

Majeure Machine Learning

Techniques

Contenu



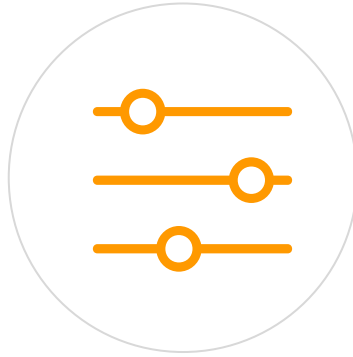
- Différents algorithmes supervisés
- Méthodes ensemblistes
- Différents algorithmes non-supervisés
- Hyperparameters Tuning

Ce que vous devrez savoir faire



- Comprendre le principe de fonctionnement de certains algorithmes de ML
- Comprendre les méthodes ensemblistes
- Etre capable de choisir les bonnes techniques
- Comprendre les différentes techniques d'Hyperparameters Tuning

Apprentissage Supervisé



Algorithmes paramétriques

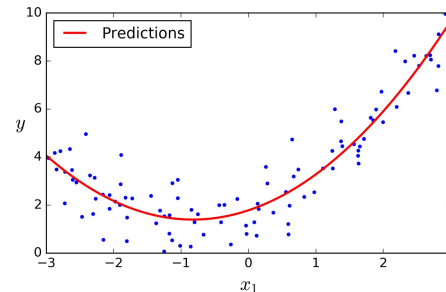
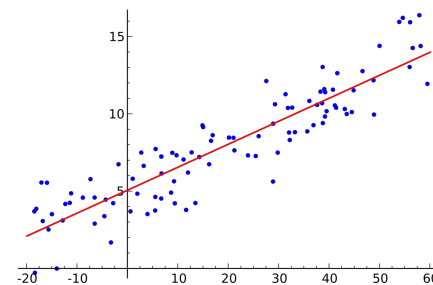
Définition : Nombre de paramètres est fixe et donc la complexité ne dépend pas du nombre de lignes des données

Régression

Régression

- Linéaire : $\hat{y}(x) = \theta_1 x_1 + \theta_0$

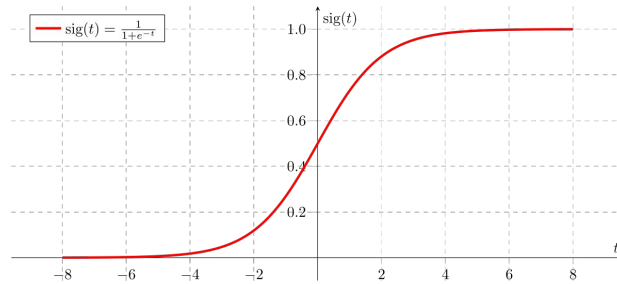
- Polynomiale : $\hat{y}(x) = \theta_2 x^2 + \theta_1 x_1 + \theta_0$



Classification

- Logistique : $\hat{y}(x) = \sigma(\theta_1 x_1 + \theta_0)$

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



Naive* Bayes

Théorème de Bayes :

$$P(A | B) = \frac{P(B | A) P(A)}{P(B)}$$

**Naive => Fait l'hypothèse que les features sont indépendantes les unes des autres.*

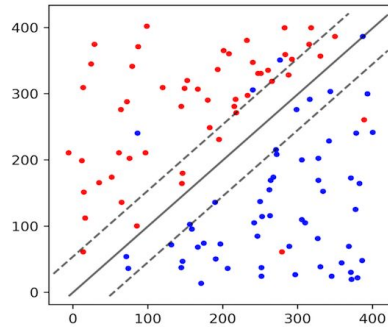
Peu probable mais marche plutôt bien quand même

- $P(X_1 \dots X_n | Y) = \prod_{i=1}^n P(X_i | Y) \Rightarrow$ Likelihood
- $P(Y) = 1/n_classes \Rightarrow$ Prior
- $P(X_n) =$ Gaussian, Bernoulli... \Rightarrow Evidence

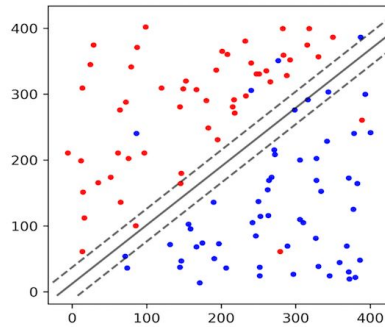
SVM (Support Vector Machine)

