

Machine Learning (part II)

Single Layer Neural Network

Angelo Ciaramella

Introduction

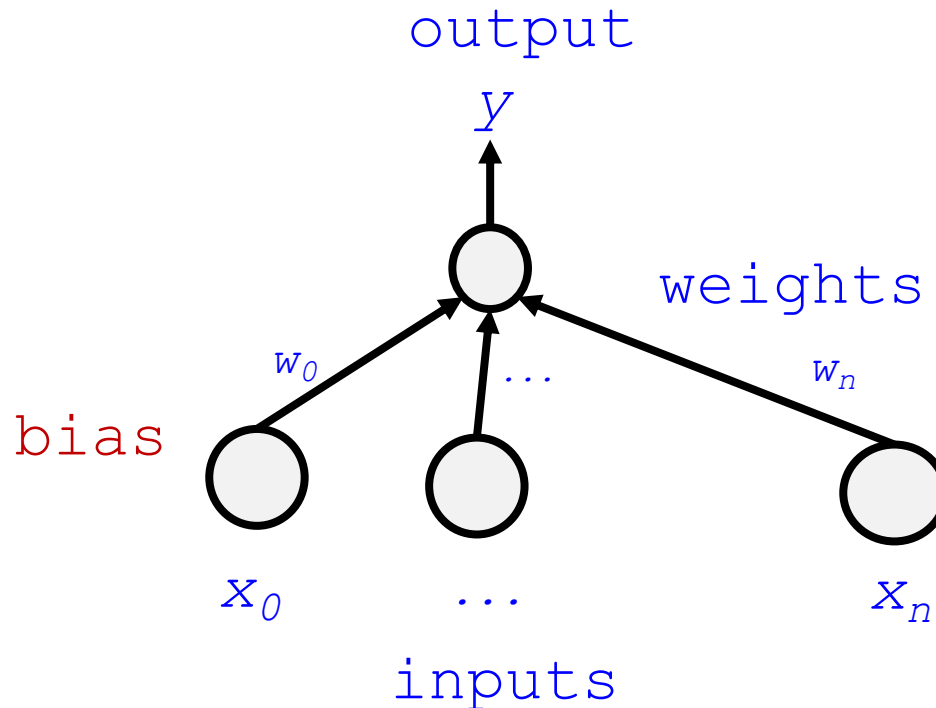
- Linear discriminant functions
 - Linear functions of the input variables
- Generalization
 - Consider a *non-linear function*



Linear discriminant function

activation function

$$y(\mathbf{x}) = \sigma(\mathbf{w}^T \mathbf{x} + w_0)$$



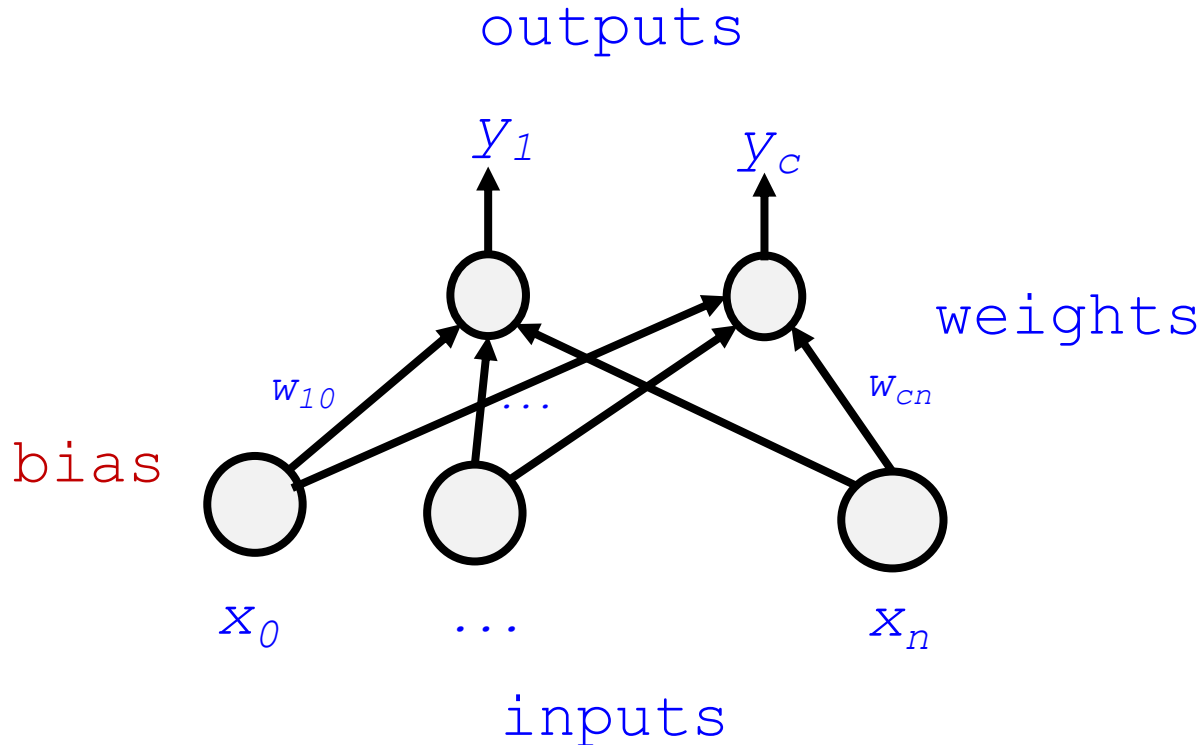
- bias is a **threshold**



Multiple linear discriminant

- $K > 2$ classes

$$y_k(\mathbf{x}) = \sigma(\mathbf{w}_k^T \mathbf{x} + w_{k0})$$



- bias is a **threshold**



The rules of Probability

- Sum rule

$$p(X) = \sum_Y p(X, Y)$$

joint probability

- Product rule

$$p(X, Y) = p(Y|X)p(X)$$
$$p(X, Y) = p(Y, X)$$

- Bayes' theorem

$$p(Y|X) = \frac{p(X|Y)p(Y)}{p(X)}$$



Probabilistic view

■ Classification

- probabilistic view of classification
- models with linear decision boundaries arise from simple assumptions about the distribution of the data

