

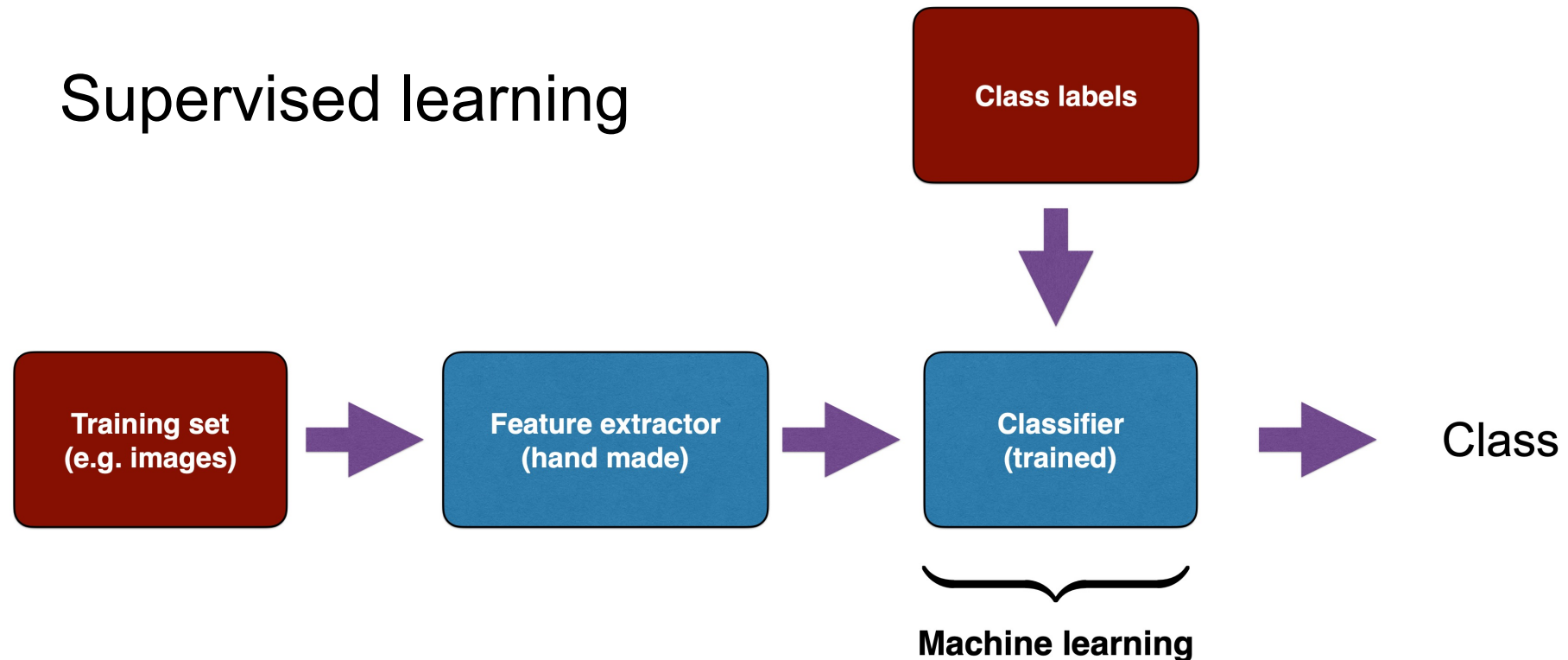
# Introduction to Machine Learning

Idar Dyrdal



# Machine learning (Pattern recognition)

- Recognition of individuals (instance recognition)
- Discrimination between classes (pattern recognition, classification)



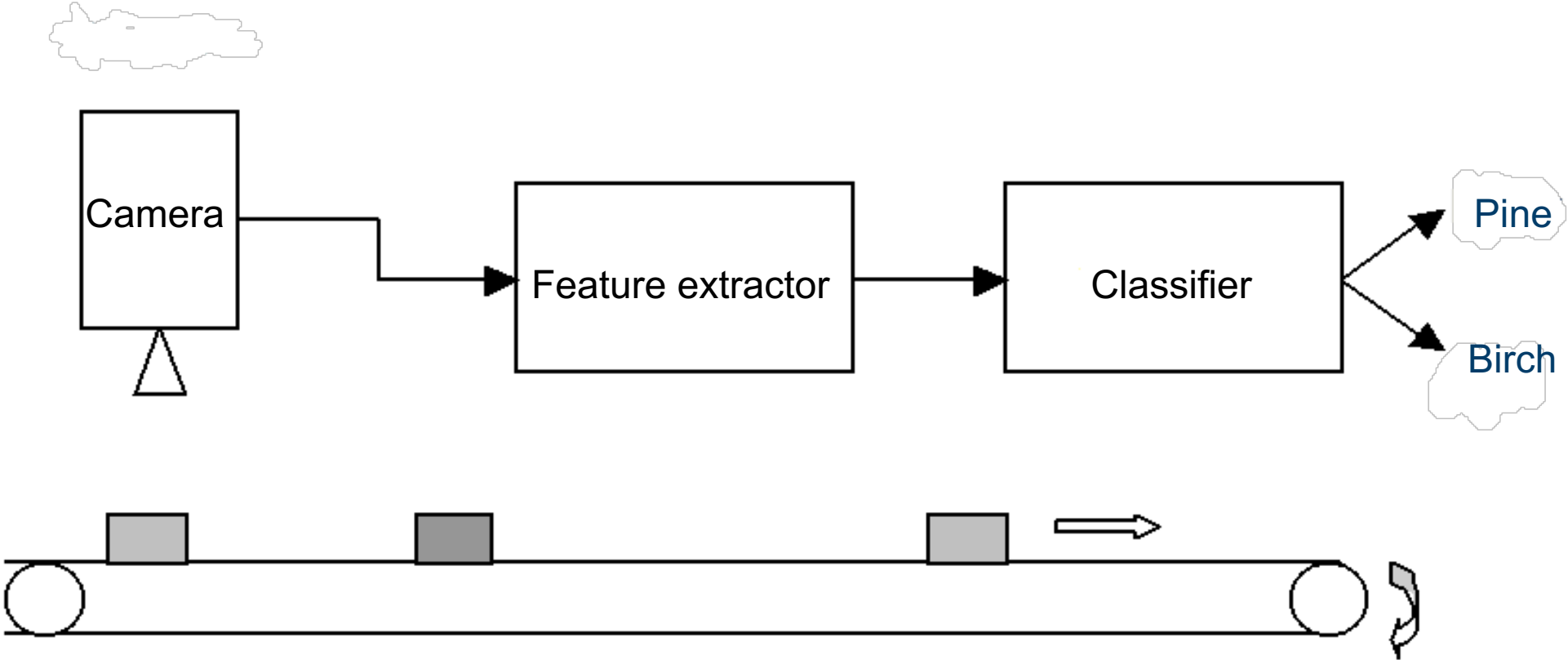
# Pattern recognition in practice

## Working applications of Image Pattern recognition:

- Reading license plates, postal codes, bar codes
- Character recognition
- Automatic diagnosis of medical samples
- Fingerprint recognition
- Face detection and recognition
- ...



# Classification system



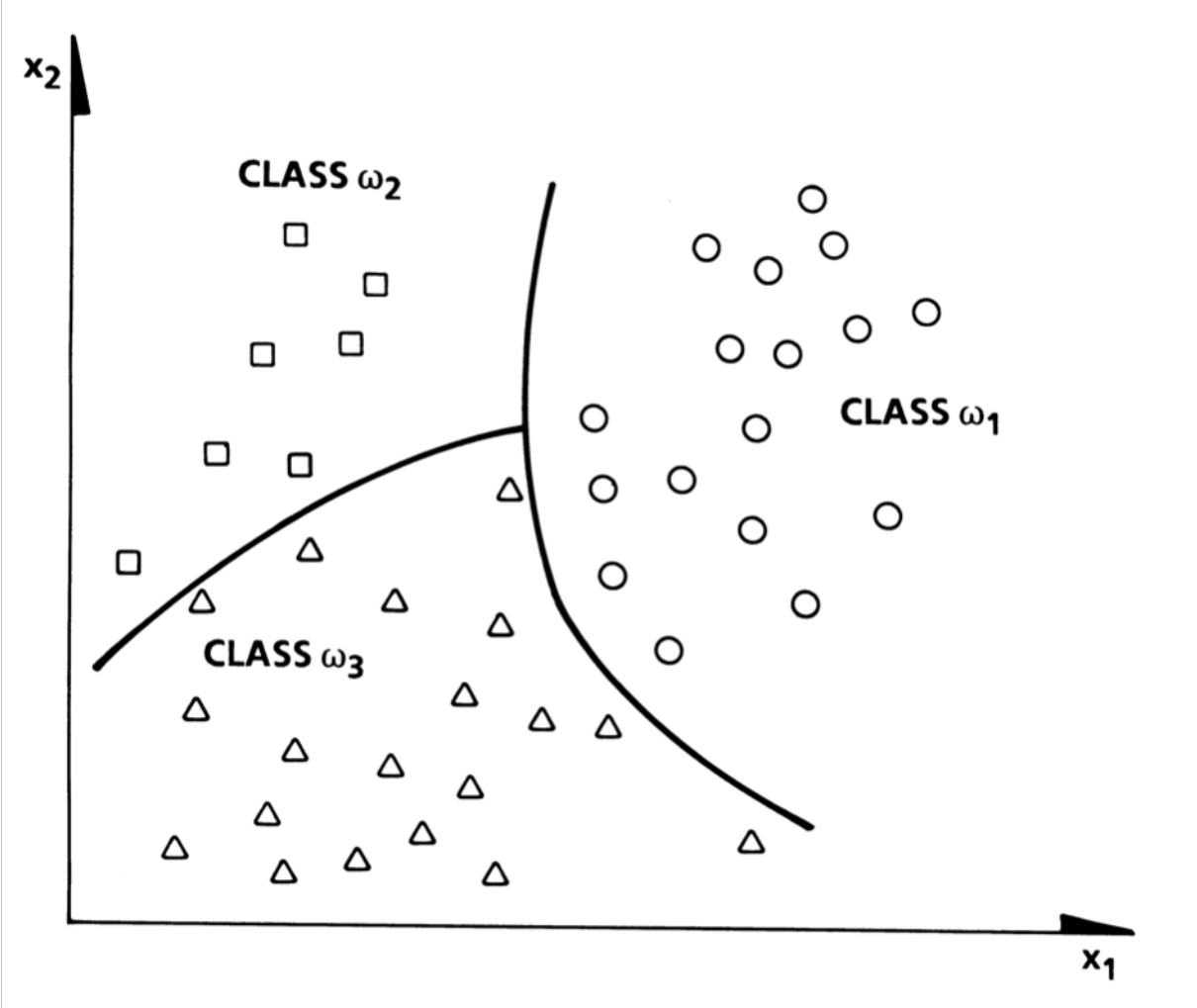
# Image features for object recognition



$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ \vdots \\ x_d \end{bmatrix}$$

