

# Classification and Pattern Recognition

## Slides 3rd Year Signal and Image Processing - Master

Jean-Yves Tournet<sup>(1)</sup>

(1) University of Toulouse, ENSEEIHT-IRIT-TéSA  
Thème 1 : Analyse et Synthèse de l'Information  
[jyt@n7.fr](mailto:jyt@n7.fr)

October 2018

## Summary

- ▶ **Chapter 1 : Introduction**
  - ▶ Examples
  - ▶ Classification Model
- ▶ **Chapter 2 : Preprocessing**
  - ▶ Signal Modeling and Representation
  - ▶ Principal Component Analysis (PCA)
  - ▶ Linear Discriminant Analysis (LDA)
- ▶ **Chapter 3 : Statistical Methods**
  - ▶ Bayesian Rule
  - ▶ Supervised Learning
  - ▶ Unsupervised Learning
- ▶ **Chapter 4 : Linear Discriminant Functions and Neural Networks**
  - ▶ Linear discriminant functions
  - ▶ Support vector machines (SVMs)
  - ▶ Neural networks
- ▶ **Chapter 5 : Decision Trees**

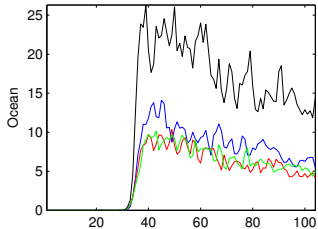
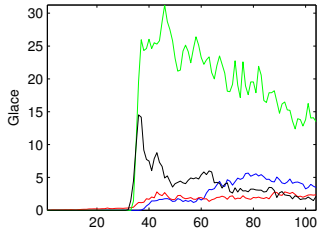
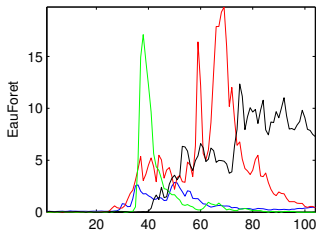
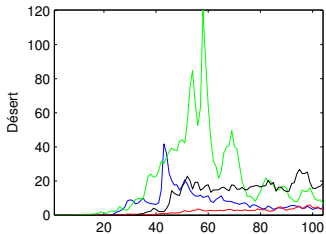
## Summary

- ▶ **Chapter 1 : Introduction**
  - ▶ **Examples**
    - ▶ Altimetry
    - ▶ Electromyography (EMG)
    - ▶ Accelerometry
  - ▶ **Classification Model**
- ▶ **Chapter 2 : Preprocessing**
- ▶ **Chapter 3 : Statistical Methods**
- ▶ **Chapter 4 : Support Vector Machines (SVMs) and Neural Networks**
- ▶ **Chapter 5 : Decision Trees**

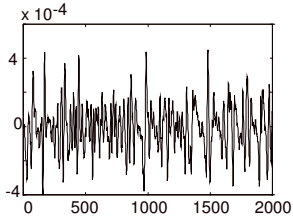
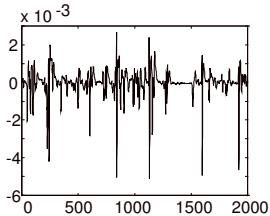
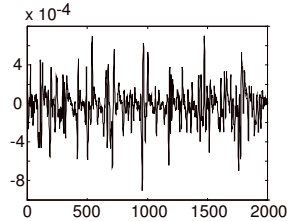
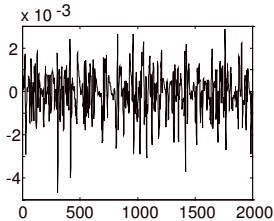
## References

- ▶ R. O. Duda, P. E. Hart and D. G. Stork, [Pattern Classification](#), 2nd edition, Wiley, 2000.
- ▶ S. Theodoridis and K. Koutroumbas, [Pattern Recognition](#), 4th edition, Academic Press, 2008.
- ▶ A. Jain, R. Duin and J. Mao, [Statistical Pattern Recognition: A Review](#), IEEE Transactions Pattern Analysis and Machine Intelligence, vol. 22, no. 1, pp. 4-37, Jan. 2000.

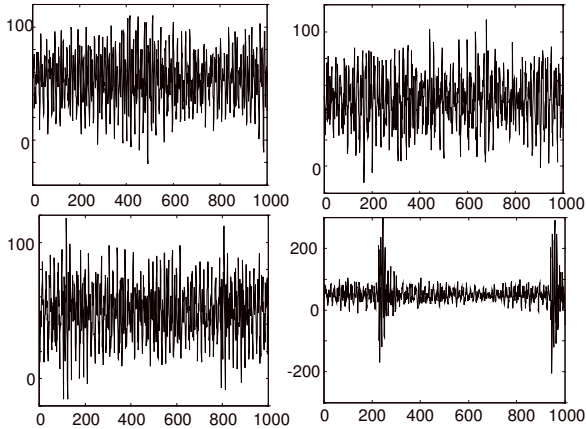
## Example 1 : Altimetric Signals



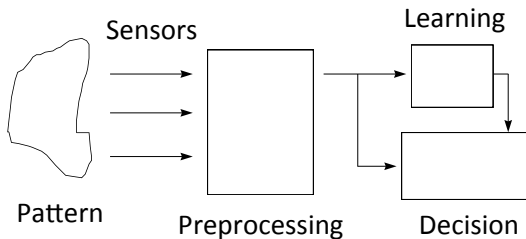
## Example 2 : EMG Signals



## Example 3 : Accelerometric Signals



## Classification Model





## Summary

- ▶ Chapter 1 : Introduction
- ▶ Chapter 2 : Preprocessing
  - ▶ Signal Modeling and Representation
  - ▶ Principal Component Analysis (PCA)
  - ▶ Linear Discriminant Analysis (LDA)
- ▶ Chapter 3 : Statistical Methods
- ▶ Chapter 4 : Support Vector Machines (SVMs) and Neural Networks
- ▶ Chapter 5 : Decision Trees

## Modeling and Representation

### Stationary Signals

- ▶ Linear models (MA, AR, ARMA)
- ▶ Non-linear models (Volterra time-series, bilinear models, ...)
- ▶ Features (Reflexion coefficients, cepstral coefficients,...)
- ▶ ...

### Non-stationary signals

- ▶ Wavelet analysis
- ▶ Time-frequency distributions
- ▶ ...

## Summary

- ▶ Chapter 1 : Introduction
- ▶ Chapter 2 : Preprocessing
  - ▶ Signal Modeling and Representation
  - ▶ Principal Component Analysis (PCA)
  - ▶ Linear Discriminant Analysis (LDA)
- ▶ Chapter 3 : Statistical Methods
- ▶ Chapter 4 : Support Vector Machines (SVMs) and Neural Networks
- ▶ Chapter 5 : Decision Trees

## PCA of samples

First step: define a norm for  $\mathbf{x} \in \mathbb{R}^p$

$$\|\mathbf{x}\|_M^2 = \langle \mathbf{x}, \mathbf{x} \rangle_M = \mathbf{x}^T \mathbf{M} \mathbf{x}$$

where  $\mathbf{M}$  is a symmetric positive definite matrix of size  $p \times p$

►  $\mathbf{M} = I_p$

$$d^2(\mathbf{x}, \mathbf{y}) = \sum_{j=1}^p [\mathbf{x}(j) - \mathbf{y}(j)]^2$$

►  $\mathbf{M} = \text{diag}\left(\frac{1}{\sigma_1^2}, \dots, \frac{1}{\sigma_p^2}\right)$

$$d^2(\mathbf{x}, \mathbf{y}) = \sum_{j=1}^p [\mathbf{x}^*(j) - \mathbf{y}^*(j)]^2$$

where  $\mathbf{x}^*(j) = \frac{x(j) - m(j)}{\sigma_j}$  (centered and reduced data). **In what follows, data are centered, reduced and  $\mathbf{M} = I_p$ .**

## PCA of samples

Second step: minimize the difference between the samples and their projections

$$\text{Minimize } I_q = \frac{1}{n} \sum_{i=1}^n \|\mathbf{x}_i - \tilde{\mathbf{x}}_i\|^2 \Leftrightarrow \text{Maximize } J_q = \frac{1}{n} \sum_{i=1}^n \|\tilde{\mathbf{x}}_i\|^2$$

*Property: solutions can be found recursively with  $J(\mathbf{u}) = \mathbf{u}^t \mathbf{T} \mathbf{u}$*

- ▶ **Determination of principal components**
- ▶ **Number of principal components**  
 $\mathbf{T}$  of size  $p \times p$  invertible  $\implies p$  principal components
- ▶ **How to choose the number of principal components?**

$$I_q = \frac{1}{n} \sum_{i=1}^n \|\mathbf{x}_i\|^2 \left[ 1 - \frac{\sum_{j=1}^q \lambda_j}{\sum_{j=1}^p \lambda_j} \right]$$

## PCA of variables

- ▶ **Number of principal components**

Covariance matrix of size  $N \times N$  with  $p$  principal components

- ▶ **Normalization**

$$v'_j(i) = \frac{x_i(j) - m(j)}{\sqrt{n}\sigma(j)}$$

- ▶ **Interpretation**

Representation in the correlation circle

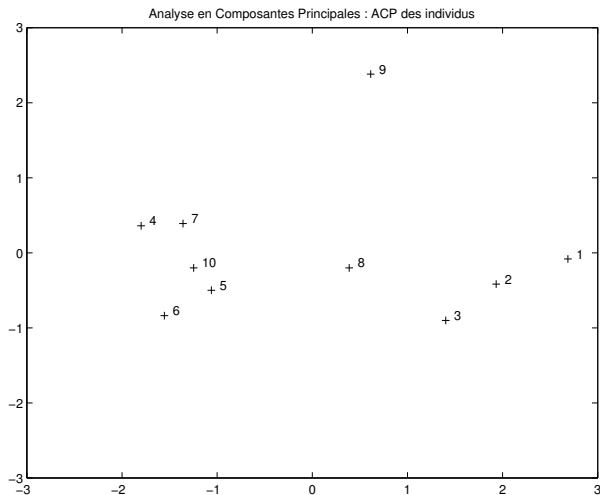
Correlation between variables

$$\|\mathbf{v}'_k - \mathbf{v}'_j\|^2 = 2(1 - r_{jk})$$

## Toy Example

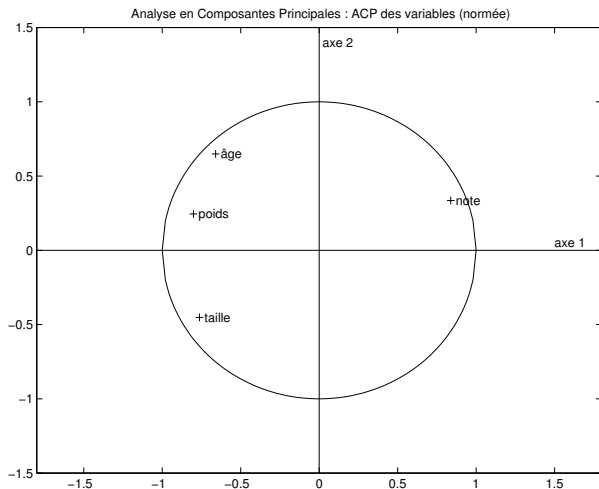
	Weight	Size	Age	Grade		Weight	Size	Age	Grade
$\mathbf{x}_1$	45	1.50	13	14	$\mathbf{x}_6$	60	1.70	14	7
$\mathbf{x}_2$	50	1.60	13	16	$\mathbf{x}_7$	70	1.60	14	8
$\mathbf{x}_3$	50	1.65	13	15	$\mathbf{x}_8$	65	1.60	13	13
$\mathbf{x}_4$	60	1.70	15	9	$\mathbf{x}_9$	60	1.55	15	17
$\mathbf{x}_5$	60	1.70	14	10	$\mathbf{x}_{10}$	65	1.70	14	11

## PCA of samples

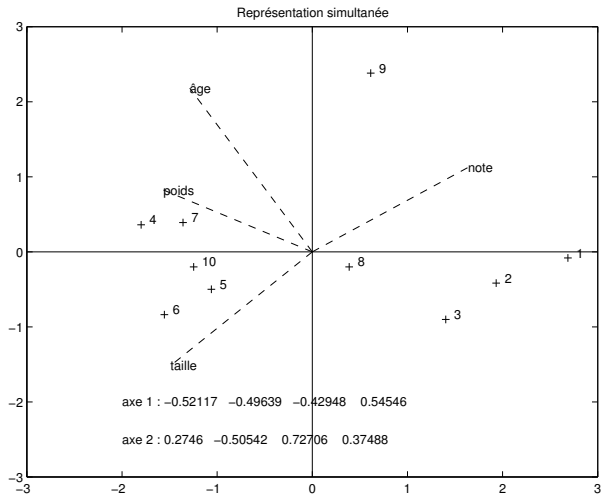




## PCA of variables



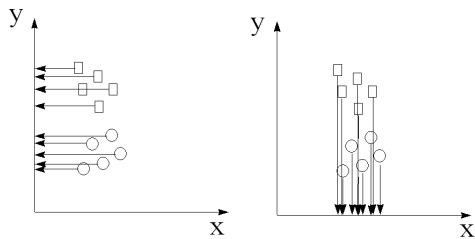
## Joint PCA



## Summary

- ▶ Chapter 1 : Introduction
- ▶ Chapter 2 : Preprocessing
  - ▶ Signal Modeling and Representation
  - ▶ Principal Component Analysis (PCA)
  - ▶ Linear Discriminant Analysis (LDA)
- ▶ Chapter 3 : Statistical Methods
- ▶ Chapter 4 : Support Vector Machines (SVMs) and Neural Networks
- ▶ Chapter 5 : Decision Trees

## Fisher criterion



## Fisher criterion

### Two classes

$$J(\mathbf{u}) = \frac{(\tilde{m}_1 - \tilde{m}_2)^2}{\tilde{s}_1^2 + \tilde{s}_2^2} = \frac{\mathbf{u}^t \mathbf{B} \mathbf{u}}{\mathbf{u}^t \mathbf{S} \mathbf{u}}$$

where  $\tilde{m}_i = \mathbf{u}^t \mathbf{m}_i$ ,  $\tilde{x} = \mathbf{u}^t \mathbf{x}$  and  $\tilde{s}_i^2 = \sum_{\mathbf{x} \in W_i} (\tilde{x} - \tilde{m}_i)^2$ .

**B** between-class scatter matrix

$$\mathbf{B} = (\mathbf{m}_1 - \mathbf{m}_2)(\mathbf{m}_1 - \mathbf{m}_2)^t$$

**S** within-class scatter matrix

$$\begin{aligned} \mathbf{S} &= \mathbf{S}_1 + \mathbf{S}_2 \\ &= \sum_{\mathbf{x} \in W_1} (\mathbf{x} - \mathbf{m}_1)(\mathbf{x} - \mathbf{m}_1)^t + \sum_{\mathbf{x} \in W_2} (\mathbf{x} - \mathbf{m}_2)(\mathbf{x} - \mathbf{m}_2)^t \end{aligned}$$

## Fisher criterion

### Two classes

- ▶ Property

$$\mathbf{T} = \mathbf{S} + \frac{n_1 n_2}{n} \mathbf{B}$$

since

$$\mathbf{T} = \sum_{x \in W_1 \cup W_2} (\mathbf{x} - \mathbf{m})(\mathbf{x} - \mathbf{m})^t$$

- ▶ Determination of the discriminant axis

$$\mathbf{u} = k \mathbf{S}^{-1} (\mathbf{m}_1 - \mathbf{m}_2)$$

## Fisher criterion

Generalization to  $K \geq 2$  classes

**S**: within-class scatter matrix

$$\mathbf{S} = \sum_{i=1}^K \mathbf{S}_i = \sum_{i=1}^K \sum_{\mathbf{x} \in W_i} (\mathbf{x} - \mathbf{m}_i)(\mathbf{x} - \mathbf{m}_i)^t$$

**T**: covariance matrix

$$\mathbf{T} = \sum_{i=1}^K \sum_{\mathbf{x} \in W_i} (\mathbf{x} - \mathbf{m})(\mathbf{x} - \mathbf{m})^t = \mathbf{S} + \sum_{i=1}^K n_i (\mathbf{m}_i - \mathbf{m})(\mathbf{m}_i - \mathbf{m})^t$$

hence the following between-class scatter matrix

$$\mathbf{B} = \sum_{i=1}^K n_i (\mathbf{m}_i - \mathbf{m})(\mathbf{m}_i - \mathbf{m})^t$$

## Properties

### Determination of the discriminant axes

Maximization of  $I(\mathbf{u}) = \frac{\mathbf{u}^t \mathbf{B} \mathbf{u}}{\mathbf{u}^t \mathbf{T} \mathbf{u}} \Leftrightarrow$  Maximization of  $J(\mathbf{u}) = \frac{\mathbf{u}^t \mathbf{B} \mathbf{u}}{\mathbf{u}^t \mathbf{S} \mathbf{u}}$

#### Remarks

- ▶ the eigenvalues of  $\mathbf{T}^{-1} \mathbf{B}$  belong to  $[0, 1]$
- ▶ limit cases  $\mu = 0$  and  $\mu = 1$
- ▶ Number of discriminant axes:  $K - 1$
- ▶ Discriminant power of variable  $\#i$ :  $J(\mathbf{e}_i)$  or  $I(\mathbf{e}_i)$