■ Posterior probability of class C_1

$$\begin{aligned} \text{a posterior probability} \quad p(C_1|\mathbf{x}) &= \frac{\text{likelihood} \quad p(\mathbf{x}|C_1)p(C_1) \quad \text{a prior probability}}{p(\mathbf{x}|C_1)p(C_1) + p(\mathbf{x}|C_2)p(C_2)} \\ &= \frac{1}{1 + \exp(-a)} = \sigma(a) \quad \text{normalization factor} \end{aligned}$$



Probabilistic view

- Making decision
 - X-ray image x
 - $P(C_1)$ probability that a person has cancer
 - $P(C_1 | x)$ probability that a person has cancer after observing information of X-ray



Minimizing the misclassification rate

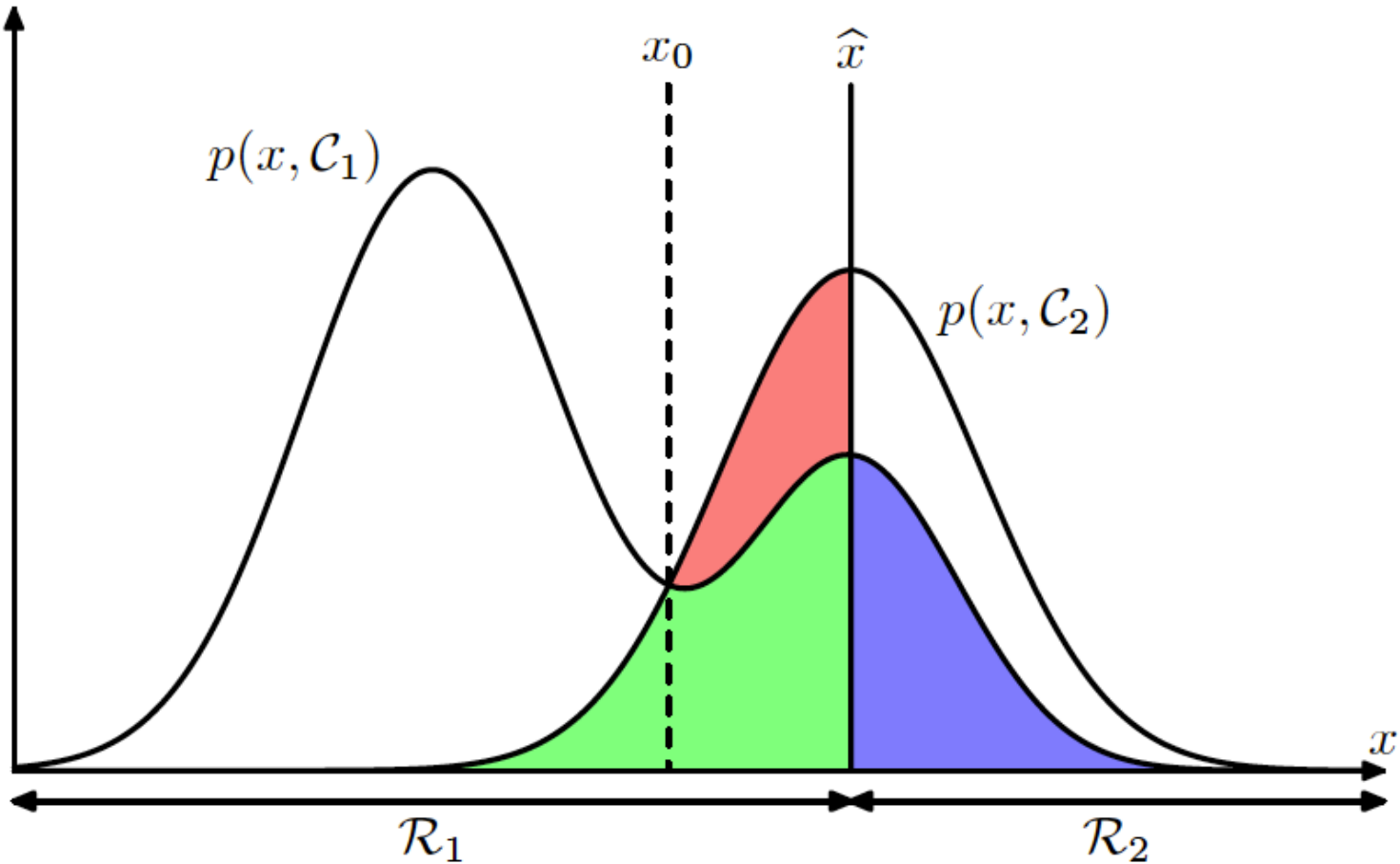
■ Decision boundaries

$$\begin{aligned} p(\text{mistake}) &= p(\mathbf{x} \in \mathcal{R}_1, \mathcal{C}_2) + p(\mathbf{x} \in \mathcal{R}_2, \mathcal{C}_1) \\ &= \int_{\mathcal{R}_1} p(\mathbf{x}, \mathcal{C}_2) d\mathbf{x} + \int_{\mathcal{R}_2} p(\mathbf{x}, \mathcal{C}_1) d\mathbf{x} \end{aligned}$$

$$\begin{aligned} p(\text{correct}) &= \sum_{k=1}^K p(\mathbf{x} \in \mathcal{R}_k, \mathcal{C}_k) \\ &= \sum_{k=1}^K \int_{\mathcal{R}_k} p(\mathbf{x}, \mathcal{C}_k) d\mathbf{x} \end{aligned}$$



