

Model-based Statistical Learning



Pr. Charles BOUVEYRON

Professor of Statistics
Chair of the Institut 3IA Côte d'Azur
Université Côte d'Azur & Inria

charles.bouveyron@univ-cotedazur.fr
[@cbouveyron](https://twitter.com/cbouveyron)

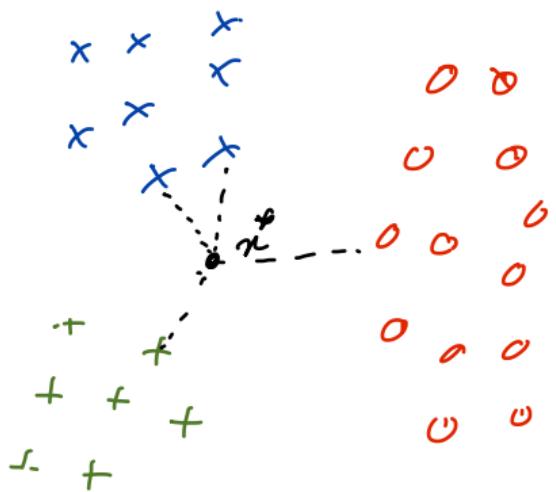
k-nearest neighbors (kNN):

k-NN is a very basic but still reference algorithm for classification. The spirit of k-NN is to compare any new observation x^* to the learning data and classify it according to its neighbors.

Algorithm: input x^* , k

(i) find in the learning dataset (x, y) the k nearest neighbors of x^*

(ii) assign x^* to the class which is the most represented among the neighbors.



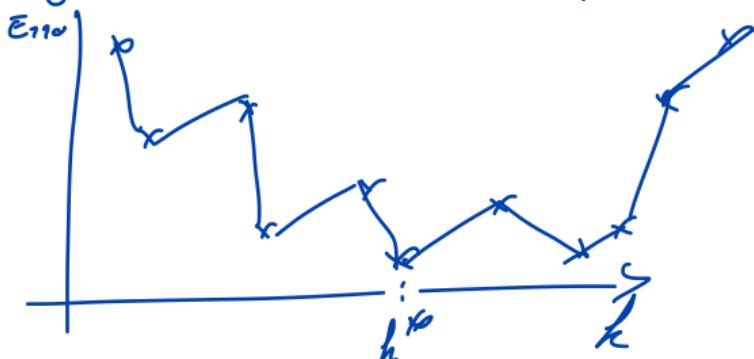
$$k=4$$

$y_{k\text{-NN}}^{\text{blue}}$ = "blue".

Rule: the choice of k in $k\text{-NN}$ is a very sensitive question since it may end up with different classifications for different values of k .

\Rightarrow the appropriate way to choose k is to rely on cross-validation.

Rmk: even with CV, the choice of h may be quite difficult because the curve of the classif. error according to h will be quite irregular.



Rmk: h -NN can be applied to any types of data if we are able to define a distance

- continuous data:
 - euclidean
 - Manhattan

- graph:
 - shortest path
 - average random walk

- text:
 - ...

Summary of k-NN :

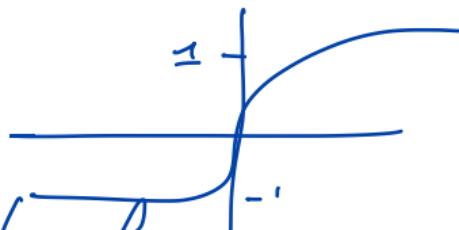
- ⊕ simple and easy to understand method
- ⊕ it can be applied to different data types
- ⊖ the choice of k is sensitive and not easy even with CV.
- ⊖ need to keep the whole learning data set and to find the neighbors in this large data set (which may be long!).

Logistic regression:

Logistic regression is, as LDA, one of the reference method for classification and is still very popular in many domains (even GAFAs). This method is however limited to the binary classification case.

The natural of logistic regression is to transform a classification problem as a regression one thanks to the Logistic function.

Remark: Logistic regression can also be seen as a very simple Neural Network with only an activation function.



The statistical model of Log. reg is:

$$\text{Logit} \left(P(Y=1 | X=x) \right) = \alpha + \underline{\beta^t} X + \epsilon \quad (*)$$

$$\text{and Logit}(p) = \log \left(\frac{p}{1-p} \right) \quad \text{with } \epsilon \sim N(0, \sigma^2)$$

⇒ Learning this model from data resumes to estimating the model parameters $\theta = \{\alpha, \beta, \sigma^2\}$. This can be done thanks to the maximum likelihood (ML) technique.

Unfortunately, the direct optimization of the likelihood is not possible and we have to rely to a numerical optimization procedure (stochastic gradient, ...) to get $\hat{\Theta}_{RL}$

\Rightarrow Once the model is learned, it is possible to use it to make a prediction for any new observation x^* .

