

# Linear models for classification (supervised learning)

The goal in classification is to take an **input vector  $\mathbf{x}$**  and to assign it to **one of  $K$  discrete classes  $C_k$**  ( $k = 1, \dots, K$ )

The simplest models for classification are **linear in the input space**, i.e. the decision boundaries that separate the classes are linear functions of the input vector  $\mathbf{x}$  – they are  **$(D - 1)$ -dimensional hyperplanes within the  $D$ -dimensional input space  $\mathbf{x}$** .

- **Least squares:**  $y_k(\mathbf{x}) = \mathbf{w}_k^T \mathbf{x} + w_{k0}$

$$\text{Tr}[(\mathbf{XW} - \mathbf{T})^T(\mathbf{XW} - \mathbf{T})] = \sum_{\text{samples}} (\text{data}_{\text{sample}} * \text{weights} - \text{label}_{\text{sample}})^2$$

Not robust to outliers and can fail with multiple classes.

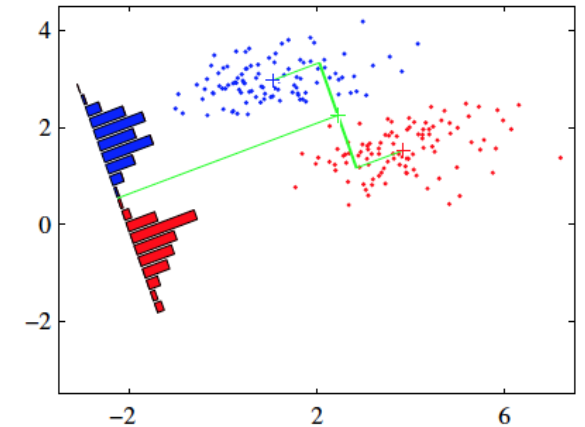
- **Fisher linear discriminant:** Minimize the ratio:

$$\frac{(\text{difference of means of projections})^2}{(\text{within class variance of projections})} = \frac{(m_2 - m_1)^2}{s_1^2 + s_2^2} = \frac{\text{between-class variance}}{\text{total within-class variance}}$$

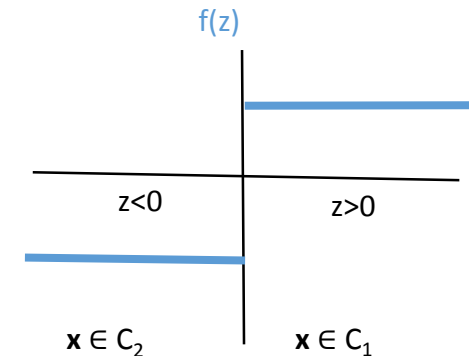
- **Perceptron algorithm.**  $y(\mathbf{x}) = f(\mathbf{w}^T * \boldsymbol{\varphi}(\mathbf{x}))$  Can be minimized with stochastic gradient descent: an iterative algorithm where, at iteration  $\tau$ ,

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

You don't know in advance the number of iterations needed, and it's hard to distinguish between very-slowly-converging solution and no possible solution.