Minimizing the misclassification rate



Minimizing the misclassification rate



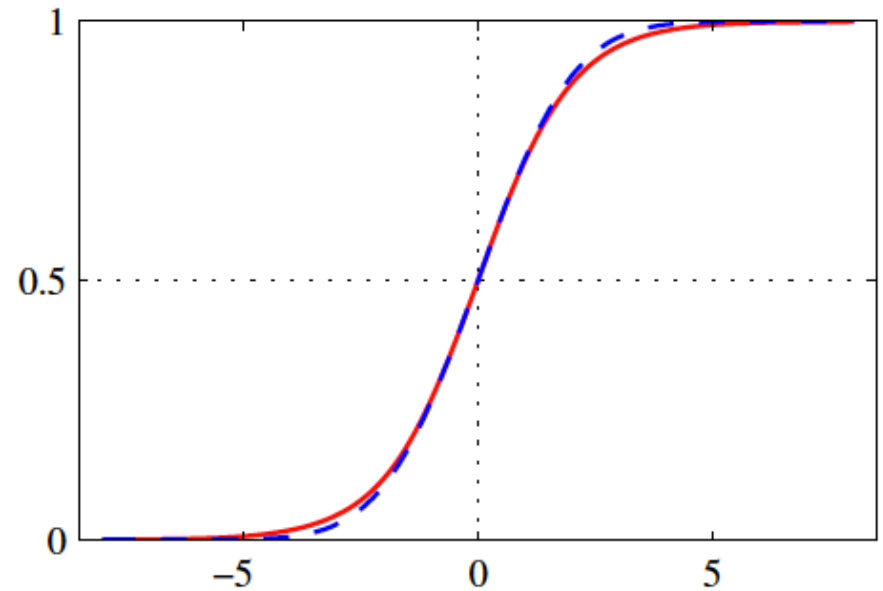
Probabilistic view

■ where

$$a = \ln \frac{p(\mathbf{x}|\mathcal{C}_1)p(\mathcal{C}_1)}{p(\mathbf{x}|\mathcal{C}_2)p(\mathcal{C}_2)}$$

■ and logistic sigmoid function

$$\sigma(a) = \frac{1}{1 + \exp(-a)}$$



Logistic sigmoid

- Symmetry property

$$\sigma(-a) = 1 - \sigma(a)$$

- Inverse (logit)

$$a = \ln \left(\frac{\sigma}{1 - \sigma} \right)$$



Probabilistic view

- $K > 2$

$$\begin{aligned} p(\mathcal{C}_k | \mathbf{x}) &= \frac{p(\mathbf{x} | \mathcal{C}_k) p(\mathcal{C}_k)}{\sum_j p(\mathbf{x} | \mathcal{C}_j) p(\mathcal{C}_j)} \\ &= \frac{\exp(a_k)}{\sum_j \exp(a_j)} \end{aligned}$$

normalized exponential
or softmax function

$$a_k = \ln p(\mathbf{x} | \mathcal{C}_k) p(\mathcal{C}_k).$$

$$\begin{aligned} a_k \gg a_j & \quad p(\mathcal{C}_k | \mathbf{x}) \simeq 1 \\ & \quad p(\mathcal{C}_j | \mathbf{x}) \simeq 0 \end{aligned}$$



Probabilistic view

- Gaussian class-conditional density

$$p(\mathbf{x}|\mathcal{C}_k) = \frac{1}{(2\pi)^{D/2}} \frac{1}{|\boldsymbol{\Sigma}|^{1/2}} \exp \left\{ -\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_k)^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu}_k) \right\}$$

- Posterior probability

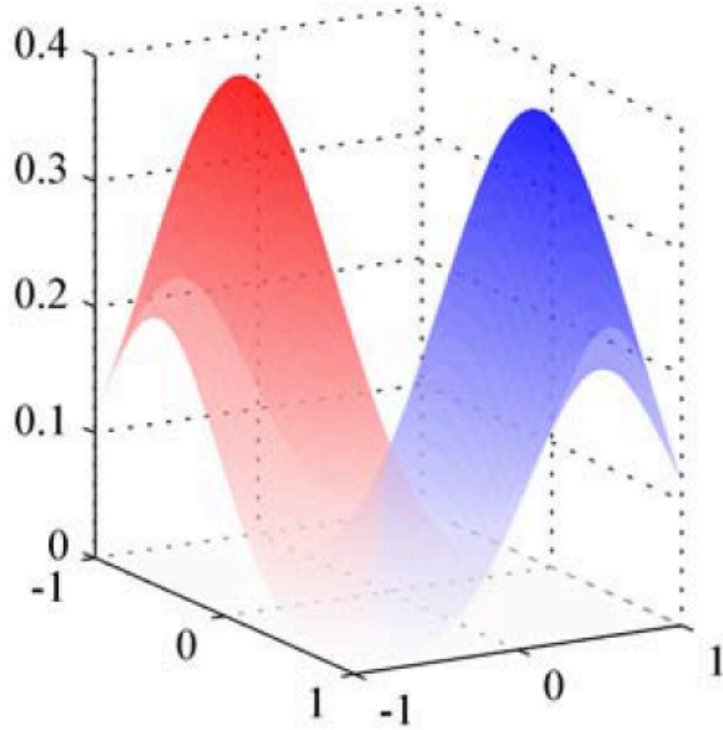
$$p(\mathcal{C}_1|\mathbf{x}) = \sigma(\mathbf{w}^T \mathbf{x} + w_0)$$

