

Pattern recognition

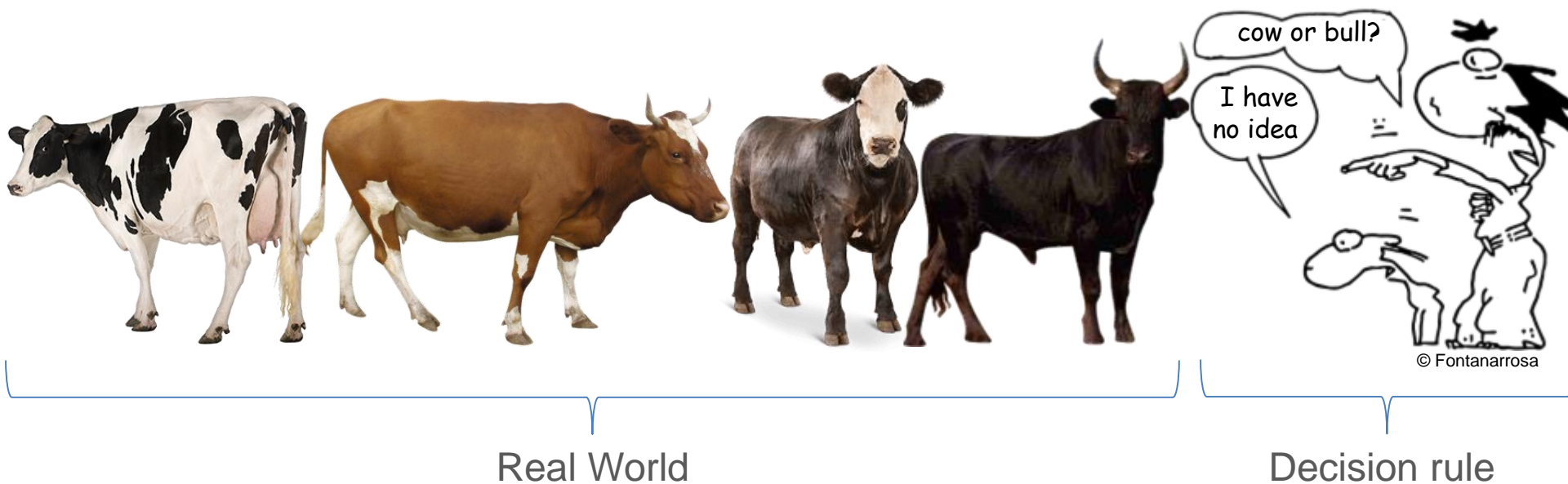
Juan Pablo Martín

Universidad Tecnológica Nacional

jpmartin@gicom.com.ar



Objective





Detection Problem

The general problem can be written as:

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

The aim is to build a classifier d , or a decision rule d which minimizes a criteria (e.g. the error probability)

$$P_e(d) = P(d(X) \neq Y)$$

where X is an observation and Y the associated hypothesis.

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



© Fontanarrosa

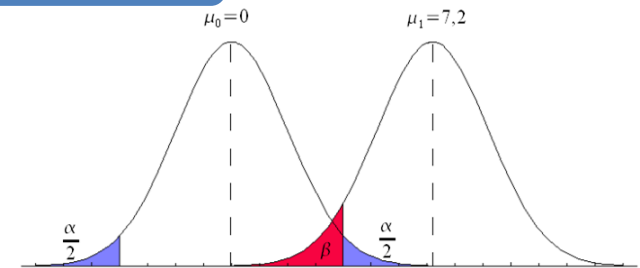
Resolution Approaches

Rule based

- If X is grey and $\text{weight}(X) > 1000$, then
 - X is an elephant
- else
 - X is a mouse
- Needs and expert to translate knowledge into **rules**; are complex, long, not reliable.

Hypothesis testing

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



Bayes test

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

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

minimize

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

$$\frac{P(x|\omega_1)}{P(x|\omega_0)} \underset{q(x)}{\overset{D_1}{>}} \underset{threshold}{\underset{D_0}{<}} \frac{P(\omega_0)}{P(\omega_1)}$$

Likelihood
ratio

Bayes test – Example: Gaussian case

Linear

function!

$$\underbrace{(M_1 - V^T)^T \Sigma^{-1} X + \frac{1}{2} M_0^T \Sigma^{-1} M_0}_{q(x)} \underbrace{U - \frac{1}{2} M_1^T \Sigma^{-1} M_1}_{\text{constant threshold}} \begin{matrix} D_1 \\ > 0 \\ < 0 \\ D_0 \end{matrix}$$

Real World Problems

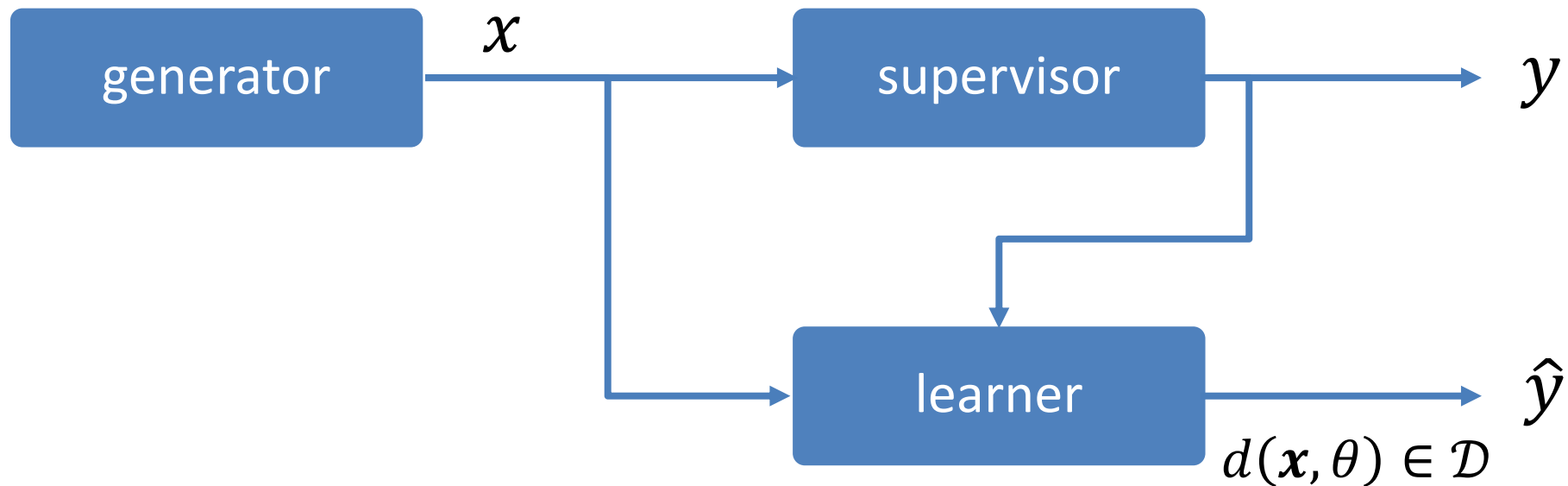
Knowing the distribution function

- $p(X|\omega_0)$ and $p(X|\omega_1)$ has to be known ...

