_W w \end{aligned}$$

$$\therefore J(w) = \frac{w^T S_B w}{w^T S_W w}$$

$S_B$  is the between class covariance, and  $S_W$  is the within class covariance. Thus

$$\frac{\partial J}{\partial w} = 0 = \frac{S_B w}{w^T S_W w} - \frac{w^T S_B w S_W w}{(w^T S_W w)^2}$$

$$(w^T S_B w) S_W w = (w^T S_W w) S_B w$$

$$S_B w = (m_2 - m_1)(m_2 - m_1)^T w \propto m_2 - m_1$$

J

Therefore

$$w \propto S_w^{-1} (m_2 - m_1) \quad \text{⑦}$$

This is the Fisher's linear discriminant. We build a classifier from this by deciding:

$$\begin{cases} x \in C_1 & \text{if } y(x) = w^T x \geq y_0 \\ x \in C_2 & \text{otherwise} \end{cases}$$

for some threshold  $y_0$ . Since  $y = w^T x$  is a sum of random variables we can assume

$$p(y|C_k) = N(\mu_k, \sigma_k^2)$$

From this we can estimate  $\mu_k, \sigma_k$  (MLE), and using decision theory we can find the optimal  $y_0$ .

Now let us discuss the relation of LDA to least squares. We use a different target scheme. If  $x_m \in C_1$  then  $t_m = \frac{N}{N_1}$ , and if  $x_m \in C_2$  then  $t_m = -\frac{N}{N_2}$ .

Consider

$$E(w) = \frac{1}{2} \sum_{n=1}^N (w^T x_n + w_0 - t_n)^2$$

$$\frac{\partial E}{\partial w_0} = 0 = \sum_{n=1}^N (w^T x_n + w_0 - t_n), \quad w_0 = \frac{1}{N} \sum_{n=1}^N (t_n - w^T x_n)$$

$$w_0 = \frac{1}{N} \left( N_1 \frac{N}{N_1} - N_2 \frac{N}{N_2} \right) - w^T \left( \frac{1}{N} \sum_{n=1}^N x_n \right)$$

$$w_0 = -w^T m \quad \text{where } m = \frac{1}{N} \sum_{n=1}^N x_n.$$

$$\frac{\partial E}{\partial w} = 0 = \sum_{n=1}^N (w^T x_n + w_0 - t_n) x_n$$

$$0 = \sum_{n=1}^N x_n x_n^T w + w_0 \sum_{n=1}^N x_n - \sum_{n=1}^N t_n x_n$$

$$\sum_{n=1}^N t_n x_n = \sum_{x_n \in C_1} N x_n - \sum_{x_n \in C_2} N x_n = N(m_1 - m_2)$$

$$w_0 \sum_{n=1}^N x_n = -w^T m (N_1 m_1 + N_2 m_2)$$

$$= -(\bar{N}_1 m_1 + \bar{N}_2 m_2) \frac{(m_1 - m_2)^T w}{N}$$

$$\sum_{n=1}^N x_n x_n^T = \sum_{x_n \in C_1} x_n x_n^T + \sum_{x_n \in C_2} x_n x_n^T$$

$$= \sum_{x_n \in C_1} (x_n - m_1)(x_n - m_1)^T + N_1 m_1 m_1^T$$

$$+ \sum_{x_n \in C_2} (x_n - m_2)(x_n - m_2)^T + N_2 m_2 m_2^T$$

Therefore

$$(S_w + N_1 m_1 m_1^T + N_2 m_2 m_2^T - \frac{1}{N} (N_1 m_1 + N_2 m_2)(N_1 m_1 + N_2 m_2)^T) w = N(m_1 - m_2)$$

$$(S_w + \frac{N_1 N_2}{N} (m_1 m_1^T + m_2 m_2^T - 2 m_1 m_2^T)) w = N(m_1 - m_2)$$

Finally,

$$(S_w + \frac{N_1 N_2}{N} S_B) w = N(m_1 - m_2)$$

Since  $S_B w \propto m_2 - m_1$ , the solution of this equation is  $w \propto S_w^{-1}(m_2 - m_1)$ , which is the same as (1)A.

The classification will be  $x \in C_1$  if  $y(x) > 0$  and  $x \in C_2$  otherwise, where  $y = w^T x + w_b$ . 9

Fisher's LDA can be generalized for multiple classes, but we are not going to consider it here.