Marge

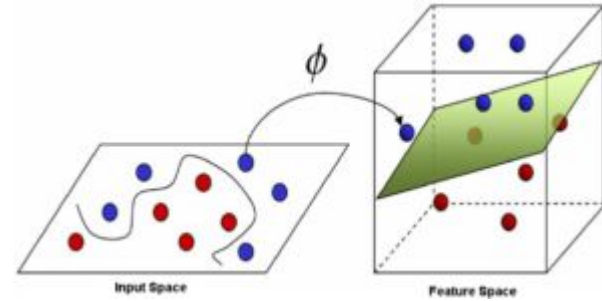
SVM Parameter C



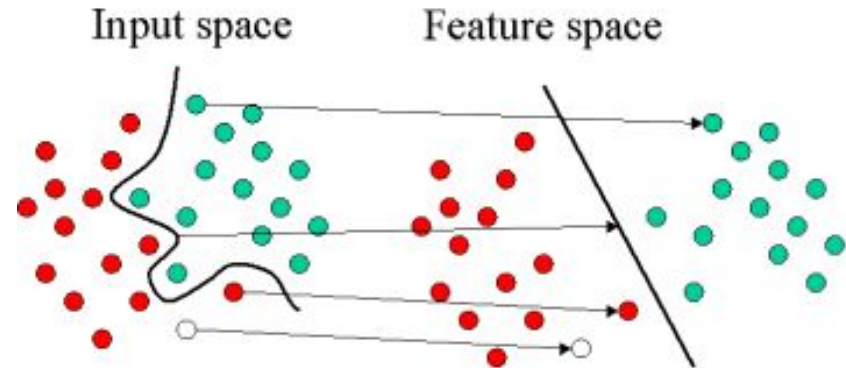
C = 1



C = 100



Kernel



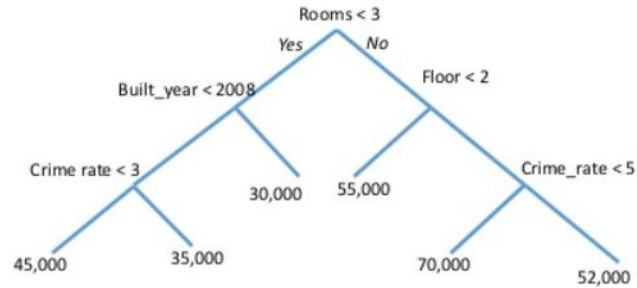


Algorithmes non-paramétriques

Définition : Nombre de paramètres et donc complexité du modèle dépendante du nombre de lignes

Arbre de décision (CART)

Régression

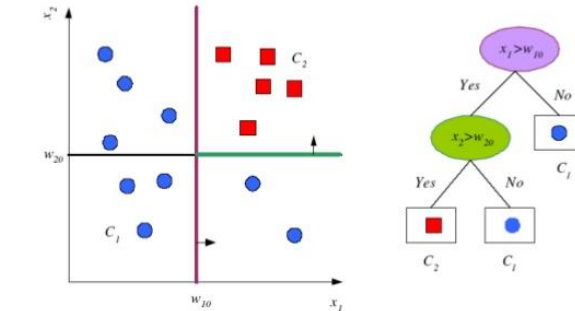


Tester toutes les séparations possibles
Trouver la feature et la valeur qui sépare le mieux les données*
Séparer par cette valeur
Entrer dans une des séparations
Si $Nb_element > min_element_par_feuille$:
Recommencer le process

Sinon :
Si $nb_elements > 0$:
Prédiction = $mean(elements) / majorité(elements)$
Sinon :
Prédiction = élément
Changer de séparation ou remonter

Classification

Decision Tree

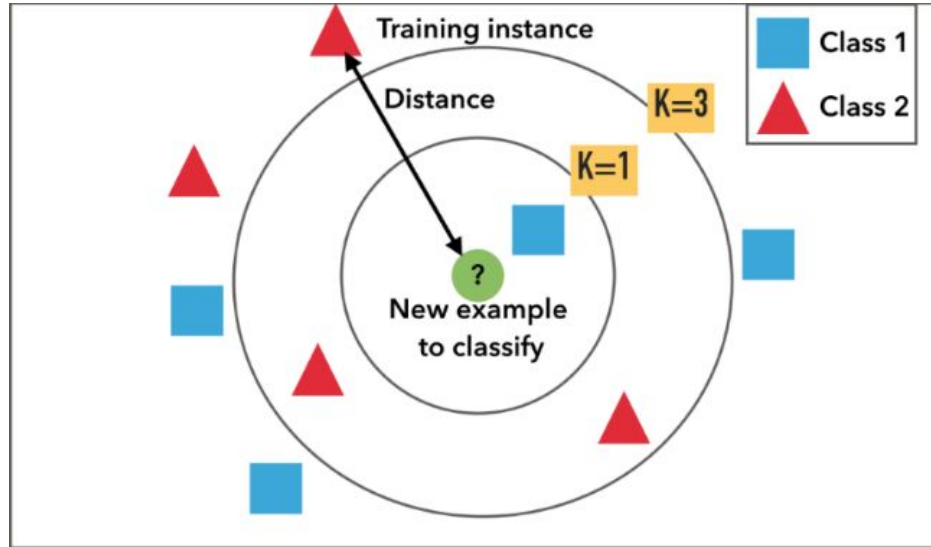


► 10

Lec 2: Decision Trees - Nearest Neighbors

- * Cost function :
- Classification : Gini criterion
 - Régression : Somme des erreurs au carré