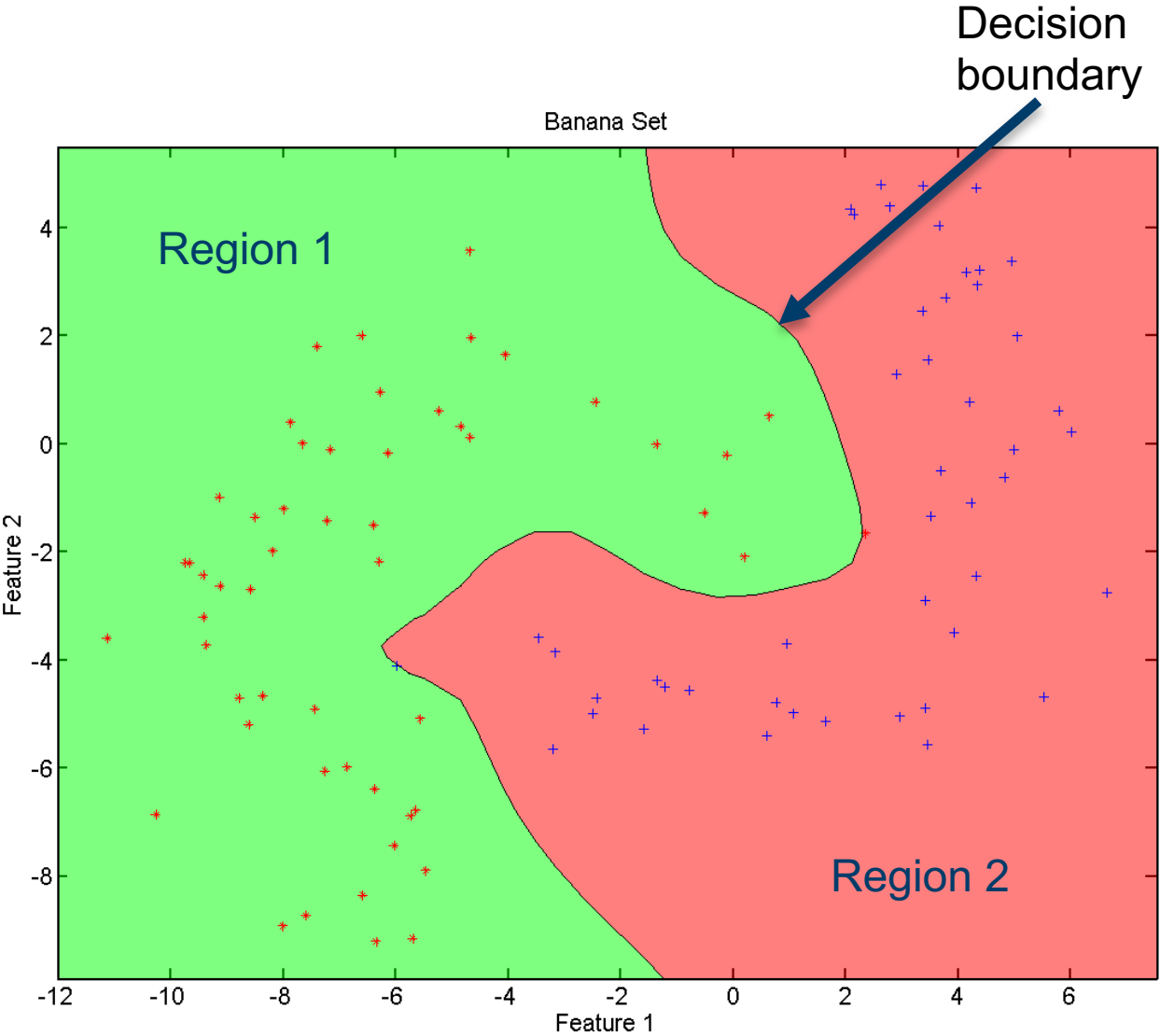
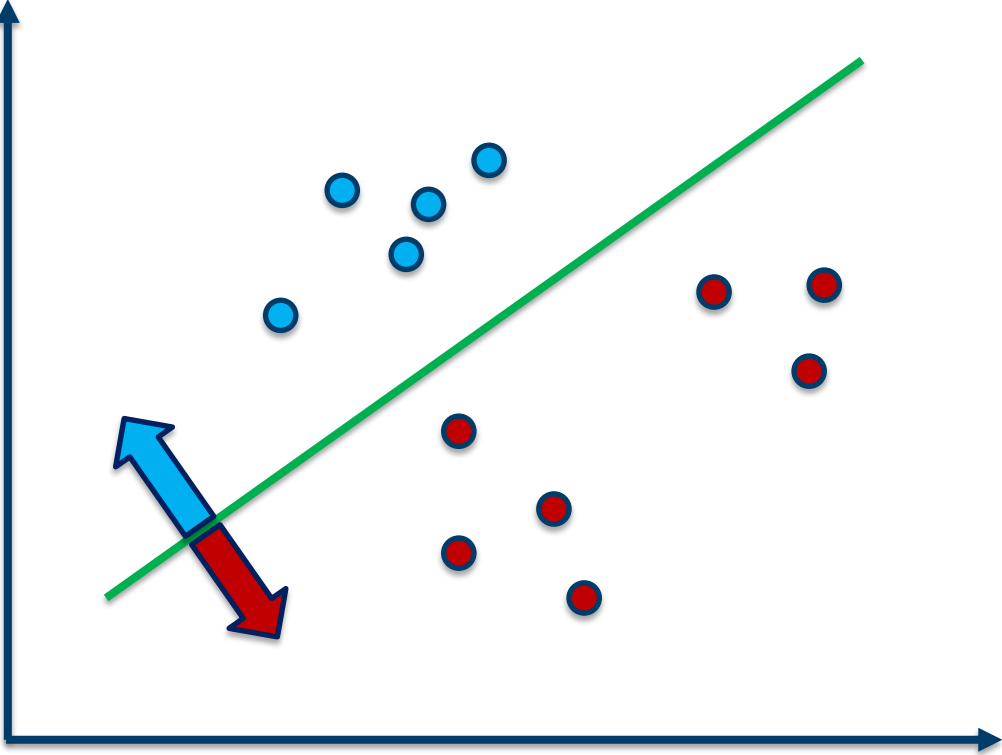
# Feature vector and feature space

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ \vdots \\ x_d \end{bmatrix}$$



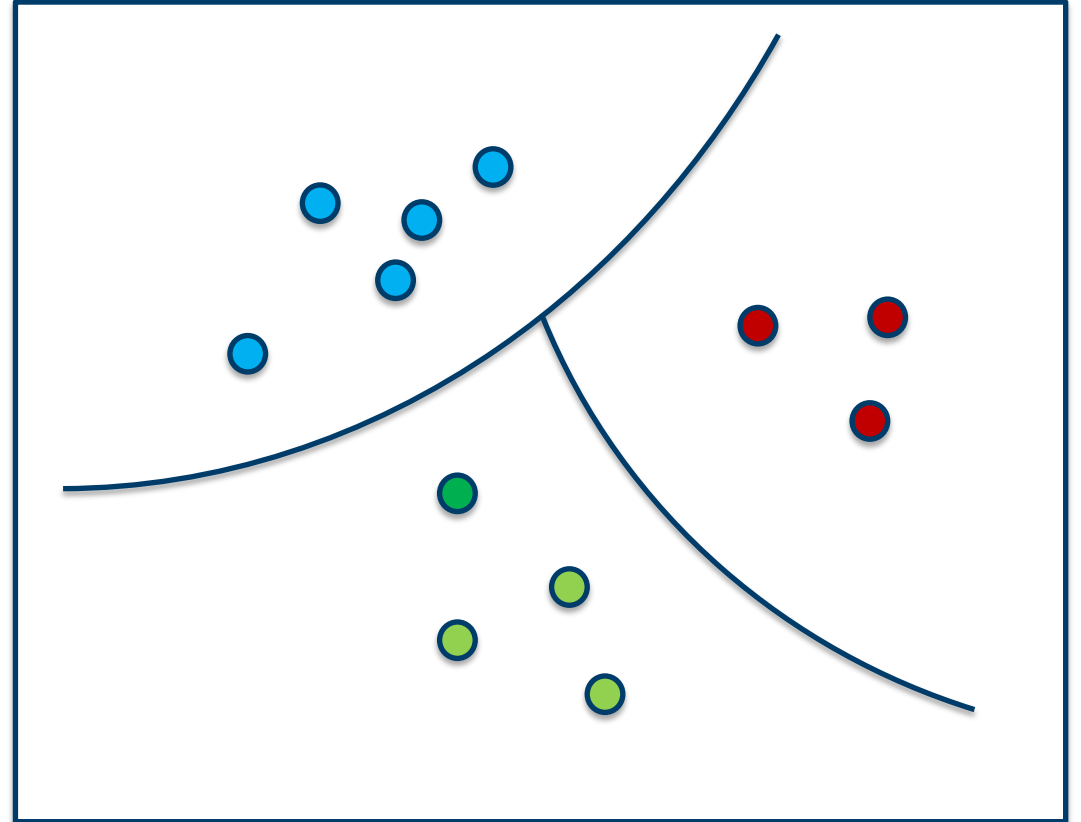
# Training of classifiers

Learn a function to predict the class from the given features



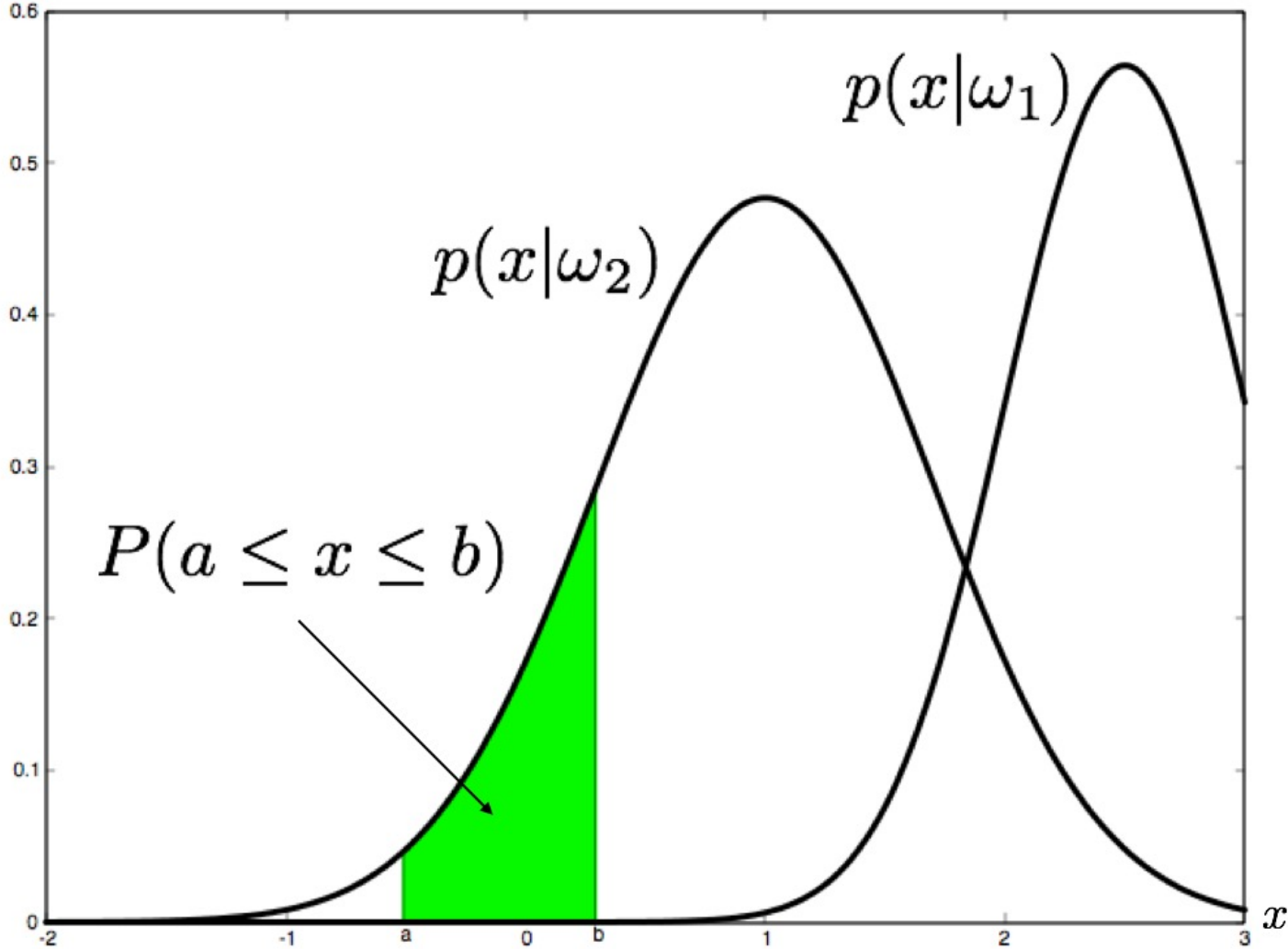
# Classifiers and training methods

- Bayes (parametric) classifier
- Nearest-neighbors and K-nearest-neighbors
- Parzen windows
- Linear and higher order discriminant functions
- Neural nets
- Support Vector Machines (SVM)
- Decision trees
- Random forest
- ...