## Probabilistic Generative Models

Here we model  $p(x|C_i)$  and  $p(C_i)$  to finally employ Bayes theorem and use  $p(C_i|x)$  for classification. Let us consider a two-class problem:

$$p(C_1|x) = \frac{p(x|C_1)p(C_1)}{p(x|C_1)p(C_1) + p(x|C_2)p(C_2)} = \frac{1}{1 + e^{-\alpha}} \equiv \sigma(\alpha)$$

where

$$\alpha = \log \frac{p(x|C_1)p(C_1)}{p(x|C_2)p(C_2)} = \log \frac{p(C_1|x)}{p(C_2|x)}$$

$\sigma(\alpha)$  is called the sigmoid function, and it maps the entire real line on the  $(0, 1)$  interval



$$1 - e^{-\alpha} = \gamma, \quad e^{-\alpha} = \frac{1}{\gamma} - 1, \quad e^\alpha = \frac{1}{\gamma} - 1 = \frac{\sigma}{1-\sigma}$$

]

The inverse of  $\sigma$  is the logit function

$$a = \log \frac{\sigma}{1-\sigma} = \log \frac{p(c_1|x)}{p(c_2|x)}$$



For  $k > 2$  we have

$$p(c_k|x) = \frac{p(x|c_k)p(c_k)}{\sum_i p(x|c_i)p(c_i)} = \frac{e^{a_k}}{\sum_j e^{a_j}}$$

where  $a_k = \log p(x|c_k)p(c_k)$ . The above function is known as Normalized exponential or softmax function since if  $a_k \gg a_j$  for any  $j \neq k$  then  $p(c_k|x) \approx 1$  and  $p(c_j|x) \approx 0$ .

Now we assume some form of distribution for the class conditional. Let us assume equal covariances:

$$p(x|c_k) = \frac{1}{(2\pi)^{D/2}} \frac{1}{|\Sigma|^k} e^{-\frac{1}{2}(x-\mu_k)^T \Sigma^{-1}(x-\mu_k)}$$

Consider  $k=2$ . Then  $p(c_1|x) = \sigma(a)$  where

$$\begin{aligned} a &= \log \frac{p(x|c_1)}{p(x|c_2)} + \log \frac{p(c_1)}{p(c_2)} \\ &= -\frac{1}{2}(x-\mu_1)^T \Sigma^{-1}(x-\mu_1) + \frac{1}{2}(x-\mu_2)^T \Sigma^{-1}(x-\mu_2) + \log \frac{p(c_1)}{p(c_2)} \\ &= x^T \Sigma^{-1} \mu_1 - \frac{1}{2} \mu_1^T \Sigma^{-1} \mu_1 - x^T \Sigma^{-1} \mu_2 + \frac{1}{2} \mu_2^T \Sigma^{-1} \mu_2 + \log \frac{p(c_1)}{p(c_2)} \\ &= (\bar{\Sigma}^{-1}(\mu_1 - \mu_2))^T x - \frac{1}{2} \mu_1^T \Sigma^{-1} \mu_1 + \frac{1}{2} \mu_2^T \Sigma^{-1} \mu_2 + \log \frac{p(c_1)}{p(c_2)} \end{aligned}$$

(11)

Define

$$\omega = \Sigma^{-1}(\mu_1 - \mu_2)$$

$$w_0 = \frac{1}{2} \mu_2^T \Sigma^{-1} \mu_2 - \frac{1}{2} \mu_1^T \Sigma^{-1} \mu_1 + \log \frac{p(C_1)}{p(C_2)}$$

Then

$$Q = \omega^T X + w_0 //$$

The decision boundary is a linear function in  $x$ -space. The priors  $p(C_k)$  only enter the bias  $w_0$  and their effect is only to translate the decision boundary.

Now let  $k=2$ . Recall that  $p(C_k|x) = \frac{e^{Q_k}}{\sum_j e^{Q_j}}$  where

$$\begin{aligned} Q_k &= \log p(x|C_k) + \log p(C_k) \\ &= -\frac{D}{2} \log \pi - \frac{1}{2} \log |\Sigma| - \frac{1}{2} (x - \mu_k)^T \Sigma^{-1} (x - \mu_k) \\ &= -\frac{D}{2} \log \pi - \frac{1}{2} \log |\Sigma| - \frac{1}{2} x^T \Sigma^{-1} x + \mu_k^T \Sigma^{-1} x - \frac{1}{2} \mu_k^T \Sigma^{-1} \mu_k \\ &\quad + \log p(C_k) \end{aligned}$$

this will be the same for all  $e^{Q_j}$ , so it cancels with the denominator.

