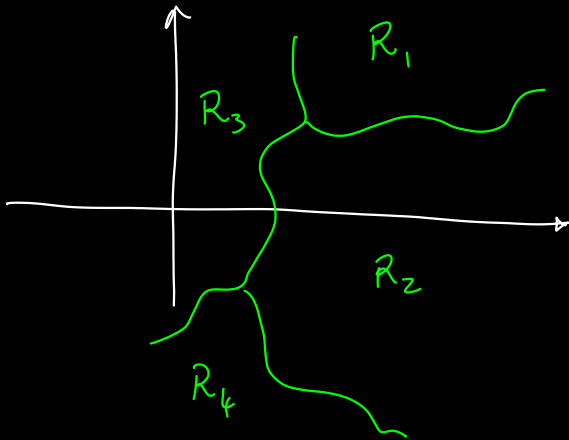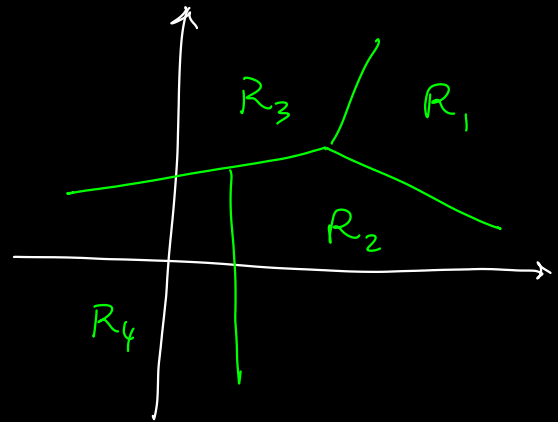Linear and Probabilistic Models for Classification

## Decision regions



## Decision regions (linear models)



## Example: 2 classes

$$\begin{cases} y(x) = w^T x + w_0 \\[2mm] x \in R_1 \text{ iff } y(x) \geq 0 \\ x \in R_2 \text{ if } y(x) < 0 \end{cases} \quad (\text{sign}[y(x)])$$

Example $\overbrace{\qquad}^{w}$

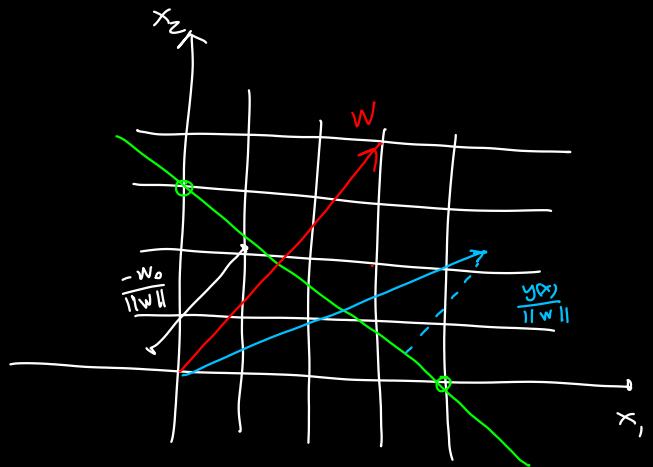$w_0 = -12, \; w_1 = 3, \; w_2 = 4$

decision boundary : $(x_1, x_2)$ s.t. $y(x) = 0$

$$3x_1 + 4x_2 - 12 = 0$$



$$\|w\| = \sqrt{3^2 + 4^2} = \sqrt{25} = 5$$

$$\frac{-w_0}{\|w\|} = \frac{12}{5} \qquad \frac{y(x)}{\|w\|}$$

## Least squares

check derivations on the book

# Fisher's linear discriminant

linear classification viewed as dimensionality reduction

$$y = w^T x \qquad \text{from } D \text{ to } 1 \text{ dimension}$$

we can select the projection that maximizes class separation

Example 2-class

$$C_1, N_1 \quad, \quad C_2, N_2$$

$$\underline{m_1} = \frac{1}{N_1} \sum_{n \in C_1} x_n \qquad \underline{m_2} = \frac{1}{N_2} \sum_{n \in C_2} x_n$$

Simple measure of separation:

$$m_2 - m_1 = w^T (\underline{m_2} - \underline{m_1})$$

but this can be arbitrarily large (increasing $\|w\|$)
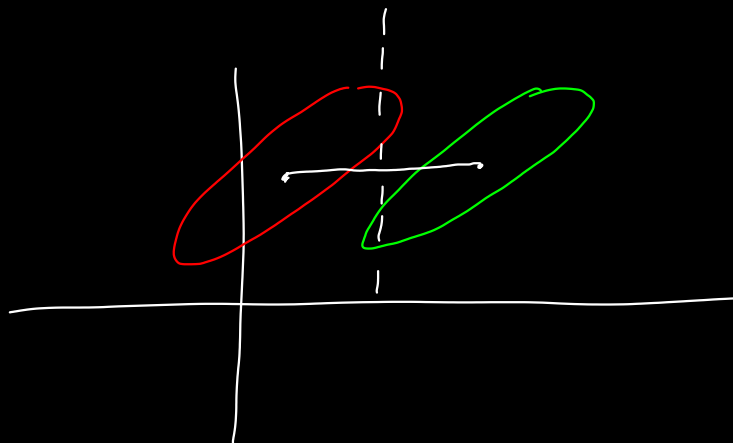
we can constrain $\sum_i w_i = 1$