K-Nearest-Neighbors



**Prédiction = moyenne / majorité des
k plus proches voisins**



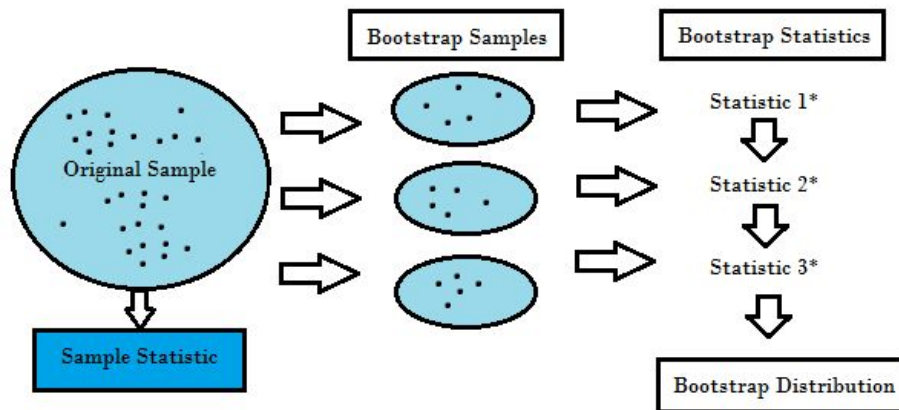
Algorithmes ensembliste

Même avec un seul algorithme on peut améliorer sa performance à l'aide d'un ensemble

Bootstrapping

Objectif : Pouvoir générer plusieurs ensembles d'entraînement

```
 $D_{bootstrap} \leftarrow \{\}$   
pour  $N$  itérations  
    choisir aléatoirement et uniformément un entier  $n$  parmi  $\{1, \dots, N\}$   
     $D_{bootstrap} \leftarrow D_{bootstrap} \cup \{(x_n, t_n)\}$   
retourner  $D_{bootstrap}$ 
```



Combinaisons de modèles - Bagging (bootstrap aggregaation)

A utiliser lorsque les modèles ont une forte capacité

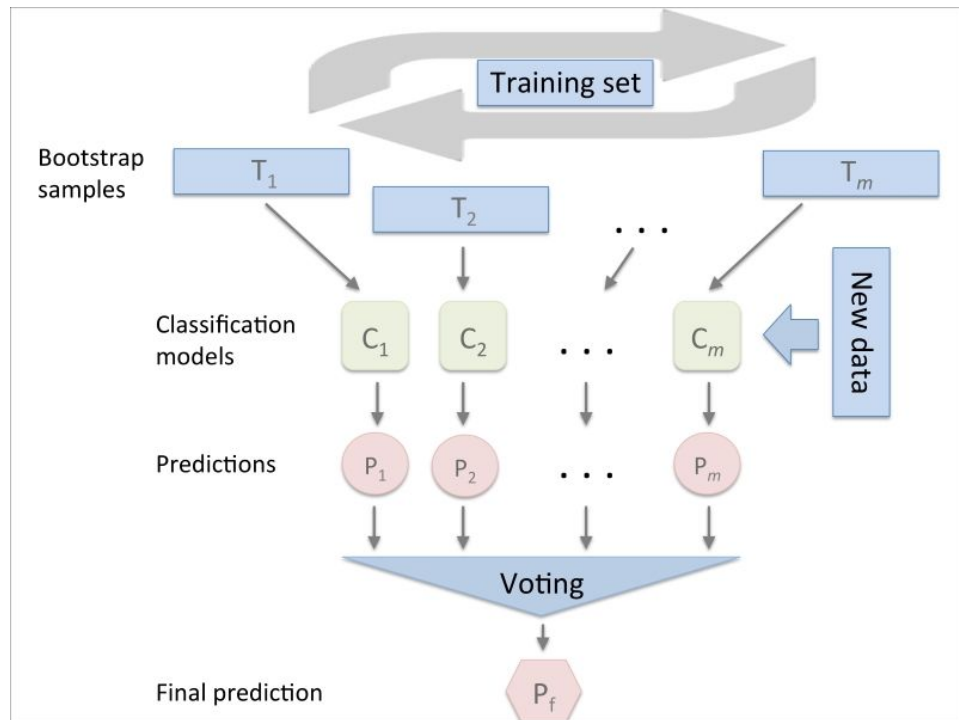
- **Bagging** : entraîne M modèles avec un algorithme donné, sur M ensembles de données *bootstrap*