Activation function

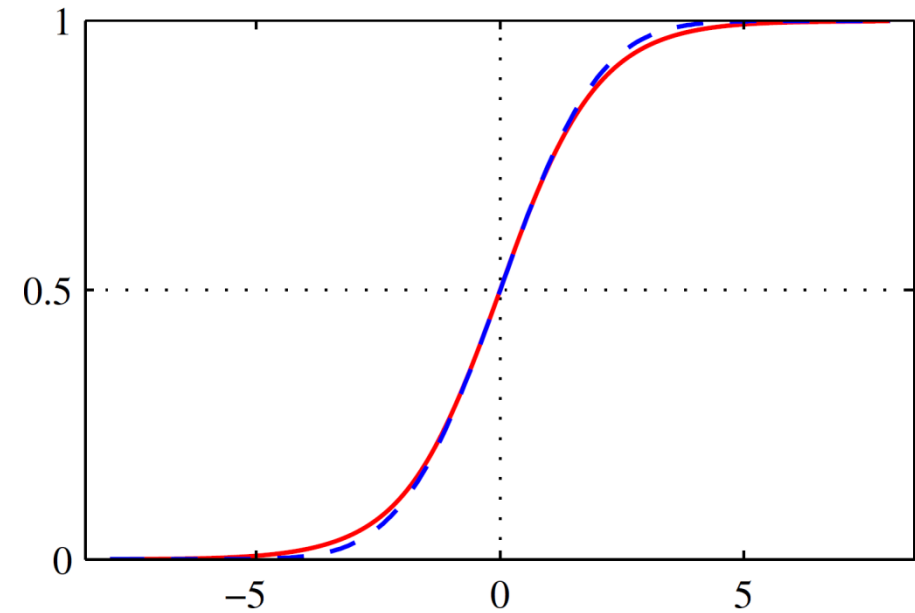


# Probabilistic Generative Models

Here the approach is to **find the parameters of a generalized linear model**, by fitting class-conditional densities and class priors separately and then applying Bayes' theorem to get  $P(C_k|\mathbf{x})$ . This represents an example of **generative modelling**, because we could take such a model and generate synthetic data by drawing values of  $\mathbf{x}$  from the marginal distribution  $p(\mathbf{x})$ .

Logistic function:

$$\begin{aligned} p(C_1|\mathbf{x}) &= \frac{p(\mathbf{x}|C_1)p(C_1)}{p(\mathbf{x}|C_1)p(C_1) + p(\mathbf{x}|C_2)p(C_2)} \\ &= \frac{1}{1 + \exp(-a)} = \sigma(a) \\ a &= \ln \frac{p(\mathbf{x}|C_1)p(C_1)}{p(\mathbf{x}|C_2)p(C_2)} \end{aligned}$$



Softmax function (for more than 2 classes):

$$p(C_k|\mathbf{x}) = \frac{p(\mathbf{x}|C_k)p(C_k)}{\sum_j p(\mathbf{x}|C_j)p(C_j)} = \frac{\exp(a_k)}{\sum_j \exp(a_j)} \quad a_k = \ln p(\mathbf{x}|C_k)p(C_k)$$

# Example

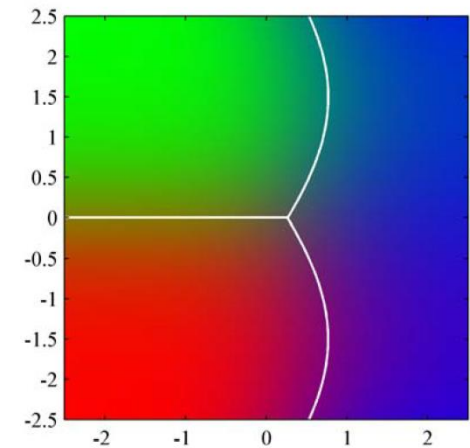
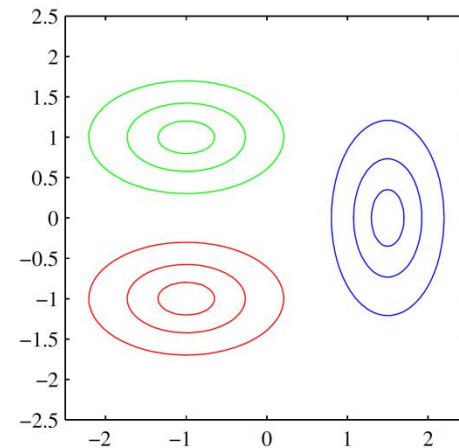
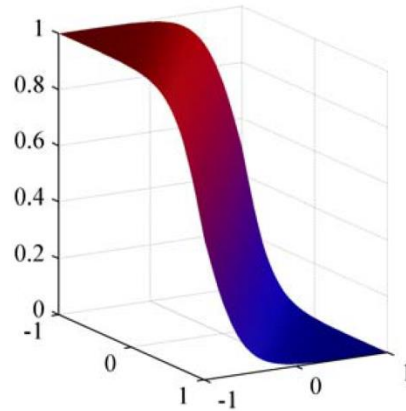
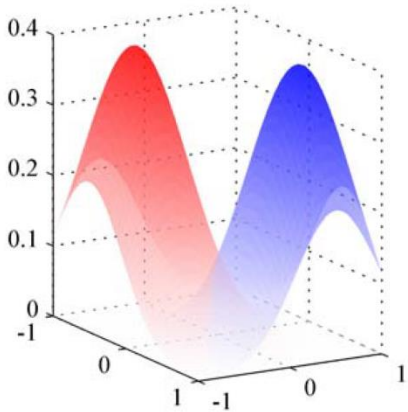
Assume that the **class-conditional densities are Gaussian** and then explore the resulting form for the posterior probabilities.