If that is not the case

- Can assume that
 - $p(X|\omega_0)$ and $p(X|\omega_1) \in \mathcal{F}_\theta$
- then estimate θ based on data and plugin the estimator in the decision rule

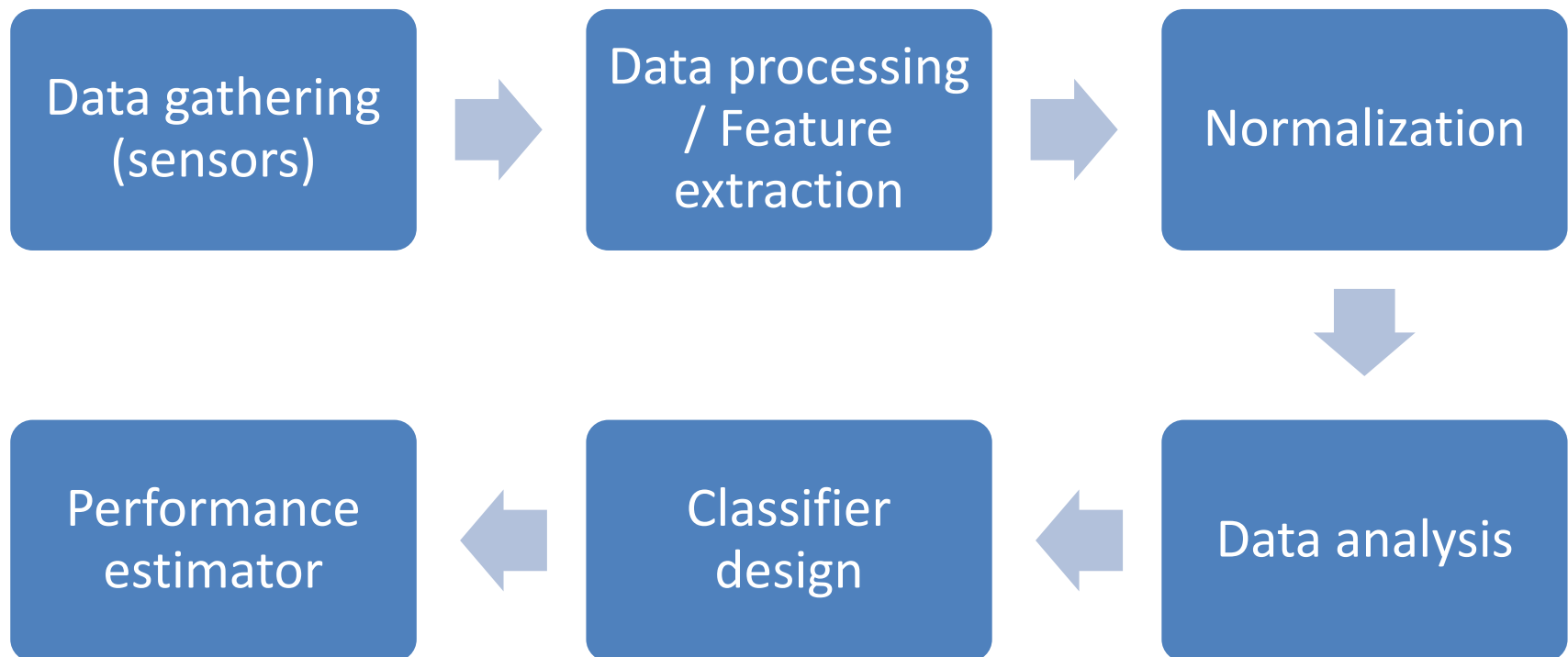
Learning Model



The knowledge of the probabilistic model is replaced by

$$A_n = \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_n, y_n)\} \quad \leftarrow \text{training set of data}$$