## Example (Matlab)

▶  $\mathbf{X} = [10; 11; -1 - 1; 1 - 1; 0 - 1; 00; \mathbf{0\ 1; -1\ 0; -1\ 1; 0\ 2}]$

▶  $\mathbf{T} = \text{cov}(\mathbf{X}) * (n - 1) \quad \rightarrow \mathbf{T} = \begin{pmatrix} 6.0000 & 0 \\ 0 & 9.6000 \end{pmatrix}$

▶  $\mathbf{m} = [5/6; -4/3]$

▶  $\mathbf{B} = \mathbf{m} * \mathbf{m}'$

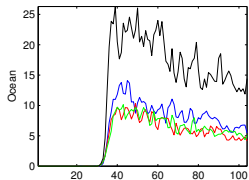
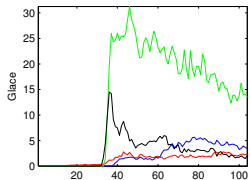
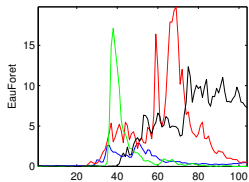
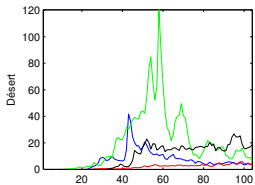
▶  $[\mathbf{U}\ \mathbf{V}] = \text{eig}(\text{inv}(\mathbf{T}) * \mathbf{B})$

$$\rightarrow \mathbf{U} = \begin{pmatrix} -0.8480 & 0.7071 \\ -0.5300 & -0.7071 \end{pmatrix}, \mathbf{V} = \begin{pmatrix} 0 & 0 \\ 0 & 0.7222 \end{pmatrix}$$

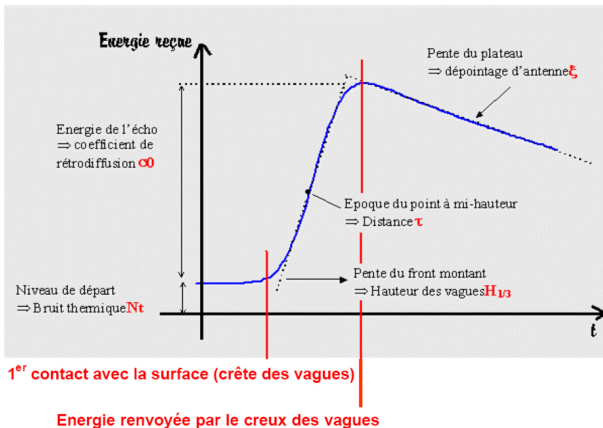
Thus, the only discriminant axis is defined by the vector  $(1, -1)$ :

**Straight line of equation  $y = -x$ .**

## Waveforms



## Analytical model



## Classification strategy

### Waveform-based classification

- ▶ **Preprocessing:** linear discriminant analysis
- ▶ **Classification:** Bayesian classifier, centroid distance rule,  $k$  nearest neighbors,...

### Parameter-based classification

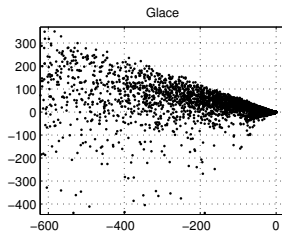
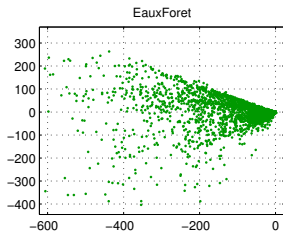
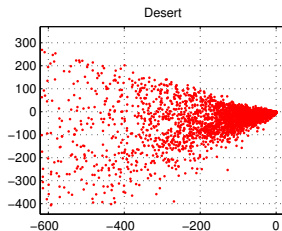
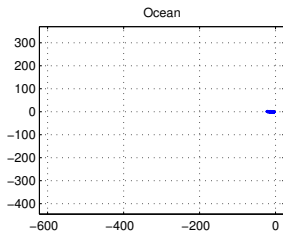
- ▶ **Preprocessing:** linear discriminant analysis
- ▶ **Classification:** Bayesian classifier, centroid distance rule,  $k$  nearest neighbors,...

## Real Data

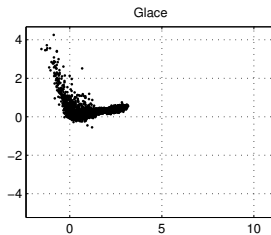
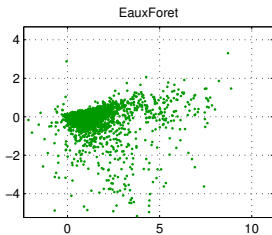
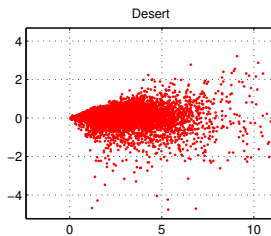
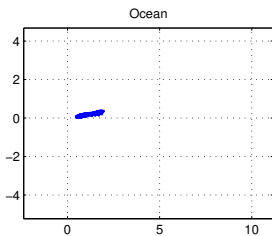
### Waveform-based classification

- ▶ **Class  $\omega_1$** : Oceans
- ▶ **Class  $\omega_2$** : Deserts (Algeria, Lybia, South Africa)
- ▶ **Class  $\omega_3$** : Water and Forests (Amazonia, Canada, Congo, Russia)
- ▶ **Class  $\omega_4$** : Ice (Arctic continental ice, Groenland continental ice, Antarctic sea ice, Arctic sea ice)

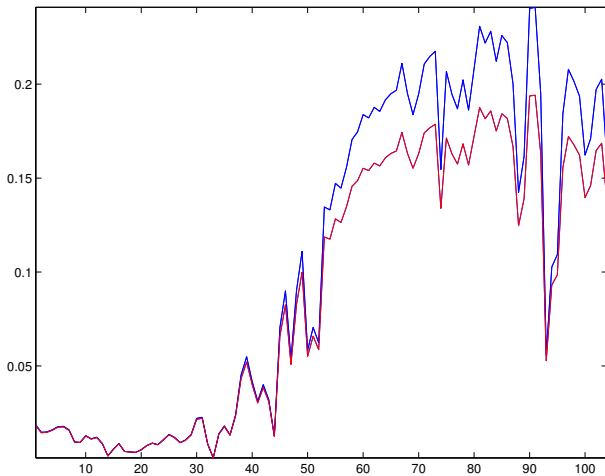
## PCA of Samples (Waveforms)



## Discriminant Analysis (Waveforms)



## Fisher criterion





## Features

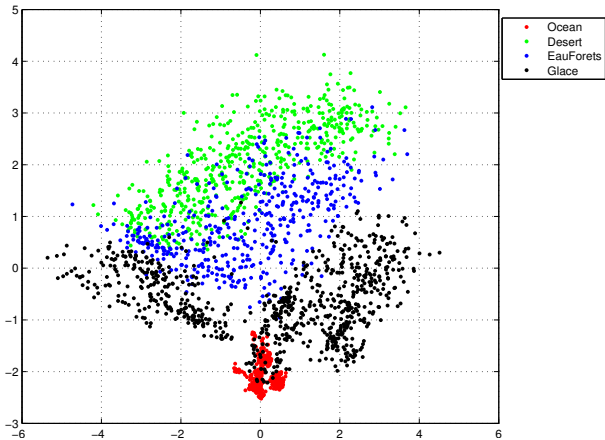
### Parameters estimated from waveforms (Retracking)

- **ku\_swh\_1Hz** : Significant Wave Height of ku band
- **ku\_std\_swh\_1Hz** : Standard deviation of swh 18 Hz
- **ku\_sigma0\_1Hz** : Backscatter coefficient ku band
- **s\_sigma0\_1Hz** : Backscatter coefficient s band
- **elevation\_1Hz** : Elevation above the reference ellipsoid
- **ku\_peakiness\_1Hz** : Peakiness of waveforms in ku band
- **ku\_ice2\_width\_1Hz** : Width of the leading edge

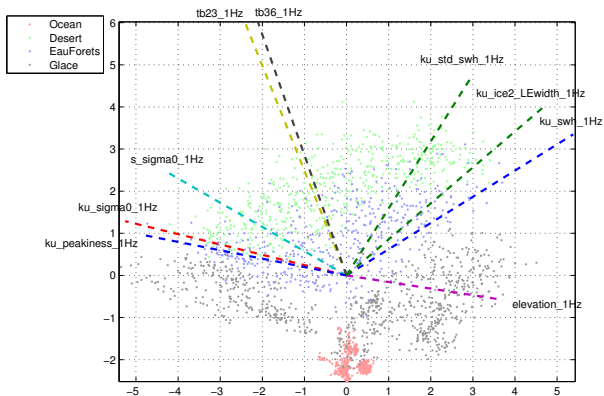
### Parameters from radiometers

- **tb23\_1Hz** : Brightness Temperature at 23 GHz
- **tb36\_1Hz** : Brightness Temperature at 36 GHz

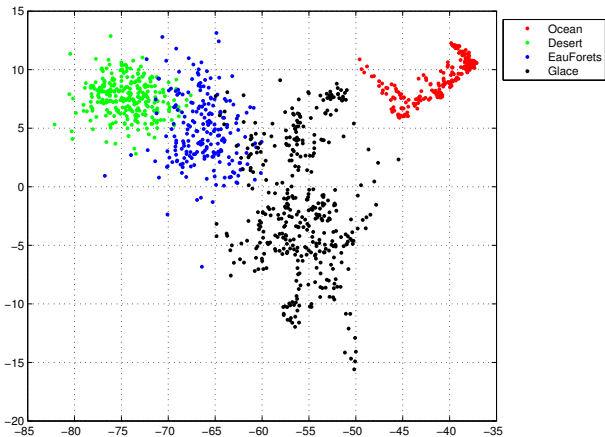
## PCA of Samples (Features)



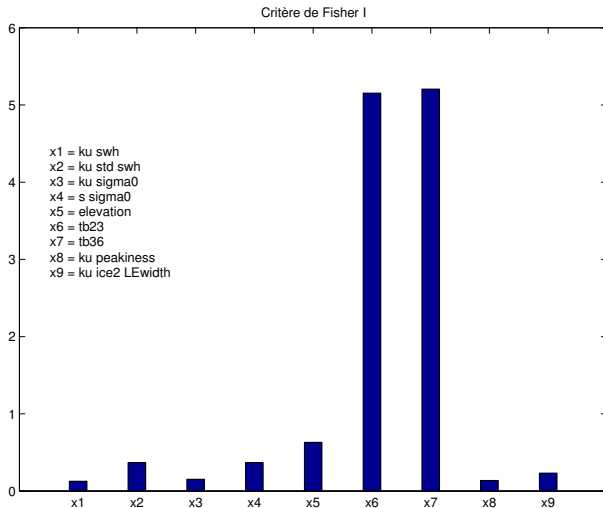
## Joint PCA (features)



## Discriminant Analysis (features)



## Fisher Criterion (features)



Discriminant power of each parameter for classification

## Summary

- ▶ Chapter 1 : Introduction
- ▶ Chapter 2 : Preprocessing
- ▶ Chapter 3 : Statistical Methods
  - ▶ Bayesian rule
    - ▶ Theory
    - ▶ Two classes
    - ▶ Uninformative cost function
    - ▶ Gaussian densities
  - ▶ Supervised learning
    - ▶ Parametric methods
    - ▶ Non-parametric methods
  - ▶ Unsupervised learning
    - ▶ Optimization methods
    - ▶ Hierarchical classification
- ▶ Chapter 4 : Support Vector Machines (SVMs) and Neural Networks
- ▶ Chapter 5 : Decision Trees

## Classification

### Notations

- ▶  $K$  classes  $\omega_1, \dots, \omega_K$
- ▶  $\mathbf{x} = [x(1), \dots, x(p)]^T$  **measurements**  $\in X = \mathbb{R}^p$
- ▶  $A$ : set of possible **actions**  $a_1, \dots, a_q$  where  $a_i =$  "assign the vector  $\mathbf{x}$  to the class  $\omega_i$ ",  $\forall i = 1, \dots, K$

### Definition

$$d: \begin{array}{l} X \rightarrow A \\ \mathbf{x} \mapsto d(\mathbf{x}) \end{array}$$

### Remark

Classification **with reject option**:  $A = \{a_0, a_1, \dots, a_K\}$  where  $a_0 =$  "do not classify the vector  $\mathbf{x}$ "

## Bayesian Rule

### Hypothesis: Probabilistic Model

- ▶ *A priori* probability of class  $\omega_i$

$$P(\omega_i)$$

- ▶ Probability density function of the observation vector  $\mathbf{x}$  conditionally to the class  $\omega_i$

$$f(\mathbf{x} | \omega_i)$$

### Conclusion

- ▶ *A posteriori* probability that  $\mathbf{x}$  belongs to class  $\omega_i$

$$P(\omega_i | \mathbf{x}) = \frac{f(\mathbf{x} | \omega_i) P(\omega_i)}{f(\mathbf{x})}$$

with  $f(\mathbf{x}) = \sum_{i=1}^K f(\mathbf{x} | \omega_i) P(\omega_i)$ .



## MAP Classifier

### Definition

$$d^*(\mathbf{x}) = a_j \Leftrightarrow P(\omega_j | \mathbf{x}) \geq P(\omega_k | \mathbf{x}), \forall k \in \{1, \dots, K\}$$

### Equiprobable Classes: Maximum Likelihood Classifier

$$d^*(\mathbf{x}) = a_j \Leftrightarrow f(\mathbf{x} | \omega_j) \geq f(\mathbf{x} | \omega_k), \forall k \in \{1, \dots, K\}$$

### Property

The MAP classifier minimizes the probability of error

## Proof (2 classes)

$$\begin{aligned}P_e &= P[d(\mathbf{x}) = a_1 \cap \mathbf{x} \in \omega_2] + P[d(\mathbf{x}) = a_2 \cap \mathbf{x} \in \omega_1] \\ &= P[d(\mathbf{x}) = a_1 \mid \mathbf{x} \in \omega_2] P(\omega_2) + P[d(\mathbf{x}) = a_2 \mid \mathbf{x} \in \omega_1] P(\omega_1)\end{aligned}$$

Let  $R_i = \{\mathbf{x} \in \mathbb{R}^p / d(\mathbf{x}) = a_i\}$  be the acceptance region for class  $\omega_i$

$$\begin{aligned}P_e &= \int_{R_2} P(\omega_1) f(\mathbf{x} \mid \omega_1) d\mathbf{x} + \int_{R_1} P(\omega_2) f(\mathbf{x} \mid \omega_2) d\mathbf{x} \\ &= P(\omega_2) \left[ 1 - \int_{R_2} f(\mathbf{x} \mid \omega_2) d\mathbf{x} \right] + \int_{R_2} P(\omega_1) f(\mathbf{x} \mid \omega_1) d\mathbf{x} \\ &= P(\omega_2) + \int_{R_2} [P(\omega_1) f(\mathbf{x} \mid \omega_1) - P(\omega_2) f(\mathbf{x} \mid \omega_2)] d\mathbf{x} \\ &= P(\omega_2) - \int_{R_2} [P(\omega_2 \mid \mathbf{x}) - P(\omega_1 \mid \mathbf{x})] f(\mathbf{x}) d\mathbf{x}\end{aligned}$$

$P_e$  is minimum when  $R_2 = \{\mathbf{x} / P(\omega_2 \mid \mathbf{x}) > P(\omega_1 \mid \mathbf{x})\}$

## Probability of Error

- ▶ Definition ( $K$  classes)

$$P_e = \sum_{i=1}^K P[d(x) = a_i \cap x \notin \omega_i]$$

- ▶ Proof  
admitted

## Cost Functions

### Definitions

- ▶ Cost of making the **decision**  $a_j$  knowing that  $x$  belongs to **class**  $\omega_i$

$$c(a_j, \omega_i) \geq 0$$

- ▶ **Average** cost of making the **decision**  $a_j$  given  $x$

$$R(a_j | x) = \sum_{i=1}^K c(a_j, \omega_i) P(\omega_i | x)$$

- ▶ **Average** cost of making the **decision**  $d(x)$  given  $x$

$$R_d(x) = \sum_{i=1}^K c(d(x), \omega_i) P(\omega_i | x)$$

- ▶ **Average overall** cost associated with the **decision rule**  $d$

$$R_d = \int_{\mathbb{R}^p} R_d(x) f(x) dx$$

## Bayesian classification

### Definition

Determine a decision function  $d$  which minimizes the average overall cost  $R_d$ . Since  $f(\mathbf{x}) \geq 0$ , the function  $d$  which minimizes  $R_d$  denoted as  $d^*$  is defined by

$$R_{d^*}(\mathbf{x}) \leq R_d(\mathbf{x}) \quad \forall \mathbf{x} \in \mathbb{R}^p$$

### Classification without reject option

$$d^*(\mathbf{x}) = a_j \iff R_j(\mathbf{x}) \leq R_k(\mathbf{x}) \quad \forall k$$

with

$$R_j(\mathbf{x}) = \sum_{i=1}^K c_{ji} P(\omega_i | \mathbf{x}) = R(a_j | \mathbf{x}) \quad \text{and} \quad c_{ji} = c(a_j, \omega_i)$$

## Particular case: 2 classes

### Definition of costs

$$R_1(\mathbf{x}) = c_{11}P(\omega_1 | \mathbf{x}) + c_{12}P(\omega_2 | \mathbf{x})$$

$$R_2(\mathbf{x}) = c_{21}P(\omega_1 | \mathbf{x}) + c_{22}P(\omega_2 | \mathbf{x})$$

with  $c_{11} < c_{12}$  and  $c_{22} < c_{21}$ .