Problem:

Fig 4.6

Idea: 1) maximize distance of projected means
2) minimize within-class projected variance

$$S_K^2 = \sum_{n \in C_K} (y_n - m_K)^2$$

$$y(x_n) = w^T x_n$$

total within-class variance:   $S_1^2 + S_2^2$

$$J(w) = \frac{(m_2 - m_1)^2}{S_1^2 + S_2^2} \qquad \text{Fisher's criterion}$$

$$= \frac{\left(w^T \underline{m}_2 - w^T \underline{m}_1\right)^2}{\sum\limits_{n \in C_1} \left(w^T x_n - w^T \underline{m}_1\right)^2 - \sum\limits_{n \in C_2} \left(w^T x_n - w^T \underline{m}_2\right)^2} =$$

$$w^T (x_n - \underline{m}_1)\left[w^T (x_n - \underline{m}_1)\right]^T$$

$$w^T (x_n - \underline{m}_1)(x_n - \underline{m}_1)^T w$$

$$= \frac{w^T (\underline{m}_2 - \underline{m}_1)(\underline{m}_2 - \underline{m}_1)^T w}{w^T \left[\sum\limits_{n \in C_1} (x_n - \underline{m}_1)(x_n - \underline{m}_1)^T + \sum\limits_{n \in C_2} (x_n - \underline{m}_2)(x_n - \underline{m}_2)^T\right] w}$$

$$= \frac{w^T S_B w}{w^T S_W w}$$

$S_B \leftarrow$   between-class covariance

$S_W \leftarrow$   within-class covariance

$$\text{differentiating} \Rightarrow w \propto S_W^{-1}(\underline{m}_2 - \underline{m}_1)$$

$y = w^T x$ is roughly gaussian because of the central limit theorem.

least-squares $\longrightarrow$ minimize error w.r.t. $t$

fisher $\longrightarrow$ maximize class separation in output space

If the target for class $C_1$ is $\dfrac{N}{N_1}$ (reciprocal of prior)
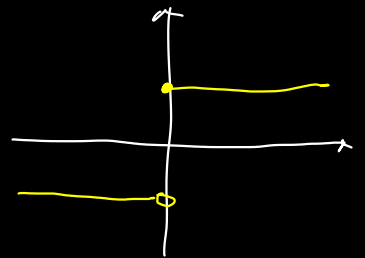
$$C_2 \text{ is } -\frac{N}{N_2}$$

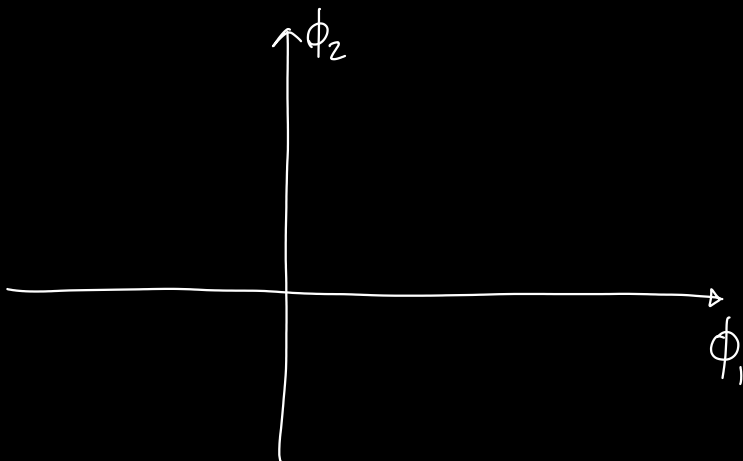$\Longrightarrow$ least squares $\equiv$ fisher

## Perception

Rosenblatt 1962

$$y(x) = f\left(w^T \phi(x)\right) \qquad f(a) = \begin{cases} +1, & a \geq 0 \\ -1, & a < 0 \end{cases}$$



optimization of $w$ through error minimization



$t \in \{-1, +1\}$

$w^T \phi(x_n) > 0 \longrightarrow C_1$

$w^T \phi(x_n) < 0 \longrightarrow C_2$

we would like $w^T \phi(x_n) t_n > 0$ for all patterns
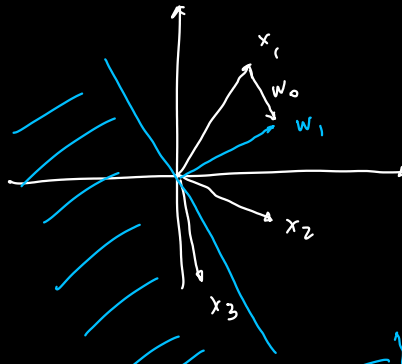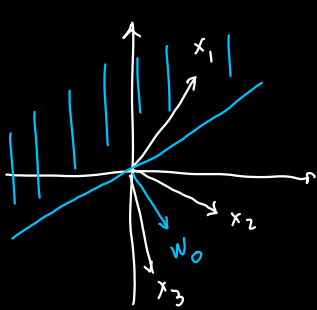
$$E_p(w) = - \sum_{n \in M} w^T \phi(x_n) t_n \qquad M = \text{misclassified}$$