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

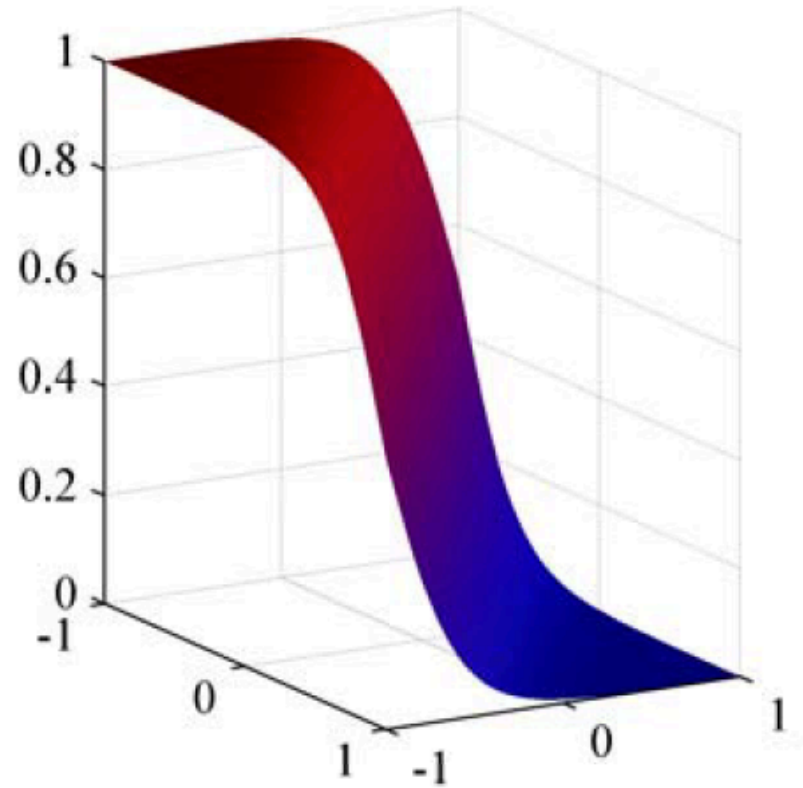
$$w_0 = -\frac{1}{2}\boldsymbol{\mu}_1^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_1 + \frac{1}{2}\boldsymbol{\mu}_2^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_2 + \ln \frac{p(\mathcal{C}_1)}{p(\mathcal{C}_2)}$$



Probabilistic view



Class-conditional densities
for two classes



$p(C_1|\mathbf{x})$



Probabilistic view

- $K > 2$

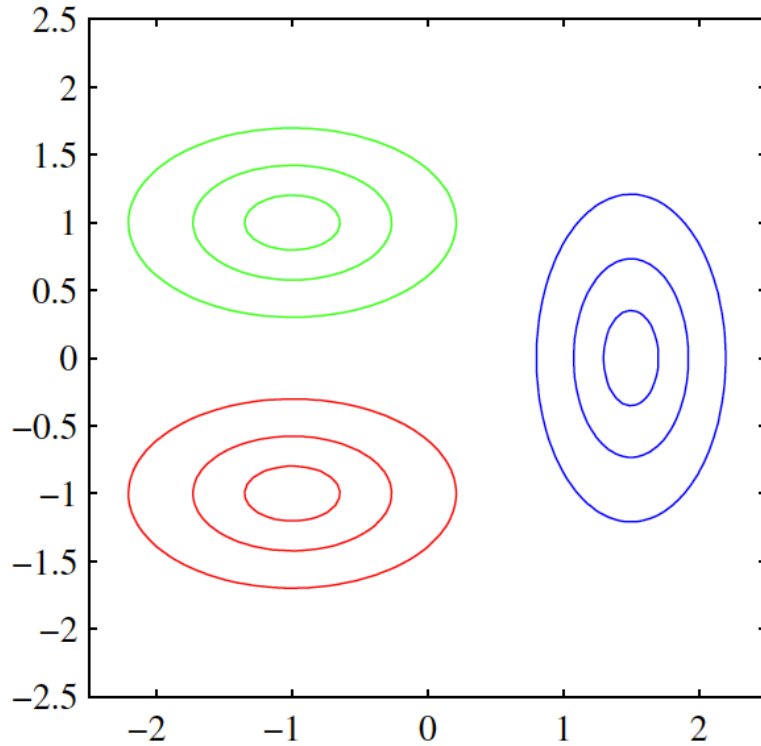
$$a_k(\mathbf{x}) = \mathbf{w}_k^T \mathbf{x} + w_{k0}$$

- where

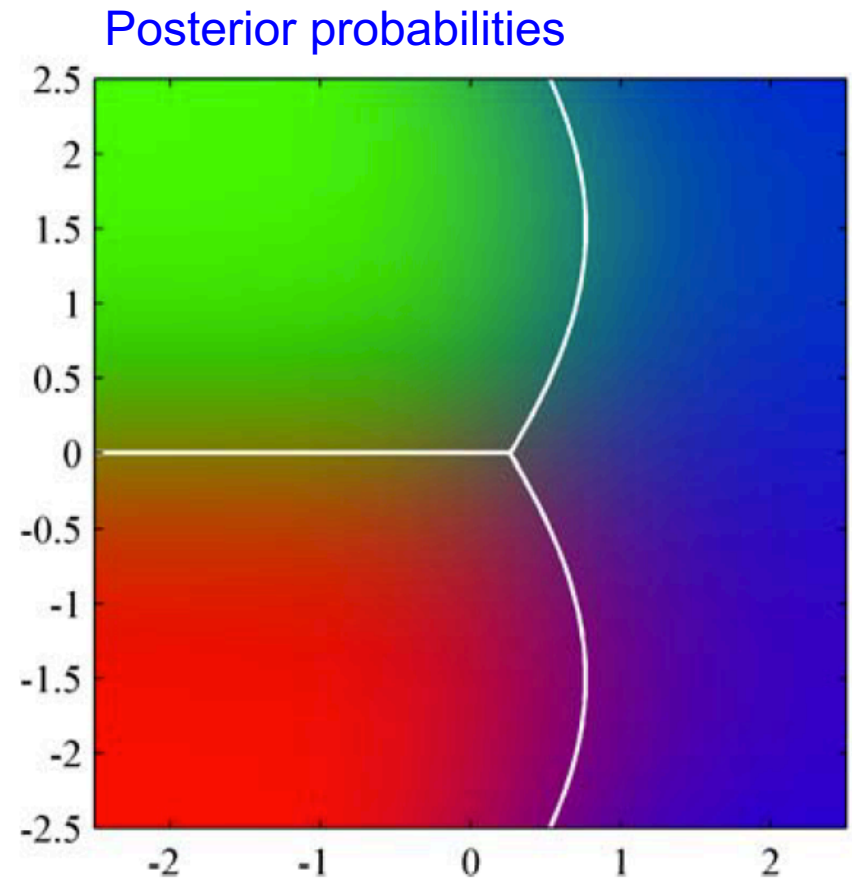
$$\begin{aligned}\mathbf{w}_k &= \Sigma^{-1} \boldsymbol{\mu}_k \\ w_{k0} &= -\frac{1}{2} \boldsymbol{\mu}_k^T \Sigma^{-1} \boldsymbol{\mu}_k + \ln p(\mathcal{C}_k)\end{aligned}$$