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

With:

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

$$\begin{aligned} \mathbf{w} &= \boldsymbol{\Sigma}^{-1}(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2) \\ w_0 &= -\frac{1}{2}\boldsymbol{\mu}_1^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_1 + \frac{1}{2}\boldsymbol{\mu}_2^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_2 + \ln \frac{p(\mathcal{C}_1)}{p(\mathcal{C}_2)}. \end{aligned}$$



the boundary between the red and green classes, which have the same covariance matrix, is linear, whereas those between the other pairs of classes are quadratic

# Maximum likelihood solution

## Continuous case:

Once we have specified a parametric functional form for the class-conditional densities  $p(\mathbf{x}/C_k)$ , we can then **determine the values of the parameters**, together with the prior class probabilities  $p(C_k)$ , using maximum likelihood.

Consider first the case of **two classes**, each having a **Gaussian class-conditional** density with a shared covariance matrix. The likelihood function is given by:

$$p(\mathbf{t}|\pi, \mu_1, \mu_2, \Sigma) = \prod_{n=1}^N [\pi \mathcal{N}(\mathbf{x}_n|\mu_1, \Sigma)]^{t_n} [(1 - \pi) \mathcal{N}(\mathbf{x}_n|\mu_2, \Sigma)]^{1-t_n}$$

Then we just need to set the derivative with respect to  $\pi, \mu_1, \mu_2$  and  $\Sigma$  to zero to determine the parameters.

# Maximum likelihood solution

## Discrete case:

Assume values  $x_i \in \{0, 1\} \rightarrow$  If there are  $D$  inputs, then a general distribution would correspond to  $2^D - 1$  independent variables, which grows exponentially with the number of features. Therefore we make the **naive Bayes assumption** in which the feature values are treated as independent, conditioned on the class  $C_k$ .