pour $m = 1, \dots, M$

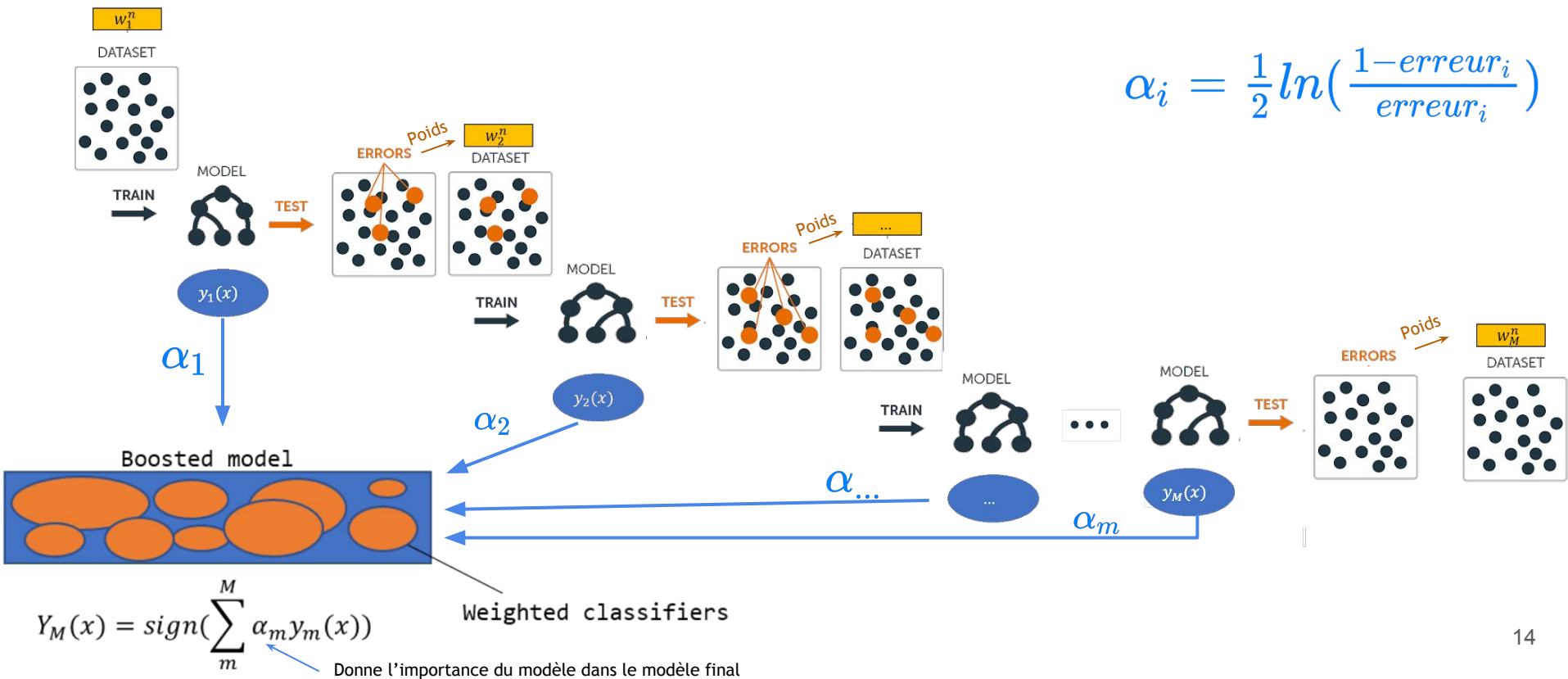
- génère un ensemble de données *bootstrap* $\mathcal{D}_{\text{bootstrap}}$ à partir de \mathcal{D}
- entraîner un modèle $y_m(\mathbf{x})$ sur $\mathcal{D}_{\text{bootstrap}}$

retourner le modèle ensemble (comité)