### Bayesian Classifier

$$d^*(\mathbf{x}) = a_1 \Leftrightarrow R_1(\mathbf{x}) \leq R_2(\mathbf{x})$$

$$d^*(\mathbf{x}) = a_1 \Leftrightarrow \frac{f(\mathbf{x} | \omega_1)}{f(\mathbf{x} | \omega_2)} \geq \frac{c_{12} - c_{22}}{c_{21} - c_{11}} \frac{P(\omega_2)}{P(\omega_1)}$$

### Vocabulary

$\frac{f(\mathbf{x} | \omega_1)}{f(\mathbf{x} | \omega_2)}$  is the **likelihood ratio**.

## Non-informative cost function

### Definition

$$c_{ji} = 1 - \delta_{ij} = \begin{cases} 0 & \text{si } i = j \\ 1 & \text{si } i \neq j \end{cases}$$

### Conditional Risks

$$R_j(\mathbf{x}) = \sum_{i=1}^K c_{ji} P(\omega_i | \mathbf{x}) = \sum_{i \neq j}^K P(\omega_i | \mathbf{x}) = 1 - P(\omega_j | \mathbf{x})$$

### MAP Classifier

$$d^*(\mathbf{x}) = a_j \Leftrightarrow P(\omega_j | \mathbf{x}) \geq P(\omega_k | \mathbf{x}), \forall k \in \{1, \dots, K\}$$

### Maximum Likelihood Classifier

$$d^*(\mathbf{x}) = a_j \Leftrightarrow f(\mathbf{x} | \omega_j) \geq f(\mathbf{x} | \omega_k), \forall k \in \{1, \dots, K\}$$

## Gaussian Case

### Densities

$$f(\mathbf{x} | \omega_i) = \frac{1}{(2\pi)^{p/2} \sqrt{\det \Sigma_i}} \exp \left[ -\frac{1}{2} (\mathbf{x} - \mathbf{m}_i)^t \Sigma_i^{-1} (\mathbf{x} - \mathbf{m}_i) \right]$$

### General Case

$$d^*(\mathbf{x}) = a_i \Leftrightarrow g_i(\mathbf{x}) \geq g_k(\mathbf{x}) \quad \forall k$$

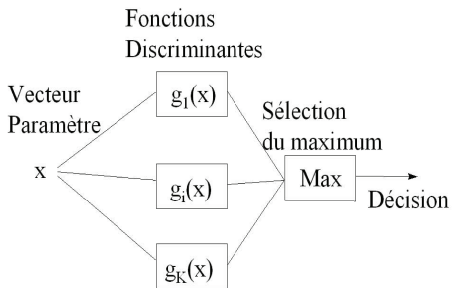
with

$$g_i(\mathbf{x}) = -(\mathbf{x} - \mathbf{m}_i)^t \Sigma_i^{-1} (\mathbf{x} - \mathbf{m}_i) - \ln \det \Sigma_i + 2 \ln P(\omega_i)$$



## Gaussian Case

### Classifier



**Discriminant functions**

Quadratic term + linear term  
+ constant

## Identical covariance matrices ( $\Sigma_i = \Sigma$ )

### Centroid Distance Rule

$$d^*(\mathbf{x}) = a_i \Leftrightarrow d_M(\mathbf{x}, \mathbf{m}_i) \leq d_M(\mathbf{x}, \mathbf{m}_k) \quad \forall k$$

where  $d_M$  is the **Mahalanobis distance**

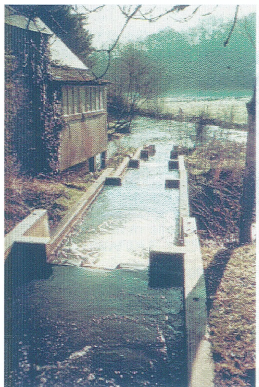
$$d_M(\mathbf{x}, \mathbf{m}_k) = \sqrt{(\mathbf{x} - \mathbf{m}_k)^t \Sigma^{-1} (\mathbf{x} - \mathbf{m}_k)}$$

### Affine discriminant functions

$$d^*(\mathbf{x}) = a_i \Leftrightarrow \left( \mathbf{x} - \frac{1}{2} (\mathbf{m}_i + \mathbf{m}_k) \right)^t \Sigma^{-1} (\mathbf{m}_i - \mathbf{m}_k) \geq \ln \frac{P(\omega_k)}{P(\omega_i)}$$

# Contexte de l'étude

## *Les passes à poissons*



à bassins successifs



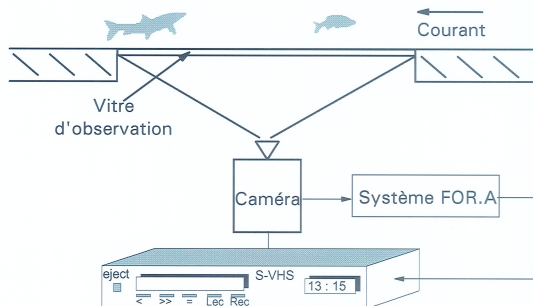
à ralentisseurs



ascenseur

## *Suivi des passes à poissons*

- Dispositif vidéo actuel de comptage par espèces :



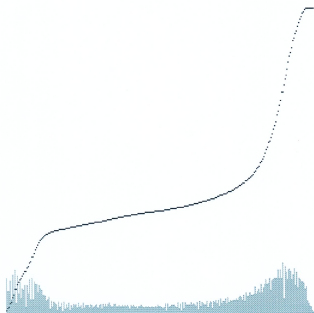
- Dépouillement des bandes vidéo enregistrées par un opérateur humain

# Acquisition des images

## *Nouvelles conditions d'acquisition*



Image



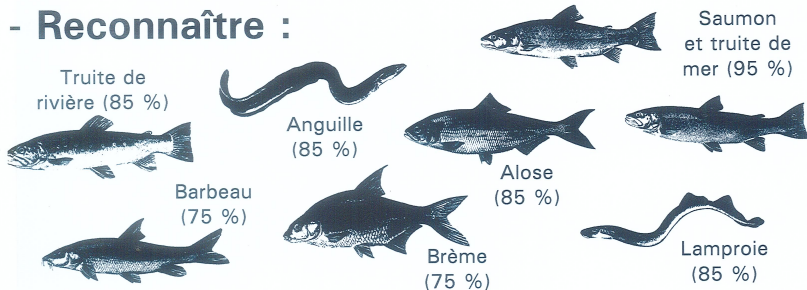
Histogramme : bimodal

## *Cahier des charges*

### *du système de vision automatique*

- Compter les poissons, toutes espèces confondues, au delà d'une certaine taille

- Reconnaître :

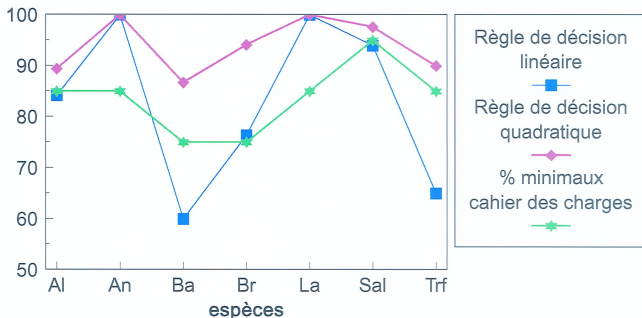


- Fournir : La comptabilisation des différentes espèces  
Les dates et heures de passage

# Mesure d'efficacité de la règle de classement

## Résultats obtenus sur la BDT

% de reconnaissance dynamiques sur la base de données test  
(comparaison des règles de décision quadratique et linéaire)



## Summary

- ▶ Chapter 1 : Introduction
- ▶ Chapter 2 : Preprocessing
- ▶ Chapter 3 : Statistical Methods
  - ▶ Bayesian rule
    - ▶ Theory
    - ▶ Two classes
    - ▶ Uninformative cost function
    - ▶ Gaussian densities
  - ▶ Supervised learning
    - ▶ Parametric methods
    - ▶ Non-parametric methods
  - ▶ Unsupervised learning
    - ▶ Optimization methods
    - ▶ Hierarchical classification
- ▶ Chapter 4 : Support Vector Machines (SVMs) and Neural Networks
- ▶ Chapter 5 : Decision Trees



## MMSE Estimator

$$\theta \in \mathbb{R}$$

$$E \left[ \left( \theta - \hat{\theta}_{MMSE}(\mathbf{x}) \right)^2 \right] = \min_{\pi} E \left[ (\theta - \pi(\mathbf{x}))^2 \right]$$

with  $\mathbf{x} = (x_1, \dots, x_n)$ , hence the **MMSE estimator**

$$\hat{\theta}_{MMSE}(\mathbf{x}) = E[\theta | \mathbf{x}]$$

$$\theta \in \mathbb{R}^p$$

The **MMSE estimator**  $\hat{\boldsymbol{\theta}}_{MMSE}(\mathbf{x}) = E[\boldsymbol{\theta} | \mathbf{x}]$  minimizes the quadratic cost

$$E \left[ (\boldsymbol{\theta} - \pi(\mathbf{x}))^t \mathbf{Q} (\boldsymbol{\theta} - \pi(\mathbf{x})) \right]$$

for any symmetric definite positive matrix  $\mathbf{Q}$  (and in particular for  $\mathbf{Q} = I_p$ , the identity matrix).

## MAP Estimator

$$\theta \in \mathbb{R}$$

The **MAP estimator**  $\hat{\theta}_{\text{MAP}}(\mathbf{x})$  minimizes the average of a “uniform” cost function

$$c(\theta - \pi(\mathbf{x})) = \begin{cases} 0 & \text{si } |\theta - \pi(\mathbf{x})| \leq \frac{\Delta}{2} \\ 1 & \text{si } |\theta - \pi(\mathbf{x})| > \frac{\Delta}{2} \end{cases}$$

and is defined by

$$c(\theta - \hat{\theta}_{\text{MAP}}(\mathbf{x})) = \min_{\pi} c(\theta - \pi(\mathbf{x}))$$

If  $\Delta$  is arbitrary small,  $\hat{\theta}_{\text{MAP}}(\mathbf{x})$  is the value of  $\pi(\mathbf{x})$  which maximizes the posterior  $p(\theta|\mathbf{x})$  hence its name **MAP estimator**. The MAP estimator is computed by setting to zero the derivative of  $p(\theta|\mathbf{x})$  (or of its logarithm) with respect to  $\theta$ .

$$\theta \in \mathbb{R}^p$$

Determine the values of  $\theta_i$  which make the partial derivatives of  $p(\boldsymbol{\theta}|\mathbf{x})$  (or of its logarithm) with respect to  $\theta_i$  equal to zero.

## Summary

- ▶ Chapter 1 : Introduction
- ▶ Chapter 2 : Preprocessing
- ▶ Chapter 3 : Statistical Methods
  - ▶ Bayesian rule
    - ▶ Theory
    - ▶ Two classes
    - ▶ Uninformative cost function
    - ▶ Gaussian densities
  - ▶ Supervised learning
    - ▶ Parametric methods
    - ▶ Non-parametric methods
  - ▶ Unsupervised learning
    - ▶ Optimization methods
    - ▶ Hierarchical classification
- ▶ Chapter 4 : Support Vector Machines (SVMs) and Neural Networks
- ▶ Chapter 5 : Decision Trees

## Histogram

### Density estimation

- **Probability:**  $P_R$

Let  $\mathbf{x}$  be a random vector of  $\mathbb{R}^p$

$$P_R = P[\mathbf{x} \in R] = \int_R f(\mathbf{u}) d\mathbf{u}$$

- **Number of vectors  $\mathbf{x}_i$  belonging to the region  $R$ :**  $N_R$

$$P[N_R = k] = \mathcal{C}_n^k P_R^k (1 - P_R)^{n-k} \quad k \in \{0, \dots, n\}$$

$$E\left(\frac{N_R}{n}\right) = P_R, \quad \text{Var}\left(\frac{N_R}{n}\right) = \frac{P_R(1 - P_R)}{n}$$

$N_R/n$  is an unbiased and convergent estimator of  $P_R$ .

- **Estimator of the probability density function:**  $f(\mathbf{x})$

If  $f$  is approximately constant on  $R$

$$\int_R f(\mathbf{u}) d\mathbf{u} \simeq f(\mathbf{x}) V_R \Rightarrow \boxed{\hat{f}(\mathbf{x}) = \frac{N_R}{nV_R}}$$

## Remarks

- ▶ If the volume  $V_R$  is fixed and the number of data  $n$  increases,  $N_R/n$  converges to  $P_R$  but  $\frac{N_R}{nV_R}$  converges to the mean of  $f$  on  $R$  and not to  $f(\mathbf{x})$ .
- ▶ If the volume  $V$  decreases ( $n$  being fixed), the region  $R$  can contain zero observation, hence  $\hat{f}(\mathbf{x}) \simeq 0$ .

⇒  $n$  and  $R$  should vary simultaneously

We build a sequence of regions  $R_k$  containing  $\mathbf{x}$  and  $N_k$  observations

- ▶ let  $V_k$  be the “volume” of the region  $R_k$
- ▶ let  $\hat{f}_k(\mathbf{x}) = \frac{N_k}{kV_k}$  be the  $k$ th estimation of  $f(\mathbf{x})$ .

How can we define the regions  $R_k$  in order to ensure  $\hat{f}_k(\mathbf{x})$  is a good estimator of  $f(\mathbf{x})$ ?

## Convergence

To ensure the **convergence** of  $\hat{f}_k(\mathbf{x})$  to  $f(\mathbf{x})$  when  $k \rightarrow \infty$ , the following three conditions have to be satisfied:

$$\begin{aligned}\lim_{k \rightarrow \infty} V_k &= 0 \\ \lim_{k \rightarrow \infty} N_k &= \infty \\ \lim_{k \rightarrow \infty} \frac{N_k}{k} &= 0.\end{aligned}$$

The problem is to choose the sequence of pairs  $(N_k, V_k)$  which satisfies these conditions. **Two methods**:

- ▶ **Parzen windows**
- ▶  **$k$  nearest neighbors**

## Parzen windows

### Definition of a kernel $\phi$

Let  $R_k$  be an hypercube of side  $h_k$  and dimension  $p$ . Thus  $V_k = h_k^p$  and

$$N_k = \sum_{i=1}^k \phi \left( \frac{\mathbf{x} - \mathbf{x}_i}{h_k} \right)$$

with

$$\phi(\mathbf{u}) = 1 \text{ if } |u(j)| \leq 1/2 \quad \forall j \in \{1, \dots, p\}$$

$$\phi(\mathbf{u}) = 0 \text{ else}$$

Indeed,  $\phi \left( \frac{\mathbf{x} - \mathbf{x}_i}{h_k} \right) = 1$  if  $\mathbf{x}_i$  is in the region  $R_k$ .

### Density estimator

$$\hat{f}_k(x) = \frac{1}{kh_k^p} \sum_{i=1}^k \phi \left( \frac{\mathbf{x} - \mathbf{x}_i}{h_k} \right)$$

## Effect of parameter $h_k$

Define

$$\delta_k(\mathbf{x}) = \frac{1}{V_k} \phi\left(\frac{\mathbf{x}}{h_k}\right)$$

- ▶ When  $h_k$  is “large”,  $\delta_k(\mathbf{x})$  varies slowly
  - ▶  $\widehat{f}_k(\mathbf{x})$  is the superimposition of  $k$  slowly varying functions
  - ▶  $\widehat{f}_k(\mathbf{x})$  is an estimator of  $f(\mathbf{x})$  with limited resolution.
- ▶ When  $h_k$  is “small”,  $\delta_k(\mathbf{x})$  varies rapidly
  - ▶  $\widehat{f}_k(\mathbf{x})$  is the superimposition of  $k$  narrow functions centered on the observed samples
  - ▶  $\widehat{f}_k(\mathbf{x})$  is a noisy estimation of  $f(\mathbf{x})$ .