Probabilistic view



Class-conditional densities
for three classes

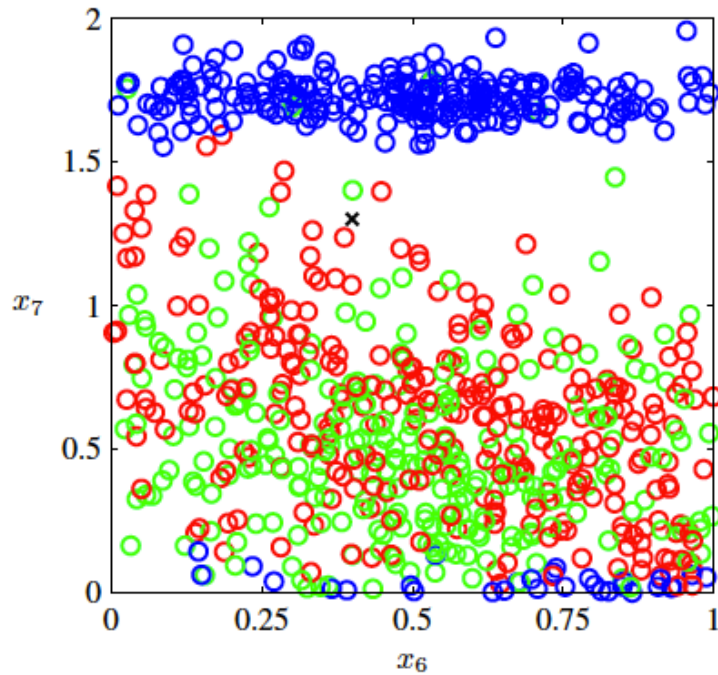


Fisher's linear discriminant

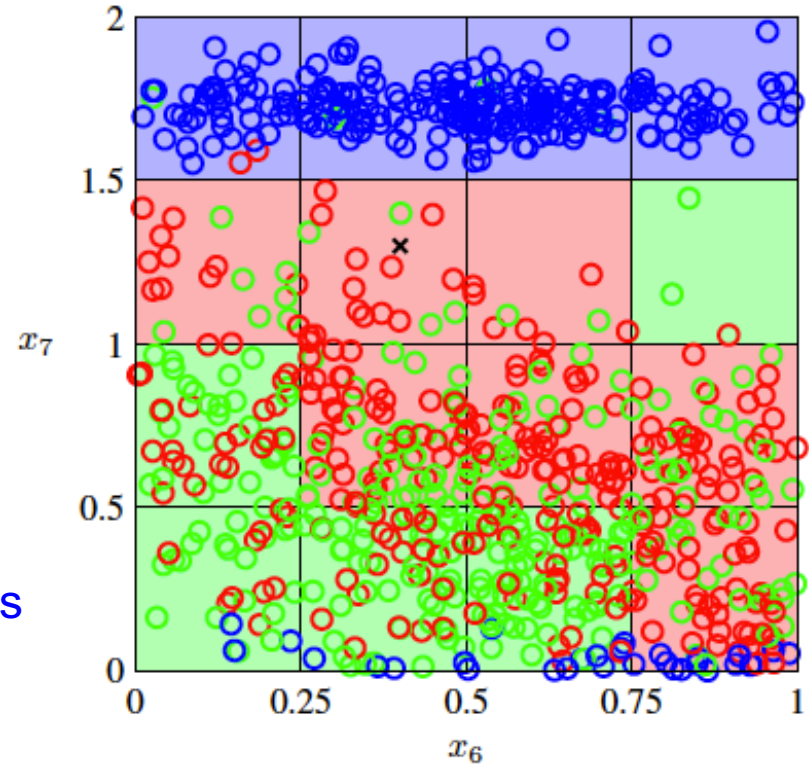
- Curse of dimensionality
 - The design of a good classifier becomes rapidly more difficult as the dimensionality of the input space increases
 - Pre-processing
 - To reduce its dimensionality
 - Fisher discriminant aims to achieve an optimal linear dimensionality reduction