$$y_{\text{COM}}(\mathbf{x}) = \frac{1}{M} \sum_{m=1}^M y_m(\mathbf{x})$$

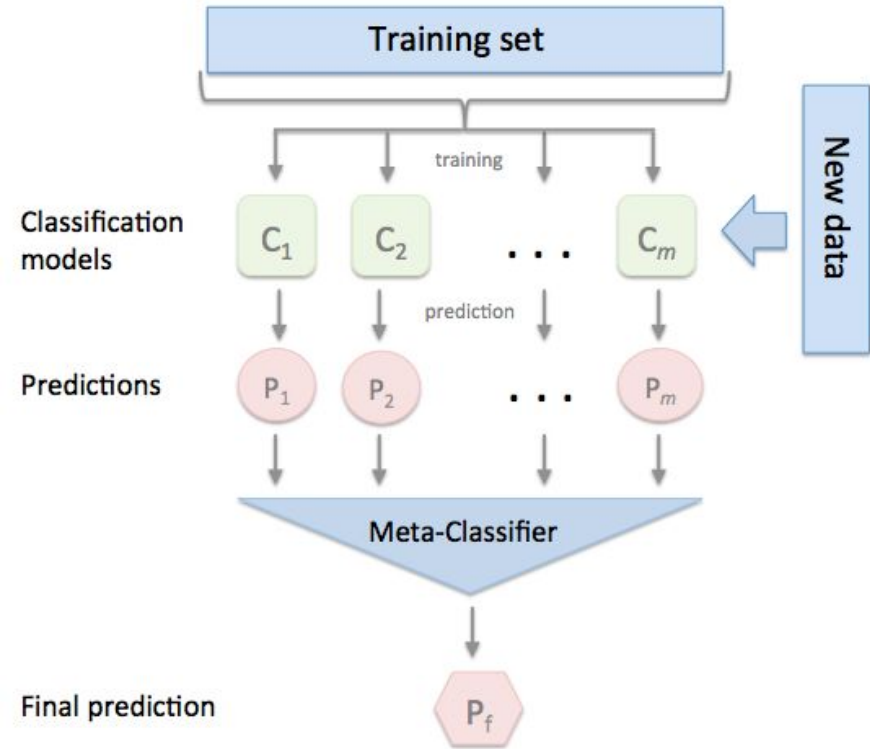


Combinaisons de modèles - Boosting



$$\alpha_i = \frac{1}{2} \ln\left(\frac{1 - \text{erreur}_i}{\text{erreur}_i}\right)$$

Combinaisons de modèles - Stacking

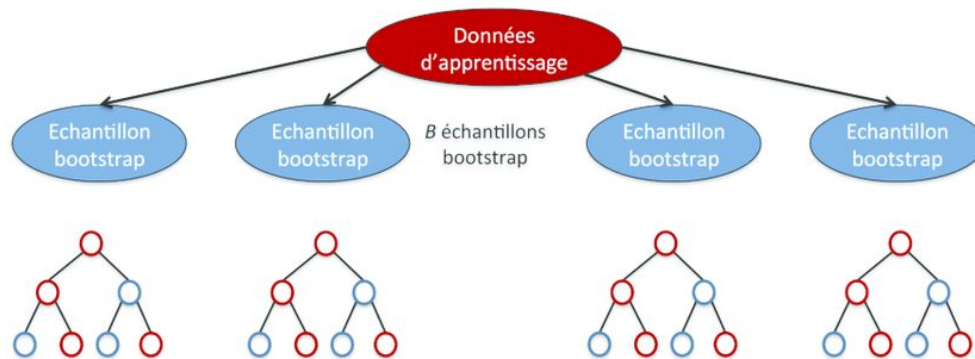


Meta-Classifier : c'est un modèle qui va prendre en entrée les prédictions des sous-modèles (peu importe les types de modèles) et qui va apprendre à les combiner pour en obtenir un résultat optimum

Random forest

Random forest = tree bagging + feature sampling

Feature sampling => donner accès à différentes informations selon les arbres (explication jeu du “qui est-ce”)



A chaque nœud : tirage aléatoire de m variables parmi les p variables

Classification : $m = \sqrt{p}$

Régression : $m = p/3$

Prédiction

Classification : classe majoritaire prédite par les B arbres

Régression : moyenne des valeurs prédites par les B arbres

Gradient Boosting

Gradient Boosting = Descente de Gradient + Boosting

Dans le **boosting** l'importance du vote de chaque classifieur est calculé par (selon sa performance) :

$$\alpha_i = \frac{1}{2} \ln\left(\frac{1 - \text{erreur}_i}{\text{erreur}_i}\right)$$

Le **Gradient Boosting** calcule ce poids différemment. On applique alors une descente de gradient afin d'optimiser ce poids pour chaque arbre afin d'obtenir le meilleur "meta-model"