## Generalization to other kernels

$\widehat{f}_k(\mathbf{x})$  is a **probability density function** if  $\phi$  satisfies the following conditions

$$\begin{aligned}\phi(\mathbf{u}) &\geq 0 \quad \forall \mathbf{u} \\ \int_{\mathbb{R}^p} \phi(\mathbf{u}) d\mathbf{u} &= 1.\end{aligned}$$

In order to ensure that  $\widehat{f}_k(\mathbf{x})$  is an **unbiased and convergent** estimator of  $f(\mathbf{x})$ ,  $\phi$  has to satisfy the following (necessary and sufficient) conditions

$$\begin{aligned}\sup_{\mathbf{u}} \phi(\mathbf{u}) &< +\infty, \quad \lim_{h_k \rightarrow \infty} \frac{1}{V_k} \phi\left(\frac{\mathbf{u}}{h_k}\right) = \delta(u) \\ \lim_{k \rightarrow \infty} V_k &= 0, \quad \lim_{k \rightarrow \infty} kV_k = +\infty\end{aligned}$$

These conditions are for instance verified for  $V_k = \frac{V_1}{\sqrt{k}}$  or  $V_k = \frac{V_1}{\ln k}$

## $k$ -nearest neighbor method

### Assumptions

$N_k$  is fixed (e.g.,  $N_k$  is the integer part of  $\sqrt{k}$ )  
 $V_k$  is estimated

### Density estimator

$$\hat{f}_k(\mathbf{x}) = \frac{N_k}{kV_k}$$

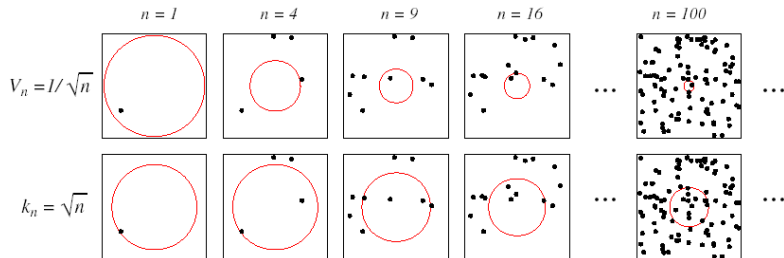
### Convergence

The estimator  $\hat{f}_k(\mathbf{x})$  converges to  $f(\mathbf{x})$  if and only if

$$\begin{aligned}\lim_{k \rightarrow \infty} N_k &= \infty \\ \lim_{k \rightarrow \infty} \frac{N_k}{k} &= 0\end{aligned}$$

## Fenêtres de Parzen vs. k-PPV (2)

- Illustration des méthodes
  - des fenêtres de Parzen
  - des k plus proches voisins



## The $k$ -nearest neighbor rule

### The nearest neighbor rule

$$d(x) = a_j \text{ if the nearest neighbor of } x \text{ belongs to } \omega_j$$

The observed vector  $x$  is affected to the class of its nearest neighbor.

### Inequality of Cover and Hart

$$P^* \leq P_1 \leq P^* \left( 2 - \frac{K}{K-1} P^* \right)$$

### The $k$ -nearest neighbor rule

$x$  is assigned to the class most common amongst its  $k$ -nearest neighbors (with a given distance measure)

## Probability of error

### Inequalities

$$P^* \leq P_k \leq P^* + \frac{1}{\sqrt{ke}} \quad \text{or} \quad P^* \leq P_k \leq P^* + \sqrt{\frac{2P_1}{k}}$$

### Approximations

When  $P^*$  is small, the following results can be obtained

$$P_1 \approx 2P^* \quad \text{and} \quad P_3 \approx P^* + 3(P^*)^2$$

## Summary

- ▶ **Chapter 1** : Introduction
- ▶ **Chapter 2** : Preprocessing
- ▶ **Chapter 3** : Statistical Methods
  - ▶ Bayesian rule
    - ▶ Theory
    - ▶ Two classes
    - ▶ Uninformative cost function
    - ▶ Gaussian densities
  - ▶ Supervised learning
    - ▶ Parametric methods
    - ▶ Non-parametric methods
  - ▶ **Unsupervised learning**
    - ▶ **Optimization methods**
    - ▶ Hierarchical classification
    - ▶ Density-based methods
    - ▶ Mixture models
    - ▶ Spectral methods
- ▶ **Chapter 4** : Support Vector Machines (SVMs) and Neural Networks
- ▶ **Chapter 5** : Decision Trees

## Unsupervised learning

$N$  unlabelled data vectors of  $\mathbb{R}^p$  denoted as  $\mathbf{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$  which should be split into  $K$  classes  $\omega_1, \dots, \omega_K$ .

### Motivations

- ▶ Supervised learning is **costly**
- ▶ The classes can **change** with time
- ▶ Provide some **information** about the data structure

### Optimal solution

Number of partitions of  $\mathbf{X}$  in  $K$  subsets

$$P(N, K) = \frac{1}{K!} \sum_{i=0}^K i^N (-1)^{K-i} C_K^i \quad K < N$$

*Example:*  $P(100, 5) \approx 10^{68}$  !

## Partition with minimum mean square error

### Mean square error of a partition

Mean square error (MSE) of the class  $\omega_i$  and of the partition  $\mathcal{X}$

$$E_i^2 = \sum_{k=1}^{N_i} d^2(\mathbf{x}_k, \mathbf{g}_i), \quad E^2 = \sum_{i=1}^K E_i^2$$

where  $\mathbf{g}_i = \frac{1}{N_i} \sum_{k=1}^{N_i} \mathbf{x}_k$  is the centroid of the class  $\omega_i$

### Properties

$$\sum_{k=1}^{N_i} d^2(\mathbf{x}_k, \mathbf{y}) = \sum_{k=1}^{N_i} d^2(\mathbf{x}_k, \mathbf{g}_i) + N_i d^2(\mathbf{g}_i, \mathbf{y})$$

In particular, for  $\mathbf{y} = \mathbf{g}$  (data centroid), we obtain

$$\sum_{i=1}^K \sum_{k=1}^{N_i} d^2(\mathbf{x}_k, \mathbf{g}) = \underbrace{\sum_{i=1}^K \sum_{k=1}^{N_i} d^2(\mathbf{x}_k, \mathbf{g}_i)}_{E^2} + \sum_{i=1}^K N_i d^2(\mathbf{g}_i, \mathbf{g})$$

MSE of  $\mathcal{X}$  = within-class MSE + between-class MSE



## ISODATA algorithm

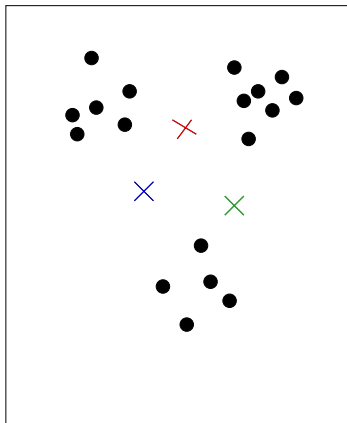
Search a partition of  $\mathbf{X}$  ensuring a **local minimum** of  $E^2$

1. Initial **choice** of the number of classes and the class centroids
2. Assign each vector  $\mathbf{x}_i$  to  $\omega_j$  (using the **centroid distance rule**) such that

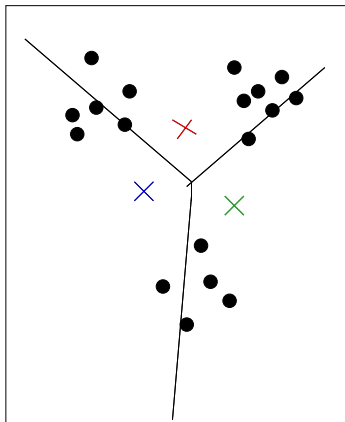
$$d(\mathbf{x}_i, \mathbf{g}_j) = \inf_k d(\mathbf{x}_i, \mathbf{g}_k)$$

3. **Compute the centroids**  $\mathbf{g}_k^*$  of the new classes  $\omega_k^*$
  4. **Repeat steps 2) and 3)** until convergence
- ▶ *Improved version of ISODATA*
    - ▶ Two classes are merged if their centroids are close
    - ▶ A class is split if it contains too many vectors  $\mathbf{x}_i$  or if its mean square error is too large
  - ▶ *Convergence*: see notes or textbooks

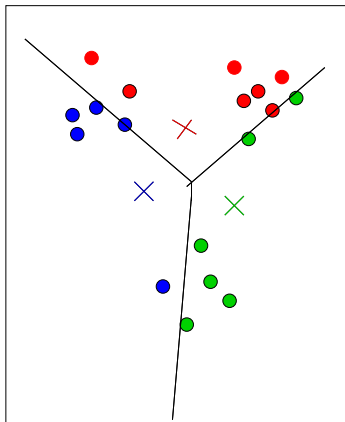
## K-means



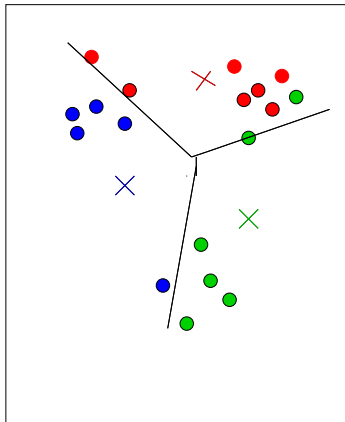
## K-means



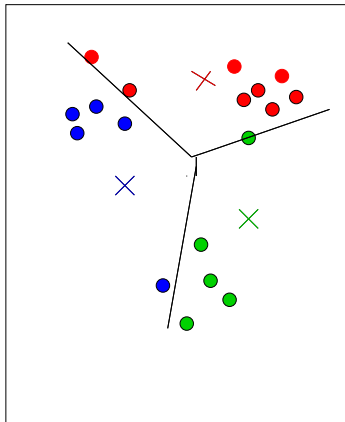
## K-means



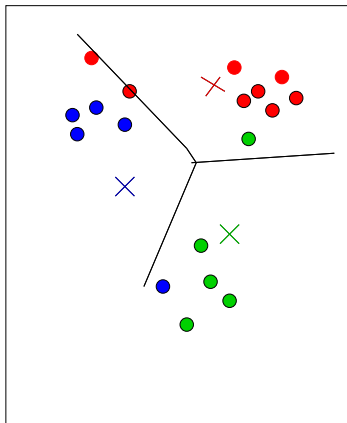
## K-means



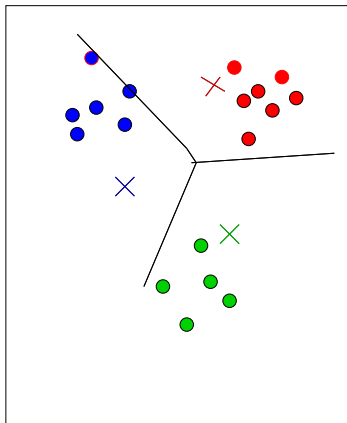
## K-means



## K-means

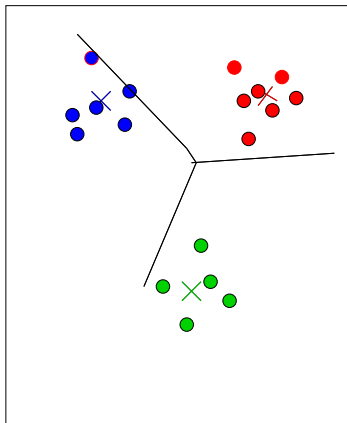


## K-means

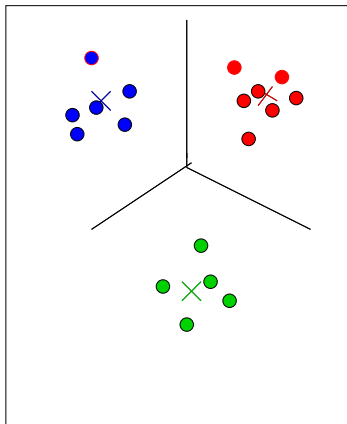




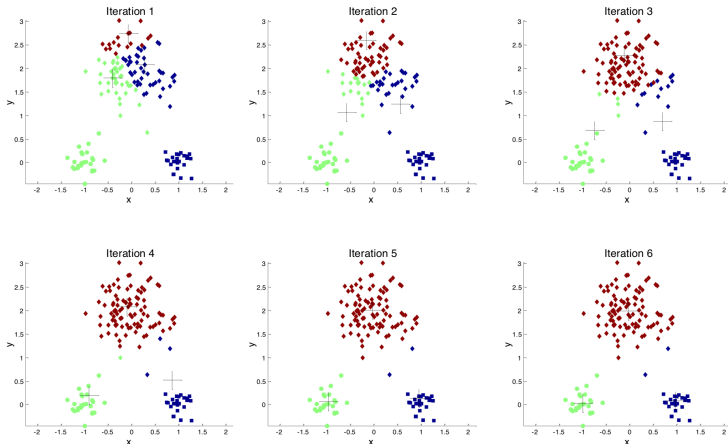
## K-means



## K-means

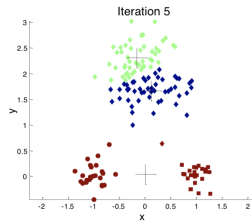
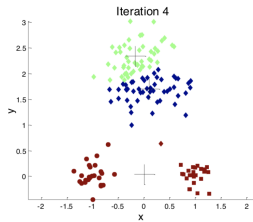
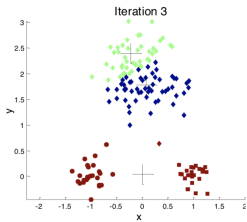
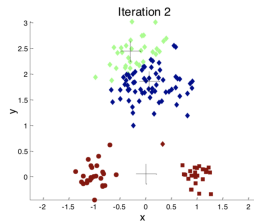
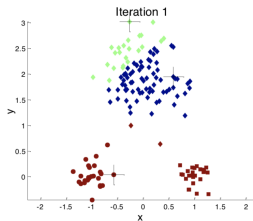


## Initialization Problems

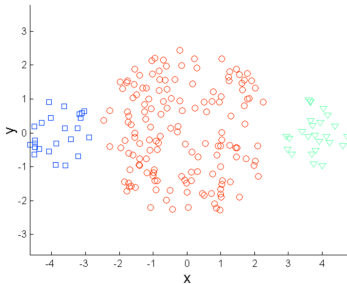


Thanks to Jing Gao from SUNY Buffalo university for her slides!!

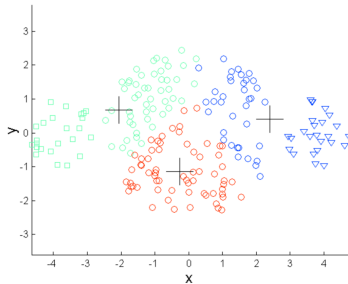
## Initialization Problems



## Limitations of K-means: Differing Sizes

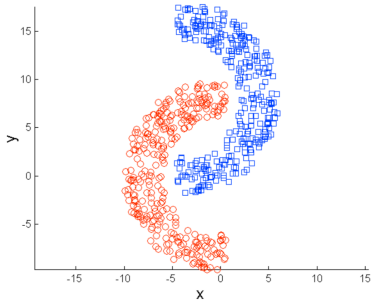


Original Points

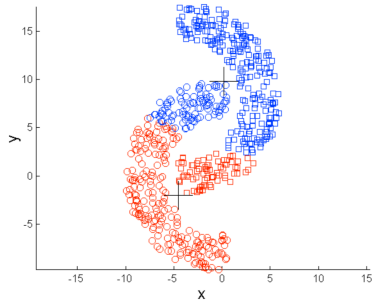


K-means (3 Clusters)

# Limitations of K-means: Irregular Shapes

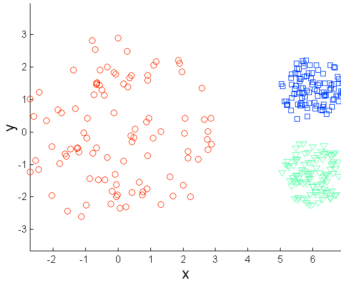


Original Points

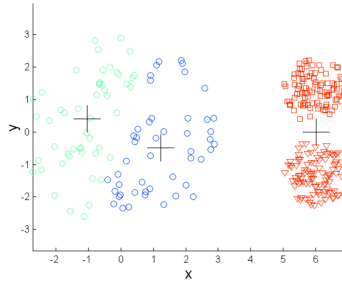


K-means (2 Clusters)

## Limitations of K-means: Differing Density



Original Points



K-means (3 Clusters)

## Summary

- ▶ Chapter 1 : Introduction
- ▶ Chapter 2 : Preprocessing
- ▶ Chapter 3 : Statistical Methods
  - ▶ Bayesian rule
    - ▶ Theory
    - ▶ Two classes
    - ▶ Uninformative cost function
    - ▶ Gaussian densities
  - ▶ Supervised learning
    - ▶ Parametric methods
    - ▶ Non-parametric methods
  - ▶ Unsupervised learning
    - ▶ Optimization methods
    - ▶ Hierarchical classification
    - ▶ Density-based methods
    - ▶ Mixture models
    - ▶ Spectral methods
- ▶ Chapter 4 : Support Vector Machines (SVMs) and Neural Networks
- ▶ Chapter 5 : Decision Trees



## Hierarchical classification

**Ascending hierarchy:** method of distances

- ▶ Distance *Min* (single linkage algorithm)

$$d(X_i, X_j) = \min_{\mathbf{x} \in X_i, \mathbf{y} \in X_j} d(\mathbf{x}, \mathbf{y})$$

This distance favors elongated classes

- ▶ Distance *Max* (complete linkage algorithm)

$$d(X_i, X_j) = \max_{\mathbf{x} \in X_i, \mathbf{y} \in X_j} d(\mathbf{x}, \mathbf{y})$$

- ▶ Average linkage algorithm

$$d(X_i, X_j) = \frac{1}{N_i N_j} \sum_{\mathbf{x} \in X_i, \mathbf{y} \in X_j} d(\mathbf{x}, \mathbf{y})$$

