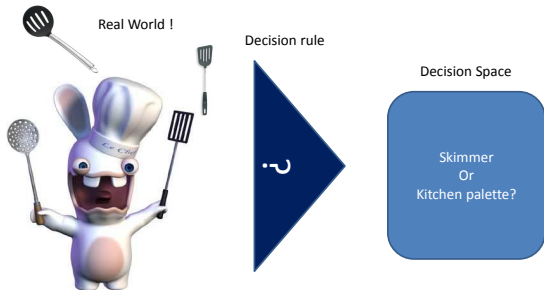


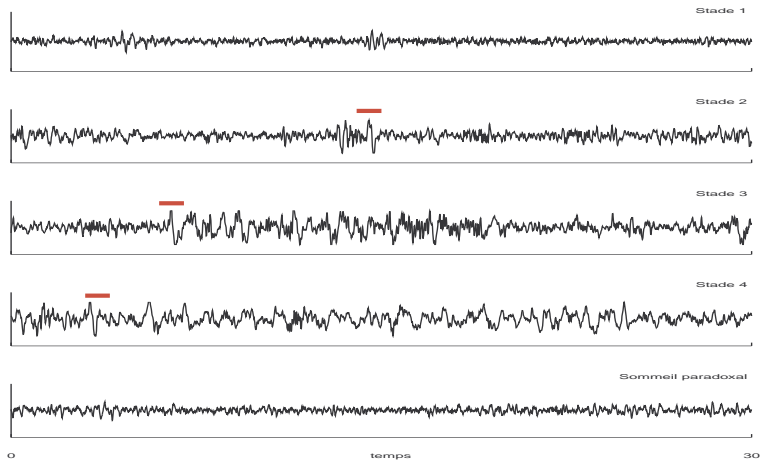
Pattern recognition and applications to monitoring

Pierre Beuseroy - Edith Grall - Paul Honeine - Régis Lengellé

– introduction –



- Real world
- Data gathering
- Decision rule
- Decision space



Example : Complex K detection in EEG sleep signal.

The problem can be written as follow :

$$\begin{cases} \omega_0 : \mathbf{x} = \mathbf{b} & \text{hypothesis "noise alone"} \\ \omega_1 : \mathbf{x} = \mathbf{b} + \mathbf{s} & \text{hypothesis "signal plus noise"} \end{cases}$$

The aim is to build a detector d , or a decision rule d which minimizes a criteria - the error probability for example

$$P_e(d) = p(d(\mathbf{X}) \neq Y),$$

where \mathbf{X} is an observation and Y its associate hypothesis.

The strategy to design a solution to this problem depends on the nature of the available information on the problem

Rule based approach

If X is grey and $weight(X) > 1000$ then

X is an elephant

Else

X is a mouse ...

⇒ Need an expert

⇒ Need to translate expert thought to rules - complex, long, most of the time not reliable.

Hypothesis testing approach

Express the problem as an hypothesis testing problem :

$$\begin{cases} H_0 : X \in \omega_0; & X \sim p(X | \omega_0) \\ H_1 : X \in \omega_1; & X \sim p(X | \omega_1) \end{cases}$$

Deduce the test.

Data driven approach

Use a set of data \mathcal{A}

2 classes

$$\begin{cases} H_0 : X \in \omega_0; & X \sim p(X|\omega_0) \\ H_1 : X \in \omega_1; & X \sim p(X|\omega_1) \end{cases}$$

1 class

$$\begin{cases} H_0 : X \in \omega_0 \\ H_1 : X \notin \omega_0 \end{cases}$$

Multiclass

$$\begin{cases} H_0 : X \in \omega_0 \\ H_1 : X \in \omega_1 \\ H_2 : X \in \omega_2 \\ \dots \\ H_K : X \in \omega_K \end{cases}$$

$$\begin{cases} H_0 : X \in \omega_0; & X \sim p(X | \omega_0) \\ H_1 : X \in \omega_1; & X \sim p(X | \omega_1) \end{cases}$$