Fonction de coût

Comparatif

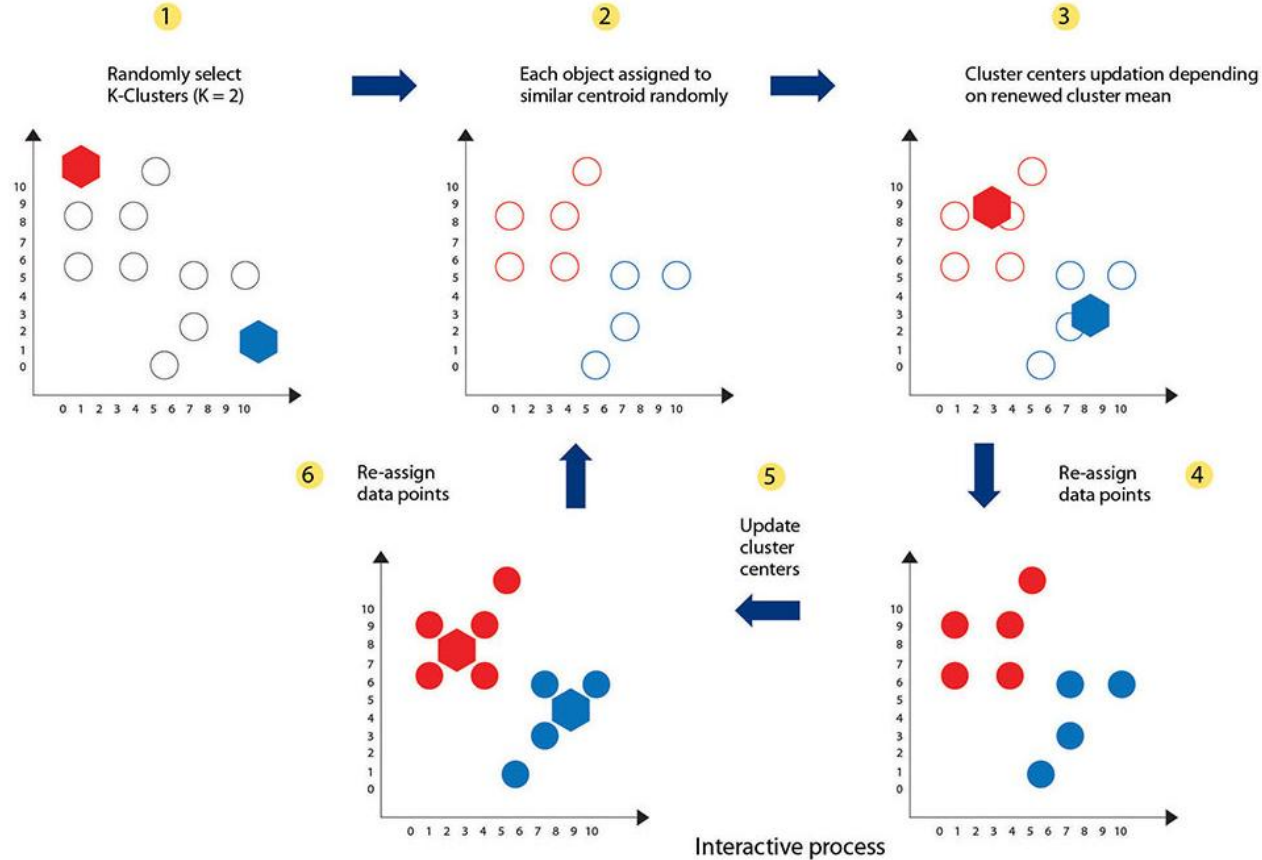
Algorithme	Hyperparamètres	Avantages	Inconvénients
Régression	<ul style="list-style-type: none"> - Degré si polynomiale - Type de régularisation - Coefficient de régularisation 	<ul style="list-style-type: none"> - Rapide - Explicable 	<ul style="list-style-type: none"> - Peu adaptée aux problèmes complexes et non linéaires
Arbre de décision	<ul style="list-style-type: none"> - Nombre minimum d'éléments de feuille - Profondeur maximum 	<ul style="list-style-type: none"> - Non-linéaire - Explicable - Rapide 	<ul style="list-style-type: none"> - Enclin à overfitting - Performances moyennes
Naive Bayes	<ul style="list-style-type: none"> - Loi de la probabilité d'évidence 	<ul style="list-style-type: none"> - Rapide - Nécessite peu d'exemples - Efficace 	<ul style="list-style-type: none"> - Hypothèse d'indépendance forte - Performances moyennes
SVM	<ul style="list-style-type: none"> - Kernel - Coefficient de régularisation 	<ul style="list-style-type: none"> - Très performant - Applicable à tous les problèmes 	<ul style="list-style-type: none"> - Long et coûteux
KNN	<ul style="list-style-type: none"> - K 	<ul style="list-style-type: none"> - Très rapide 	<ul style="list-style-type: none"> - Performances moyennes
Random Forest	<ul style="list-style-type: none"> - Nombre d'arbres + Hyperparamètres des arbres 	<ul style="list-style-type: none"> - Très performant - Applicable à tous les problèmes 	<ul style="list-style-type: none"> - Long et coûteux - Enclin à overfitting
Gradient Boosting	<ul style="list-style-type: none"> - Fonction de coût + Hyperparamètres des arbres 	<ul style="list-style-type: none"> - Très (Très) performant - Applicable à tous les problèmes 	<ul style="list-style-type: none"> - Très long et coûteux