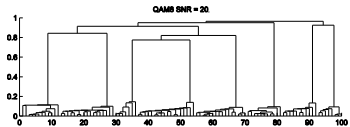
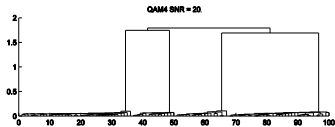
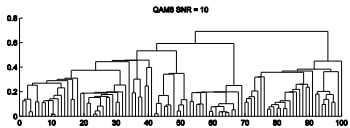
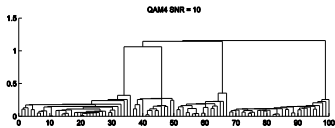
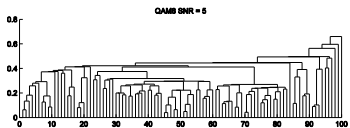
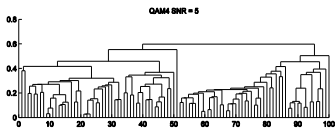
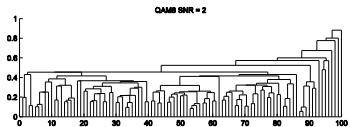
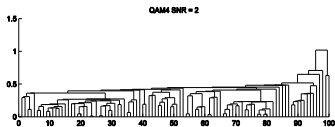
- ▶ Distance between the means

$$d(X_i, X_j) = d(\mathbf{g}_i, \mathbf{g}_j)$$

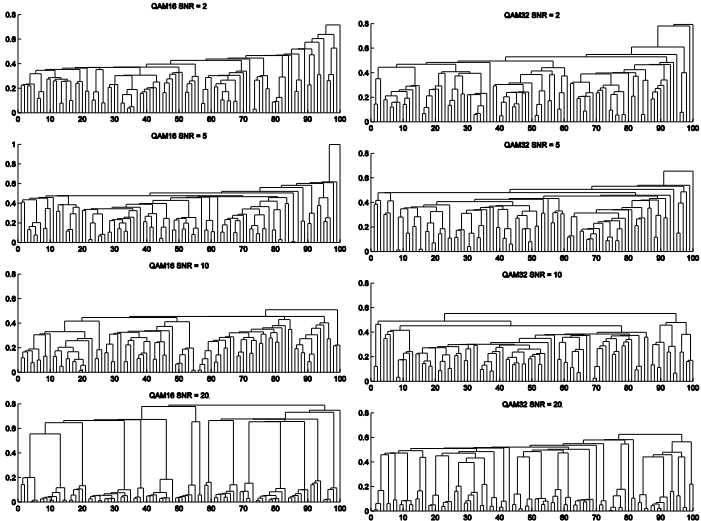
where  $X_i$  and  $X_j$  have cardinals  $N_i$  and  $N_j$  and centroids  $\mathbf{g}_i$  and  $\mathbf{g}_j$ .

**Representation using a tree whose nodes indicate the different groups**

# Classification of Modulations



# Classification of Modulations



## Summary

- ▶ Chapter 1 : Introduction
- ▶ Chapter 2 : Preprocessing
- ▶ Chapter 3 : Statistical Methods
  - ▶ Bayesian rule
    - ▶ Theory
    - ▶ Two classes
    - ▶ Uninformative cost function
    - ▶ Gaussian densities
  - ▶ Supervised learning
    - ▶ Parametric methods
    - ▶ Non-parametric methods
  - ▶ Unsupervised learning
    - ▶ Optimization methods
    - ▶ Hierarchical classification
    - ▶ Density-based methods
    - ▶ Mixture models
    - ▶ Spectral methods
- ▶ Chapter 4 : Linear Discriminant Functions and Neural Networks
- ▶ Chapter 5 : Decision Trees

## Density-based clustering methods

### Definitions

- ▶ **Neighborhood** of a point  $p \in \mathbf{X}$

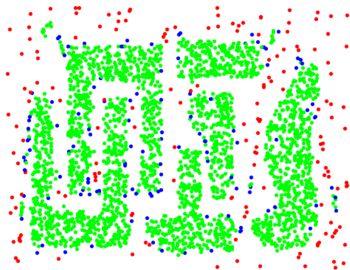
$$N_\epsilon(p) = \{q \in \mathbf{X} \mid d(p, q) \leq \epsilon\}$$

- ▶ **Core point**: a point  $p$  is a core point if it has more than MinPts neighbors in  $N_\epsilon(p)$ .
- ▶ **Border point**: a point  $p$  is a border point if it has less than MinPts neighbors in  $N_\epsilon(p)$  and if it is in the neighborhood of a core point
- ▶ **Noise point**: a noise point is a point that is not a core point nor a border point

## Core, border and noise points



Original Points



Point types: **core**,  
**border** and **outliers**

$\epsilon = 10$ , MinPts = 4

Thanks to Jing Gao from SUNY Buffalo university for her slides!!

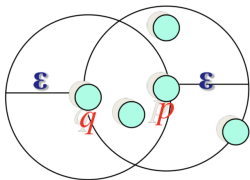
## Relationships between different points

### Directly density-reachable points

A point  $q$  is directly reachable from  $p$  for  $(\epsilon, \text{MinPts})$  if

- ▶  $q \in N_\epsilon(p)$
- ▶  $p$  is a core point, i.e.,  $N_\epsilon(p)$  contains more than  $\text{MinPts}$  points.

### Example



- ▶  $q$  is directly density-reachable from  $p$
- ▶  $p$  is not directly density-reachable from  $q$  ( $q$  is not a core point)

It is an asymmetric relationship!

$\text{MinPts} = 4$

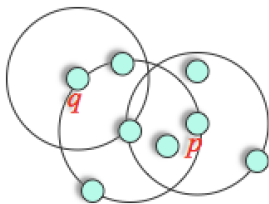
## Relationships between different points

### Density-reachable points

A point  $q$  is density-reachable from  $p$  for  $(\epsilon, \text{MinPts})$  if there is chain of points  $p_i$  such that

- ▶  $p_1 = p$  and  $p_N = q$
- ▶ each point  $p_{i+1}$  is directly density-reachable from  $p_i$ .

### Example



- ▶  $q$  is density-reachable from  $p$
- ▶  $p$  is not density-reachable from  $q$  ( $q$  is not a core point)

It is an asymmetric relationship!

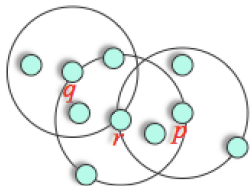


## Relationships between different points

### Density-connectivity

Points  $p$  and  $q$  are density-connected for  $(\epsilon, \text{MinPts})$  if there is an object  $r$  such that both  $p$  and  $q$  are density reachable from  $r$ .

### Example



It is a symmetric relationship!

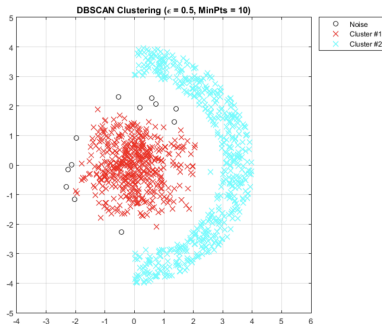
## Density-Based Spatial Clustering of Applications with Noise (DBSCAN)

### Cluster definition

A **cluster**  $C$  is defined as a maximal set of density-connected points

- ▶ **Maximality**: if  $p \in C$  and if  $q$  is density-reachable from  $p$ , then  $q \in C$ .
- ▶ **Connectivity**: for all  $(p, q) \in C$ ,  $p$  and  $q$  are density-connected.

### Example



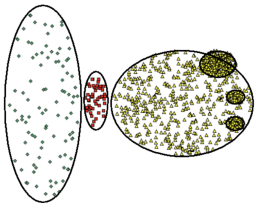
## DBSCAN

### A simple algorithm

- ▶ Select a point  $p$
- ▶ Determine all density-reachable points from  $p$  for  $(\epsilon, \text{MinPts})$ .
- ▶ If  $p$  is a core point, i.e., if the cardinal of  $N_\epsilon(p)$  is larger than  $\text{MinPts}$ , a cluster is formed
- ▶ If  $p$  is a border point, DBSCAN visits the next point
- ▶ Continue the procedure until all points have been visited

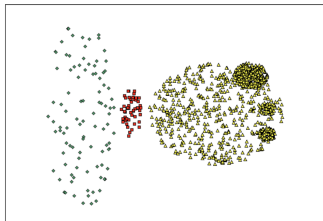
The DBSCAN algorithm generally provides a result **independent of the order** the points have been processed.

## How can we choose $\epsilon$ and MinPts?

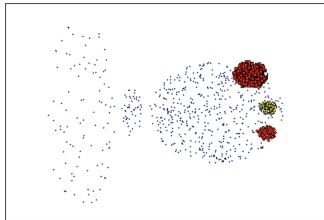


Original Points

- Cannot handle varying densities
- sensitive to parameters—hard to determine the correct set of parameters



(MinPts=4, Eps=9.92).



(MinPts=4, Eps=9.75)

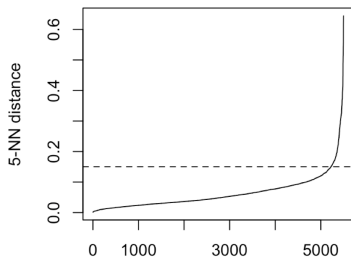
Thanks to Jing Gao from SUNY Buffalo university for her slides!!

## How can we choose $\epsilon$ and MinPts?

### Some ideas

- ▶ Points belonging to a cluster have a distance to their  $k$ -nearest neighbor (denoted as  $k$ -dist) smaller than noise points
- ▶ Compute all the  $k$ -dist values and sort them in increasing order
- ▶ Choose  $\epsilon$  as the value of  $k$ -dist associated with a sharp change in the curve (this value does not vary significantly with the value of  $k$ )

### Example



## Conclusions

### Pros

- ▶ Clusters of arbitrary shapes
- ▶ Robustness to noise

### Cons

- ▶ Problems with clusters of different densities
- ▶ Parameter determination can be difficult

### Some references (related to DBSCAN)

- ▶ M. Ester, H.-P. Kriegel, J. Sander and X. Xu, "A density-based algorithm for discovering clusters in large spatial databases with noise," in Proc. Int. Conf. Knowledge Discovery and Data Mining (KDD'96), Portland, Oregon, Aug. 1996.
- ▶ J. Gao, Slides on Data Mining and Bioinformatics, University at Buffalo, <https://www.cse.buffalo.edu/~jing/cse601/fa13/>.

## Summary

- ▶ Chapter 1 : Introduction
- ▶ Chapter 2 : Preprocessing
- ▶ Chapter 3 : Statistical Methods
  - ▶ Bayesian rule
    - ▶ Theory
    - ▶ Two classes
    - ▶ Uninformative cost function
    - ▶ Gaussian densities
  - ▶ Supervised learning
    - ▶ Parametric methods
    - ▶ Non-parametric methods
  - ▶ Unsupervised learning
    - ▶ Optimization methods
    - ▶ Hierarchical classification
    - ▶ Density-based methods
    - ▶ Mixture models
    - ▶ Spectral methods
- ▶ Chapter 4 : Linear Discriminant Functions and Neural Networks
- ▶ Chapter 5 : Decision Trees

## Gaussian mixture (unsupervised learning)

**Idea:**  $N$  unlabelled data vectors of  $\mathbb{R}^P$  denoted as  $\mathbf{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$  emerge from  $K$  Gaussian components/classes denoted as  $\omega_1, \dots, \omega_K$ .

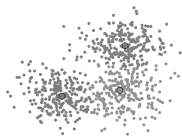
### Gaussian mixture model (GMM)

- ▶ Definition

$$p(\mathbf{x}|\theta) = \sum_{j=1}^K \pi_j p(\mathbf{x}|\boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j) \quad (1)$$

- ▶  $\boldsymbol{\theta} = (\pi_1, \dots, \pi_j, \dots, \pi_K, \boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_j, \dots, \boldsymbol{\mu}_K, \boldsymbol{\Sigma}_1, \dots, \boldsymbol{\Sigma}_j, \dots, \boldsymbol{\Sigma}_K)$  contains all the parameters of the mixture model.

Example:  $K = 3$  seems reasonable ( $\boldsymbol{\theta}$  is unknown)





## Gaussian mixture (supervised learning)

When the  $N$  data  $\mathbf{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$  are *complete* (i.e. labelled, assigned to classes), the parameter estimation problem is straightforward (each Gaussian can be estimated separately). Besides the data, we also have their labels:  $\{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\}$  where  $y_i \in \{\omega_1, \dots, \omega_K\}$ .

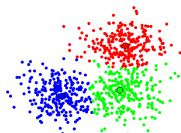
### Gaussian mixture model (GMM)

- ▶ Definition

$$p(\mathbf{x}|\boldsymbol{\theta}) = \sum_{j=1}^K \pi_j p(\mathbf{x}|\boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j) \quad (2)$$

- ▶ **Assignments:** binary variable  $\delta(j|i)$  assigns data  $\mathbf{x}_i$  to the  $j$ th Gaussian (class  $\omega_j$ ) if  $\delta(j|i) = 1$  ( $\delta(j|i) = 0$  otherwise).

Example:  $K = 3$  is now given (along with data assignments).



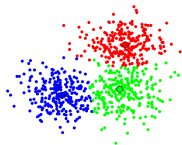
## Gaussian mixture (supervised learning)

When the  $N$  data  $\mathbf{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$  are *complete* (i.e., labelled or assigned to classes), the parameter estimation problem is straightforward (each Gaussian density can be estimated separately).

### Estimation of the Gaussian mixture model

- ▶  $\hat{\pi}_j \leftarrow \frac{\hat{n}_j}{n}$       with       $\hat{n}_j = \sum_{i=1}^n \delta(j|i)$
- ▶  $\hat{\boldsymbol{\mu}}_j \leftarrow \frac{1}{\hat{n}_j} \sum_{i=1}^n \delta(j|i) \mathbf{x}_i$
- ▶  $\hat{\boldsymbol{\Sigma}}_j \leftarrow \frac{1}{\hat{n}_j} \sum_{i=1}^n \delta(j|i) (\mathbf{x}_i - \hat{\boldsymbol{\mu}}_j)(\mathbf{x}_i - \hat{\boldsymbol{\mu}}_j)^T$   
where  $\delta(j|i) = 1$  assigns data  $\mathbf{x}_i$  to the  $j$ th Gaussian density

Example:  $K = 3$  is now given (along with data assignments).



## Gaussian mixture (back to unsupervised case)

Without assignment (unsupervised learning), we start from an initial setting of the parameters  $\theta : \theta^0 = (\pi_1^0, \dots, \pi_K^0, \mu_1^0, \Sigma_1^0, \dots, \mu_K^0, \Sigma_K^0)$  and compute

$$P(y_i = \omega_j | \mathbf{x}_i, \theta^0) = \frac{p(\mathbf{x}_i | y_i = \omega_j, \theta^0) P(\omega_j)}{\sum_j p(\mathbf{x}_i | y_i = \omega_j, \theta^0) P(\omega_j)} \quad (3)$$

$$P(y_i = \omega_j | \mathbf{x}_i, \theta^0) = \frac{\pi_j^0 p(\mathbf{x}_i | \mu_j^0, \Sigma_j^0)}{\sum_j \pi_j^0 p(\mathbf{x}_i | \mu_j^0, \Sigma_j^0)} \quad (4)$$

### Soft assignment

- ▶  $\hat{\delta}(j|i) = P(y_i = \omega_j | \mathbf{x}_i, \theta^0)$
- ▶  $\sum_{j=1}^K \hat{\delta}(j|i) = 1.$

## Expectation-Maximization (EM-Algorithm for GMM estimation)

Without assignment (unsupervised case), we start from  $\theta^0$  and iterate soft-assignments that “complete the incomplete data” (Expectation-step) and parameter refinement (Maximisation-step). We can show that the new setting of the parameters  $\theta^{(k+1)}$  maximises the log-likelihood of the “completed” data.