# Class conditional probability density functions



# Bayesian decision theory

## Overview

Class conditional densities:

$$p(\mathbf{x}|\omega_i), \text{ for each class } \omega_1, \omega_2, \dots, \omega_c$$

Prior probabilities:

$$P(\omega_1), P(\omega_2), \dots, P(\omega_c)$$

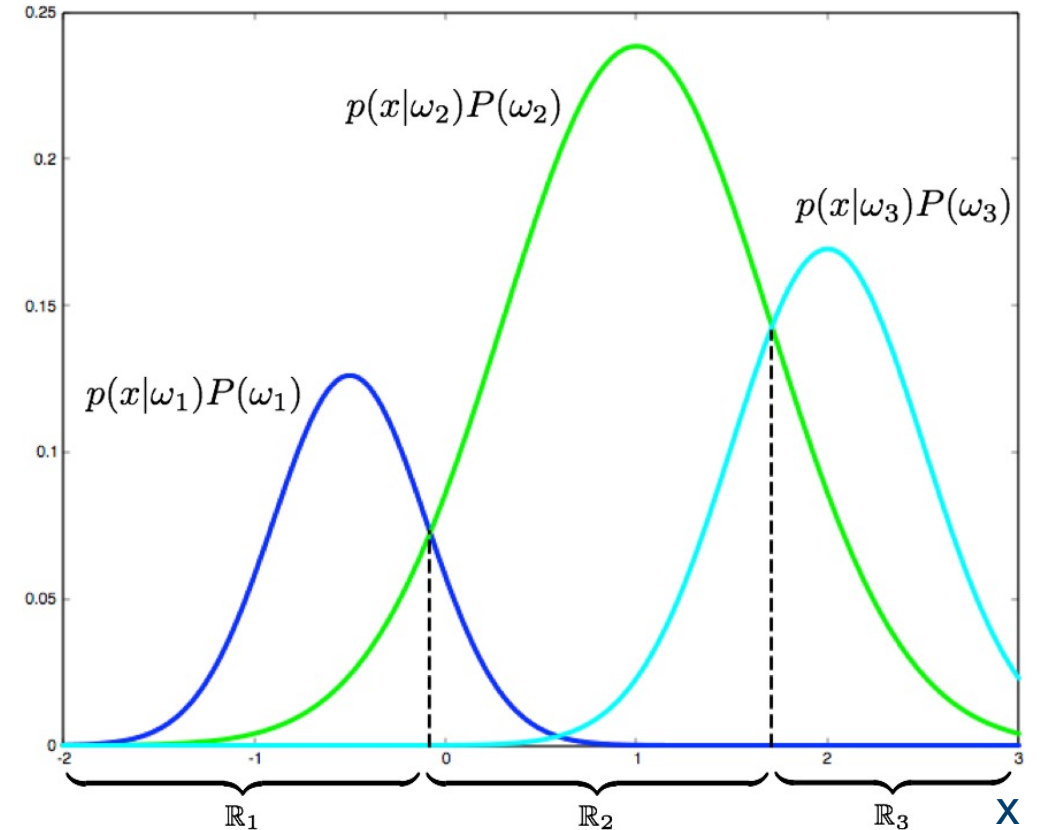
Posterior probabilities given by Bayes rule:

$$P(\omega_i|\mathbf{x}) = \frac{p(\mathbf{x}|\omega_i)P(\omega_i)}{\sum_{j=1}^c p(\mathbf{x}|\omega_j)P(\omega_j)}, i = 1, \dots, c$$

(a function of the measured feature vector  $\mathbf{x} = [x_1, x_2, \dots, x_d]^t$ ).

Minimum error rate classification:

*Assign the unknown object to the class with maximum posterior probability!*



# Density estimation

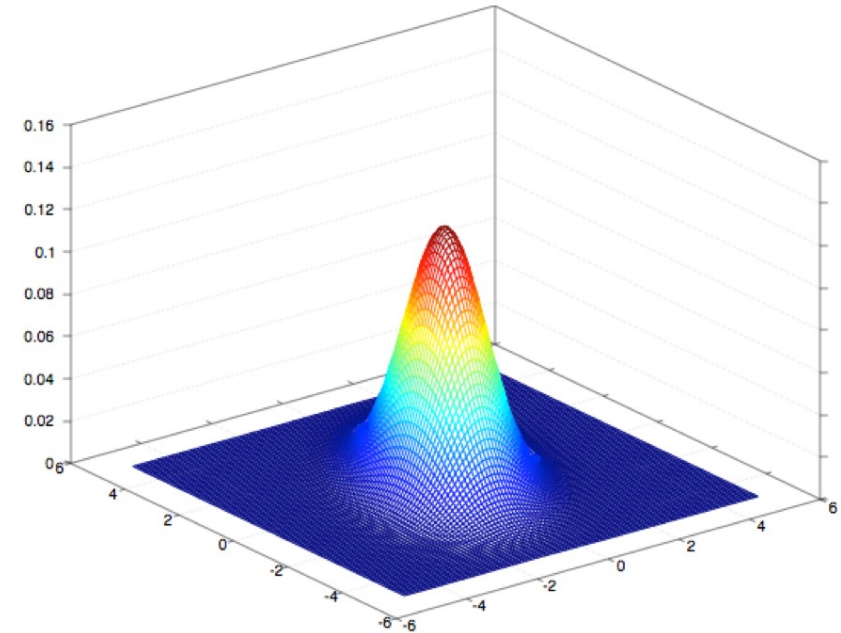
## Example – Gaussian distribution:

### Parametric methods:

- Assume a given shape of the density function
- Use the training set to estimate the unknown parameters.

### Non-parametric (distribution free) methods:

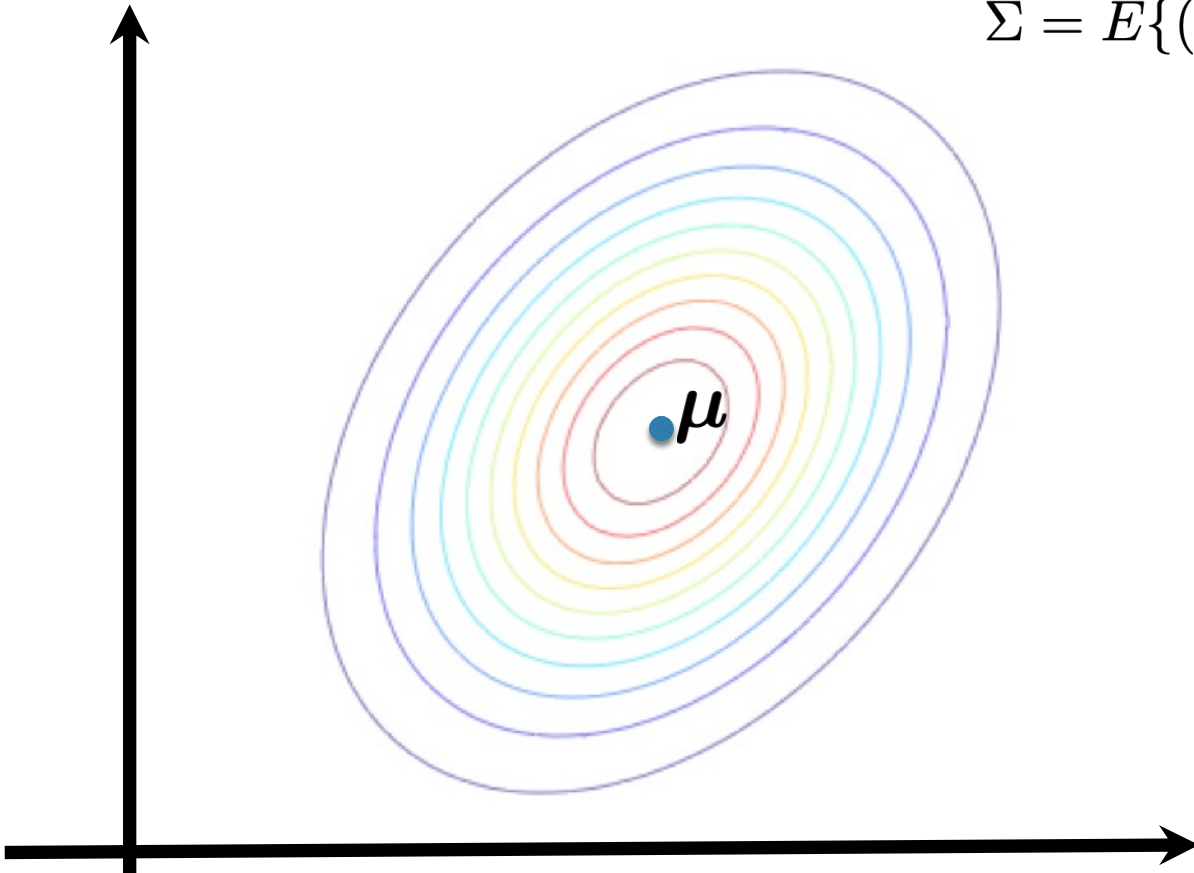
- Point estimation of the density using the training set directly
- Parzen windows
- Nearest neighbor estimation (leads directly to the nearest-neighbor and k-nearest-neighbor classifiers).



$$p(\mathbf{x}|\omega_i) = \frac{1}{(2\pi)^{\frac{d}{2}} |\Sigma_i|^{\frac{1}{2}}} \exp \left[ -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_i)^t \Sigma_i^{-1} (\mathbf{x} - \boldsymbol{\mu}_i) \right]$$

Parameters:  $\boldsymbol{\mu}_i$  and  $\Sigma_i$

# Parameter estimation



$$\Sigma = E\{(\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^t\} = \begin{bmatrix} \sigma_{11} & \sigma_{12} & \dots & \sigma_{1d} \\ \vdots & \vdots & & \vdots \\ \sigma_{d1} & \sigma_{d2} & \dots & \sigma_{dd} \end{bmatrix}$$

Parameter estimates:

$$\hat{\boldsymbol{\mu}} = \mathbf{m} = \frac{1}{n} \sum_{k=1}^n \mathbf{x}_k$$

$$\hat{\Sigma} = \frac{1}{n} \sum_{k=1}^n (\mathbf{x}_k - \mathbf{m})(\mathbf{x}_k - \mathbf{m})^t$$

# Discriminant functions

Estimate of the density in a given point:

$$\hat{p}(\mathbf{x}|\omega_i) = \frac{1}{(2\pi)^{\frac{d}{2}} |\hat{\Sigma}_i|^{\frac{1}{2}}} \exp \left[ -\frac{1}{2} (\mathbf{x} - \hat{\boldsymbol{\mu}}_i)^t \hat{\Sigma}_i^{-1} (\mathbf{x} - \hat{\boldsymbol{\mu}}_i) \right]$$

From Bayes rule:

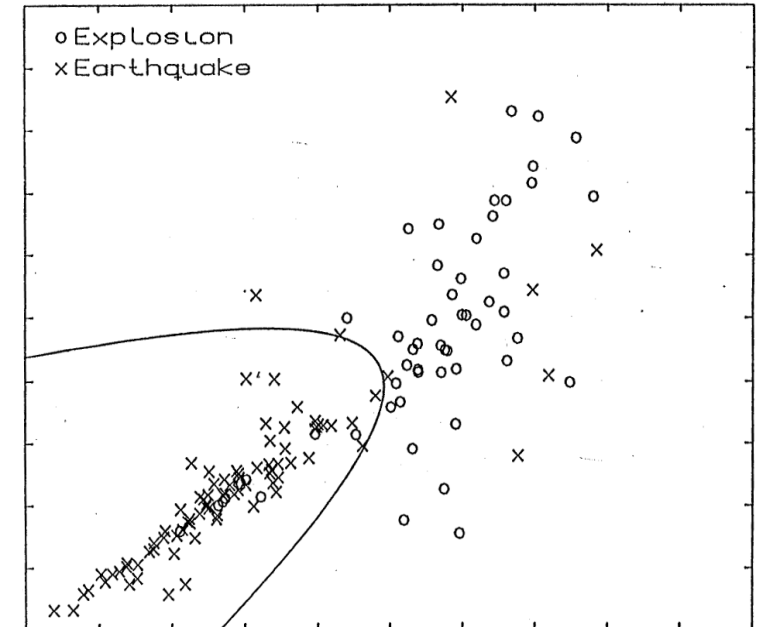
$$\hat{P}(\omega_i|\mathbf{x}) = \frac{\hat{p}(\mathbf{x}|\omega_i)P(\omega_i)}{\sum_{j=1}^c \hat{p}(\mathbf{x}|\omega_j)P(\omega_j)}$$

Examples of discriminant functions:

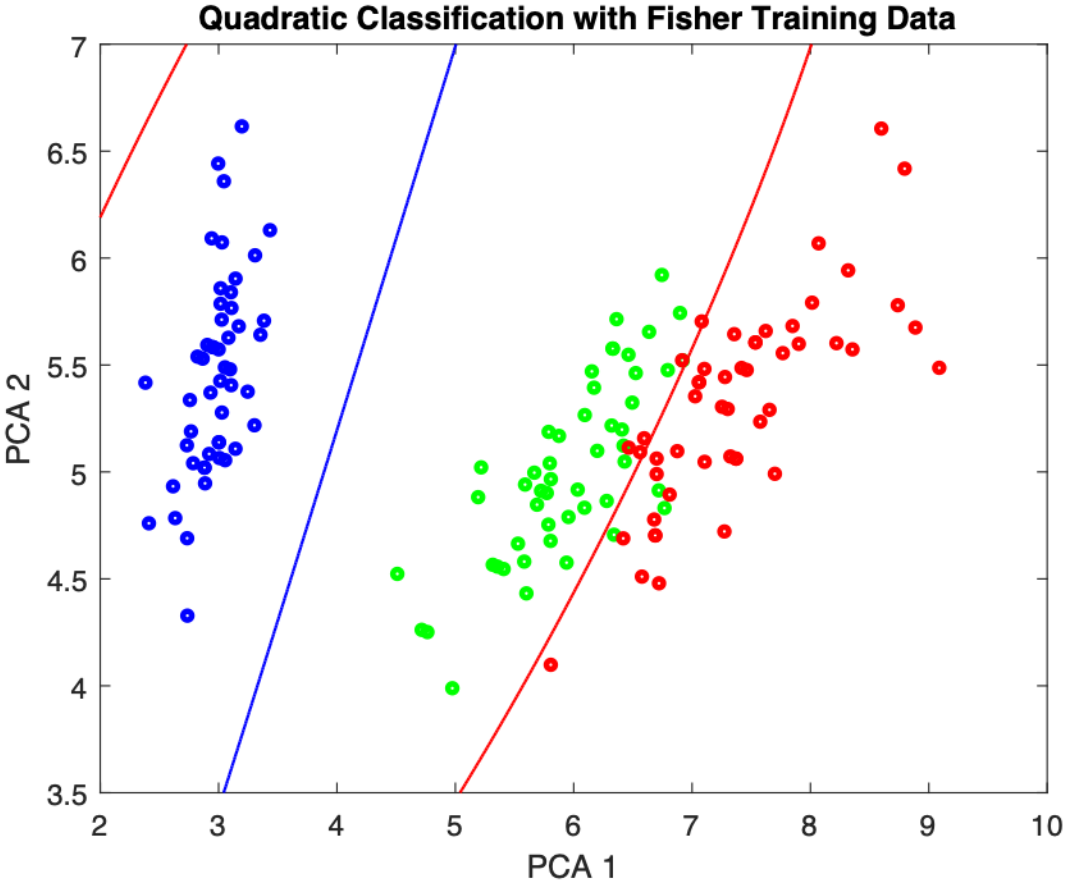
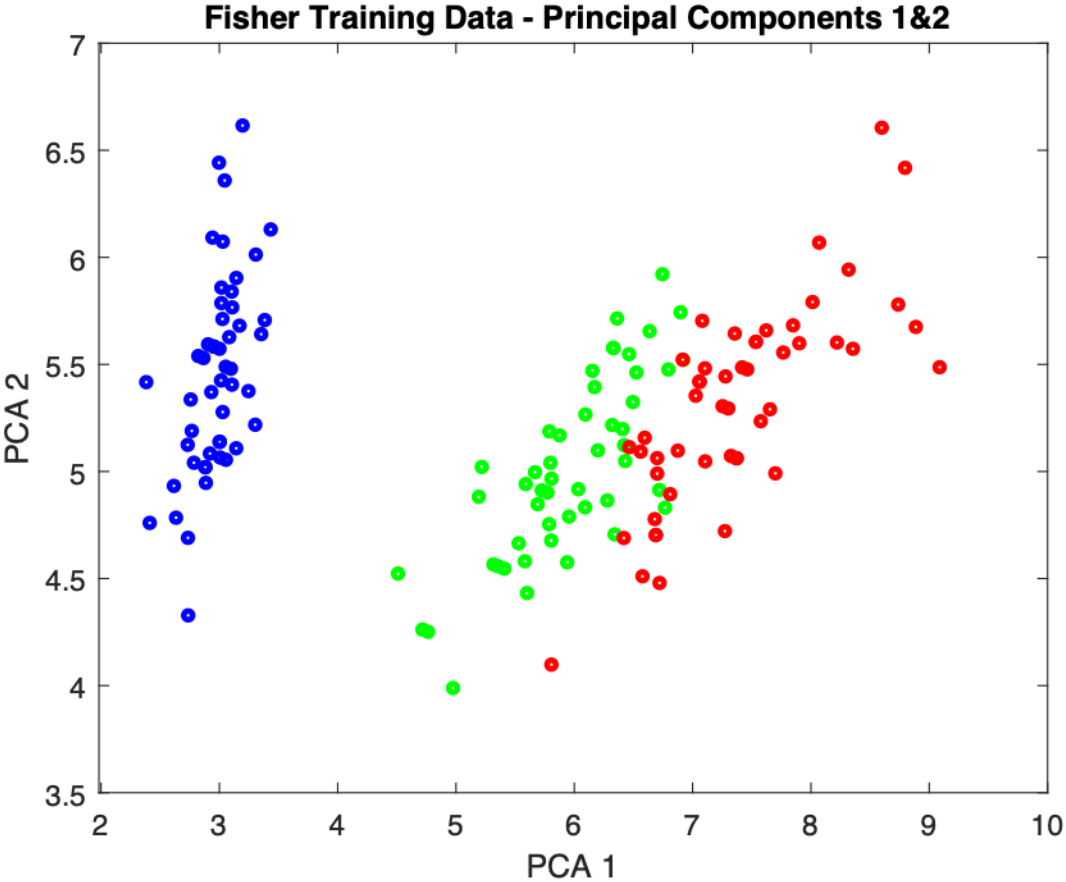
$$g_i(\mathbf{x}) = \ln \hat{P}(\omega_i|\mathbf{x}) \quad \text{or} \quad g_i(\mathbf{x}) = \ln \hat{p}(\mathbf{x}|\omega_i) + \ln P(\omega_i)$$

Decision rule:

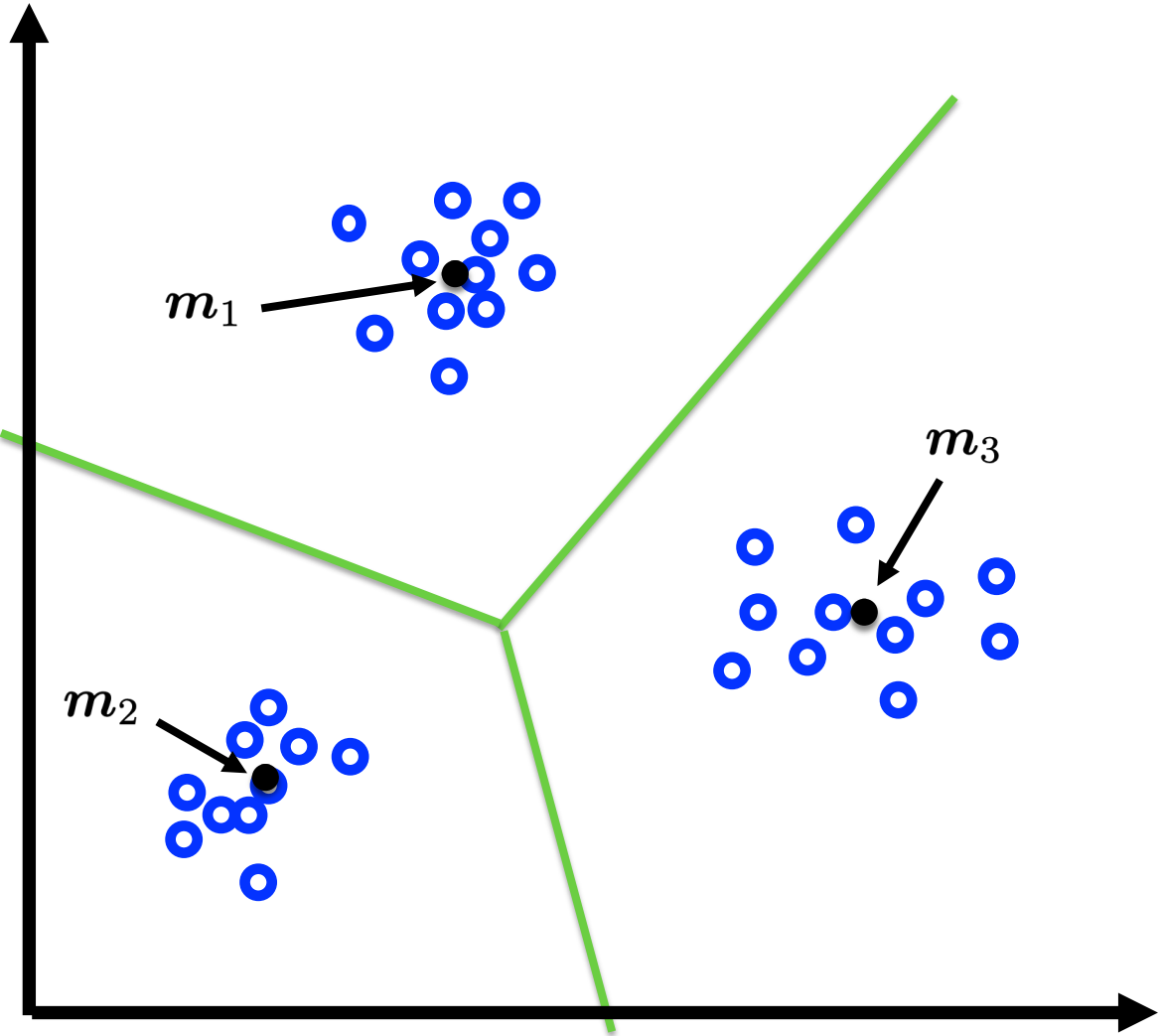
*Choose the class with maximum discriminant function value.*



# Quadratic classifier - example



# Linear classifier



Example:

Uncorrelated features and  
common covariance matrices



Linear decision boundaries

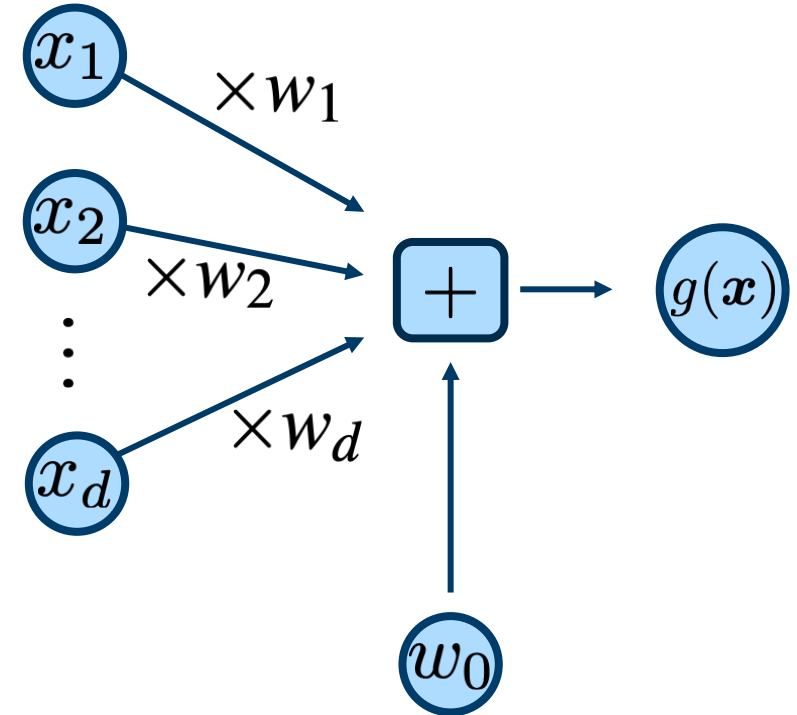
# Linear classifier (contd.)