Non Parametric Methods



Non Parametric Methods

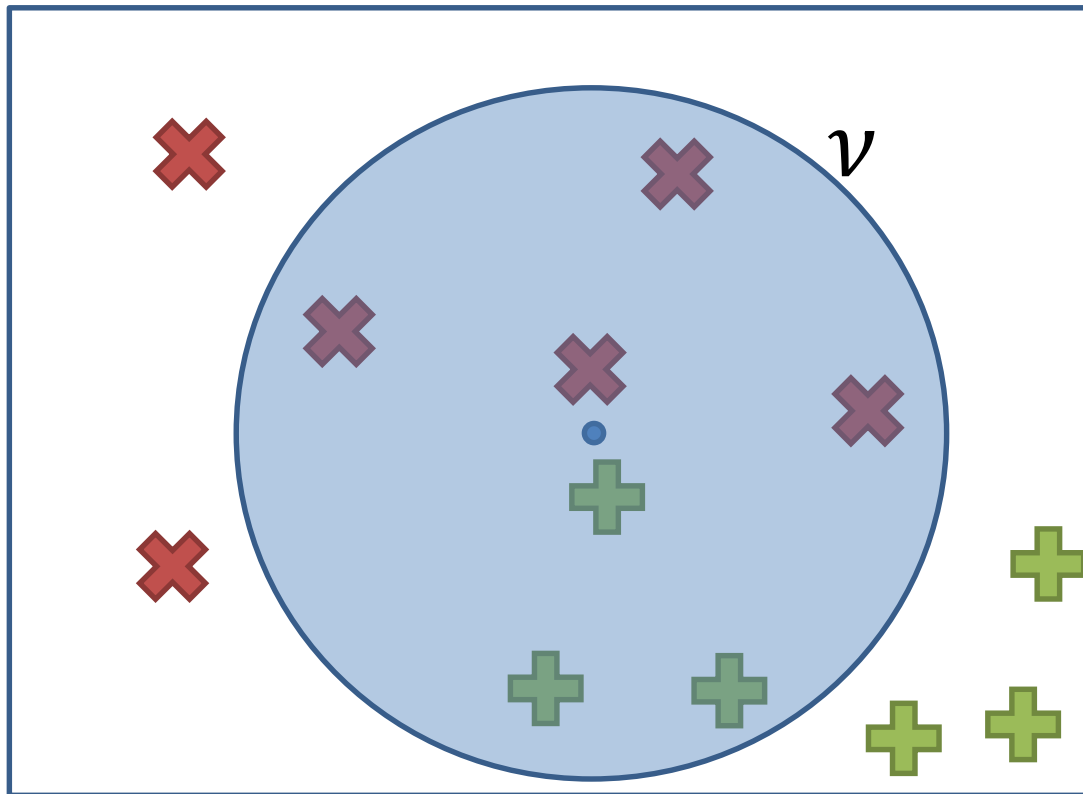


Parzen density
estimator



k-nearest neighbor
estimator

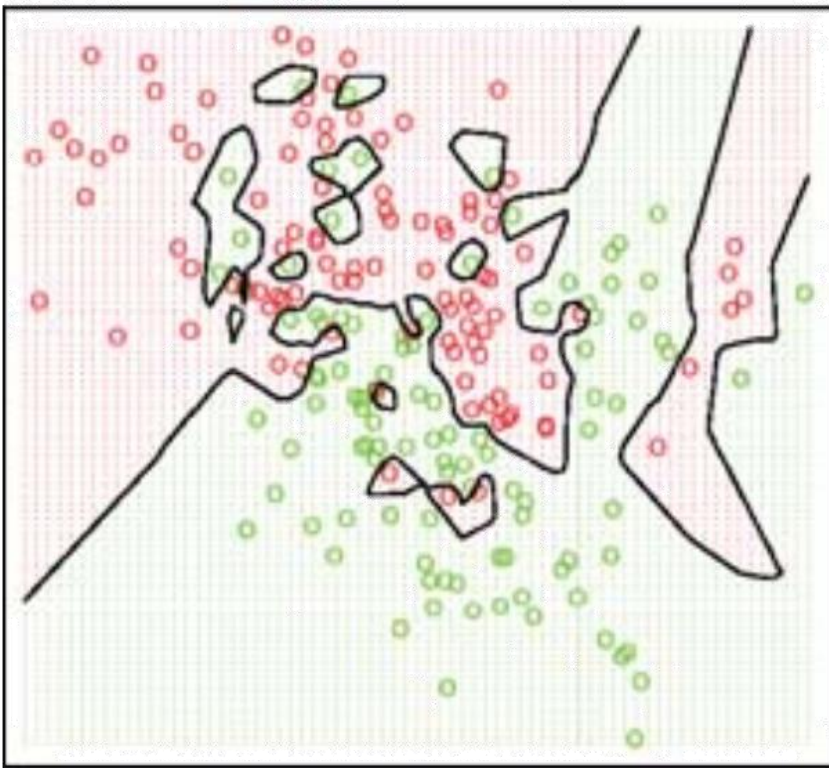
kNN Method



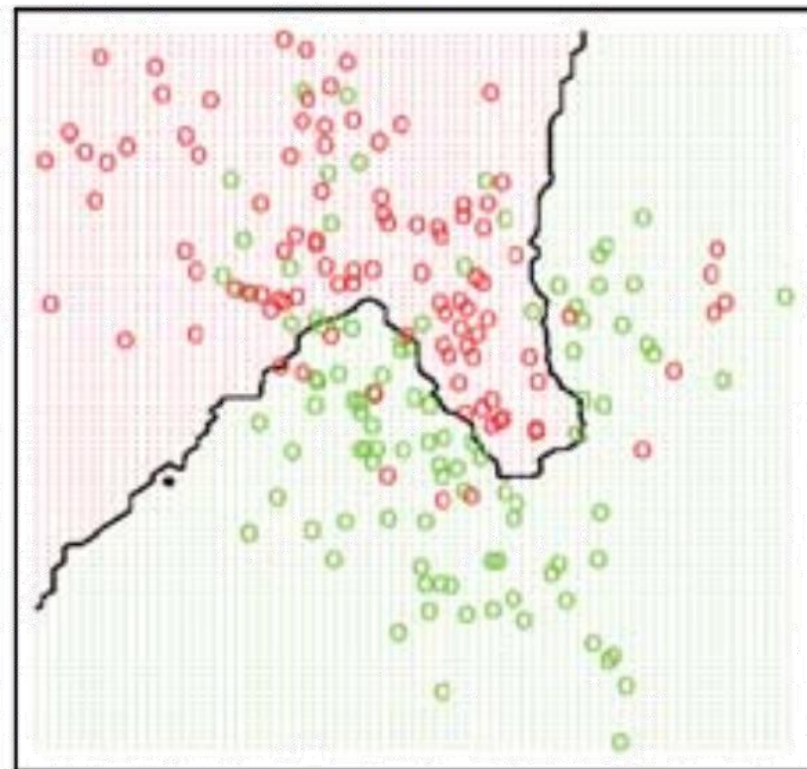
- Centered at the unlabeled sample a spherical volumen ν is enlarged
 - when k samples fall inside the volumen
 - the proportion is counted
 - the new sample **is labeled** with the most occurent class.
- ✓ For the two classes, k should be odd.

Examples

1NN



15NN



Support Vector Machines

Problem of Functional Learning

Find within $\mathcal{D} = \{d(\mathbf{x}, \theta) : \theta \in \Theta\}$, the function which gives the best approximation of y according to a risk functional

$$J(d) = \int Q(\underbrace{d(\mathbf{x}, \theta)}_{\hat{y}}, y) p(\mathbf{x}, y) d\mathbf{x} dy$$

where Q expresses the cost associated with each couple (\mathbf{x}, y) .