### EM-algorithm:

1. **Initialization** Specify  $\theta^{(k=0)}$ .
2. **Repeat**

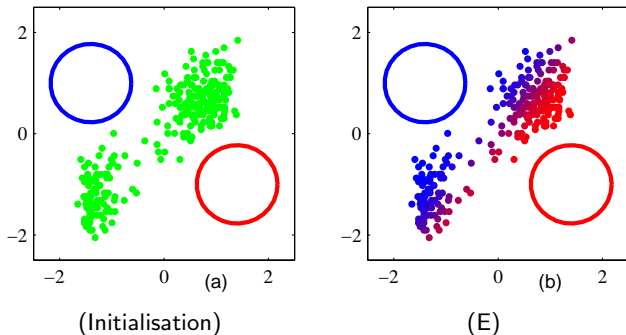
(E-step) soft-assignments of  $\mathbf{x}_i \forall i$

$$\widehat{\delta}(j|i) \leftarrow P(y_i = \omega_j | \mathbf{x}_i, \theta^{(k)}) \quad (5)$$

(M-step) Refine  $\theta^{(k+1)}$ :

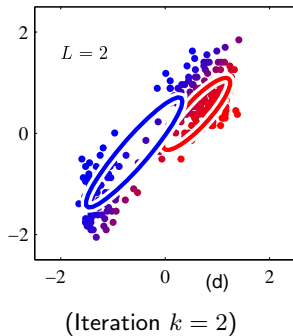
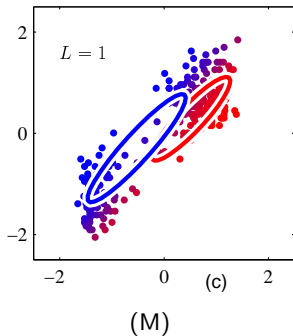
- ▶  $\widehat{\pi}_j \leftarrow \frac{\widehat{n}_j}{n}$  with  $\widehat{n}_j = \sum_{i=1}^n \widehat{\delta}(j|i)$
- ▶  $\widehat{\boldsymbol{\mu}}_j \leftarrow \frac{1}{\widehat{n}_j} \sum_{i=1}^n \widehat{\delta}(j|i) \mathbf{x}_i$
- ▶  $\widehat{\boldsymbol{\Sigma}}_j \leftarrow \frac{1}{\widehat{n}_j} \sum_{i=1}^n \widehat{\delta}(j|i) (\mathbf{x}_i - \widehat{\boldsymbol{\mu}}_j)(\mathbf{x}_i - \widehat{\boldsymbol{\mu}}_j)^T$
- ▶  $k \leftarrow k + 1$

## EM-Algorithm for GMM estimation (demo from book by Bishop\*)

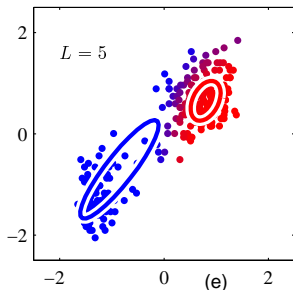


\* Christopher Bishop, *Pattern Recognition and Machine Learning*. New-York: Springer Verlag, 2006.

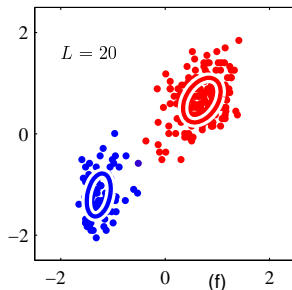
## EM-Algorithm for GMM estimation (demo 2/3)



## EM-Algorithm for GMM estimation (demo 3/3)



(Iteration  $k = 5$ )



(Iteration  $k = 20$ )

### Convergence

- ▶ The EM algorithm monotonically increases the log-likelihood of the data.
- ▶ We have  $l(\theta^0) < l(\theta^1) < \dots < l(\theta^k)$  with  $l(\theta^k) = \sum_{i=1}^n \ln p(\mathbf{x}_i | \theta^k)$

## Mélange de lois AGGD

Vraisemblance du mélange

$$p(x|\theta) = \sum_{j=1}^M p_j p(x|\theta_j)$$

avec

$$p(x|\theta_j) = \begin{cases} \frac{\delta_j^{1/\lambda_j}}{\gamma_j^{1/\lambda_j} \Gamma(1+1/\lambda_j)} \exp\left(-\frac{\delta_j}{\gamma_j} \left(\frac{\mu_j - x}{\alpha_j}\right)^{\lambda_j}\right) & \text{si } x < \mu_j \\ \frac{\delta_j^{1/\lambda_j}}{\gamma_j^{1/\lambda_j} \Gamma(1+1/\lambda_j)} \exp\left(-\frac{\delta_j}{\gamma_j} \left(\frac{x - \mu_j}{1 - \alpha_j}\right)^{\lambda_j}\right) & \text{si } x \geq \mu_j \end{cases}$$

- $\alpha_j$  : paramètre d'**asymétrie**
- $\mu_j$  : paramètre de **décalage**
- $\gamma_j$  : paramètre d'**échelle**
- $\lambda_j$  : paramètre de **forme**
- $\delta_j = \frac{2\alpha_j^{\lambda_j}(1-\alpha_j)^{\lambda_j}}{\alpha_j^{\lambda_j} + (1-\alpha_j)^{\lambda_j}}$
- $M$  : nombre de composantes



## Algorithme EM + Binarisation

**Augmentation des données** avec des labels :  $\mathbf{x} \rightarrow (\mathbf{x}, \mathbf{Z})$  où  
 $\mathbf{x} = (x^1, \dots, x^N)$  et  $\mathbf{Z} = (z^1, \dots, z^N)$  avec  $z^i$  de dimension  $M$   
(correspondant aux  $M$  classes)  
⇒ Log-vraisemblance augmentée

$$L(\mathbf{x}, \mathbf{Z}; \theta) = \sum_{i=1}^N \sum_{m=1}^M z_m^i \log(p_m p(x^i; \theta_m))$$

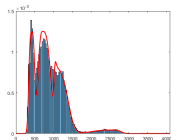
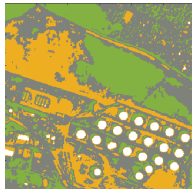
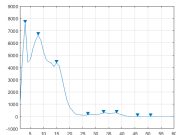
BUT : estimer  $\mathbf{p} = (p_1, \dots, p_M)$  et  $\theta$

- **Étape E** *Expectation* : calcul d'une espérance conditionnelle de la log-vraisemblance (fonction  $Q$ ) afin de mettre à jour les labels  $\mathbf{Z}$ ,
- **Étape B** *Binarization* : seuillage de la matrice des labels  $\mathbf{Z}$  en utilisant la loi multinomiale,
- **Étape M** *Maximization* : mise à jour des paramètres par maximisation de la vraisemblance de chaque classe (ML).

Mélange de lois gaussiennes généralisées asymétriques (AGGD)  
Estimation des paramètres du mélange  
Sélection de modèle  
Conclusions et perspectives

Algorithme EMB + Minimum message length  
Fusion de modes proches

## Algorithme 2



19 / 78

## Summary

- ▶ Chapter 1 : Introduction
- ▶ Chapter 2 : Preprocessing
- ▶ Chapter 3 : Statistical Methods
- ▶ Chapter 4 : Support Vector Machines (SVMs) and Neural Networks
  - ▶ Linear discriminant functions
    - ▶ The perceptron algorithm
    - ▶ The optimal Wiener-Hopf filter
    - ▶ The LMS algorithm
  - ▶ Support vector machines (SVMs)
    - ▶ Introduction
    - ▶ Optimal separating hyperplane
    - ▶ Optimization problem
    - ▶ The “soft-margin SVM” classifier
    - ▶ Non-separable case: the “ $\nu$ -SVM” classifier
    - ▶ Non-linear preprocessing - kernels
  - ▶ Neural networks
    - ▶ The multi-layer perceptron
    - ▶ The backpropagation algorithm
- ▶ Chapter 5 : Decision Trees

## Linear discriminant functions

### Bayesian classifier (2 classes)

$$d(\mathbf{x}) = a_1 \Leftrightarrow \hat{f}(\mathbf{x}|\omega_1) P(\omega_1) \geq \hat{f}(\mathbf{x}|\omega_2) P(\omega_2)$$

i.e.

$$d(\mathbf{x}) = a_1 \Leftrightarrow g(\mathbf{x}) > 0$$

$$d(\mathbf{x}) = a_2 \Leftrightarrow g(\mathbf{x}) < 0$$

where  $g(\mathbf{x}) = \hat{f}(\mathbf{x}|\omega_1) P(\omega_1) - \hat{f}(\mathbf{x}|\omega_2) P(\omega_2)$  is a **discriminant function** (in general it is non-linear)

### Linear discriminant functions (2 classes)

$$d(\mathbf{x}) = a_1 \Leftrightarrow g(\mathbf{x}) = \mathbf{w}^t \mathbf{x} > 0$$

$$d(\mathbf{x}) = a_2 \Leftrightarrow g(\mathbf{x}) = \mathbf{w}^t \mathbf{x} < 0$$

with  $\mathbf{w} = (w_1, \dots, w_{p-1}, w_p)^t$  and  $\mathbf{x} = (x_1, \dots, x_{p-1}, 1)^t$

*Pb: how can we choose  $\mathbf{w}$ ?*

## The perceptron algorithm

### Hypotheses

Two linearly separable classes: there exists  $\mathbf{w}^*$  such that

$$\begin{aligned}(\mathbf{w}^*)^t \mathbf{x} &> 0 & \forall \mathbf{x} \in \omega_1 \\ (\mathbf{w}^*)^t \mathbf{x} &< 0 & \forall \mathbf{x} \in \omega_2\end{aligned}$$

### Determination of $\mathbf{w}^*$

Optimization of a cost function

$$J(\mathbf{w}) = - \sum_{\mathbf{x} \in Y} \delta_{\mathbf{x}} \mathbf{w}^t \mathbf{x}$$

with

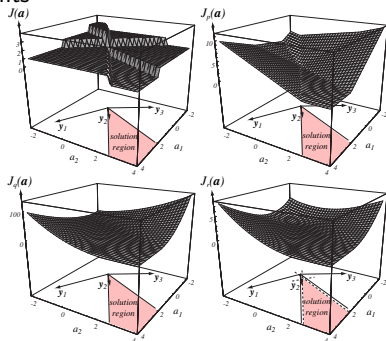
$$\delta_{\mathbf{x}} = 1 \text{ if } \mathbf{x} \in \omega_1, \quad \delta_{\mathbf{x}} = -1 \text{ if } \mathbf{x} \in \omega_2$$

and  $Y$  is the set of vectors of the training set that are misclassified by the following rule

$$d(\mathbf{x}) = a_1 \text{ if } \mathbf{w}^t \mathbf{x} > 0, \quad d(\mathbf{x}) = a_2 \text{ if } \mathbf{w}^t \mathbf{x} < 0$$

## Properties

- ▶  $J(\mathbf{w}) \geq 0$  and  $J(\mathbf{w}) = 0$  if there is no error
- ▶ The cost function is continuous **piecewise linear** and is not defined at discontinuity points



**FIGURE 5.11.** Four learning criteria as a function of weights in a linear classifier. At the upper left is the total number of patterns misclassified, which is piecewise constant and hence unacceptable for gradient descent procedures. At the upper right is the Perceptron criterion (Eq. 16), which is piecewise linear and acceptable for gradient descent. The lower left is squared error (Eq. 32), which has nice analytic properties and is useful even when the patterns are not linearly separable. The lower right is the square error with margin (Eq. 33). A designer may adjust the margin  $b$  in order to force the solution vector to lie toward the middle of the  $b = 0$  solution region in hopes of improving generalization of the resulting classifier. From: Richard O. Duda, Peter E. Hart, and David G. Stork, *Pattern Classification*. Copyright © 2001 by John Wiley & Sons, Inc.

## The steepest descent algorithm

### Definition

$\mathbf{w}_k(n)$ : value of the weight  $\mathbf{w}_k$  at the  $n$ th iteration

$$\mathbf{w}_k(n+1) = \mathbf{w}_k(n) - \mu \left. \frac{\partial J(\mathbf{w})}{\partial \mathbf{w}_k} \right|_{\mathbf{w}_k = \mathbf{w}_k(n)}$$

$\mu > 0$ : learning parameter

Update of the weights  $\mathbf{w}_k$

$$\mathbf{w}(n+1) = \mathbf{w}(n) + \mu \sum_{\mathbf{x} \in Y} \delta_{\mathbf{x}} \mathbf{x}$$

### Remarks

- ▶ The rule is repeated until convergence
- ▶ Convergence after a finite number of iterations  $n_0$

$$\hat{\mathbf{w}}_k(n_0) = \hat{\mathbf{w}}_k(n_0+1) = \hat{\mathbf{w}}_k(n_0+2) = \dots$$

when the learning classes  $\omega_1$  and  $\omega_2$  are linearly separable.

## Convergence results

### Fixed gain $\mu$

$$\mu \in ]0, 2]$$

### Variable gain $\mu_n$

$$\sum_{n=0}^{+\infty} \mu(n) = +\infty \text{ and } \sum_{n=0}^{+\infty} \mu^2(n) < +\infty$$

- ▶ **Fast convergence** if  $\mu(n)$  is close to 2
- ▶ **Small residual error** if  $\mu(n)$  is close to 0

*Example:*  $\mu(n) = \frac{c}{n}$ .



## Sequential algorithm

### The perceptron algorithm (Rosenblatt, 1950)

1. Initialize the weights:  $\mathbf{w}_k(0)$
2. Process a new vector from the learning set

We choose an input vector  $\mathbf{x}(n)$  with known class and we define the desired output  $d(n)$  as follows

$$d(n) = 1 \text{ if } \mathbf{x}(n) \in \omega_1, \quad d(n) = -1 \text{ if } \mathbf{x}(n) \in \omega_2$$

3. Compute the network output ( $f_h$  is the "sign" function)

$$y(n) = f_h \left( \sum_{i=1}^p \mathbf{w}_i(n) \mathbf{x}_i(n) \right)$$

4. Weight update

$$\hat{\mathbf{w}}_k(n+1) = \hat{\mathbf{w}}_k(n) + \mu \mathbf{x}_k(n) [d(n) - y(n)]$$

where  $e(n) = d(n) - y(n) \in \{-2, 0, 2\}$  is the output error

5. Back to 2) until convergence

## Convergence

- ▶ Convergence after  $n_0$  iterations ( $e(n) = 0$ )

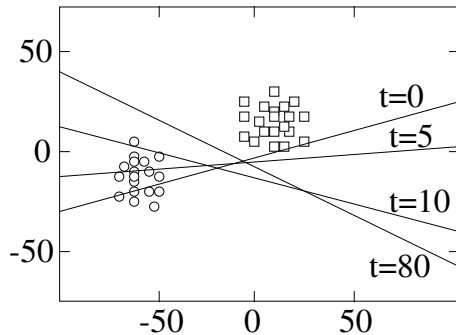
$$\hat{\mathbf{w}}_k(n_0) = \hat{\mathbf{w}}_k(n_0 + 1) = \hat{\mathbf{w}}_k(n_0 + 2) = \dots$$

when the two learning classes  $\omega_1$  and  $\omega_2$  are **linearly separable** (separated by a hyperplane)

- ▶ Variable learning parameter:  $\mu(n)$
- ▶ Generalization to  $K > 2$  classes: **Kesler's construction**

## Toy example

**Example:** 1 layer, 80 input samples,  $\mu = 0.01$



## Other algorithms

### Optimal Wiener-Hopf filter

$$J(\mathbf{w}) = \frac{1}{2} E [e^2(n)] = \frac{1}{2} E [(d(n) - \mathbf{w}^t \mathbf{x}(n))^2]$$

### Steepest descent method

$$\mathbf{w}_k(n+1) = \mathbf{w}_k(n) - \mu \left. \frac{\partial J(\mathbf{w})}{\partial \mathbf{w}_k} \right|_{\mathbf{w}_k = \mathbf{w}_k(n)}$$

### The LMS algorithm

$$\hat{\mathbf{w}}_k(n+1) = \hat{\mathbf{w}}_k(n) + \mu \mathbf{x}_k(n) [d(n) - y(n)]$$

## Summary

- ▶ Chapter 1 : Introduction
- ▶ Chapter 2 : Preprocessing
- ▶ Chapter 3 : Statistical Methods
- ▶ Chapter 4 : Support Vector Machines (SVMs) and Neural Networks
  - ▶ Linear discriminant functions
    - ▶ The perceptron algorithm
    - ▶ The optimal Wiener-Hopf filter
    - ▶ The LMS algorithm
  - ▶ Support vector machines (SVMs)
    - ▶ Introduction
    - ▶ Optimal separating hyperplane
    - ▶ Optimization problem
    - ▶ The “soft-margin SVM” classifier
    - ▶ Non-separable case: the “ $\nu$ -SVM” classifier
    - ▶ Non-linear preprocessing - kernels
  - ▶ Neural networks
    - ▶ The multi-layer perceptron
    - ▶ The backpropagation algorithm
- ▶ Chapter 5 : Decision Trees

## Support vector machines (SVMs)

### Learning set

$$\mathcal{B} = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\}$$

where  $\mathbf{x}_1, \dots, \mathbf{x}_n$  are  $n$  vectors of  $\mathbb{R}^p$  and  $y_1, \dots, y_n$  are binary variables

$$y_i = 1 \text{ if } \mathbf{x}_i \in \omega_1, \quad y_i = -1 \text{ if } \mathbf{x}_i \in \omega_2$$

### Hyperplane definition

$$g_{\mathbf{w},b}(\mathbf{x}) = \mathbf{w}^T \mathbf{x} - b = 0$$

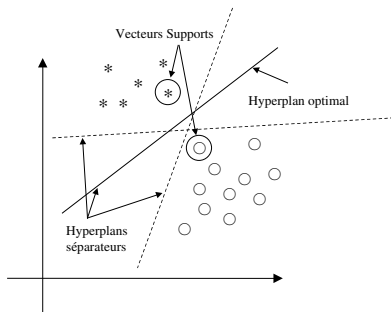
with

$$g_{\mathbf{w},b}(\mathbf{x}_i) > 0 \text{ if } \mathbf{x}_i \in \omega_1, \quad g_{\mathbf{w},b}(\mathbf{x}_i) < 0 \text{ if } \mathbf{x}_i \in \omega_2$$

### Classification rule

$$f(\mathbf{x}) = \text{sign}[g_{\mathbf{w},b}(\mathbf{x})] \tag{6}$$

# Illustration



## Problem formulation

Margin of  $\mathbf{x}_i$  with label  $y_i$  (algebraic distance to the hyperplane)

$$\gamma_i(\tilde{\mathbf{w}}) = \frac{y_i (\mathbf{w}^T \mathbf{x}_i - b)}{\|\mathbf{w}\|}$$

with  $\tilde{\mathbf{w}} = (\mathbf{w}, b)$  ( $\mathbf{x}_i$  is correctly classified by (6) if  $\gamma_i(\tilde{\mathbf{w}}) > 0$ )

Margin of the learning set

$$\gamma_{\mathcal{B}}(\tilde{\mathbf{w}}) = \min_{i \in \{1, \dots, n\}} \frac{y_i (\mathbf{w}^T \mathbf{x}_i - b)}{\|\mathbf{w}\|}$$

Since  $\gamma_{\mathcal{B}}(a\tilde{\mathbf{w}}) = \gamma_{\mathcal{B}}(\tilde{\mathbf{w}})$ ,  $\forall a > 0$ ,  $\tilde{\mathbf{w}}$  is not unique!