Apprentissage Non-Supervisé



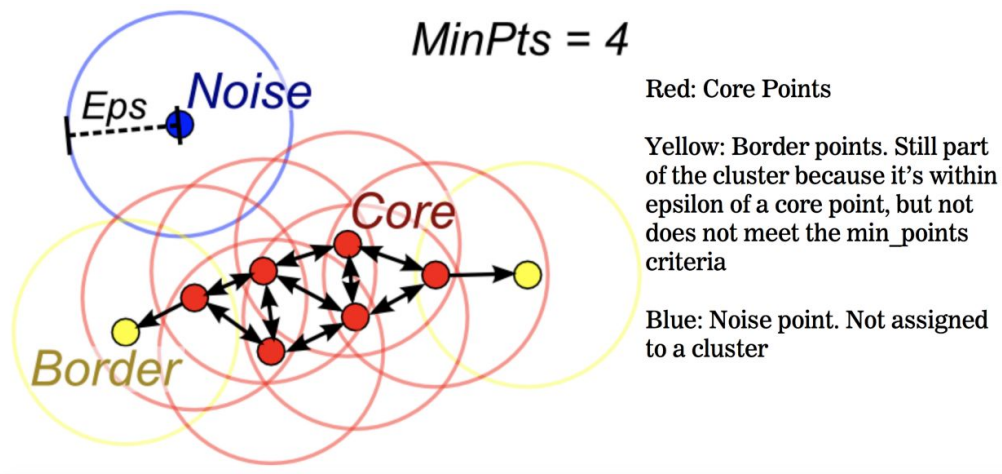
Clustering

K-Means

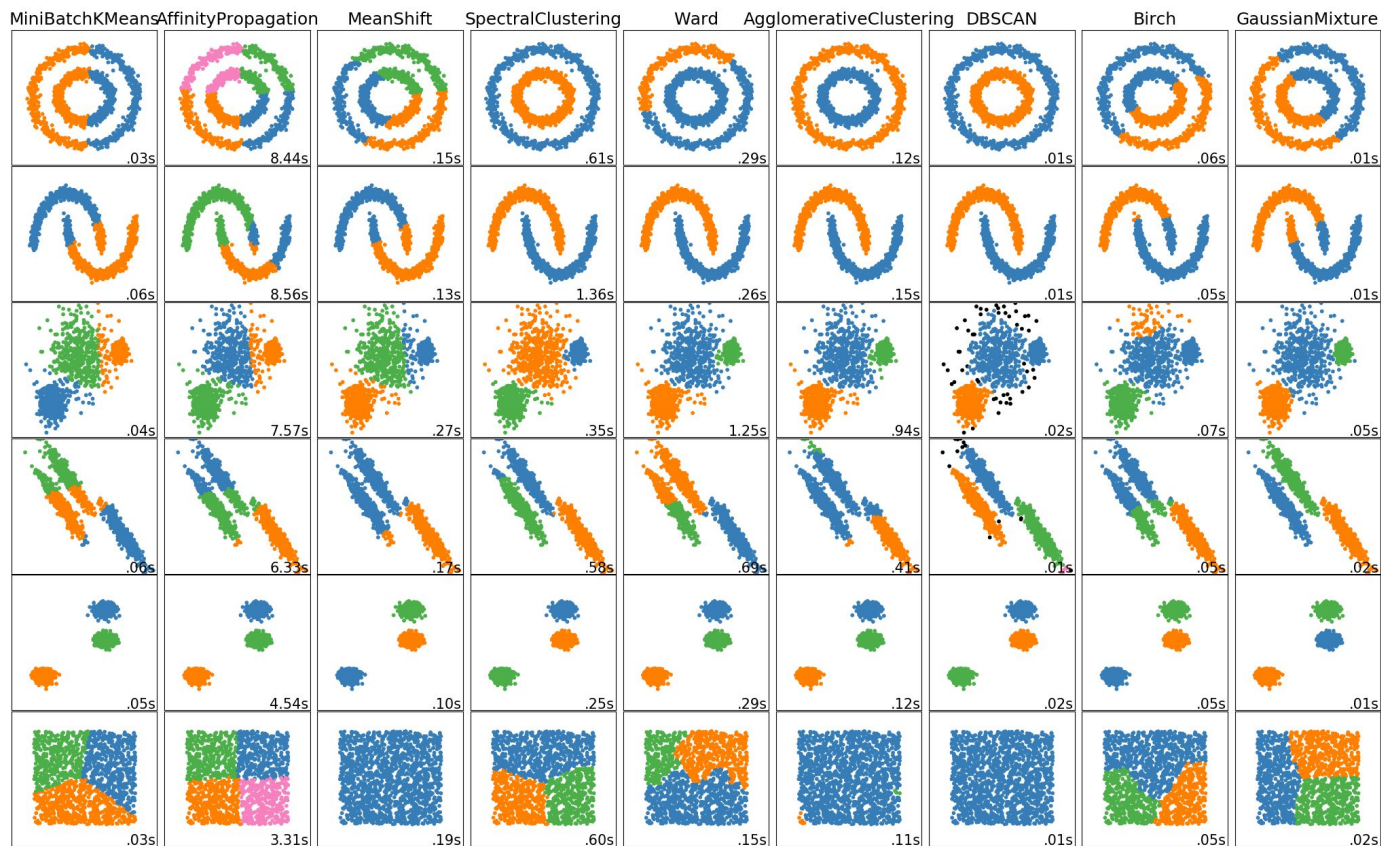


DBSCAN (Density-Based Spatial Clustering of Applications with noise)

- Détecte automatiquement les clusters
- Se base sur la densité
- Et sur un nombre minimum d'éléments
- Détecte des outliers
- Peut trouver des formes de clusters non-circulaires