Example of a cost function: error probability

$$P_e(d) = \int \mathbb{1}_{d(\mathbf{x}, \theta) \neq y} p(\mathbf{x}, y) d\mathbf{x} dy$$

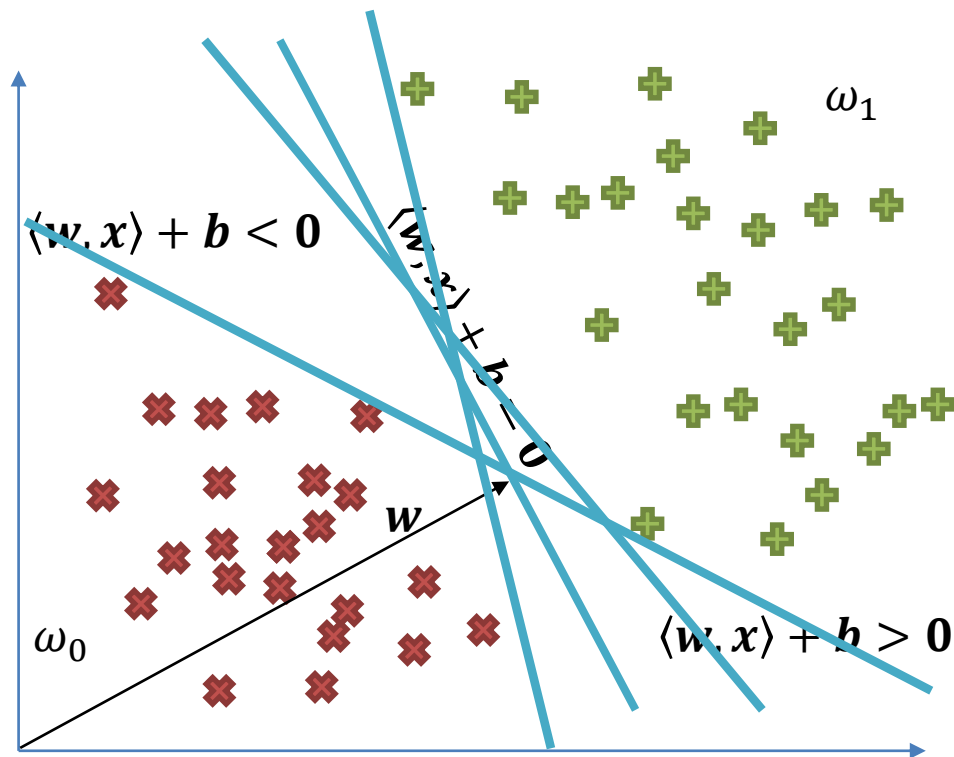
Problem of Functional Learning

As the density $p(\mathbf{x}, y)$ is unknown, the minimization of $J(d)$ is done by plugging an estimator; the empirical risk

$$J_{emp}(d) = \frac{1}{n} \sum_{k=1}^n Q(d(\mathbf{x}_k, \theta), y_k)$$

$$P_{emp}(d) = \frac{1}{2n} \sum_{k=1}^n |y_k - \underbrace{d(\mathbf{x}_k, \mathbf{w}, \mathbf{b})}_{\hat{y}}|$$

Induction principle

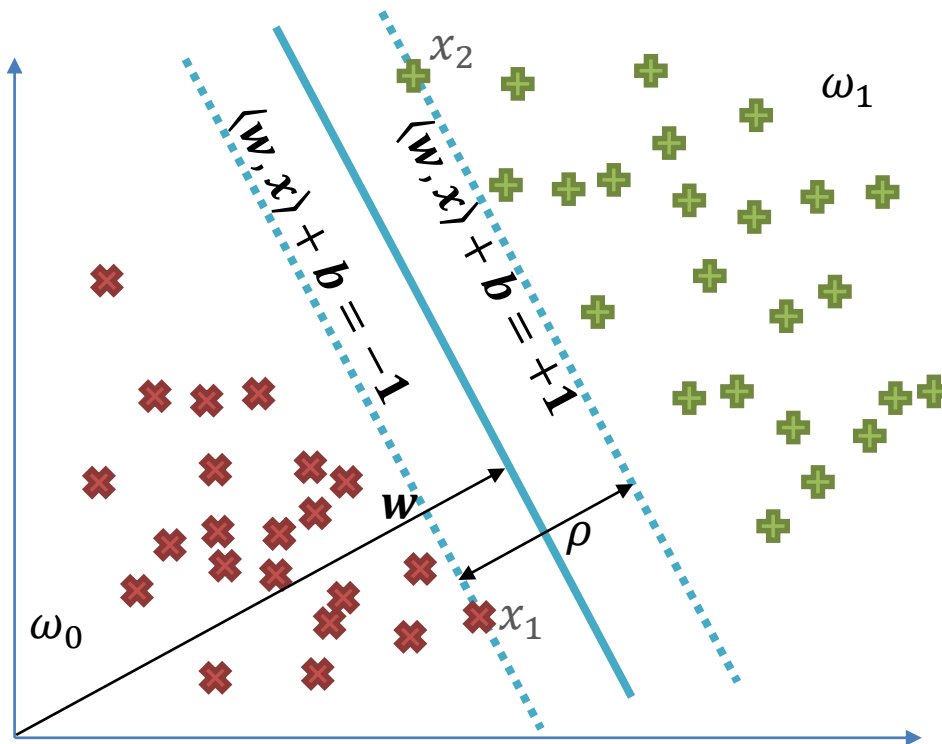


Which separating hyperplane is the best choice?

Should be the one that maximizes the margin!

- Vapnik (1965, 1992)

Margin Calculation



$$\rho = \left\langle \frac{w}{\|w\|}, x_2 - x_1 \right\rangle = \frac{2}{\|w\|}$$

Maximizing the margin, is equivalent to:

- minimizing $\frac{1}{2} \|w\|^2$
- under the constraints

$$y_i(\langle w, x_i \rangle + b) \geq 1, \\ 1 \leq i \leq n$$

Only valid for linearly separable classes

Minimization (Lagrange Multipliers)

Minimizing a convex function $f(x)$ under the constraints $g_i(x) \leq 0, i = 1, \dots, n$ is equivalent to finding the saddle point of the Lagrangian:

$$L(\mathbf{x}, \alpha) = f(\mathbf{x}) + \sum_{i=1}^n \alpha_i g_i(\mathbf{x})$$

Optimality conditions made with respect to the Lagrangian:

$$L(\mathbf{w}, b, \alpha) = \frac{1}{2} \|\mathbf{w}\|^2 + \sum_{i=1}^n \alpha_i \{y_i (\langle \mathbf{w}, \mathbf{x}_i \rangle + b) - 1\}, \quad \alpha_i \geq 0$$

Results in null derivatives with respect to the primal and dual variables:

$$\frac{\partial}{\partial \mathbf{w}} L(\mathbf{w}, b, \alpha) = 0 \quad \text{and} \quad \frac{\partial}{\partial b} L(\mathbf{w}, b, \alpha) = 0$$

Dual problem to solve:

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

Support Vectors

The optimum separator is:

$$\mathbf{w}^* = \sum_{i=1}^n \alpha_i^* y_i x_i$$

According to the Karush-Kuhn-Tucker conditions:

$$\alpha_i^* \{y_i (\langle \mathbf{w}^*, x_i \rangle + b^*) - 1\} = 0, \quad \forall i$$

Two possible cases:

1. $y_i (\langle \mathbf{w}^*, x_i \rangle + b^*) > 1$

Then $\alpha_i^* = 0$, meaning that x_i is not used to calculate \mathbf{w}^* .