**Constraints for the hyperplane:** one forces the training samples that are the closest to the hyperplane to satisfy

$$y_i (\mathbf{w}^T \mathbf{x}_i - b) = 1 \Rightarrow \min_{i \in \{1, \dots, n\}} y_i (\mathbf{w}^T \mathbf{x}_i - b) = 1$$

The vectors  $\mathbf{x}_i$  satisfying  $y_i (\mathbf{w}^T \mathbf{x}_i - b) = 1$  are called **support vectors**.



## Problem formulation

### Canonical hyperplane

$$y_i (\mathbf{w}^T \mathbf{x}_i - b) \geq 1, \quad \forall i = 1, \dots, n$$

### Classifier margin for a canonical hyperplane

$$\gamma_{\mathcal{B}}(\tilde{\mathbf{w}}) = \frac{1}{\|\mathbf{w}\|}$$

### New formulation

**Minimize**  $\frac{1}{2} \|\mathbf{w}\|^2$  **with the constraints**  $y_i (\mathbf{w}^T \mathbf{x}_i - b) \geq 1, \forall i$

Simple problem since the cost function to optimize is quadratic and the constraints are linear!

## Optimization

### Lagrangian

$$L(\tilde{\mathbf{w}}, \alpha) = \frac{1}{2} \mathbf{w}^T \mathbf{w} - \sum_{i=1}^n \alpha_i \left[ y_i (\mathbf{w}^T \mathbf{x}_i - b) - 1 \right]$$

Set to zero the partial derivatives of  $L$  with respect to  $b$  and  $w$

$$\sum_{i=1}^n \alpha_i y_i = 0 \text{ and } \mathbf{w} = \sum_{i=1}^n \alpha_i y_i \mathbf{x}_i$$

### Kühn and Tucker multipliers

For a convex optimization problem (convex function  $f(x)$  to optimize and convex constraints  $G_i(\mathbf{x}) \leq 0$ ), an optimality condition is the existence of parameters  $\alpha_i \geq 0$  such that the Lagrangian derivative is zero, i.e.,

$$L'(\mathbf{x}) = f'(\mathbf{x}) + \sum_{i=1}^n \alpha_i G'_i(\mathbf{x}) = 0$$

with  $\alpha_i = 0$  if  $G_i(\mathbf{x}) < 0$  (i.e.,  $\alpha_i G_i(\mathbf{x}) = 0$ ).

## Dual problem

Solve  $L'(\mathbf{x}) = 0$

$$\mathbf{w} = \sum_{\text{Support vectors}} \alpha_i y_i \mathbf{x}_i = \mathbf{x}^T \mathbf{Y} \alpha \quad (7)$$

with  $\alpha = (\alpha_1, \dots, \alpha_n)^T$ ,  $\mathbf{x} = (x_1, \dots, x_n)^T$ ,  $\mathbf{Y} = \text{diag}(y_1, \dots, y_n)$  and

$$\begin{cases} \alpha_i = 0 & \text{if the constraint is a strict inequality} \\ \alpha_i > 0 & \text{if the constraint is an equality} \end{cases}$$

After replacing the expression of  $w$  in the Lagrangian, we obtain

$$U(\alpha) = -\frac{1}{2} \alpha^T \mathbf{Y} (\mathbf{x} \mathbf{x}^T) \mathbf{Y} \alpha + \sum_{i=1}^n \alpha_i$$

that has to be maximized in the domain defined by  $\alpha_i \geq 0, \forall i$  and  $\sum_{i=1}^n \alpha_i y_i = 0$ .

## Remarks

### Simple optimization problem

Quadratic (hence convex) function to optimize and linear constraints

Norm of the solution  $w_0$

$$\|w_0\|^2 = \sum_{\text{support vectors}} \alpha_i^0 (1 + y_i b) = \sum_{\text{support vectors}} \alpha_i^0$$

because of the constraint  $\sum_{i=1}^n \alpha_i y_i = 0$ .

Classifier margin

$$\gamma = \frac{1}{\|w_0\|} = (\sum \alpha_i^0)^{-1/2}$$

Classification rule for a vector  $x$

$$f(x) = \text{sign} \left( \sum_{z_i \text{ support vectors}} \alpha_i^0 y_i x_i^T x - b_0 \right), \quad b_0 = \frac{1}{2} (w_0^T z^+ + w_0^T z^-),$$

where  $z^+$  (resp.  $z^-$ ) is a support vector belonging to the 1st (resp. 2nd) class.

## Extensions

### The “soft-margin” SVM classifier

$$\text{minimize } \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^n \xi_i$$

$$\text{with the constraints } y_i (\mathbf{w}^T \mathbf{x}_i - b) \geq 1 - \xi_i, \forall i$$

### The $\nu$ -SVM classifier

$$\text{Minimize } \frac{1}{2} \|\mathbf{w}\|^2 + \frac{1}{n} \sum_{i=1}^n \xi_i - \nu \gamma$$

$$\text{with the constraints } y_i (\mathbf{w}^T \mathbf{x}_i - b) \geq \gamma - \xi_i, \forall i$$

## Non-linear preprocessing

Search a non-linear transformation  $\phi : \mathbb{R}^p \rightarrow \mathbb{R}^q$  ensuring linear separability.

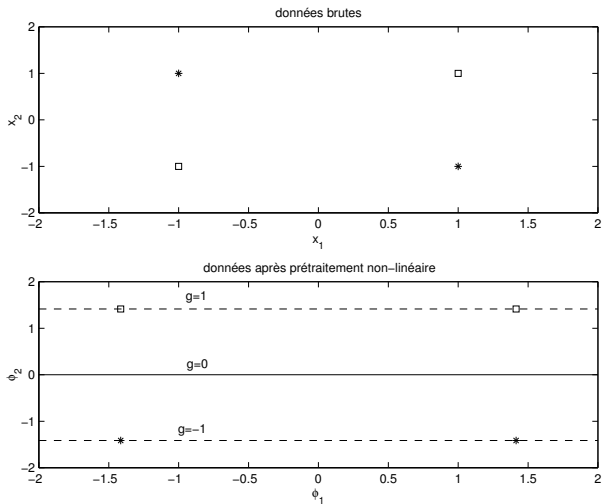
### Classical example

The class classes  $\chi_1 = \{(1, 1), (-1, -1)\}$  and  $\chi_2 = \{(1, -1), (-1, 1)\}$  are not linearly separable. Consider the application  $\phi$  defined by

$$\phi : \mathbb{R}^2 \rightarrow \mathbb{R}^6 \\ (x_1, x_2)^T \mapsto (\sqrt{2}x_1, \sqrt{2}x_1x_2, 1, \sqrt{2}x_2, x_1^2, x_2^2)^T$$

The data are separable in the plane  $(\phi_1, \phi_2)$

## Example of separability



## Non-linear SVM classifier

### Decision rule after preprocessing

$$\begin{aligned}
 f(\tilde{\mathbf{x}}) &= \text{sign}(g_{\mathbf{w},b}(\tilde{\mathbf{x}})) = \text{sign}(\mathbf{w}^T \phi(\mathbf{x}) - b) \\
 &= \text{sign}\left(\sum_{\text{support vectors}} \alpha_i^0 y_i \phi(\mathbf{x}_i)^T \phi(\mathbf{x}) - b_0\right)
 \end{aligned}$$

### Cost function to optimize

$$U(\boldsymbol{\alpha}) = -\frac{1}{2} \boldsymbol{\alpha}^T \mathbf{Y} \mathbf{G} \mathbf{Y}^T \boldsymbol{\alpha} + \sum_{i=1}^n \alpha_i$$

where  $\mathbf{G}$  is the Gram matrix defined by  $G_{ij} = \phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j)$ . The cost function  $U(\boldsymbol{\alpha})$  has to be maximized under the constraints

$$0 \leq \alpha_i \leq \frac{1}{n}, \forall i \quad \text{and} \quad \sum_{i=1}^n \alpha_i \geq \nu$$

**Conclusions:** the cost function and the decision rule only depend on the inner products  $\phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j)$  and  $\phi(\mathbf{x}_i)^T \phi(\mathbf{x})$ .



## Kernels

A kernel  $k(\mathbf{x}, \mathbf{y}) = \phi(\mathbf{x})^T \phi(\mathbf{y})$  allows inner products  $\phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j)$  and  $\phi(\mathbf{x}_i)^T \phi(\mathbf{x})$  to be computed with a reduced computational cost.

**Example:**  $\mathbf{x} = (x_1, x_2)^T$  and  $\phi(\mathbf{x}) = (x_1^2, x_2^2, \sqrt{2}x_1x_2)$

$$\begin{aligned}\phi(\mathbf{x})^T \phi(\mathbf{y}) &= \begin{pmatrix} x_1^2 \\ x_2^2 \\ \sqrt{2}x_1x_2 \end{pmatrix} \cdot \begin{pmatrix} y_1^2 \\ y_2^2 \\ \sqrt{2}y_1y_2 \end{pmatrix} \\ &= x_1^2y_1^2 + x_2^2y_2^2 + 2x_1x_2y_1y_2 \\ &= (\mathbf{x}^T \mathbf{y})^2\end{aligned}$$

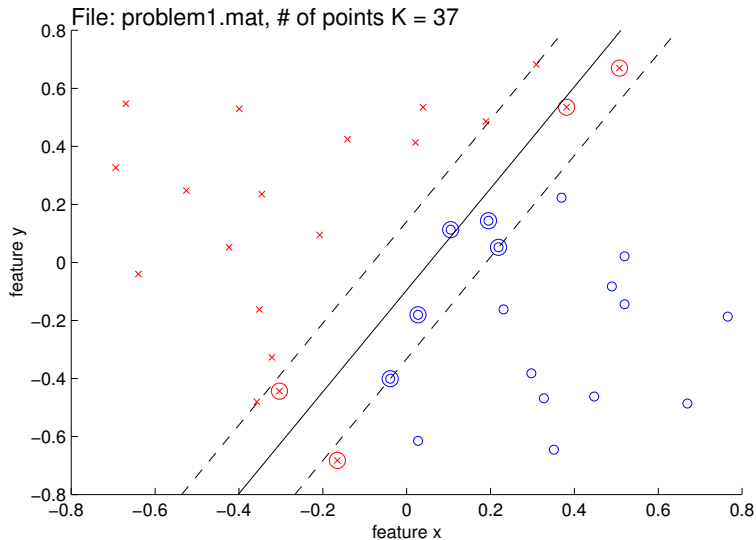
**Mercer kernels** can be expressed as inner products

$$k(\mathbf{x}, \mathbf{y}) = \phi(\mathbf{x})^T \phi(\mathbf{y})$$

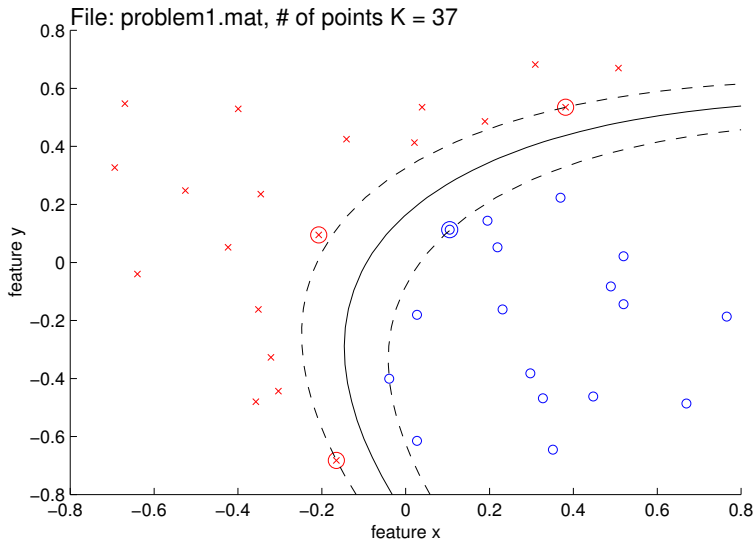
## Classical kernels

Kernel	Expression
Polynomial (of degree $p$ )	$k(\mathbf{x}, \mathbf{y}) = (\langle \mathbf{x}, \mathbf{y} \rangle)^q$ $q \in \mathbb{N}^+$
Full polynomial	$k(\mathbf{x}, \mathbf{y}) = (\langle \mathbf{x}, \mathbf{y} \rangle + c)^q$ $c \in \mathbb{R}^+, q \in \mathbb{N}^+$
RBF	$k(\mathbf{x}, \mathbf{y}) = \exp\left(-\frac{\ \mathbf{x} - \mathbf{y}\ ^2}{2\sigma^2}\right)$ $\sigma \in \mathbb{R}^+$
Mahalanobis	$k(\mathbf{x}, \mathbf{y}) = \exp\left[-(\mathbf{x} - \mathbf{y})^T \Sigma (\mathbf{x} - \mathbf{y})\right]$ $\Sigma = \text{diag}\left(\frac{1}{\sigma_1^2}, \dots, \frac{1}{\sigma_p^2}\right), \sigma_i \in \mathbb{R}^+$

## Example of linear SVMs



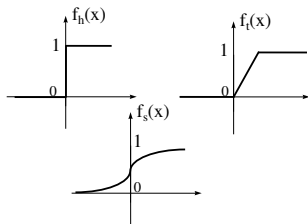
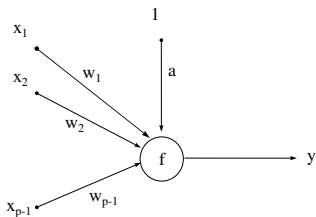
## Example of non-linear SVMs



## Summary

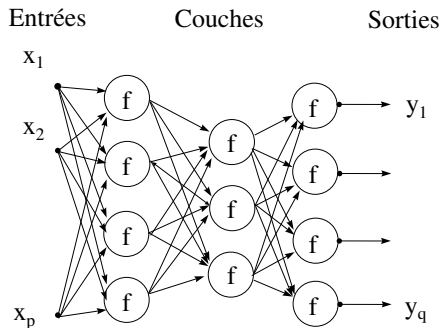
- ▶ Chapter 1 : Introduction
- ▶ Chapter 2 : Preprocessing
- ▶ Chapter 3 : Statistical Methods
- ▶ Chapter 4 : Support Vector Machines (SVMs) and Neural Networks
  - ▶ Linear discriminant functions
    - ▶ The perceptron algorithm
    - ▶ The optimal Wiener-Hopf filter
    - ▶ The LMS algorithm
  - ▶ Support vector machines (SVMs)
    - ▶ Introduction
    - ▶ Optimal separating hyperplane
    - ▶ Optimization problem
    - ▶ The “soft-margin SVM” classifier
    - ▶ Non-separable case: the “ $\nu$ -SVM” classifier
    - ▶ Non-linear preprocessing - kernels
  - ▶ Neural networks
    - ▶ The multi-layer perceptron
    - ▶ The backpropagation algorithm
- ▶ Chapter 5 : Decision Trees

## Neural networks



- ▶ **inputs:**  $x_1, \dots, x_{p-1}$
- ▶ **weights:**  $w_1, \dots, w_{p-1}$
- ▶ **bias:**  $a$  (or offset)
- ▶ **output :**  $y = f\left(\sum_{i=1}^{p-1} w_i x_i - a\right)$
- ▶ **Non-linearity**
  - ▶ Heaviside:  $f_h(x)$
  - ▶ Logical threshold:  $f_t(x)$
  - ▶ Sigmoid:  $f_s(x)$
- ▶ **Equivalent definition:**  $x_p = -1, w_p = a$  et  $y = f\left(\sum_{i=1}^p w_i x_i\right)$

## Multi-layer perceptron



► **Definition of the network structure**

- 👉 Number of outputs?
- 👉 Number of layers?
- 👉 Number of nodes per layer?