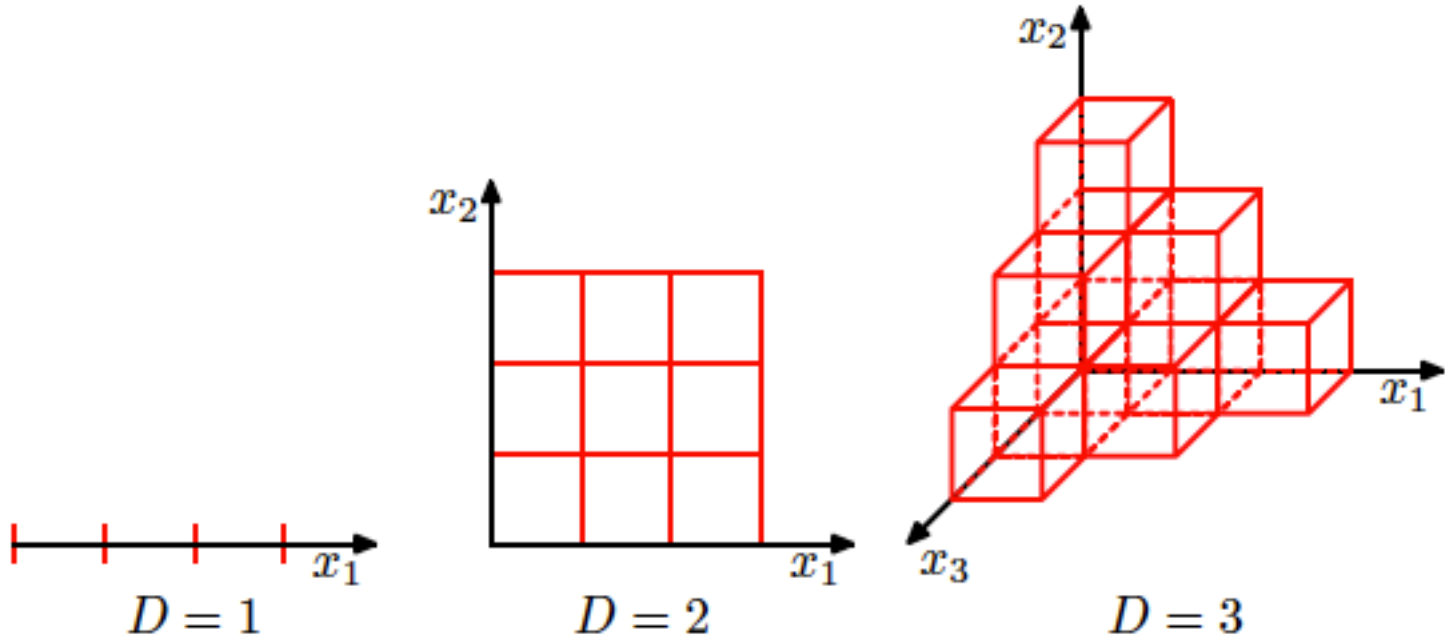
Fisher's linear discriminant



Classification by voting in the cells



Fisher's linear discriminant



The number of regions of a regular grid grows exponentially with the dimensionality of D



Fisher's linear discriminant

- Projection

$$y = \mathbf{w}^T \mathbf{x}, \quad y \geq -w_0 \text{ as class } \mathcal{C}_1$$

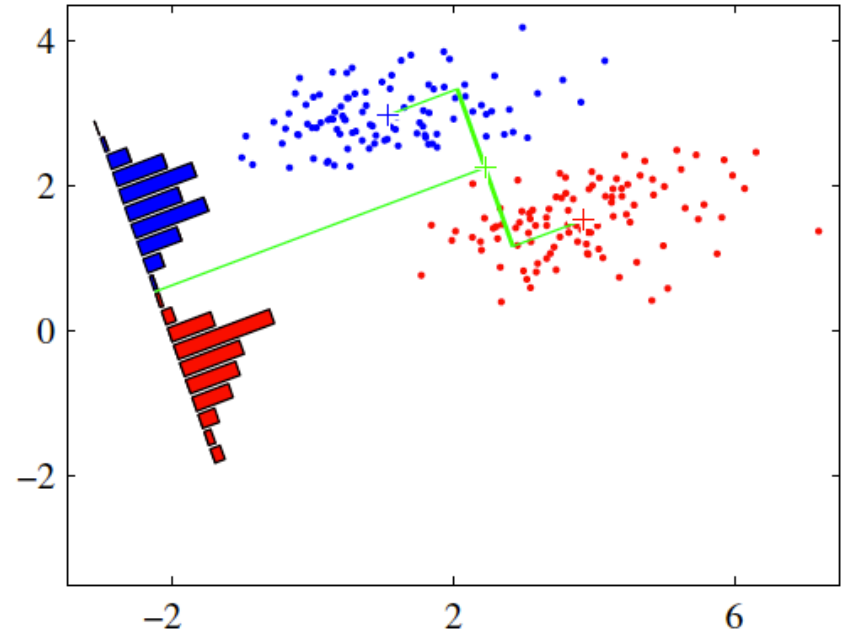
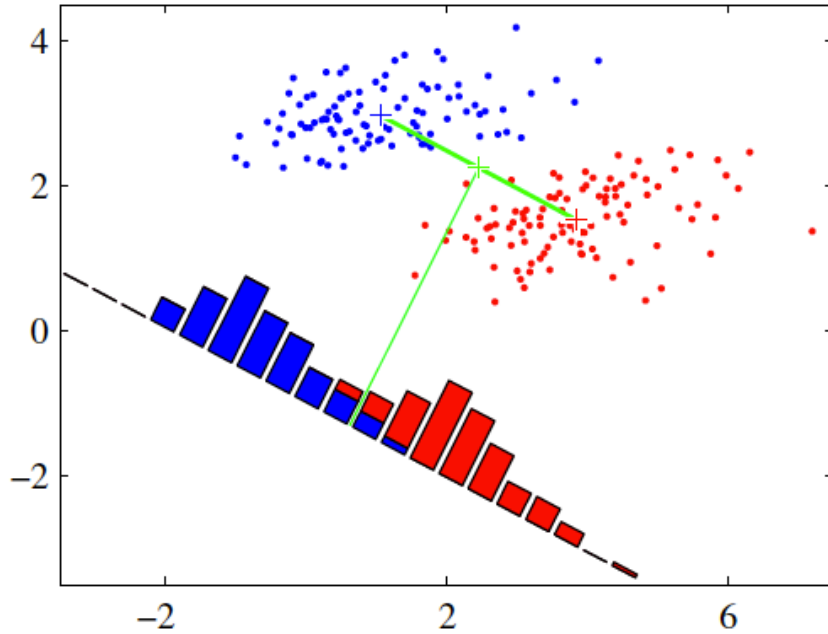
- Projection that maximizes the class separation