We want to determine $d(\mathbf{x}) = \omega_0$ if $q(\mathbf{x}) < s$ so that $\min_{q,s} P_e$

with :

$$P_e = P(D_0 | \omega_1) P(\omega_1) + P(D_1 | \omega_0) P(\omega_0)$$

$$P(D_0 | \omega_1) = \int_{\mathcal{X}} \mathbb{1}_{(q(x) < s)} p(x | \omega_1) dx$$

$$P(D_1 | \omega_0) = \int_{\mathcal{X}} \mathbb{1}_{(q(x) \geq s)} p(x | \omega_0) dx$$

For each x we wish to minimize :

$$r(x) = \mathbb{1}_{(q(x) < s)} p(x | \omega_1) P(\omega_1) + \mathbb{1}_{(q(x) \geq s)} p(x | \omega_0) P(\omega_0)$$

$$P_e = r(\mathcal{D}) = \int_{\mathcal{X}} r(x) dx$$

We choose the class ω_i such that :

$$p(x|\omega_i) P(\omega_i) \stackrel{D_i}{>} p(x|\omega_j) P(\omega_j)$$

Thus the decision rule is :

$$p(x|\omega_1) P(\omega_1) \stackrel{D_0}{<} \stackrel{D_1}{>} p(x|\omega_0) P(\omega_0)$$

or

$$\frac{p(x|\omega_1)}{p(x|\omega_0)} \stackrel{D_0}{<} \stackrel{D_1}{>} \frac{P(\omega_0)}{P(\omega_1)}$$

where

$q(x) = \frac{p(x|\omega_1)}{p(x|\omega_0)}$ is the likelihood ratio

and

$$s = \frac{P(\omega_0)}{P(\omega_1)}$$

Multiclass case

$$d(\mathbf{x}) = \arg \min_i \left(\sum_{j \neq i} P(D_i | \omega_j) P(\omega_j) \right)$$

Criteria is a cost function

We want to determine $d(\mathbf{x}) = \omega_0$ if $q(\mathbf{x}) < s$ so that $\min_{q,s} r(\mathcal{D})$

where

$$r(\mathcal{D}) = \sum_{i,j=1}^K c_{ij} P(D_i | \omega_j) P(\omega_j)$$

$$d(x) = \arg \min_i \left(\sum_j c_{ij} P(D_i | \omega_j) P(\omega_j) \right)$$

In the 2 class case :

$$\frac{p(x | \omega_1)}{p(x | \omega_0)} \underset{D_1}{\overset{D_0}{<}} \frac{(c_{10} - c_{00}) P(\omega_0)}{(c_{01} - c_{11}) P(\omega_1)}$$

We want to define α , the first order error $P(D_1|\omega_0)$. So the problem is :

$$\min_{q,s} P(D_0|\omega_1)$$

with :

$$P(D_1|\omega_0) = \alpha$$

Let assume that solving

$$\min_{q,s} \max_{\mu} P(D_0|\omega_1) + \mu(P(D_1|\omega_0) - \alpha)$$

solves the initial problem.

And :

$$\begin{aligned} & P(D_0|\omega_1) + \mu(P(D_1|\omega_0) - \alpha) \\ &= \int_{\mathcal{X}} \mathbb{1}_{(q(x) < s)} p(x|\omega_1) + \mu(\mathbb{1}_{(q(x) \geq s)} p(x|\omega_0)) dx - \mu\alpha \end{aligned}$$

Thus, for a given μ the decision rule is :

$$p(x|\omega_1) \underset{D_1}{\overset{D_0}{>}} \mu p(x|\omega_0)$$

or

$$\frac{p(x|\omega_1)}{p(x|\omega_0)} \underset{D_1}{\overset{D_0}{>}} \mu$$

which is a likelihood ratio test (as in Bayes case).

Free structure detection

Based on simple hypothesis, applying decision criterion such as Bayes' cost leads to Bayes' rule :

$$d^*(\mathbf{x}) = \begin{cases} \omega_1 & \text{si } p(\mathbf{x}|\omega_1)/p(\mathbf{x}|\omega_0) \geq \lambda_0 \\ \omega_0 & \text{sinon,} \end{cases}$$

$p(\mathbf{x}|\omega_0)$ and $p(\mathbf{x}|\omega_1)$ have to be known. The threshold λ_0 is the only parameter of the rule. It depends on the criterion one wishes to optimize.