Discriminant function:

$$g(\mathbf{x}) = \mathbf{w}^t \mathbf{x} + w_0 = w_0 + \sum_{i=1}^d w_i x_i$$

Can be rewritten as:

$$g(\mathbf{x}) = [w_0, w_1, \dots, w_d] \begin{bmatrix} 1 \\ x_1 \\ \vdots \\ x_d \end{bmatrix} = \mathbf{a}^t \mathbf{y}$$

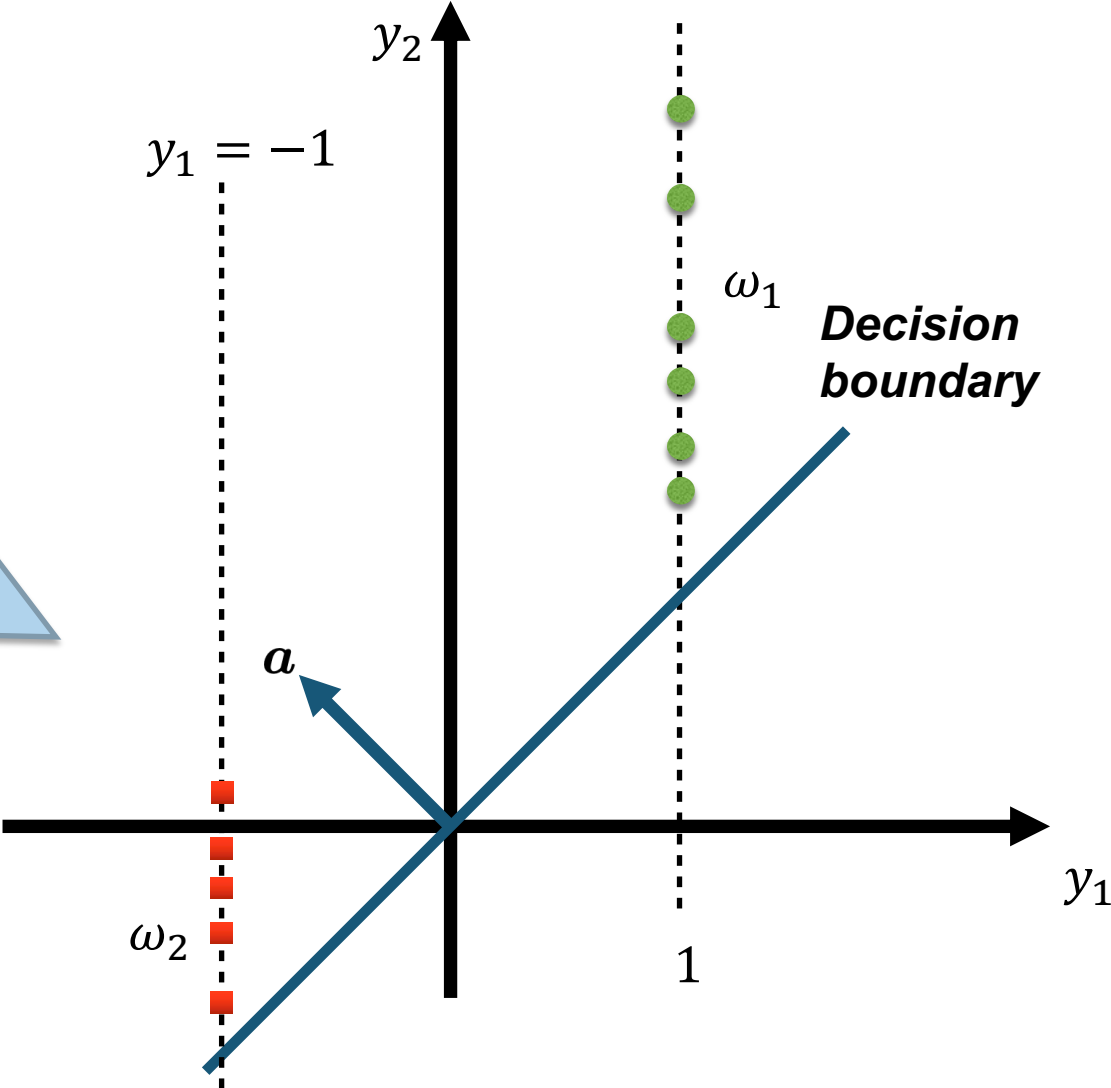
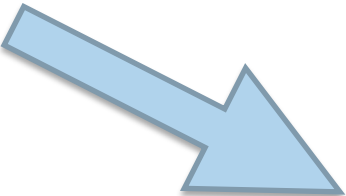
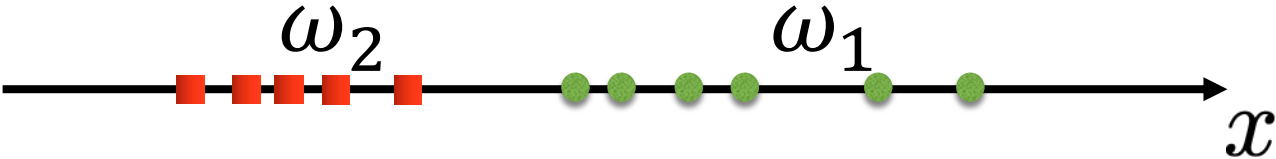


Decision rule (two-class case):

Decide  $\omega_1$  if  $\mathbf{a}^t \mathbf{y} > 0$  and  $\omega_2$  if  $\mathbf{a}^t \mathbf{y} \leq 0$



# Augmented feature space



Augmented weight vector and feature vector:

$$\mathbf{a} = \begin{bmatrix} w_0 \\ w_1 \\ \vdots \\ w_d \end{bmatrix} \quad \text{and} \quad \mathbf{y} = \begin{bmatrix} 1 \\ x_1 \\ \vdots \\ x_d \end{bmatrix}$$

Normalized samples in augmented feature space.

# Gradient descent

Find the minimum of a criterion function:

$$J(\mathbf{a}) \geq 0$$

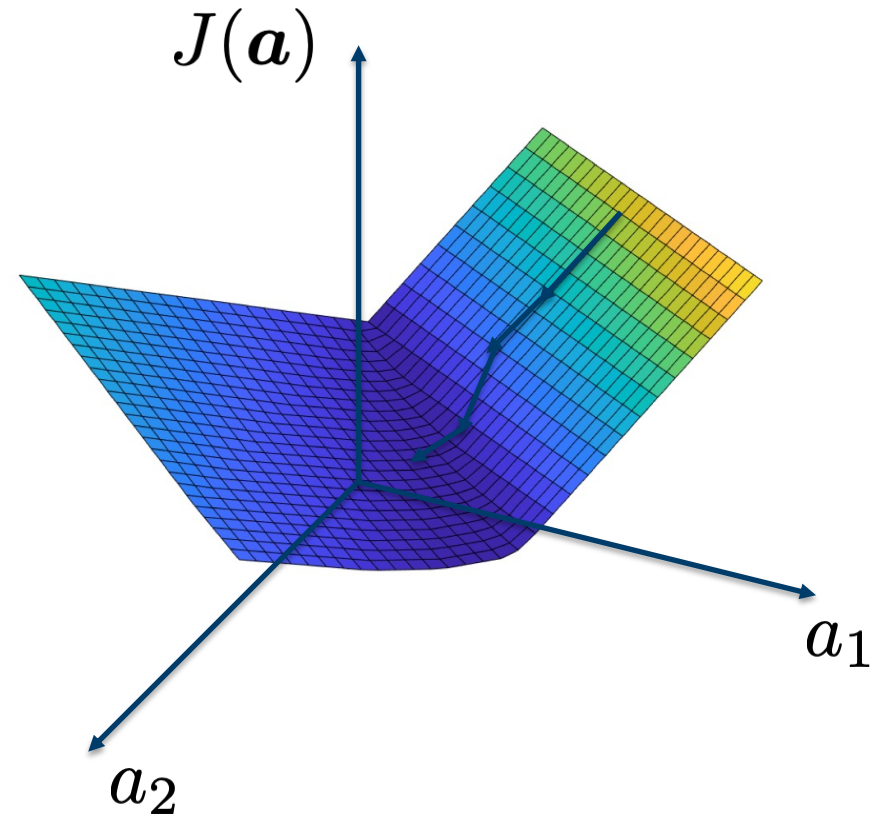
Basic algorithm:

$$\mathbf{a}_1 = \text{arbitrary}$$

$$\mathbf{a}_{k+1} = \mathbf{a}_k - \rho_k \nabla J(\mathbf{a}_k), \quad k = 1, 2, \dots$$

Optimal step length (increment):

$$\rho_k = \frac{\|\nabla J\|^2}{\nabla J^t D \nabla J} \quad \text{where} \quad D_{ij} = \frac{\partial^2 J(\mathbf{a}_k)}{\partial a_j \partial a_k} \quad \text{are the components of matrix } D$$



# Example

Perceptron criterion function:

$$J_p(\mathbf{a}) = - \sum_{\mathbf{y} \in \mathcal{Y}} \mathbf{a}^t \mathbf{y} \quad \text{where } \mathcal{Y} = \{\mathbf{y} : \mathbf{a}^t \mathbf{y} \leq 0\}$$

The goal is to satisfy a set of inequalities:

$$\mathbf{a}^t \mathbf{y}_i > 0 \quad \forall \mathbf{y}_i \in \{\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_n\}$$

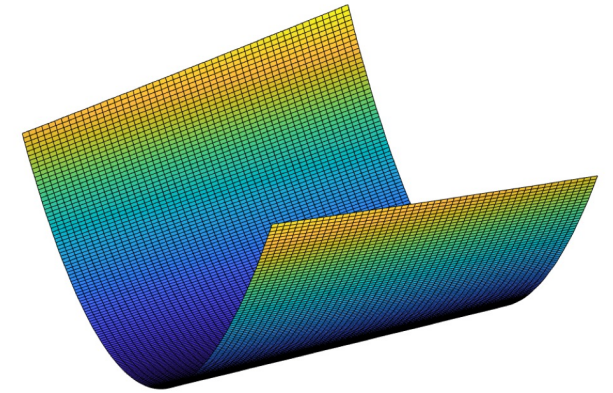
# Linear least-squares optimization

Solve a set of linear equations

$$\mathbf{a}^t \mathbf{y}_i = b_i \quad \text{where } b_i > 0, \quad i = 1, \dots, n \quad (\text{positive margins})$$