2. $y_i (\langle \mathbf{w}^*, x_i \rangle + b^*) = 1$

Then $\alpha_i^* \neq 0$ and x_i is on the margin.

Linearly Separable Solution

The optimum separator is:

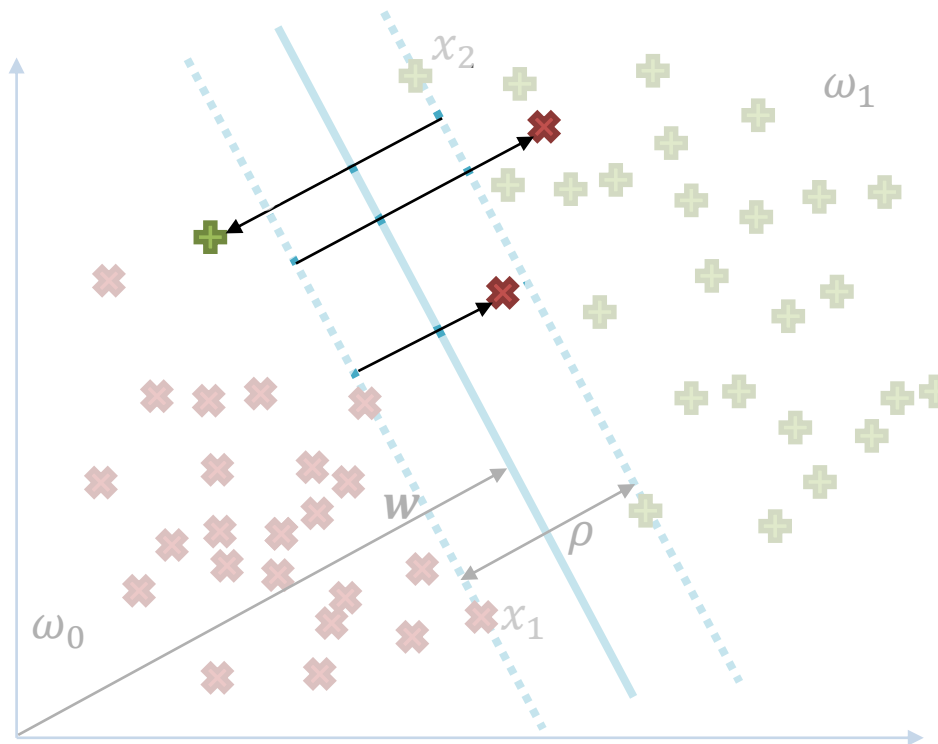
$$\mathbf{w}^* = \sum_{sv} \alpha_i^* y_i \mathbf{x}_i$$

The decision rule is:

$$d(\mathbf{x}, \boldsymbol{\alpha}^*, b^*) = \text{sign} \left(\sum_{sv} \langle \mathbf{x}, \mathbf{w}^* \rangle + b^* \right)$$

$$d(\mathbf{x}, \boldsymbol{\alpha}^*, b^*) = \text{sign} \left(\sum_{sv} \alpha_i^* y_i \langle \mathbf{x}, \mathbf{x}_i \rangle + b^* \right)$$

Classes not linearly separable



- The problem formulation has to be modified.
- Missclassified data is penalized.
 - A cost related to the distance from the sample to the margin is considered.

• A new $f(x)$ is minimized

$$\frac{1}{2} \|w\|^2 + C \underbrace{\sum_{i=1}^n \xi_i}_{\text{Penalizes missclassified samples}}, \quad C \geq 0$$

Non Separable Solution

- A new $f(x)$ is minimized

$$\frac{1}{2} \|w\|^2 + C \sum_{i=1}^n \xi_i, \quad C \geq 0$$

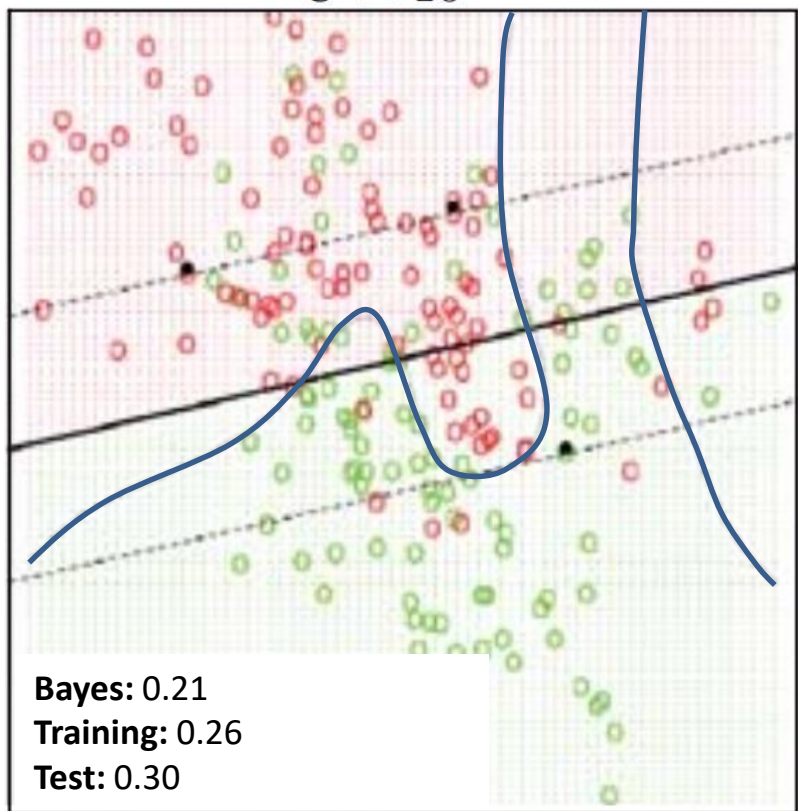
- Under the constraints

$$y_i(\langle w, x_i \rangle + b) \geq 1 - \xi_i, \quad \xi_i \geq 0, 1 \leq i \leq n$$