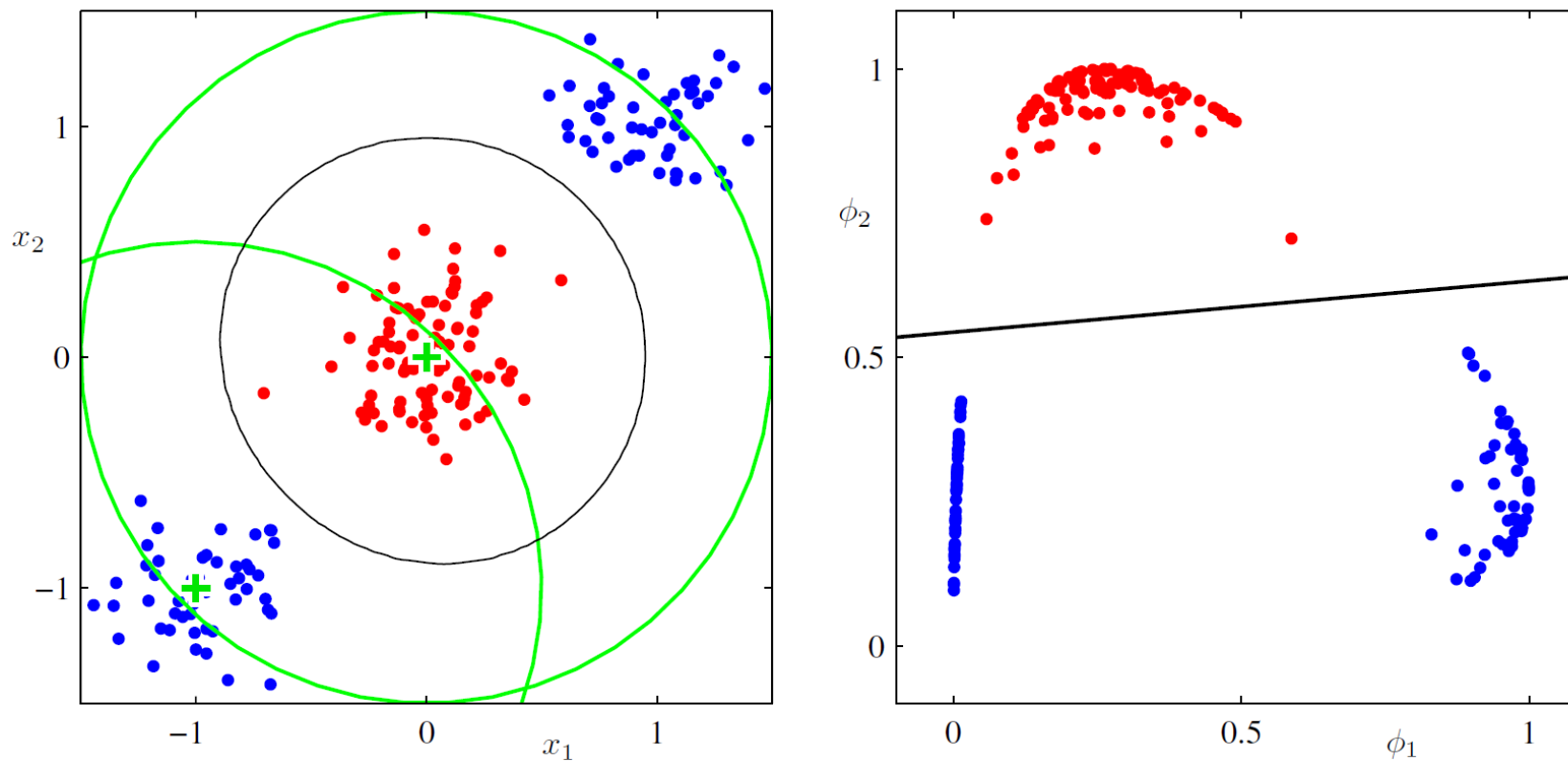
$$p(\mathbf{x}|\mathcal{C}_k) = \prod_{i=1}^D \mu_{ki}^{x_i} (1 - \mu_{ki})^{1-x_i}$$

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

$$a_k(\mathbf{x}) = \sum_{i=1}^D \{x_i \ln \mu_{ki} + (1 - x_i) \ln(1 - \mu_{ki})\} + \ln p(\mathcal{C}_k)$$

$$p(\mathcal{C}_k|\mathbf{x}) = \frac{\exp(a_k)}{\sum_j \exp(a_j)}$$

# Basis Function



$$\mathbf{x} \longrightarrow \phi(\mathbf{x})$$

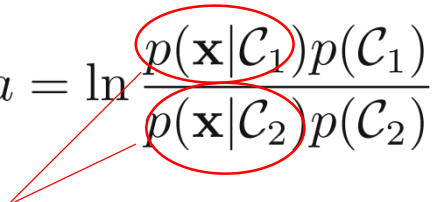
# Probabilistic Discriminative Models

An alternative approach is to **use the functional form of the generalized linear model explicitly** and to determine its parameters directly by using maximum likelihood.

In this more direct approach, we are maximizing a likelihood function defined through the conditional distribution  $p(C_k|x)$ , which represents a form of **discriminative training**. One advantage of this discriminative approach is that there will typically be **fewer adaptive parameters to be determined**, as we shall see shortly.