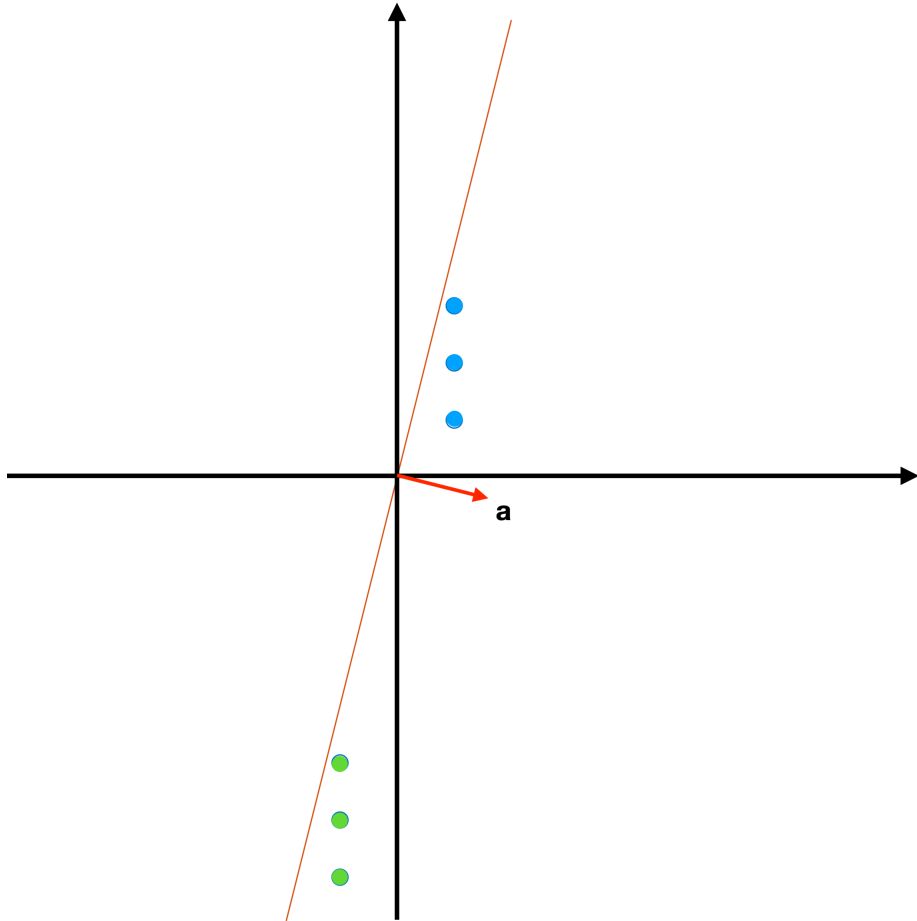
Find the minimum of a least-squares criterion function

$$J_s(\mathbf{a}) = \sum_{i=1}^n (\mathbf{a}^t \mathbf{y}_i - b_i)^2$$



Can be solved by gradient descent or by solving the *normal equations*.

# Linear least-squares optimization (contd.)



The criterion function may be written as

$$J_s(\mathbf{a}) = \|\mathbf{e}\|^2 = \|\mathbf{Y}\mathbf{a} - \mathbf{b}\|^2$$

where

$$\mathbf{Y} = \begin{bmatrix} \mathbf{y}_1^t \\ \vdots \\ \mathbf{y}_n^t \end{bmatrix} \quad \text{and} \quad \mathbf{b} = \begin{bmatrix} b_1 \\ \vdots \\ b_n \end{bmatrix}$$

Zero gradient of the criterion function leads to

$$\mathbf{a} = (\mathbf{Y}^t\mathbf{Y})^{-1}\mathbf{Y}^t\mathbf{b}$$

# Artificial Neural Network (ANN)

Used in Machine Learning and Pattern Recognition:

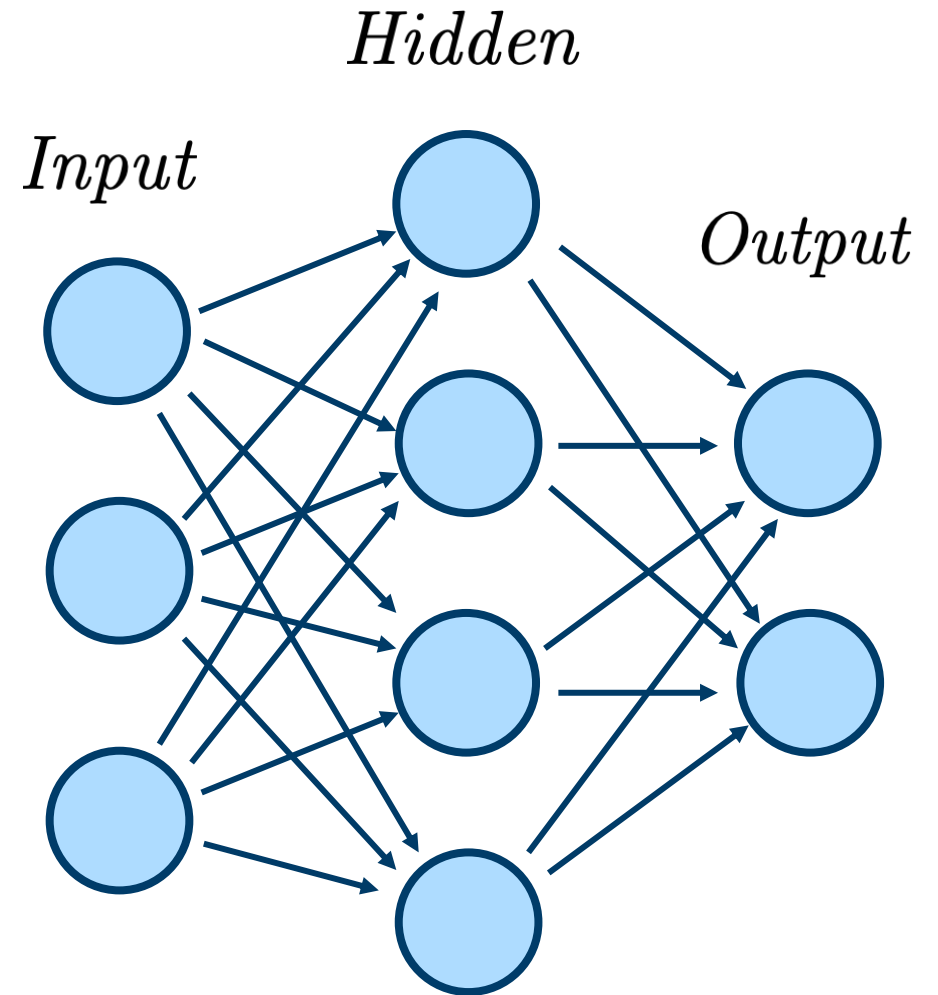
- Regression
- Classification
- Clustering
- ...

Applications:

- Speech recognition
- Recognition of handwritten text
- Image classification
- ...

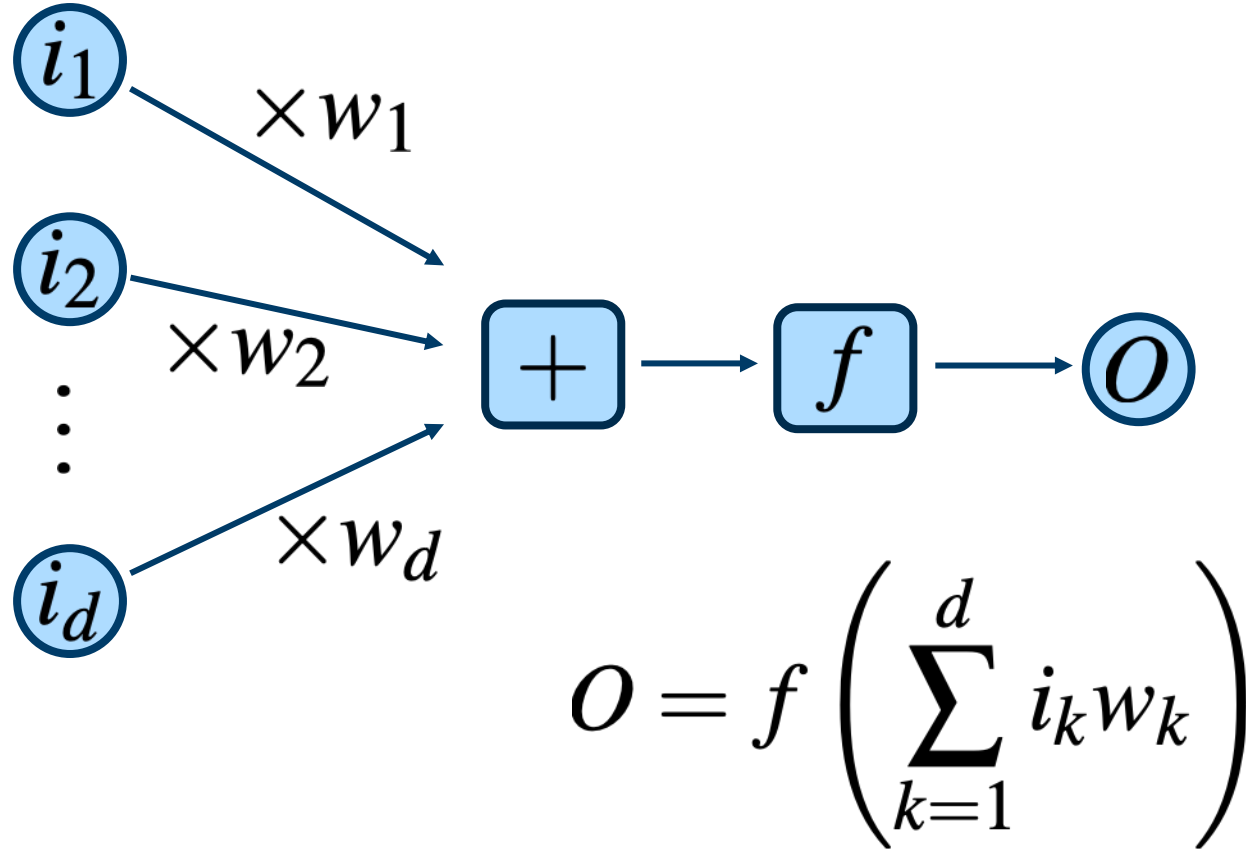
Network types:

- Feed-forward neural networks
- Recurrent neural networks (RNN)
- ...