Thus the decision rule or the detector is not subject to any structural constraints but depends only on the choice of a criteria.

Example 1

Gaussian case - same variance

$$\begin{cases} P(X|\omega_0) \sim N(M_0, \Sigma) \\ P(X|\omega_1) \sim N(M_1, \Sigma) \end{cases}$$

and $P(\omega_0) = P(\omega_1) = \frac{1}{2}$

Thus :

$$P(X|\omega_i) = (2\pi)^{-d/2} |\Sigma|^{-1/2} e^{-\frac{1}{2}(X-M_i)^T \Sigma^{-1} (X-M_i)}$$

So :

$$\log(q(X)) = \frac{1}{2} (X - M_0)^T \Sigma^{-1} (X - M_0) - \frac{1}{2} (X - M_1)^T \Sigma^{-1} (X - M_1)$$

And the decision rule is :

$$(M_1 - M_0)^T \Sigma^{-1} X + \frac{1}{2} M_0^T \Sigma^{-1} M_0 - \frac{1}{2} M_1^T \Sigma^{-1} M_1 \underset{D_1}{\overset{D_0}{<}} 0$$

Which is equivalent to

$$V^T X + U \underset{D_1}{\overset{D_0}{<}} 0 \text{ a linear function !}$$

$$\begin{cases} P(X|\omega_0) \sim N(M_0, \Sigma_0) \\ P(X|\omega_1) \sim N(M_1, \Sigma_1) \end{cases}$$

and let $P(\omega_0) = P(\omega_1) = \frac{1}{2}$

Thus :

$$\log(q(X)) = \frac{1}{2} \log\left(\frac{|\Sigma_0|}{|\Sigma_1|}\right) + \frac{1}{2} (X - M_0)^T \Sigma_0^{-1} (X - M_0) - \frac{1}{2} (X - M_1)^T \Sigma_1^{-1} (X - M_1)$$

So

$$\begin{aligned} & \frac{1}{2} X^T (\Sigma_0^{-1} - \Sigma_1^{-1}) X - M_0^T \Sigma_0^{-1} X + M_1^T \Sigma_1^{-1} X + \frac{1}{2} M_0^T \Sigma_0^{-1} M_0 \\ & - \frac{1}{2} M_1^T \Sigma_1^{-1} M_1 + \frac{1}{2} \log\left(\frac{|\Sigma_0|}{|\Sigma_1|}\right) \begin{matrix} D_0 \\ < \\ D_1 \end{matrix} 0 \end{aligned}$$

Which is equivalent to :

$$X^T T X + V^T X + U \begin{matrix} D_0 \\ < \\ D_1 \end{matrix} 0 \text{ a quadratic function !}$$

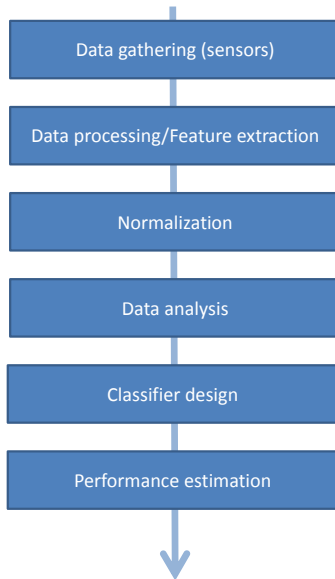
To use these approaches $p(\mathbf{x}|\omega_0)$ and $p(\mathbf{x}|\omega_1)$ have to be known...

If not the case?