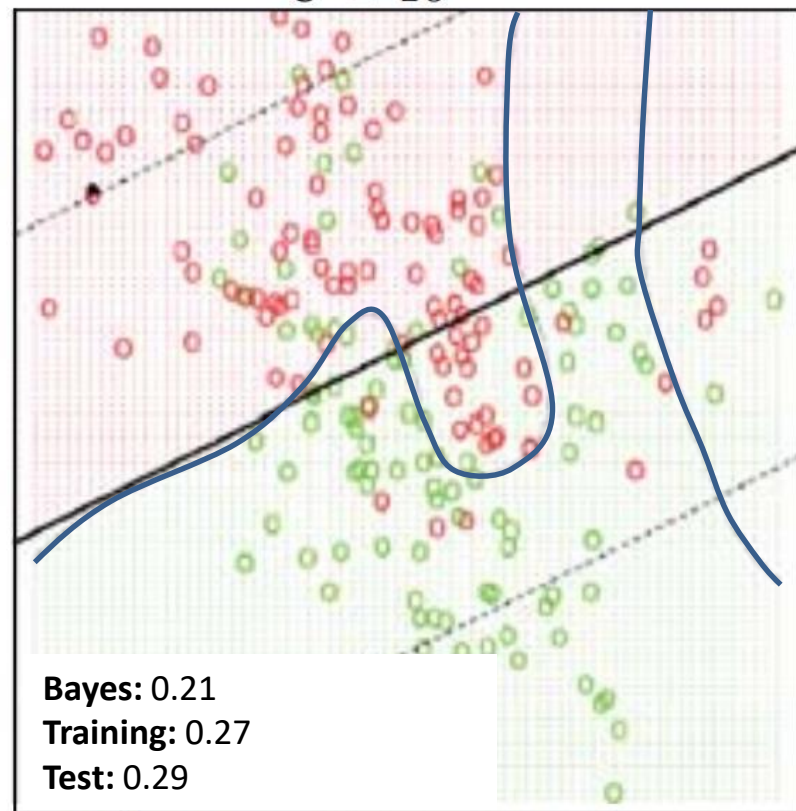
- **C large:** small margin, less training errors.
- **C small:** large margin, more training errors.

Examples

$C = 10^4$



$C = 10^{-2}$



Non-linear Classification

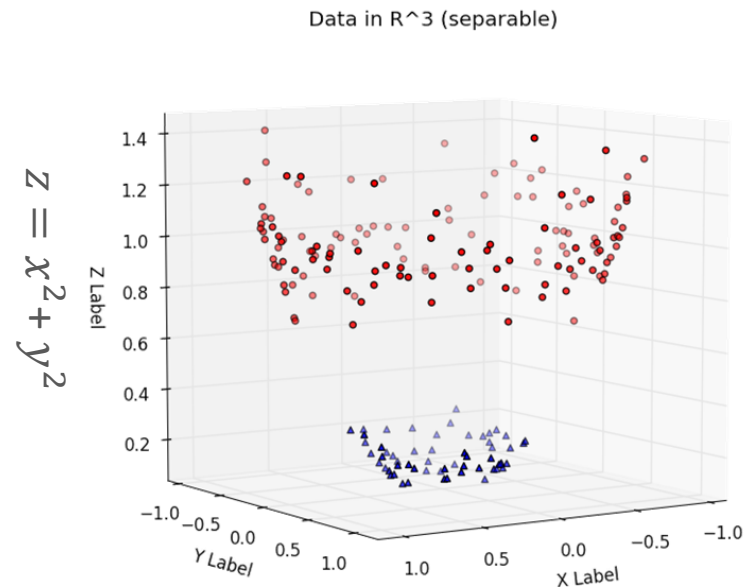
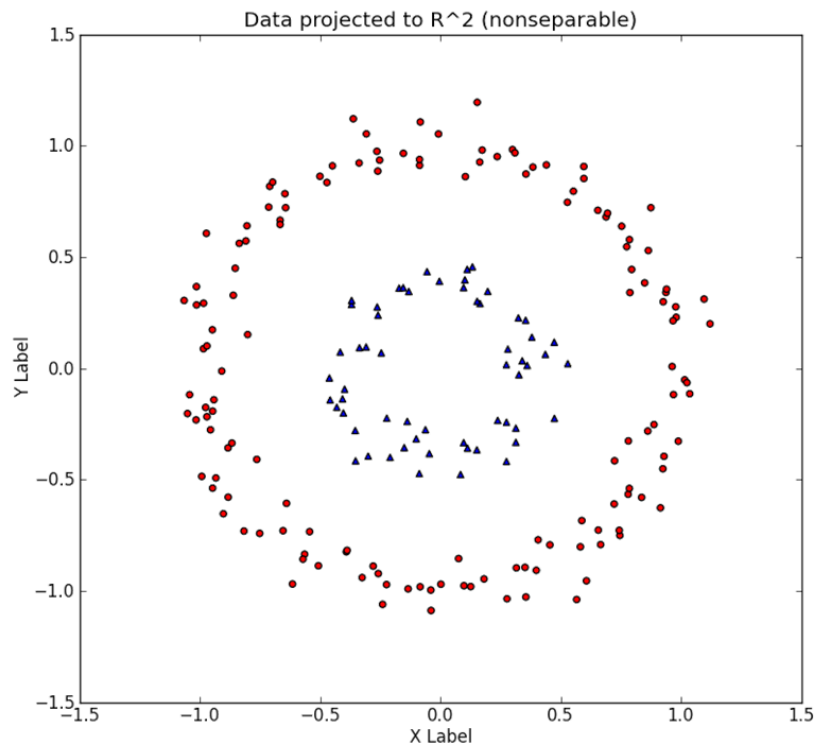
Linear classifiers have limited capabilities.