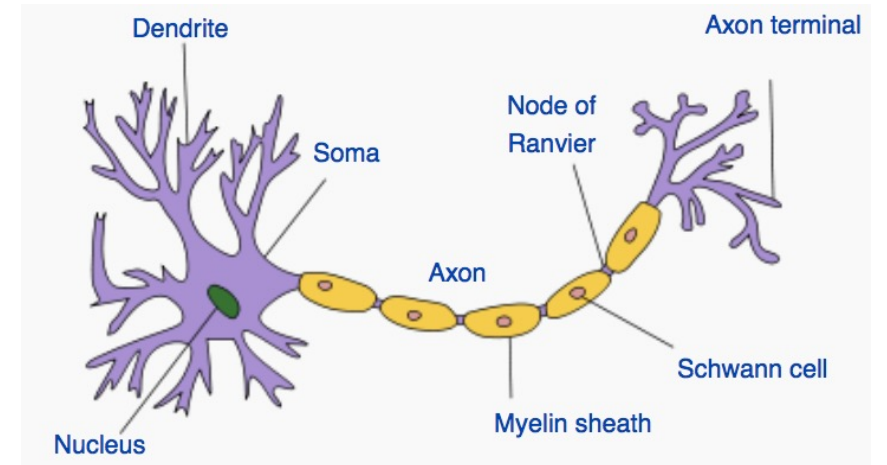


Feed-forward ANN (non-linear classifier)

# Mark 1 Perceptron (Rosenblatt, 1957-59)



## Biological neuron



(Credit: Quasar Jarosz, English Wikipedia)

# Activation functions

- Sigmoid (logistic function):

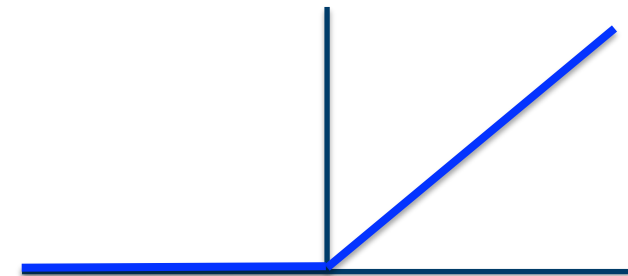
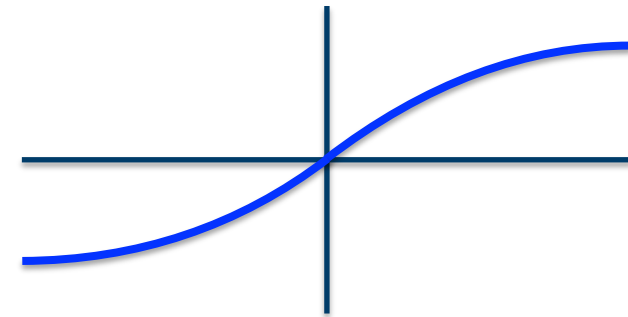
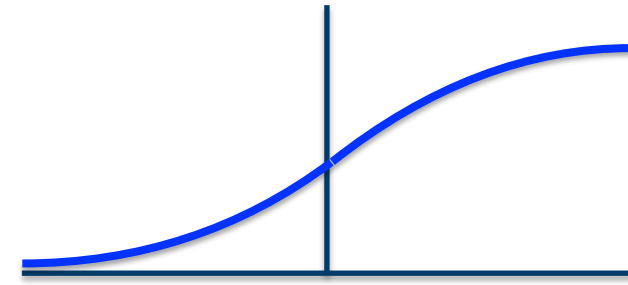
$$f(x) = \frac{1}{1 + e^{-x}}$$

- Hyperbolic tangent:

$$f(x) = \tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$$

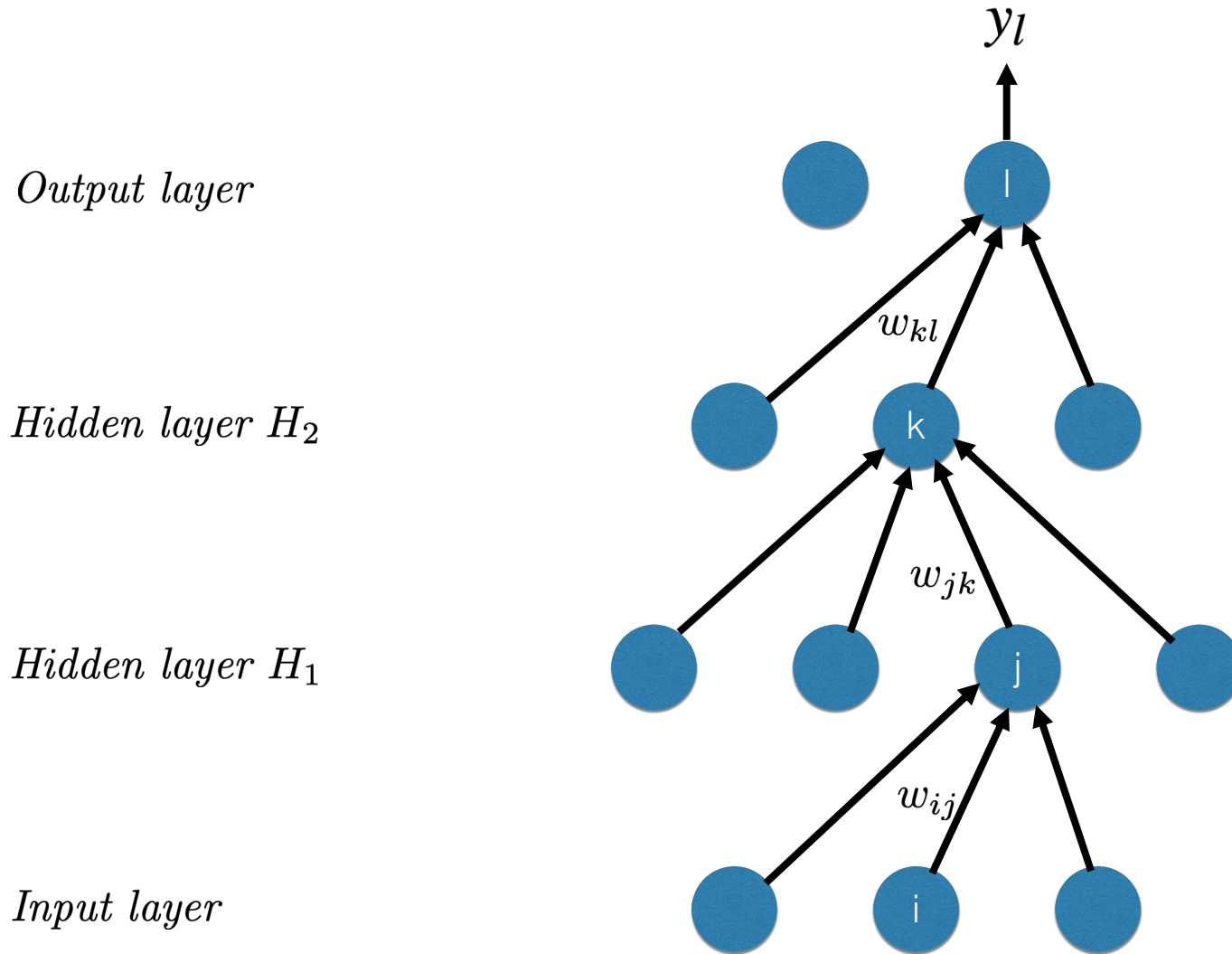
- Rectified linear unit (ReLU):

$$f(x) = \max(x, 0)$$





# Feed-forward neural network



$$y_l = f(z_l)$$

$$z_l = \sum_{k \in H_2} w_{kl} x_k$$

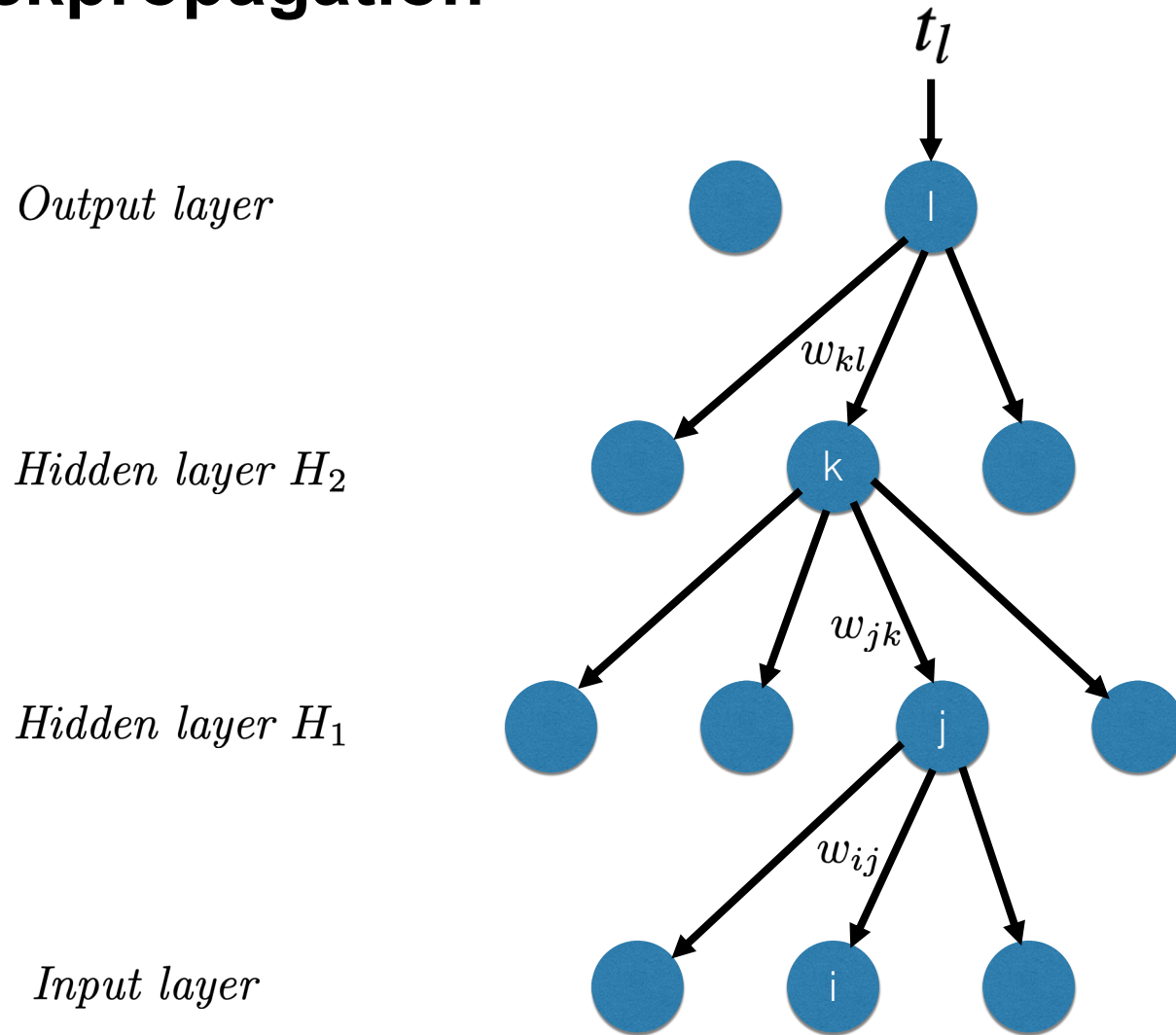
$$y_k = f(z_k)$$

$$z_k = \sum_{j \in H_1} w_{jk} x_j$$

$$y_j = f(z_j)$$

$$z_j = \sum_{i \in \text{Input}} w_{ij} x_i$$

# Backpropagation



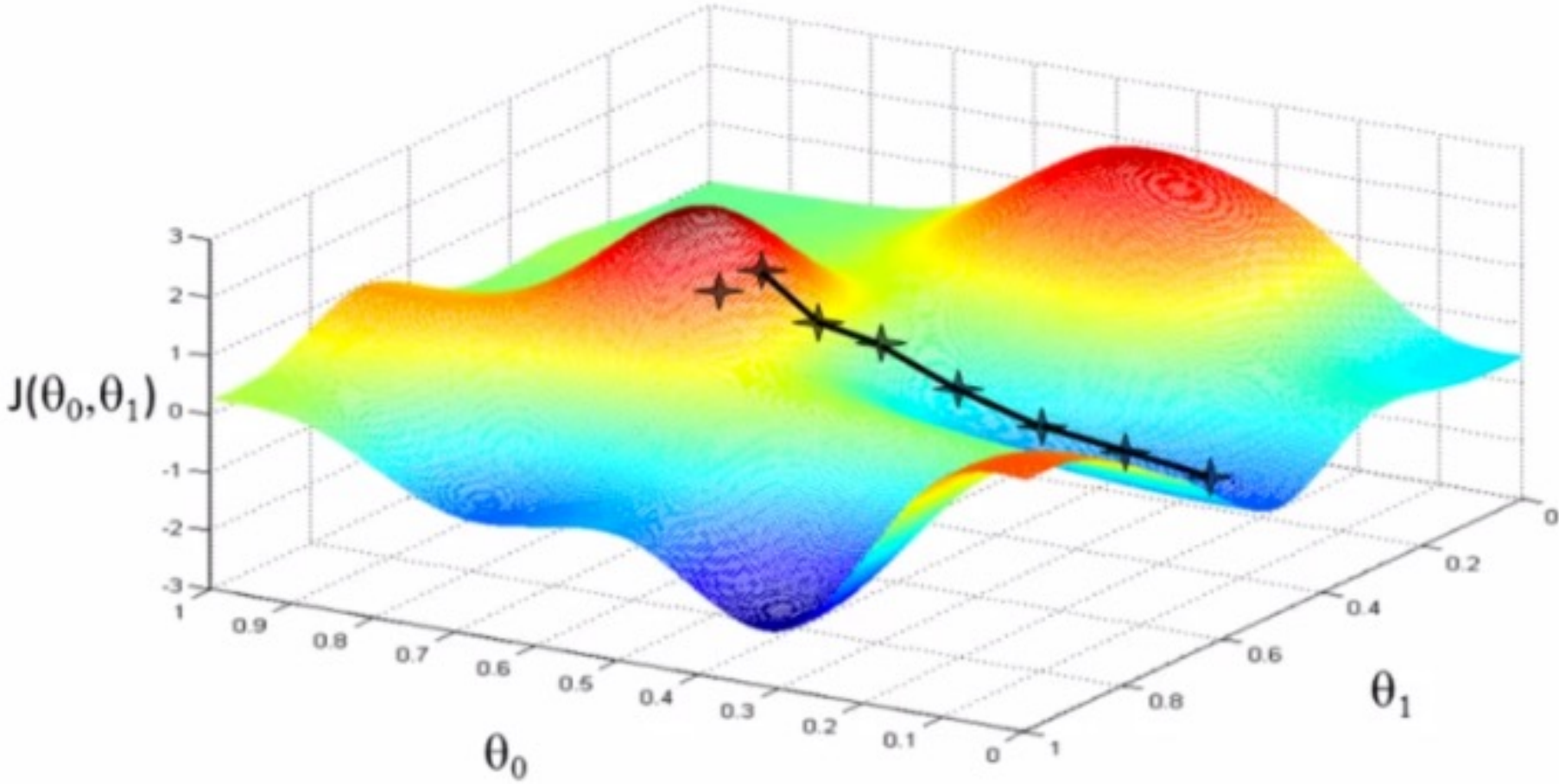
Error function (sample-by-sample measure of difference between output and target value):