► **Learning rules?**

## One-layer perceptron with $f_L(x) = x$

► **Learning rule**

- 👉 Minimization of  $J(w) = \frac{1}{2}E[e^2(n)] \Rightarrow$  Wiener-Hopf filter
- 👉 Steepest descent algorithm with instantaneous error

$$J_{\text{inst}}(w) = \frac{1}{2} \left[ d(n) - f_L \left( \sum_{i=1}^p w_i x_i(n) \right) \right]^2$$

► **LMS algorithm**

$$\hat{w}_k(n+1) = \hat{w}_k(n) + \mu e(n) x_k(n)$$

- **Properties:** convergence of the LMS algorithm to the Wiener-Hopf solution



## One-layer perceptron with sigmoidal function

### ► Learning rule

- 👉 Minimization of  $J(w) = \frac{1}{2}E[e^2(n)] \Rightarrow$  system of non-linear equations (no Wiener-Hopf solution)
- 👉 Descent algorithm with instantaneous error

$$J_{\text{inst}}(w) = \frac{1}{2} \left[ d(n) - f_s \left( \sum_{i=1}^p w_i x_i(n) \right) \right]^2$$

### ► The LMS algorithm

$$\hat{w}_k(n+1) = \hat{w}_k(n) + \mu e(n) x_k(n) y(n) [1 - y(n)]$$

### ► Properties

- 👉 Convergence to a local minimum of  $J(w)$
- 👉 Same performance obtained with  $f_h$  instead of the sigmoid
- 👉 Generalization (multi-layer network): the back-propagation algorithm

## Backpropagation

- ▶ Weight Initialization:  $w_{ij}^{(l)}(0)$  with footnotesize values
- ▶ at iteration  $n$ 
  - ☞ Compute the desired output

$$d(n) = 1 \text{ if } x(n) \in \omega_1, \quad d(n) = 0 \text{ if } x(n) \in \omega_2$$

- ☞ Compute the network output

$$y_j^{(l)}(n) = f_s \left( \sum_{i=1}^{N_{l-1}} w_{ij}^{(l)}(n) y_i^{(l-1)}(n) \right)$$

where  $y_j^{(l)}$  is the  $j$ th output of the  $l$ th layer (with  $y^{(L)}(n) = y(n)$  et  $y_j^{(0)}(n) = x_j(n)$ ),  $N_l$  is the number of nodes of the  $l$ th layer and  $L$  is the number of layers.

- ☞ Weight update

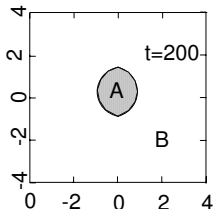
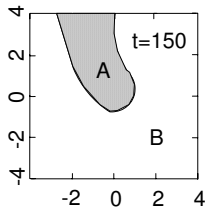
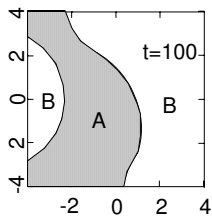
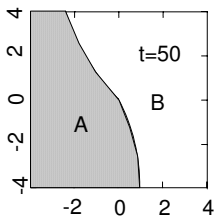
$$\hat{w}_{ij}^{(l)}(n+1) = \hat{w}_{ij}^{(l)}(n) + \mu \delta_j^{(l)}(n) y_i^{(l-1)}(n)$$

with  $\delta^{(L)}(n) = e(n)y(n)[1 - y(n)]$  and

$$\delta_j^{(l)}(n) = y_j^{(l)}(n) [1 - y_j^{(l)}(n)] \sum_k \delta_k^{(l+1)}(n) w_{jk}^{(l+1)}(n)$$

where the summation covers all the errors of the  $(l+1)$ th layer and  $e(n) = d(n) - y(n)$  is the error at the output of the network.

Class 1:  $D(0, 1)$       Class 2:  $D(0, 5) \setminus D(0, 1)$



## Remarks

- ▶ **Epoch**
- ▶ How should we choose the **parameter  $\mu$** ?
  - ▶ **small  $\mu$** : slow convergence, small residual error
  - ▶ **large  $\mu$** : fast convergence, large residual error
- ▶ Classification rule with **constant momentum** noté  $\beta \in [0, 1[$

$$\hat{w}_{ij}(n+1) = \hat{w}_{ij}(n) + \Delta \hat{w}_{ij}(n)$$

$$\Delta \hat{w}_{ij}(n) = \beta \Delta \hat{w}_{ij}(n-1) + \mu \delta_j(n) y_i(n)$$

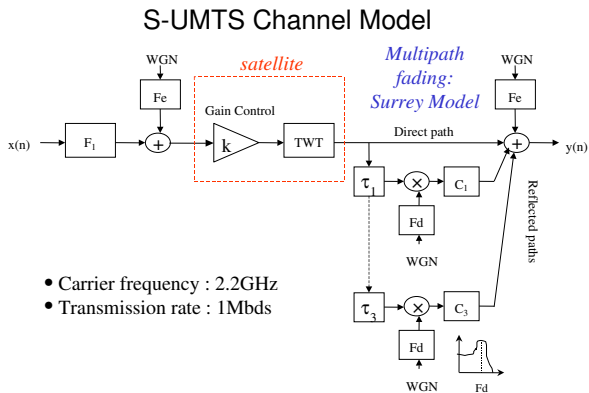
For  $\beta = 0$ : classical backpropagation rule

For  $\beta > 0$ :  $\Delta \hat{w}_{ij}(n) = \mu \sum_{k=0}^n \beta^{n-k} \delta_j(n) y_i(n)$

Convergence is ensured if  $0 \leq |\beta| \leq 1$ .

- ▶ The convergence is **accelerated** when consecutive errors have **the same sign** whereas it is **reduced** when two consecutive errors have **different signs**.
- ▶ **Stopping rules** are problem dependent. For example, in some equalization problems, we use a sequence of known samples (pilot symbols) and the learning algorithm is stopped when  $\|\text{input} - \text{output}\| < \varepsilon$
- ▶ **Universal approximation theorem** of Cybenko.

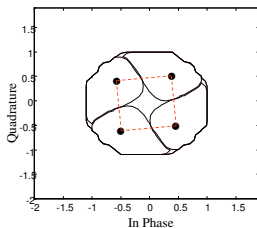
## Application to the equalization of non-linear channels



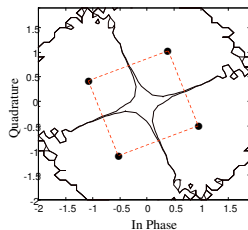
## Application to the equalization of non-linear channels

### Consequences of the non-linearity

#### Distortion of a 4-QAM signal



backoff 0dB  
no downlink noise

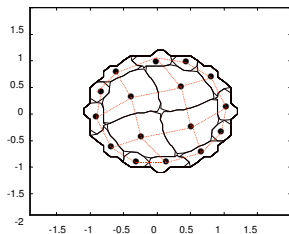


backoff 0dB  
15dB downlink SNR

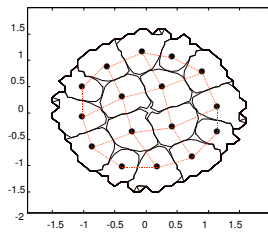
## Application to the equalization of non-linear channels

### Consequences of the non-linearity

#### Distortion of a 16-QAM signal



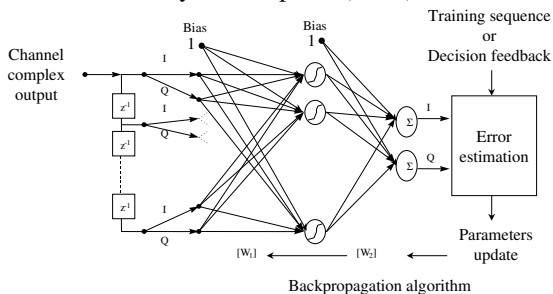
backoff -1dB  
no downlink noise



backoff -1dB  
15dB downlink noise

## Application to the equalization of non-linear channels

### Neural Network Equalizers : Multilayer Perceptron (MLP)

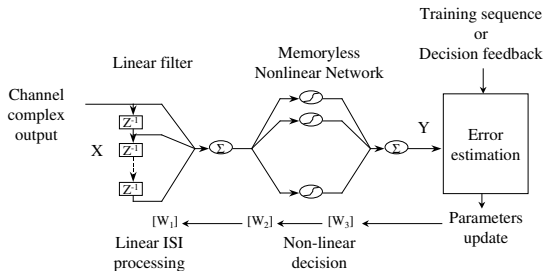


- **Advantages** : universal approximation theorem, parallelism.
- **Drawback** : slow convergence.



## Application to the equalization of non-linear channels

### Neural Network Equalizers : Linear Filter - NonLinear Network (LF-NLN)

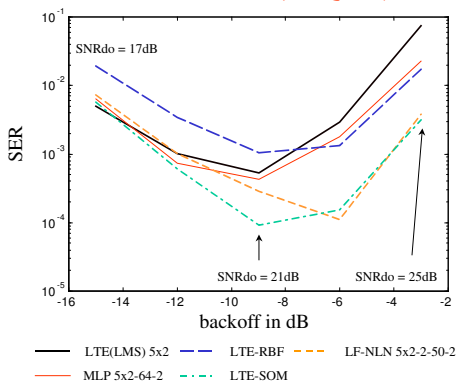


- **Advantage** : More simple than the MLP.
- **Drawback** : Non optimal.

## Application to the equalization of non-linear channels

## NN equalization of 16-QAM signals, S-UMTS

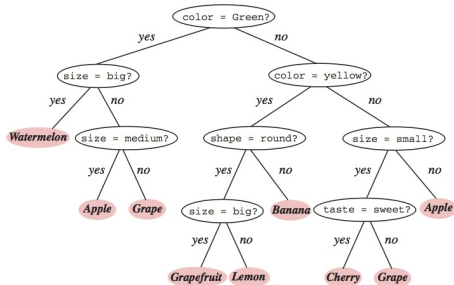
### Backoff influence (16-QAM)



## Summary

- ▶ Chapter 1 : Introduction
- ▶ Chapter 2 : Preprocessing
- ▶ Chapter 3 : Statistical Methods
- ▶ Chapter 4 : Support Vector Machines (SVMs) and Neural Networks
- ▶ Chapter 5 : Decision Trees

## Example of decision tree



**FIGURE 8.2.** A tree with arbitrary branching factor at different nodes can always be represented by a functionally equivalent binary tree—that is, one having branching factor  $B = 2$  throughout, as shown here. By convention the “yes” branch is on the left, the “no” branch on the right. This binary tree contains the same information and implements the same classification as that in Fig. 8.1. From: Richard O. Duda, Peter E. Hart, and David G. Stork, *Pattern Classification*. Copyright © 2001 by John Wiley & Sons, Inc.

- ▶ Construction of the tree
- ▶ Classification rule

## Splitting Rules

### Inhomogeneity or impurity of the data

- ▶ Entropy (Algorithm C4.5)

$$i_n = - \sum_j \frac{n_j}{n} \log_2 \left( \frac{n_j}{n} \right)$$

- ▶ Gini Index (CART)

$$i_n = \sum_j \frac{n_j}{n} \left( 1 - \frac{n_j}{n} \right)$$

### Drop of Impurity

$$\Delta i_n = i_n - P_L i_L - P_R i_R$$

where  $P_L = \frac{n_L}{n}$ ,  $P_R = \frac{n_R}{n}$  are the proportions of the sets  $D_L, D_R$ .

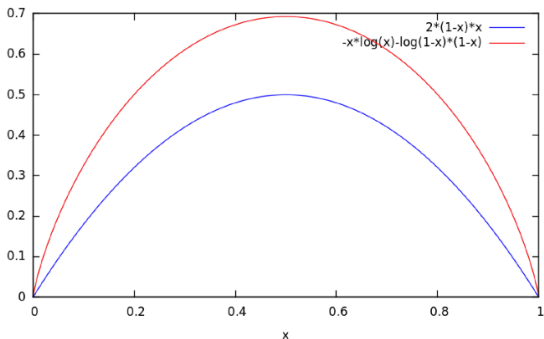
**Choose the split associated with the maximum drop of impurity!**

## Gini Index or Entropy?

Example of 2 classes ( $x = n_1/n$ )

$$\text{Gini}(S) = 2x(1 - x)$$

$$\text{Ent}(S) = -x \log x - (1 - x) \log(1 - x)$$



## Implementation by default

### Algorithm CART (`classregtree.m`)

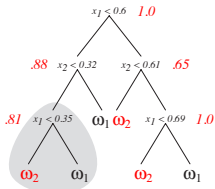
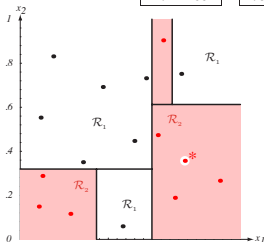
- ▶ All variables are considered for each split
- ▶ All splits are considered
- ▶ Stopping rule: pure node or number of elements less than  $n_{\min}$  (specified by the user)
- ▶ Splitting rule: Gini index

## CART algorithm (example #1 for vectors in $\mathbb{R}^2$ )

### Example 1: A simple tree classifier

Consider the following  $n = 16$  points in two dimensions for training a binary CART tree ( $B = 2$ ) using the entropy impurity (Eq. 1).

$\omega_1$ (black)		$\omega_2$ (red)	
$x_1$	$x_2$	$x_1$	$x_2$
.15	.83	.10	.29
.09	.55	.08	.15
.29	.35	.23	.16
.38	.70	.70	.19
.52	.48	.62	.47
.57	.73	.91	.27
.73	.75	.65	.90
.47	.06	.75	.36* (.32†)



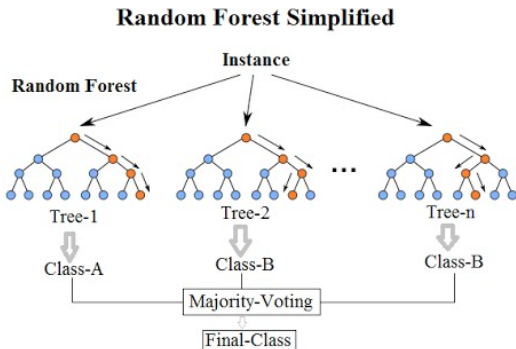


## CART algorithm (example #2 for qualitative data)

	Weight	Size	Age	Result
$x_1$	Light	Small	Young	Pass
$x_2$	Light	Small	Young	Pass
$x_3$	Light	Tall	Young	Pass
$x_4$	Light	Tall	Old	Fail
$x_5$	Light	Tall	Old	Pass
$x_6$	Light	Tall	Old	Fail
$x_7$	Heavy	Small	Old	Fail
$x_8$	Heavy	Small	Young	Fail
$x_9$	Light	Small	Old	Pass
$x_{10}$	Heavy	Tall	Old	Pass

- ▶ First branch (Gini Index)
  - ▶ **Weight**  $\Rightarrow 2(1/7 + 1/15) \sim 0.42$
  - ▶ **Size**  $\Rightarrow 12/25 \sim 0.48$
  - ▶ **Age**  $\Rightarrow 2(3/40 + 3/20) \sim 0.45$

## Random Forests

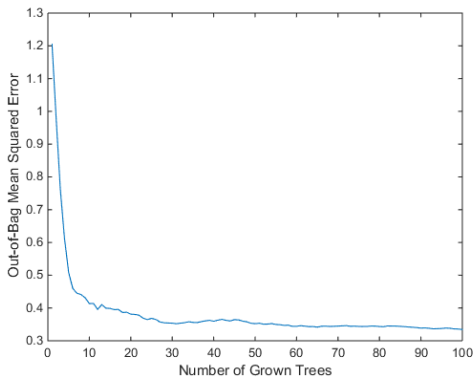


### TreeBagger.m (options by default)

- ▶ **Resampling all** the data in the training set by bootstrap (and not a subset)
- ▶ **Number** of variables to select at random for each decision split:  $\sqrt{n_{\text{var}}}$

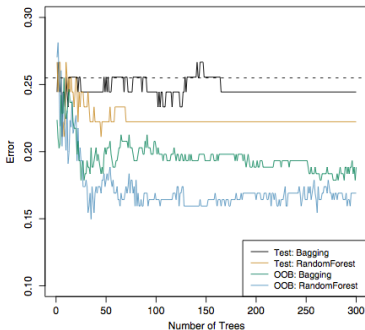
## Out-of-bag error

### Matlab example



## Comparison between CART and Random Forests

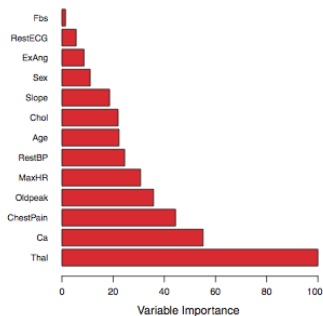
Book by G. James, D. Witten, T. Hastie and R. Tibshirani



**FIGURE 8.8.** Bagging and random forest results for the **Heart** data. The test error (black and orange) is shown as a function of  $B$ , the number of bootstrapped training sets used. Random forests were applied with  $m = \sqrt{p}$ . The dashed line indicates the test error resulting from a single classification tree. The green and blue traces show the OOB error, which in this case is considerably lower.

## Importance of the different variables

### Mean decrease in Gini index



**FIGURE 8.9.** A variable importance plot for the `Heart` data. Variable importance is computed using the mean decrease in Gini index, and expressed relative to the maximum.

## References

### CART and Random Forests

- ▶ L. Breiman, J. H. Friedman, R. A. Olshen and C. J. Stone, *Classification and Regression Trees*, Chapman & Hall, New-York, 1993.
- ▶ L. Breiman, *Random Forests*, Machine Learning, vol. 45, no. 1, pp. 5-32, 2001.
- ▶ R. O. Duda, P. E. Hart and D. G. Stork, *Pattern Classification*, 2nd edition, Wiley, 2000.
- ▶ G. James, D. Witten, T. Hastie and R. Tibshirani, *An Introduction to Statistical Learning with Applications in R*, Springer, New-York, 2013.

Thanks for your attention