Generative approach

$$\begin{aligned} p(C_1|x) &= \frac{p(x|C_1)p(C_1)}{p(x|C_1)p(C_1) + p(x|C_2)p(C_2)} \\ &= \frac{1}{1 + \exp(-a)} = \sigma(a) \end{aligned}$$

$$a = \ln \frac{p(x|C_1)p(C_1)}{p(x|C_2)p(C_2)}$$


We need to model these

Discriminative model approach

$$p(C_1|\phi) = y(\phi) = \sigma(\mathbf{w}^T \phi)$$

# Logistic regression

$$p(\mathcal{C}_1|\phi) = y(\phi) = \sigma(\mathbf{w}^T \phi)$$

Likelihood: 
$$p(\mathbf{t}|\mathbf{w}) = \prod_{n=1}^N y_n^{t_n} \{1 - y_n\}^{1-t_n}$$

$$\mathbf{t} = (t_1, \dots, t_N)^T \text{ and } y_n = p(\mathcal{C}_1|\phi_n)$$

$$\phi_n = \phi(\mathbf{x}_n)$$

We can define an error function by taking the negative logarithm of the likelihood, which gives the cross-entropy error function in the form:

$$E(\mathbf{w}) = -\ln p(\mathbf{t}|\mathbf{w}) = -\sum_{n=1}^N \{t_n \ln y_n + (1 - t_n) \ln(1 - y_n)\}$$

Taking the gradient of the error function with respect to  $\mathbf{w}$ , we obtain:

$$\nabla E(\mathbf{w}) = \sum_{n=1}^N (y_n - t_n) \phi_n$$



# Iterative reweighted least squares

For logistic regression, there is no closed-form solution as in Linear Regression, due to the nonlinearity of the logistic sigmoid function. However, the error function can be minimized by an efficient iterative technique based on the *Newton-Raphson* iterative optimization scheme, that uses a **local quadratic approximation** to the log likelihood function:

$$\mathbf{w}^{(\text{new})} = \mathbf{w}^{(\text{old})} - \mathbf{H}^{-1} \nabla E(\mathbf{w})$$

where  $\mathbf{H}$  is the Hessian matrix whose elements comprise the second derivatives of  $E(\mathbf{w})$  with respect to the components of  $\mathbf{w}$ . The update formula is (see section 4.3.3 in Bishop's book):

$$\mathbf{w}^{(\text{new})} = (\Phi^T \mathbf{R} \Phi)^{-1} \Phi^T \mathbf{R} \mathbf{z}$$

Where

$$\mathbf{z} = \Phi \mathbf{w}^{(\text{old})} - \mathbf{R}^{-1} (\mathbf{y} - \mathbf{t})$$

And  $\mathbf{R}$  is a  $N \times N$  diagonal matrix with elements:

$$R_{nn} = y_n(1 - y_n)$$

where  $y_n = \sigma(a_n)$  and  $a_n = \mathbf{w}^T \phi_n$

# Multiclass logistic regression

In the case of more than 2 classes, the posterior probabilities can be often modeled by a **softmax transformation** of linear functions of the feature variables, so that

$$p(\mathcal{C}_k|\phi) = y_k(\phi) = \frac{\exp(a_k)}{\sum_j \exp(a_j)}$$

where the ‘activations’  $a_k$  are given by

$$a_k = \mathbf{w}_k^T \phi.$$

The *cross-entropy* error function for the multiclass classification problem is given by

$$E(\mathbf{w}_1, \dots, \mathbf{w}_K) = -\ln p(\mathbf{T}|\mathbf{w}_1, \dots, \mathbf{w}_K) = -\sum_{n=1}^N \sum_{k=1}^K t_{nk} \ln y_{nk}$$

And its gradient (as in the logistic regression case) is

$$\nabla_{\mathbf{w}_j} E(\mathbf{w}_1, \dots, \mathbf{w}_K) = \sum_{n=1}^N (y_{nj} - t_{nj}) \phi_n$$