Can be implemented after a non-linear transformation,

$$\mathbf{x} \rightarrow \phi(\mathbf{x}) = [\phi_1(\mathbf{x}), \phi_2(\mathbf{x}), \dots]^T$$

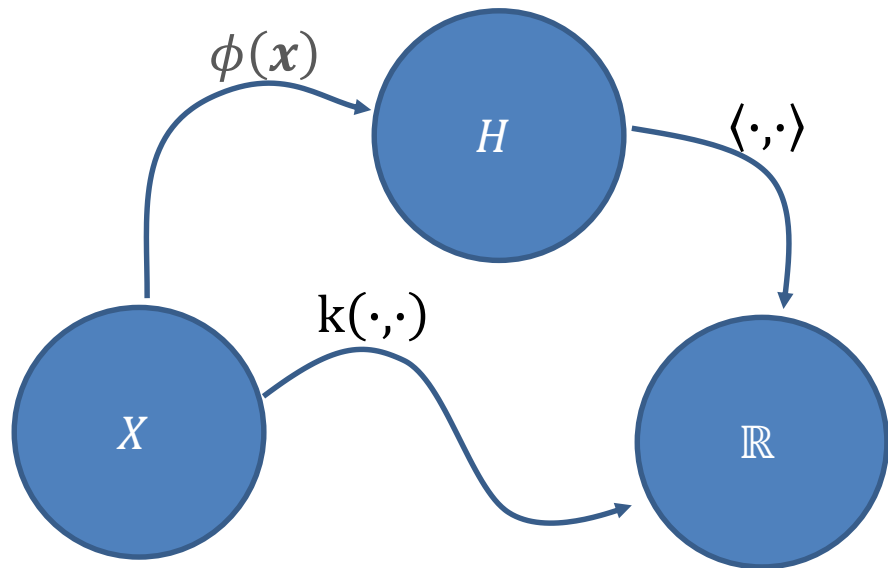
where $\phi_i(\mathbf{x})$ are non-linear functions.

A linear classifier with respect to $\phi(\mathbf{x})$ is non-linear with respect to \mathbf{x} .



$$d(\mathbf{x}, \boldsymbol{\alpha}^*, b^*) = \text{sign} \left(\sum_{sv} \alpha_i^* y_i \langle \mathbf{x}, \boldsymbol{\phi}(\mathbf{x}_i) \rangle + b^* \right) = \text{sign} \left(\sum_{sv} \alpha_i^* y_i \langle \boldsymbol{\phi}(\mathbf{x}), \boldsymbol{\phi}(\mathbf{x}_i) \rangle + b^* \right)$$

Kernel Trick



$$\forall x, x' \in X$$

$$k(x, x') = \langle \phi(x), \phi(x') \rangle_{\mathcal{H}}$$

- Almost no conditions on X
- No need for dot product
- No need to know $\phi(x)$.

Kernel Trick

$$d(\mathbf{x}, \alpha^*, b^*) = \text{sign} \left(\sum_{sv} \alpha_i^* y_i k(\mathbf{x}, \mathbf{x}_i) + b^* \right)$$

Projective kernels

Monomial of degree q	$\langle \mathbf{x}, \mathbf{x}' \rangle^q$
Polynomial of degree q	$(1 + \langle \mathbf{x}, \mathbf{x}' \rangle)^q$
Sigmoidal	$\frac{1}{\eta_0} \tanh(\beta_0 \langle \mathbf{x}, \mathbf{x}' \rangle - \alpha_0)$

... and also:

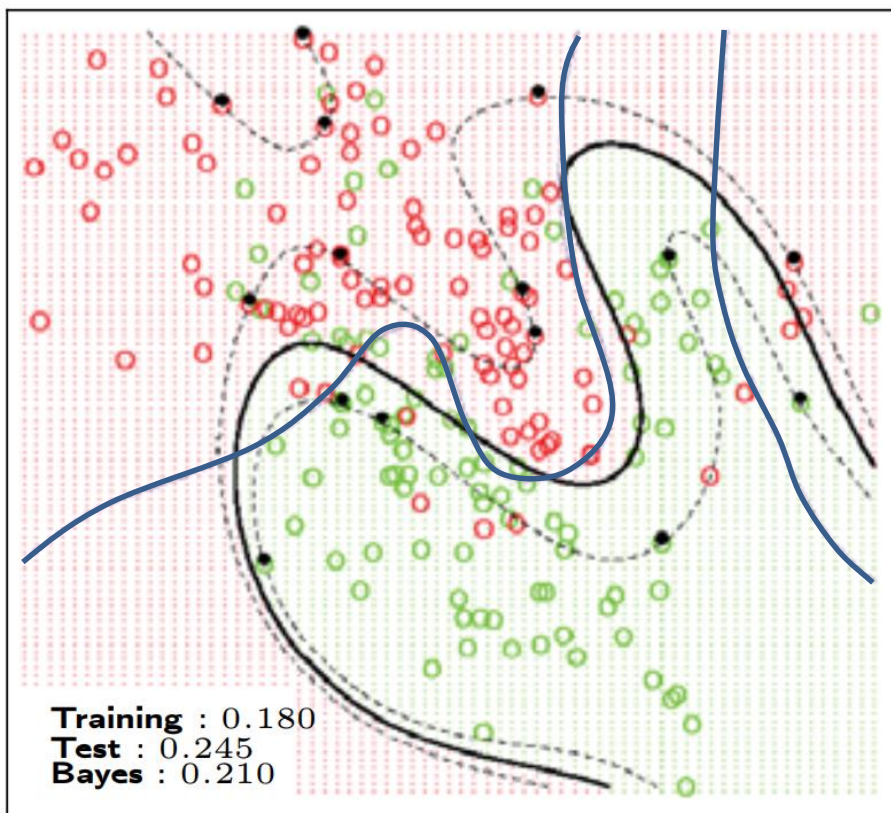
- $k_1(\mathbf{x}, \mathbf{x}') + k_2(\mathbf{x}, \mathbf{x}')$
- $k_1(\mathbf{x}, \mathbf{x}') \cdot k_2(\mathbf{x}, \mathbf{x}')$
- ...