$$\mathbf{m}_1 = \frac{1}{N_1} \sum_{n \in \mathcal{C}_1} \mathbf{x}_n, \quad \mathbf{m}_2 = \frac{1}{N_2} \sum_{n \in \mathcal{C}_2} \mathbf{x}_n.$$

$$m_2 - m_1 = \mathbf{w}^T (\mathbf{m}_2 - \mathbf{m}_1)$$



Fisher's linear discriminant



Fisher directions



Fisher's linear discriminant

- within-class variance

$$s_k^2 = \sum_{n \in \mathcal{C}_k} (y_n - m_k)^2$$

- Fisher criterion

$$J(\mathbf{w}) = \frac{(m_2 - m_1)^2}{s_1^2 + s_2^2}$$

between-class variance

within-class variance



Fisher's linear discriminant

- Fisher criterion

$$J(\mathbf{w}) = \frac{\mathbf{w}^T \mathbf{S}_B \mathbf{w}}{\mathbf{w}^T \mathbf{S}_W \mathbf{w}}$$

- Between-class covariance

$$\mathbf{S}_B = (\mathbf{m}_2 - \mathbf{m}_1)(\mathbf{m}_2 - \mathbf{m}_1)^T$$

- within-class covariance

$$\mathbf{S}_W = \sum_{n \in \mathcal{C}_1} (\mathbf{x}_n - \mathbf{m}_1)(\mathbf{x}_n - \mathbf{m}_1)^T + \sum_{n \in \mathcal{C}_2} (\mathbf{x}_n - \mathbf{m}_2)(\mathbf{x}_n - \mathbf{m}_2)^T$$



Fisher's linear discriminant

- Differentiating with respect to weights

$$\underbrace{(\mathbf{w}^T \mathbf{S}_B \mathbf{w})}_{\text{scalar factor}} \mathbf{S}_W \mathbf{w} = \underbrace{(\mathbf{w}^T \mathbf{S}_W \mathbf{w})}_{\text{scalar factor}} \mathbf{S}_B \mathbf{w}$$

direction of $(\mathbf{m}_2 - \mathbf{m}_1)$

$$\mathbf{w} \propto \mathbf{S}_W^{-1} (\mathbf{m}_2 - \mathbf{m}_1)$$

Fisher linear discriminant

Generalization for more classes

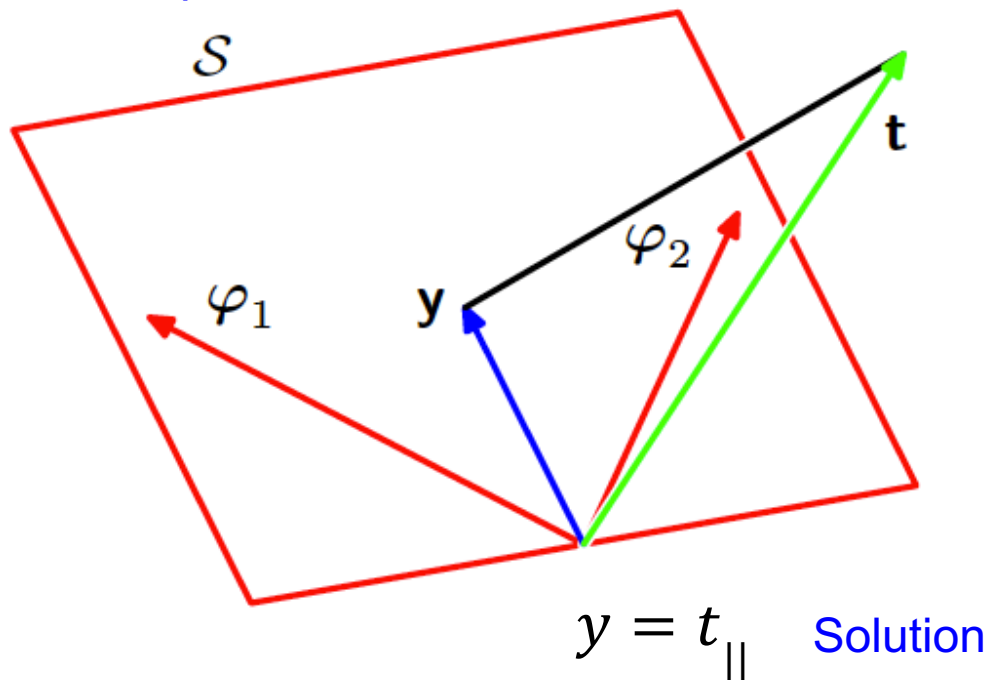


Sum-of squares error

■ Quadratic error function

$$E(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N \sum_{k=1}^c \{y_k(\mathbf{x}^n; \mathbf{w}) - t_k^n\}^2$$

Geometrical representation



$$y = \sum_{j=0}^M w_j \phi_j^n$$

$$\frac{\partial E}{\partial w_j} = 0 = \phi_j^T (y - t)$$



Gradient descent

- Error function

$$E(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N E^n(\mathbf{w})$$

- Stochastic gradient descent algorithm

$$\mathbf{w}^{(\tau+1)} = \mathbf{w}^{(\tau)} - \eta \nabla E^n$$

$$\mathbf{w}^{(\tau+1)} = \mathbf{w}^{(\tau)} - \eta \frac{\partial E^n}{\partial \mathbf{w}}$$



The perceptron algorithm

- Output

$$y(\mathbf{x}) = f(\mathbf{w}^T \phi(\mathbf{x}))$$

- Activation function

$$f(a) = \begin{cases} +1, & a \geq 0 \\ -1, & a < 0 \end{cases}$$



The perceptron algorithm

■ Classification

$$\begin{aligned} \mathbf{w}^T \phi(\mathbf{x}_n) &> 0 & \mathcal{C}_1 \\ \mathbf{w}^T \phi(\mathbf{x}_n) &< 0 & \mathcal{C}_2 \end{aligned} \quad t \in \{-1, +1\}$$