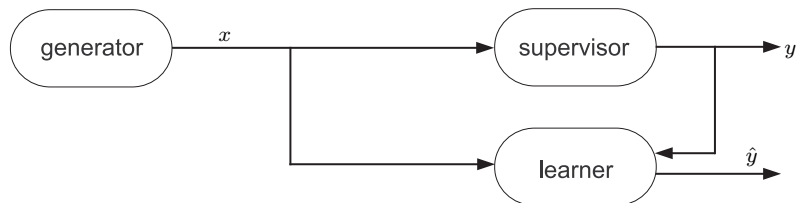
can assume that

$$p(\mathbf{x}|\omega_0) \text{ and } p(\mathbf{x}|\omega_1) \in \mathcal{F}_\theta$$

then select or estimate θ based on data
and plugging the estimator in the decision rule.



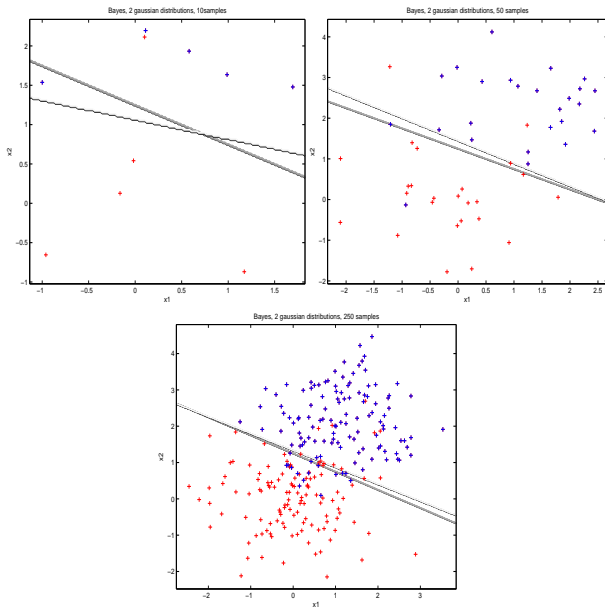
The learning model is composed of 3 elements :



- ❶ Generator : $\mathbf{X} \in \mathcal{X} \subset \mathbb{R}^l$, random variables i.i.d.
- ❷ Supervisor : $Y \in \mathcal{Y} \subset \mathbb{R}$, random variables
- ❸ Learner : represented by $d(\mathbf{x}; \theta) \in \mathcal{D}$