We thus have

$$Q_k = \omega_k^T X + w_{k0}$$

$$\omega_k = \Sigma^{-1} \mu_k$$

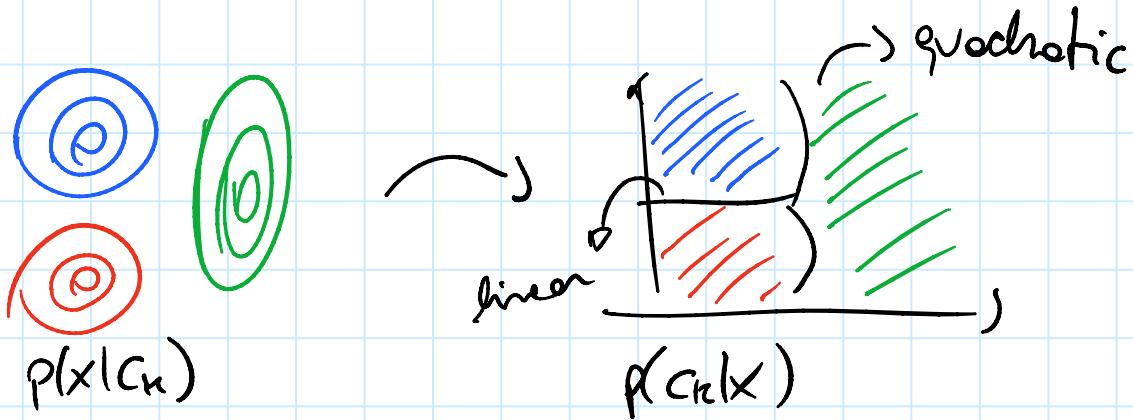
$$w_{k0} = -\frac{1}{2} \mu_k^T \Sigma^{-1} \mu_k + \log p(C_k)$$

]

The decision boundary will be given by

$p(C_k|x) = p(C_j|x)$  where these are the two largest posteriors. This gives  $a_k = a_j$  which is linear in  $x$ . (12)

If we had chosen  $\Sigma$  the same for all classes, the quadratic term would not have cancelled, and the decision boundary would be a quadratic function of  $x$ . For instance:



Once we specify a form for  $p(x|C_k)$  we can use MLE to estimate parameters and the priors. Assume we have training data  $\{(x_m, t_m)\}_{m=1}^N$ , where  $t_m = 1$  if  $x_m \in C_1$ , and  $t_m = 0$  if  $x_m \in C_2$ , where we are considering  $k=2$  problem. Assume  $x \sim N(\mu_1, \Sigma)$  for  $x \in C_1$  and  $p(C_1) = \pi$ , and also  $x \sim N(\mu_2, \Sigma)$  for  $x \in C_2$  and  $p(C_2) = 1 - \pi$ . Thus,

$$p(x, C_1) = p(x|C_1) p(C_1) = \pi N(x|\mu_1, \Sigma)$$

$$p(x, C_2) = p(x|C_2) p(C_2) = (1-\pi) N(x|\mu_2, \Sigma)$$

The date likelihood function is thus

(13)

$$P(T | \mu_1, \mu_2, \Sigma) = \prod_{n=1}^N (\pi N(x_n | \mu_1, \Sigma))^{\ell_n} ((1-\pi) N(x_n | \mu_2, \Sigma))^{1-\ell_n}$$

where  $T = (t_1, \dots, t_N)^T$

$$\begin{aligned} l = \log P &= \sum_{n=1}^N \left\{ \ell_n \left( \log \pi - \frac{1}{2} \log \Sigma - \frac{1}{2} \log |\Sigma| - \frac{1}{2} (x_n - \mu_1)^T \Sigma^{-1} (x_n - \mu_1) \right) \right. \\ &\quad \left. + (1-\ell_n) \left( \log (1-\pi) - \frac{1}{2} \log \Sigma - \frac{1}{2} \log |\Sigma| - \frac{1}{2} (x_n - \mu_2)^T \Sigma^{-1} (x_n - \mu_2) \right) \right\} \end{aligned}$$