■ Perceptron criterion

$$E_P(\mathbf{w}) = - \sum_{n \in \mathcal{M}} \mathbf{w}^T \phi_n t_n$$



The perceptron algorithm

- Gradient descent algorithm

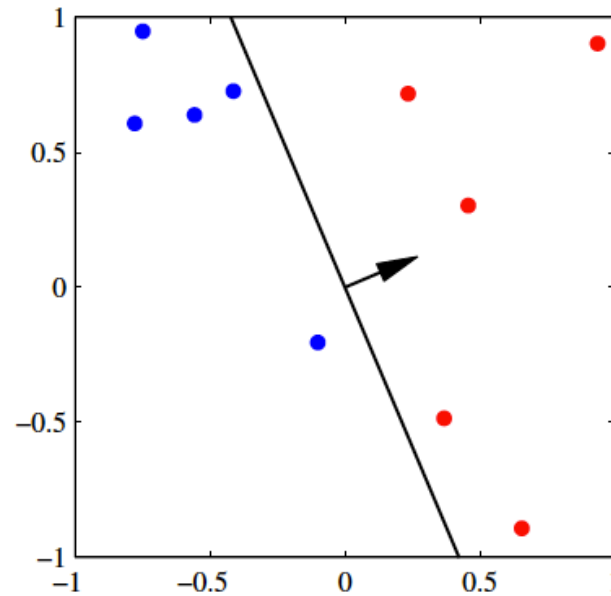
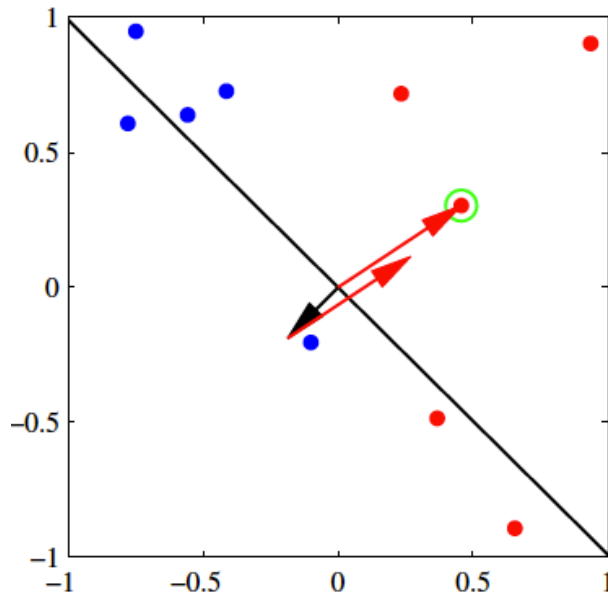
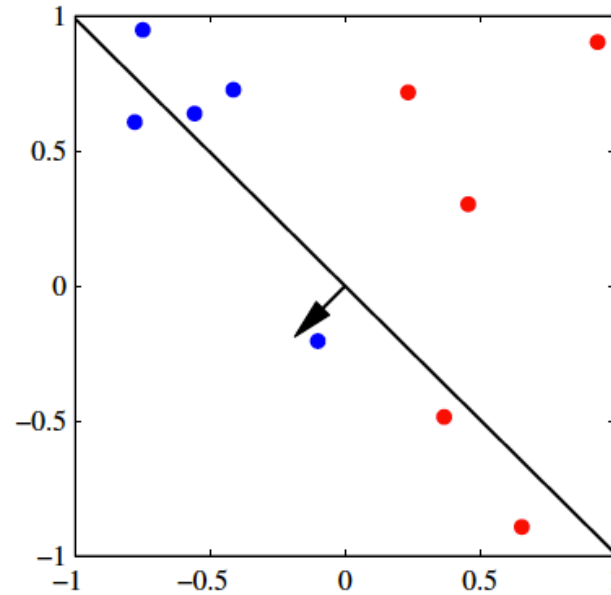
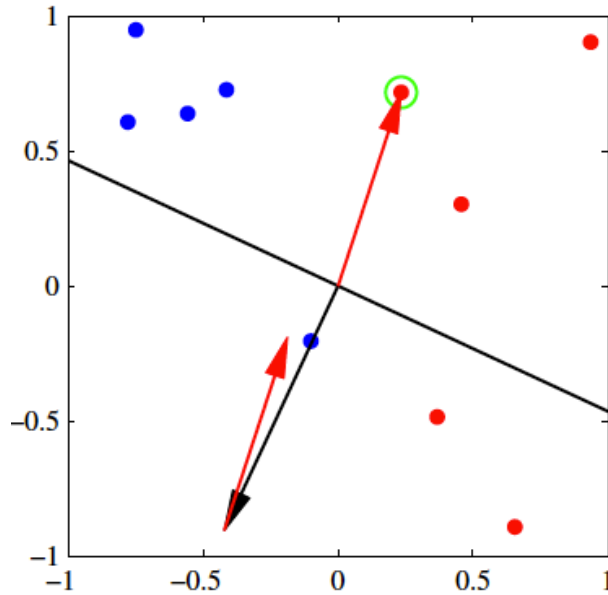
$$\mathbf{w}^{(\tau+1)} = \mathbf{w}^{(\tau)} - \eta \nabla E_P(\mathbf{w}) = \mathbf{w}^{(\tau)} + \eta \phi_n t_n$$

- Perceptron convergence theorem

- if there exists an exact solution (if the training data set is linearly separable) then the perceptron learning algorithm is guaranteed to find an exact solution in a finite number of steps



The perceptron algorithm



Linear separability

- Linearly separable
 - The points can be classified correctly by a linear decision boundary
- No linearly separable
 - exclusive-OR (XOR) problem