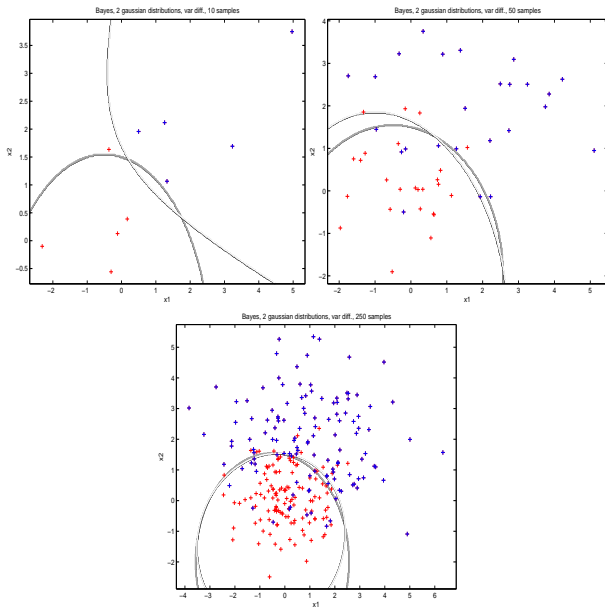
Parametric estimation : Example 1

Experimental results



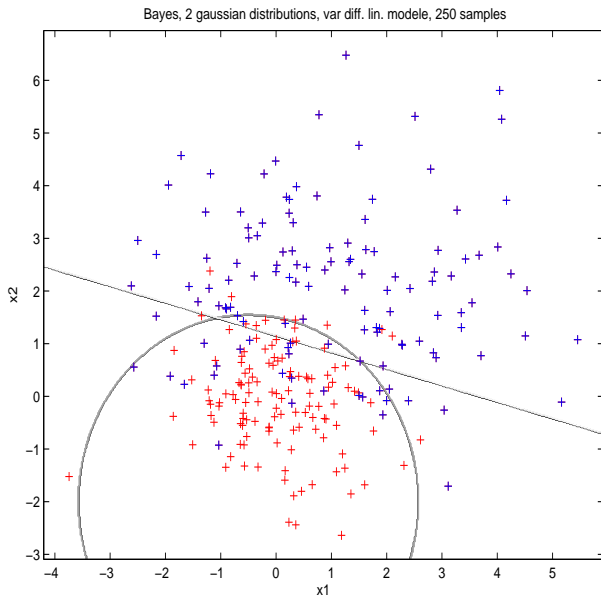
Parametric estimation : Example 2

Experimental results



Example 3

Parametric estimation : Incorrect model



Avoid to use strong assumption about data distribution.

Main methods :

- 1 Parzen density estimator,
- 2 k-nearest neighbor estimator.

Principle

Draw N samples from $p(X)$, count the proportion of them $\frac{n(X)}{N}$ that fall in a given domain of volume ν .

The volume ν is defined a priori.

Estimator

$$\hat{p}(X) = \frac{n(X)}{\nu N}$$

or

$$\hat{p}(X) = \frac{1}{N} \sum_{i=1}^N K(X - X_i) = \frac{1}{N} \sum_{i=1}^N \delta(X - X_i) * K(X)$$

with :

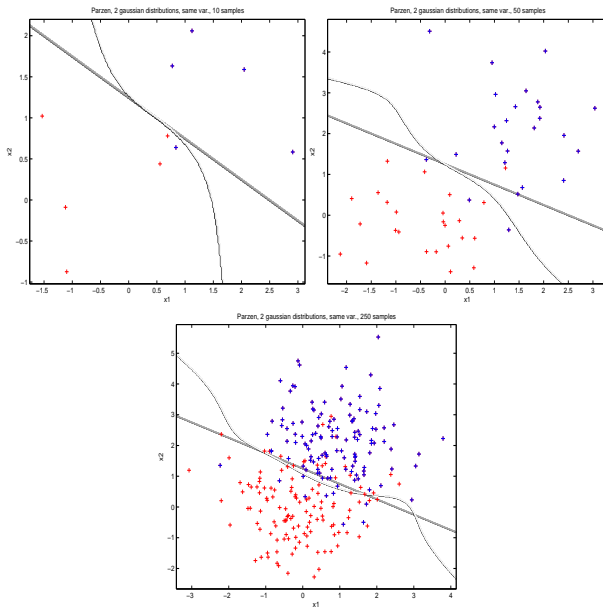
$$\int K(x) dx = 1$$

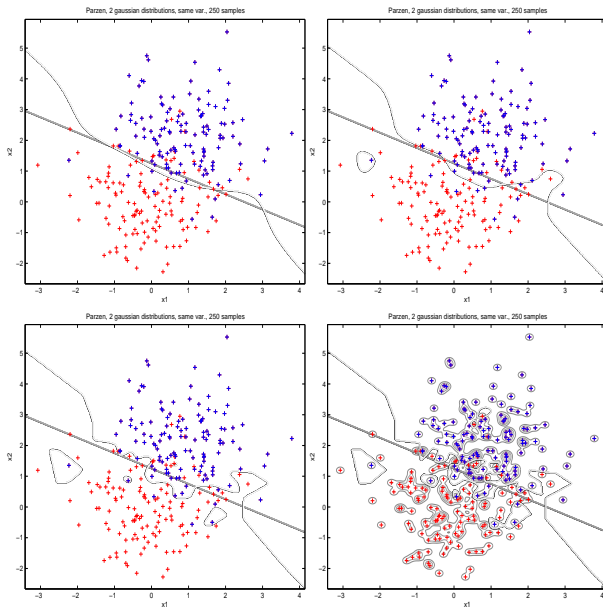
Usual kernels :

- Uniform kernel
- Gaussian kernel

$\hat{p}(X)$ depends on kernel parameters.