$$\frac{\partial l}{\partial \pi} = 0 = \sum_{n=1}^N \frac{\ell_n}{\pi} - \frac{(1-\ell_n)}{1-\pi} \left\{ = \sum_{n=1}^N \frac{(\ell_n - \pi)}{\pi(1-\pi)} \right\}$$

$$\pi = \frac{1}{N} \sum_{n=1}^N \ell_n = \frac{N_1}{N_1 + N_2} //$$

$$\frac{\partial l}{\partial \mu_1} = 0 = \sum_{n=1}^N \ell_n (\Sigma^{-1} x_n - \Sigma^{-1} \mu_1)$$

$$\sum_{n=1}^N \ell_n \mu_1 = \sum_{n=1}^N \ell_n x_n, \quad \mu_1 = \frac{1}{N_1} \sum_{n=1}^N \ell_n x_n = \frac{1}{N_1} \sum_{x_n \in C_1} x_n //$$

$$\frac{\partial l}{\partial \mu_2} = 0 = \sum_{n=1}^N (1-\ell_n) (\Sigma^{-1} x_n - \Sigma^{-1} \mu_2)$$

$$\sum_{n=1}^N (1-\ell_n) \mu_2 = \sum_{n=1}^N (1-\ell_n) x_n, \quad \mu_2 = \frac{1}{N_2} \sum_{x_n \in C_2} x_n //$$

—

(14)

The terms involving  $\Sigma$  are

$$\begin{aligned}
 &= \sum_{n=1}^N \left\{ t_n \left( -\frac{1}{2} \log |\Sigma| - \frac{1}{2} (\mathbf{x}_n - \boldsymbol{\mu}_1)^T \Sigma^{-1} (\mathbf{x}_n - \boldsymbol{\mu}_1) \right) \right. \\
 &\quad \left. + (1-t_n) \left( -\frac{1}{2} \log |\Sigma| - \frac{1}{2} (\mathbf{x}_n - \boldsymbol{\mu}_2)^T \Sigma^{-1} (\mathbf{x}_n - \boldsymbol{\mu}_2) \right) \right\} \\
 &= -\frac{N}{2} \log |\Sigma| - \frac{1}{2} \sum_{\mathbf{x}_n \in C_1} (\mathbf{x}_n - \boldsymbol{\mu}_1)^T \Sigma^{-1} (\mathbf{x}_n - \boldsymbol{\mu}_1) \\
 &\quad - \frac{1}{2} \sum_{\mathbf{x}_n \in C_2} (\mathbf{x}_n - \boldsymbol{\mu}_2)^T \Sigma^{-1} (\mathbf{x}_n - \boldsymbol{\mu}_2) \\
 &= -\frac{N}{2} \log |\Sigma| - \frac{1}{2} \ln ((N_1 S_1 + N_2 S_2) \Sigma^{-1})
 \end{aligned}$$

where we defined

$$\left\{ \begin{array}{l} S_1 = \frac{1}{N_1} \sum_{\mathbf{x}_n \in C_1} (\mathbf{x}_n - \boldsymbol{\mu}_1)(\mathbf{x}_n - \boldsymbol{\mu}_1)^T \\ S_2 = \frac{1}{N_2} \sum_{\mathbf{x}_n \in C_2} (\mathbf{x}_n - \boldsymbol{\mu}_2)(\mathbf{x}_n - \boldsymbol{\mu}_2)^T \end{array} \right.$$

Setting the  $\frac{\partial}{\partial \Sigma} (\dots) = 0$  we have

$$\begin{aligned}
 -\frac{N}{2} \Sigma^{-1} + \frac{1}{2} (\Sigma^{-1} (N_1 S_1 + N_2 S_2) \Sigma^{-1})^T &= 0 \\
 \Sigma^{-1} (N_1 S_1 + N_2 S_2)^T \Sigma^{-1} &= N \Sigma^{-1}
 \end{aligned}$$

Therefore  $\Sigma = S$ ,  $S = \frac{1}{N} (N_1 S_1 + N_2 S_2)$ .

If we have training data we can replace this into our decision rule  $a = \mathbf{w}^T \mathbf{x} + w_0 \geq 0$ . This is the plus-in classifier.

—

(15)

J

(16) 7

1