$$E(\mathbf{w}) = \sum_l (y_l - t_l)^2$$

Backpropagation:

- Compute the gradient of the error function with respect to the weights, using the chain rule.
- Adjust weights (gradient descent) level-by-level.

# Gradient descent



HOME PLOTS APPS

Design App Get More Apps Install App Package App

Signal Analyzer Classification Learner Deep Network ... Neural Net Clustering Neural Net Fitting Neural Net Pattern Re... Neural Net Time Series Regression Learner Distribution Fitter Audio Labeler Filter Builder Filter Designer

FILE APPS

Q Search Documentation

/ > Users > ffi > Desktop

Current Folder

Name
------

Details

Command Window

New to MATLAB? See resources for [Getting Started](#).

```
fx >>
```

Workspace

Name	Value
------	-------



## Welcome to the Neural Network Pattern Recognition app.

Solve a pattern-recognition problem with a two-layer feed-forward network.

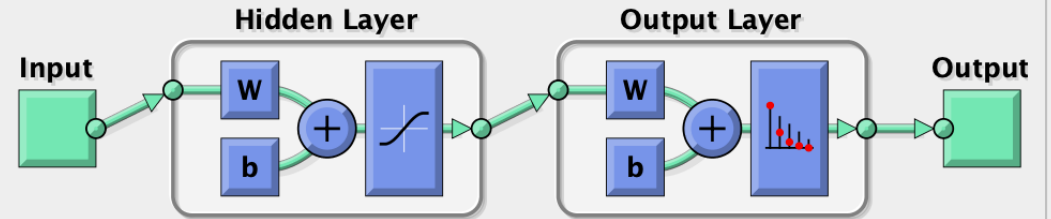
### Introduction

In pattern recognition problems, you want a neural network to classify inputs into a set of target categories.

For example, recognize the vineyard that a particular bottle of wine came from, based on chemical analysis ([wine\\_dataset](#)); or classify a tumor as benign or malignant, based on uniformity of cell size, clump thickness, mitosis ([cancer\\_dataset](#)).

The Neural Pattern Recognition app will help you select data, create and train a network, and evaluate its performance using cross-entropy and confusion matrices.

### Neural Network



A two-layer feed-forward network, with sigmoid hidden and softmax output neurons ([patternnet](#)), can classify vectors arbitrarily well, given enough neurons in its hidden layer.

The network will be trained with scaled conjugate gradient backpropagation ([trainscg](#)).

To continue, click [Next].

Neural Network Start

Welcome

Back

Next

Cancel



## Select Data

What inputs and targets define your pattern recognition problem?

### Get Data from Workspace

Input data to present to the network.

Inputs:

Target data defining desired network output.

Targets:

Samples are:  Matrix columns  Matrix rows

### Summary

No inputs selected.

No targets selected.

Want to try out this tool with an example data set?

Select inputs and targets, then click [Next].



## Select Data

What inputs and targets define your pattern recognition problem?

Get Data from

Pattern Recognition Data Set Chooser

Input data to

Select a data set:

Description

Inputs:

Simple Classes  
Iris Flowers  
Breast Cancer  
Types of Glass  
Thyroid  
Wine Vintage

Filename: [wine\\_dataset](#)

Pattern recognition is the process of training a neural network to assign the correct target classes to a set of input patterns. Once trained the network can be used to classify patterns it has not seen before.

This dataset can be used to create a neural network that classifies wines from three winerys in Italy based on constituents found through chemical analysis.

LOAD [wine\\_dataset](#).MAT loads these two variables:

wineInputs - a 13x178 matrix of thirteen attributes of 178 wines.

1. Alcohol
2. Malic acid
3. Ash
4. Alcalinity of ash
5. Magnesium
6. Total phenols
7. Flavanoids
8. Nonflavanoid phenols

Import

Cancel



Loading dataset.

Neural Network Start

Welcome

Back

Next

Cancel



## Select Data

What inputs and targets define your pattern recognition problem?

### Get Data from Workspace

Input data to present to the network.

Inputs:

wineInputs



...

Target data defining desired network output.

Targets:

wineTargets



...

Samples are:



Matrix columns



Matrix rows

Want to try out this tool with an example data set?

Load Example Data Set

### Summary

Inputs 'wineInputs' is a 13x178 matrix, representing static data: 178 samples of 13 elements.

Targets 'wineTargets' is a 3x178 matrix, representing static data: 178 samples of 3 elements.

 To continue, click [Next].

 Neural Network Start

 Welcome

 Back

 Next

 Cancel





## Validation and Test Data

Set aside some samples for validation and testing.

### Select Percentages

Randomly divide up the 178 samples:

Training: 70% 124 samples

Validation: 15% 27 samples

Testing: 15% 27 samples

Restore Defaults

### Explanation

Three Kinds of Samples:

Training:

These are presented to the network during training, and the network is adjusted according to its error.

Validation:

These are used to measure network generalization, and to halt training when generalization stops improving.

Testing:

These have no effect on training and so provide an independent measure of network performance during and after training.

Change percentages if desired, then click [Next] to continue.

Neural Network Start

Welcome

Back

Next

Cancel



## Network Architecture

Set the number of neurons in the pattern recognition network's hidden layer.

### Hidden Layer

Define a pattern recognition neural network. (patternnet)

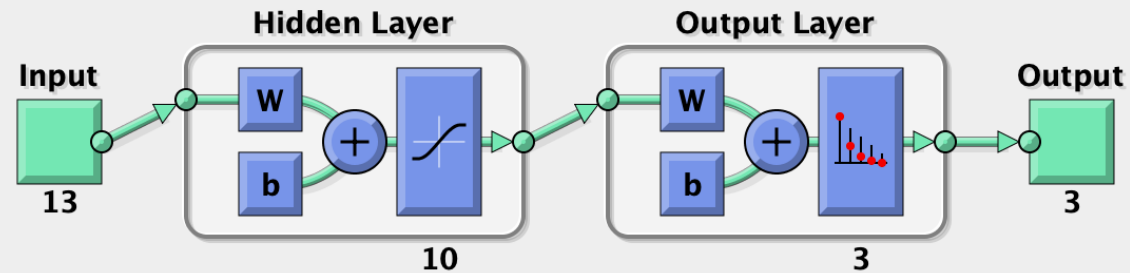
Number of Hidden Neurons:

Restore Defaults

### Recommendation

Return to this panel and change the number of neurons if the network does not perform well after training.

### Neural Network



 Change settings if desired, then click [Next] to continue.

 Neural Network Start

 Welcome

 Back

 Next

 Cancel



## Train Network

Train the network to classify the inputs according to the targets.

### Train Network

Train using scaled conjugate gradient backpropagation. (trainscg)



Training automatically stops when generalization stops improving, as indicated by an increase in the cross-entropy error of the validation samples.

### Notes

Training multiple times will generate different results due to different initial conditions and sampling.

Minimizing Cross-Entropy results in good classification. Lower values are better. Zero means no error.

Percent Error indicates the fraction of samples which are misclassified. A value of 0 means no misclassifications, 100 indicates maximum misclassifications.

### Results

	Samples	CE	%E
Training:	124	-	-
Validation:	27	-	-
Testing:	27	-	-

Plot Confusion

Plot ROC



Train network, then click [Next].

Neural Network Start

Welcome

Back

Next

Cancel



## Train Network

Train the network to classify the inputs according to the targets.

### Train Network

Train using scaled conjugate gradient backpropagation. (trainscg)

Retrain

Training automatically stops when generalization stops improving, as indicated by an increase in the cross-entropy error of the validation samples.

### Notes

Training multiple times will generate different results due to different initial conditions and sampling.

Minimizing Cross-Entropy results in good classification. Lower values are better. Zero means no error.

Percent Error indicates the fraction of samples which are misclassified. A value of 0 means no misclassifications, 100 indicates maximum misclassifications.

### Results

	Samples	CE	%E
Training:	124	9.86550e-1	0
Validation:	27	2.74688e-0	11.11111e-0
Testing:	27	2.66457e-0	0

Plot Confusion

Plot ROC

Open a plot, retrain, or click [Next] to continue.

Neural Network Start

Welcome

Back

Next

Cancel



# Summary

## Machine learning:

- Pattern classification
- Training of classifiers (supervised learning)
- Parametric and non-parametric methods
- Discriminant functions
- Quadratic and linear classifiers
- Neural Networks.

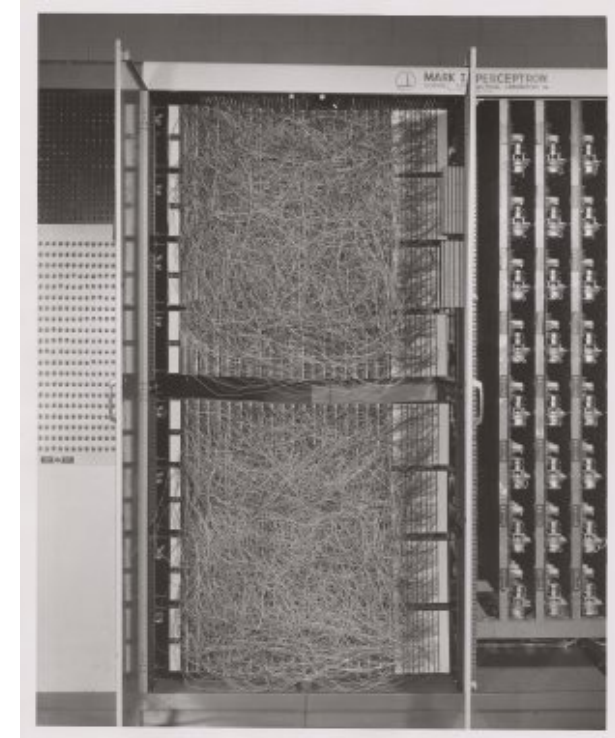
## Recommended reading:

- Szeliski 5.1 - 5.3

## Additional reading:

- Szeliski 6.1 - 6.3
- R. O. Duda, P. E. Hart, D. G. Stork (2001). *Pattern classification* (2nd ed.). Wiley, New York. ISBN 0-471-05669-3.

Mark 1 Perceptron



(Credit: Cornell Aeronautical Laboratory)