# Stochastic gradient descent
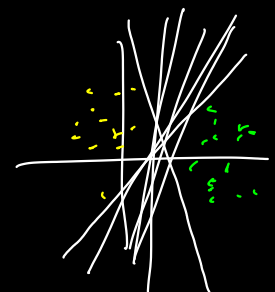
$$w^{(\tau+1)} = w^{(\tau)} - \eta \nabla E_p(w) = w^{(\tau)} + \eta \phi(x_n) t_n$$

$\uparrow$ learning rate

because  $w * constant$  does not change the solution, we can

set  $\eta = 1$

not guaranteed to reduce error at each step



Fig

# Probabilistic models

## Maximum Likelihood

In general we would like to maximize

$$\ln L = \sum_{n=1}^{N} \ln p(x_n, t_n | \theta) \cdots \quad \left[ p(x_n, t_n | \theta) = p(x_n | t_n, \theta)\, p(t_n | \theta) \right]$$

$$= \sum_{n=1}^{N} \left[ \ln p(x_n | t_n, \theta) + \ln p(t_n | \theta) \right] \cdots$$

$$= \sum_{k=1}^{K} \sum_{n:\{t_n = k\}} \left[ \ln p(x_n | C_k, \theta_k) + \ln p(C_k | \theta_F) \right]$$

$\left( \begin{array}{l} D \text{ nda makes class} \\ \text{independaca an assumption} \end{array} \right)$

When we differentiate w.r.t. the parameters for class $k$, only data points belonging to that class influence the derivative.

Example: Gaussian class-conditional likelihoods

$$p(C_k) = \pi_k$$

$$p(x | C_k) = N(x | \mu_k, \Sigma_k)$$

$$\theta = \{ \pi_1 \cdots \pi_K, \mu_1 \cdots \mu_K, \Sigma_1 \cdots \Sigma_k \}$$

$$D \rightarrow p(x, C_k) \equiv$$

$$D \left\{ \begin{array}{l} D_1 \\ D_2 \\ \vdots \\ D_k \\ D_K \end{array} \right. \longrightarrow ML \left\{ \begin{array}{l} p(C_k) \rightarrow \pi_{ML} = \frac{N_k}{N} \\ \\ p(x | C_k) \end{array} \right. \longrightarrow \left\{ \begin{array}{l} \mu_k = \frac{1}{N_k} \sum_{n:t_n = k} x_n \\ \\ \Sigma_k = \frac{1}{N_k} \sum_{n:t_n = k} (x_n - \mu_k)^T (x_n - \mu_k) \end{array} \right.$$

Special Case: equal $\Sigma$

## Example Gaussian distributions

$$p(x|C_k) = \underbrace{\frac{1}{(2\pi)^{D/2}} \frac{1}{|\Sigma_k|^{1/2}}}_{c} \exp\left\{ -\frac{1}{2}(x-\mu_k)^T \Sigma_k^{-1} (x-\mu_k) \right\}$$

$\underbrace{\phantom{\exp\left\{ -\frac{1}{2}(x-\mu_k)^T \Sigma_k^{-1} (x-\mu_k) \right\}}}_{\substack{\shortparallel \\ c}}$ quadratic

if 2-class, decision boundaries:

$$p(C_1|x) = 0.5 \Leftrightarrow \ln p(C_1|x) = c \Leftrightarrow$$

if 2-classes and $\Sigma_1 = \Sigma_2 = \Sigma \Rightarrow$

$$p(C_1|x) = \sigma(a)$$

$$a = \ln \frac{N(x|\mu_1, \Sigma) \, p(C_1)}{N(x|\mu_2, \Sigma) \, p(C_2)} = \cancel{\ln c} - \frac{1}{2}(x-\mu_1)^T \Sigma^{-1}(x-\mu_1) + \ln p(C_1)$$

$$- \cancel{\ln c} + \frac{1}{2}(x-\mu_2)^T \Sigma^{-1}(x-\mu_2) - \ln p(C_2)$$

$$= \cancel{-\frac{1}{2} x^T \Sigma^{-1} x} + \cancel{\frac{1}{2} x^T \Sigma^{-1} x} - \frac{1}{2}\mu_1^T \Sigma^{-1}\mu_1 + \frac{1}{2}\mu_2^T \Sigma^{-1}\mu_2$$

$$+ \mu_1^T \Sigma^{-1} x - \mu_2^T \Sigma^{-1} x + \ln \frac{p(C_1)}{p(C_2)}$$

$$= \underbrace{(\mu_1 - \mu_2)^T \Sigma^{-1}}_{w^T} x \underbrace{- \frac{1}{2}\mu_1^T \Sigma^{-1}\mu_1 + \frac{1}{2}\mu_2^T \Sigma^{-1}\mu_2 + \ln \frac{p(C_1)}{p(C_2)}}_{w_0}$$

$$= w^T x + w_0 \qquad \text{linear!}$$