Examples (1)





Estimator

$$\hat{p}(X) = \frac{k-1}{V(X)N}$$

Decision

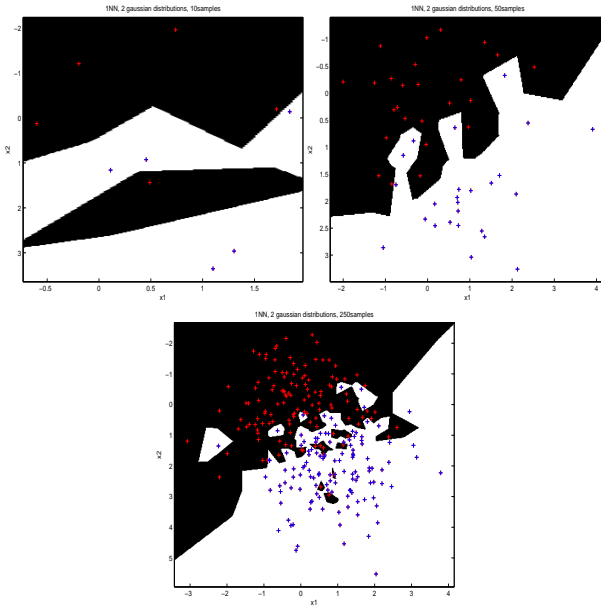
Two options :

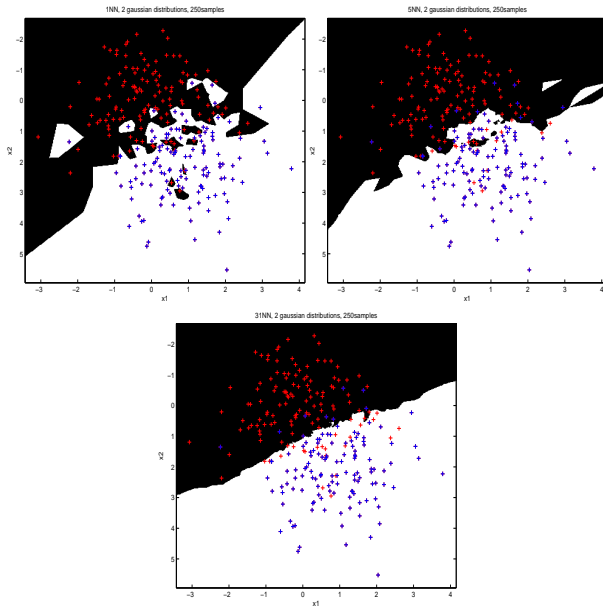
- ① Compare the volume of k NN of the different classes.
- ② Compare the number of neighbors of each class among the k NN (most popular form)

Asymptotic property

when n goes to infinity

$$r^*(X) \leq r_{kNN}(X) \leq 2r^*(X)$$





Imposed structure detection

Without knowledge on conditional distributions of the samples one has to choose another approach which can be :

- 1 define a detector class $\mathcal{D} = \{d(\mathbf{x}, \theta) : \theta \in \Theta\}$
- 2 select the detector in \mathcal{D} which is the best according to a given criteria

Simple at first sight, the implementation of this approach is based on our ability to answer properly to the following questions :

- 1 How to choose and to define the detector class \mathcal{D} ?
- 2 What criterions are pertinent to select a detector in \mathcal{D} ?
- 3 How to explore \mathcal{D} in order to find the best detector?

Chapitre 1 : *Statistical learning theory - Basic concepts*

→ *fonctionnel training, consistency, generalization capability, ...*

Chapitre 2 : *Regularization*

→ *Well posed and ill posed problems, Tikhonov regularization, ...*

Chapitre 3 : *Kernel methods*

→ *RKHS, Mercer's condition, kernelization example, ...*

Chapitre 4 : *Support Vector Classifiers*

→ *Primal problem, Lagrangian, ...*

Chapitre 5 : *Support Vector Regression*

→ *Primal problem, Lagrangian, ...*

Chapitre 6 : *Support Vector one class*

→ *Primal problem, Lagrangian, ...*

Chapitre 7 : *Support Vector multi-class*

→ *Primal problem, Lagrangian, ...*

Chapitre 8 : *Semi-supervised classification*

→ *clustering methods, Generative methods, ...*