Radial kernels

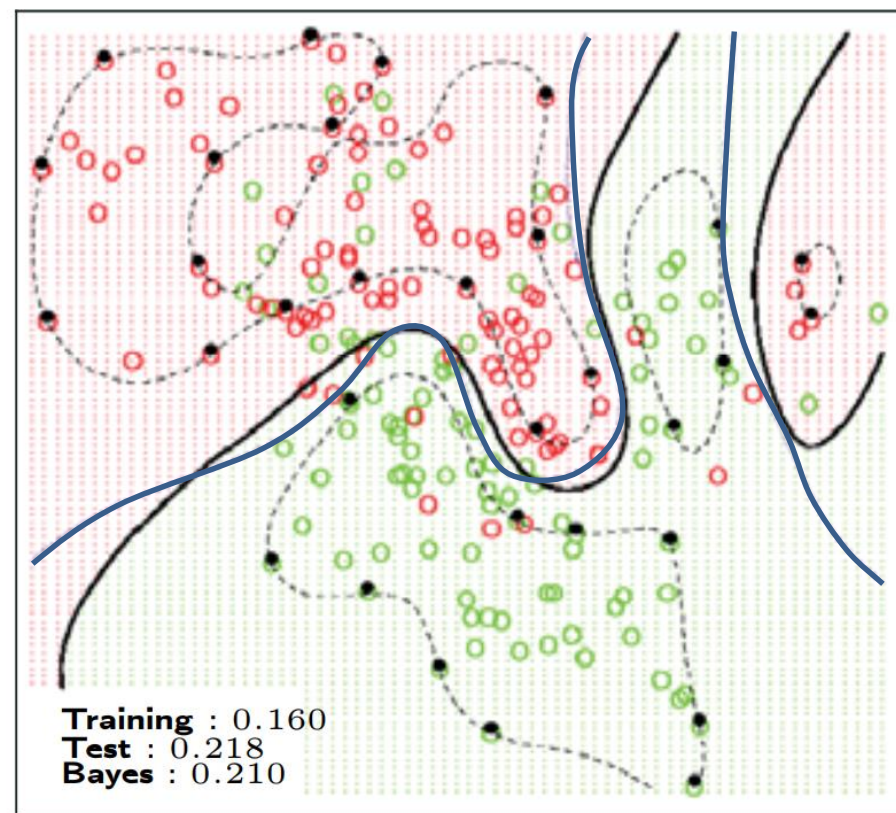
Gaussian	$\exp \left(-\frac{1}{2\sigma_0^2} \ \mathbf{x} - \mathbf{x}'\ ^2 \right)$
Exponential	$\exp \left(-\frac{1}{2\sigma_0^2} \ \mathbf{x} - \mathbf{x}'\ \right)$
Uniform	$\frac{1}{\eta_0} \mathbb{1}_{\ \mathbf{x} - \mathbf{x}'\ \leq \beta_0}$
Epanechnikov	$\frac{1}{\eta_0} (\beta_0^2 - \ \mathbf{x} - \mathbf{x}'\ ^2) \mathbb{1}_{\ \mathbf{x} - \mathbf{x}'\ \leq \beta_0}$
Cauchy	$\frac{1}{\eta_0} \frac{1}{1 + \ \mathbf{x} - \mathbf{x}'\ ^2 / \beta_0^2}$

Examples

polynome



gaussian kernel



Parameter Tuning

k-fold

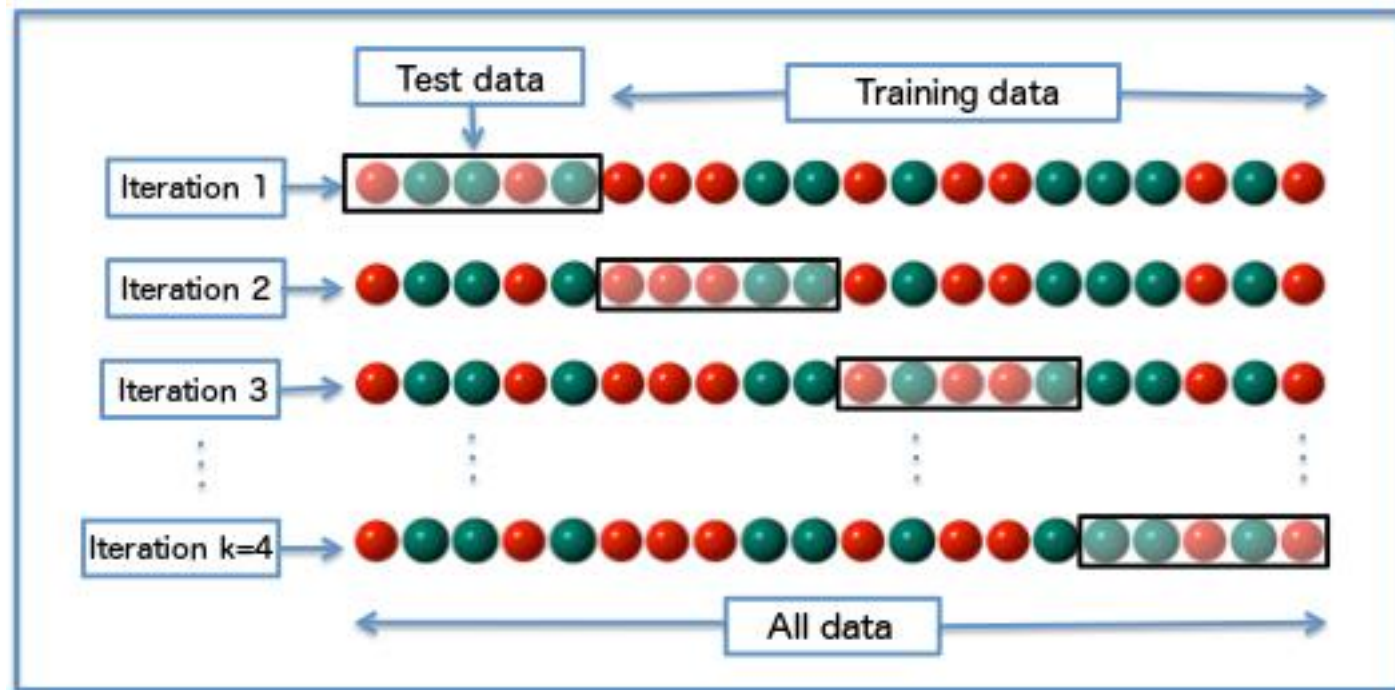
- Regular
- Random

LOOCV

- Leave-one-out

$$E = \frac{1}{k} \sum_{i=1}^k E_i$$

k-fold cross-validation



Source: www.wikipedia.org

Conclusions

Do NOT use
ML as a
black box

Focus on the concept
of the method

Aware of
the
knowledge
of the
system

Correct selection of
the method used

Every
problem is
unique

Choose wisely



ITU KALEIDOSCOPE

SANTA FE 2018

Machine learning for a 5G future

ITU KALEIDOSCOPE

SANTA FE 2018

Thank you

26-28 November
Santa Fe, Argentina