As for LDA, we will rely on the maximum a posteriori rule (MAP):

MAP rule:

$$\hat{y}^* = \arg \max_{k=1,2} P(Y=k | X=x^*)$$

where $P(Y=1 | X=x^*)$ can be obtained from
equation (*): Logit ($P(Y=1 | X=x^*)$) = $\hat{\alpha} + \hat{\beta}^t x^*$

$$\Leftrightarrow \log \left(\frac{P(Y=1)}{P(Y=-1)} \right) = \hat{\alpha} + \hat{\beta}^t x^*$$

$$\Leftrightarrow \frac{P(Y=1)}{P(Y=-1)} = \exp(\hat{\alpha} + \hat{\beta}^t x^*)$$

$$\Leftrightarrow P(Y=1) = (1 - P(Y=1)) \exp(-)$$

$$\begin{aligned} \Leftrightarrow P(Y=1) &= \exp(-) - P(Y=1) \exp(-) \\ \Leftrightarrow P(Y=1)(1 + \exp(-)) &= \exp(-) \\ \Leftrightarrow P(Y=1 | X=x^*) &= \frac{\exp(\hat{\alpha} + \hat{\beta}^t x^*)}{1 + \exp(\hat{\alpha} + \hat{\beta}^t x^*)} \end{aligned}$$

if $P(Y=1 | X=x^*) > 0.5 \Rightarrow \hat{y}^* = 1$
 otherwise $\hat{y}^* = -1$.

Summary on logistic regression:

- ⊕ a single method without tuning parameter.
- ⊕ no assumption on the class distributions
- ⊕ the method outputs probabilities to belong to the classes.
- ⊖ limited to only two classes.
- ⊖ it requires the use of an external optimization algorithm.

Outline

1. Introduction
 2. Reminder on the learning process
 3. Model-based statistical learning
 4. Linear models for classification
 5. Mixture models and the EM algorithm
- (...)

Clustering and classification

The task of clustering is the same as in supervised classification, i.e. predict a categorical target variable y , but without having access to examples of both x and y for learning.

(Sup.) classification

$$(x, y) \xrightarrow{\text{learn}} \hat{f} \xrightarrow{\text{predict}} \hat{f}(x^A) = \hat{y}$$

(Unsup.) clustering

$$x \xrightarrow{\text{learn}} \hat{f}(x) = \hat{y}$$

\Rightarrow Clustering is by nature a more challenging task compared to classification.

Definition: the goal of clustering is to form groups such that :

- data in the same group should behave in a similar way
- data in different groups should behave in different manners.

This task is simple but of course heavily combinatorial.

↳ This is why we needed to develop clustering algorithm which approximate the best solution.

Among the possible algorithms for clustering :

- k -means.
 - DB-scan
 - spectral clustering.
 - model-based clustering.
- } non model-based techniques

↳ most of those techniques rely on two ingredients:

- the mixture model
- the EM algorithm.

The mixture model

The mixture model is a family of statistical models that assume that the overall p.d.f is a mixture of several prob. density functions:

$$p(x) = \sum_{k=1}^K \pi_k p_k(x)$$



where $p_k()$ is a specific p.d.f and the weights π_k are such that $\pi_k > 0$ and $\sum_{k=1}^K \pi_k = 1$.

The mixture model

The mixture model is a class of models which extremely flexible. Some examples:

$$* p(x) = 0.6 N(x; 0, 1) + 0.4 U(0, 1)$$

$$* p(x) = 0.3 dN(x; 0, 1) + 0.3 N(x; 2, 1) + 0.4 St(x, \dots)$$

$$* p(x) = \sum_{h=1}^K \pi_h N(x; \mu_h, \Sigma_h)$$

This model is known as the Gaussian Mixture model (GMM)

$$* p(x) = \sum_{h=1}^K \pi_h P(x; \lambda_h)$$

Poisson Mixture Model

$$\# p(x) = \sum_h \pi_h B(\quad)$$

The Gaussian mixture model (GMM)

The GMM model is clearly the most popular mixture model :

$$p(x) = \sum_{k=1}^K \pi_k N(x; \mu_k, \Sigma_k)$$

where μ_k is the mean of the k^{th} component and Σ_k is the covariance matrix of this same component.

Rank: when using this model for clustering, we naturally associate the notion of mixture component to the groups that we are looking for.

GMM for supervised classification

we indeed recognize that :

- the model behind QDA is the Gaussian mixture model:

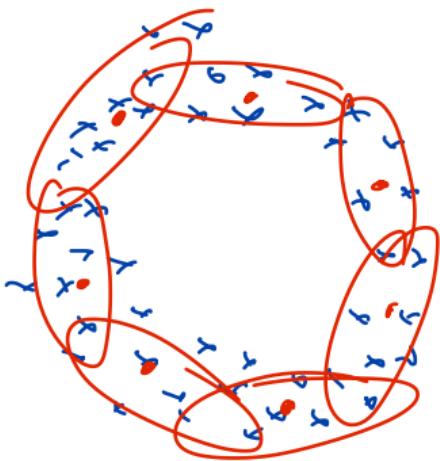
$$\begin{cases} P(Y=h) = \pi_h \\ p(x | Y=h) = \mathcal{N}(x; \mu_h, \Sigma_h) \end{cases}$$

✓ marginarialis
over \Rightarrow

$$p(x) = \sum_{h=1}^K \pi_h \mathcal{N}(x; \mu_h, \Sigma_h)$$

- LDA : $p(x) = \sum_{k=1}^K \pi_k, \mathcal{N}(x; \mu_k, \Sigma)$
 ↳ it is a restricted / constrained GNDP.

Even though the GPR model is based on a single specific density, it can be used to fit data that are very complex and extremely far to be Gaussian:



GMM for clustering

From now, we are going to focus on the specific task of clustering with a specific Mixt. model which is the GMM.

Clustering with GMM is a two-step procedure:

- 1) fit the GMM model from data with the EM algorithm. (estimate model parameters)
- 2) deduce the clustering from the fitted model.

Remark: Let's assume for the moment that K is known.