Autres





Réduction de dimensions

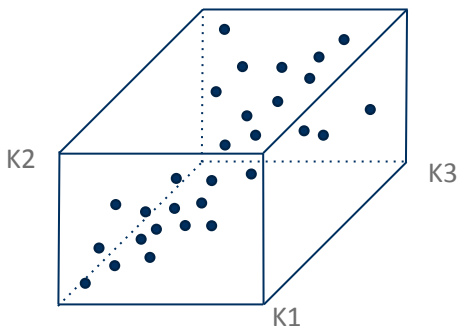
Motivation

x1	x2	x3	...	xn
v11	v12	v13	...	v1n
v21	v22	v23	...	v2n
...
vm1	vm2	vm3	...	vmn

Pas de visualisation possible



k1	k2	k3
u11	u12	u13
u21	u22	u23
...
um1	um2	um3



+

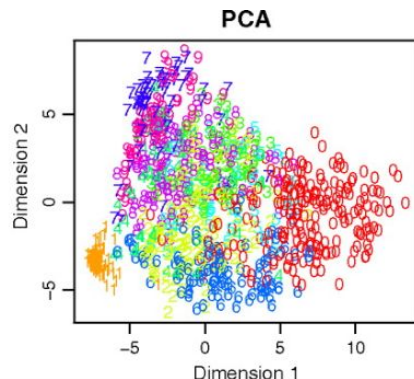
Amélioration du
temps
d'apprentissage

Deux exemples de solutions

PCA

(Principal Component Analysis)

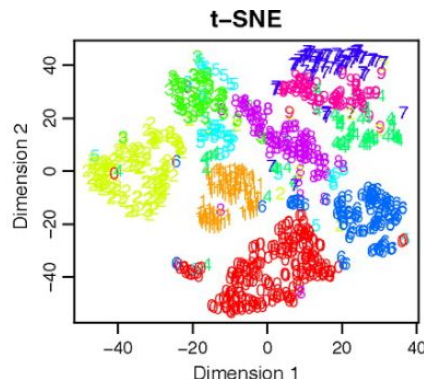
- Procédé d'algèbre linéaire
- Simple et efficace
- Peut être limité



T-SNE

(T-Distributed Stochastic Neighbor Embedding)

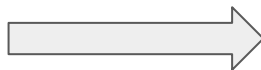
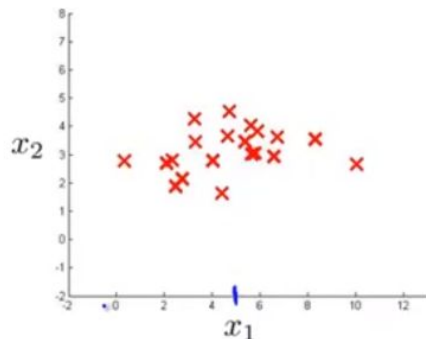
- Utilisation d'embeddings
- Adapté à quelques milliers de lignes
- Lourd



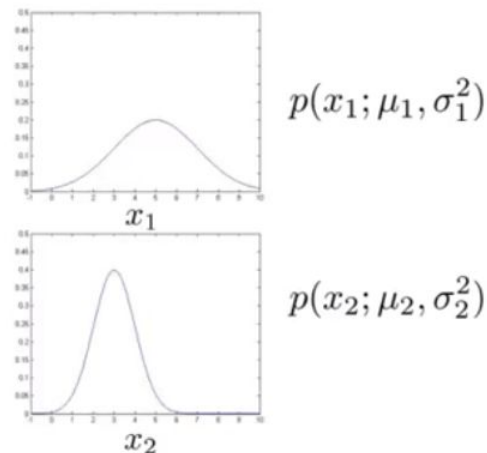


Détection d'anomalie

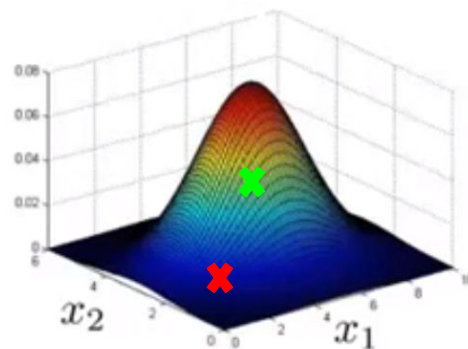
Gaussian Detection



Calcul de la moyenne et de la variance par feature



Produit des probabilités



Principe de l'algorithme:

- Pour chaque feature:
 - On calcule la moyenne et la variance
 - On en déduit une loi normale
- On calcule la probabilité finale en multipliant les lois normales précédemment calculées
- On fixe une valeur de seuil (valeur arbitraire)
- Si la probabilité de notre exemple est inférieure à la valeur de seuil alors on la détecte comme anomalie

Hyperparameters Tuning

Problématique

Exemple :

Algorithme A

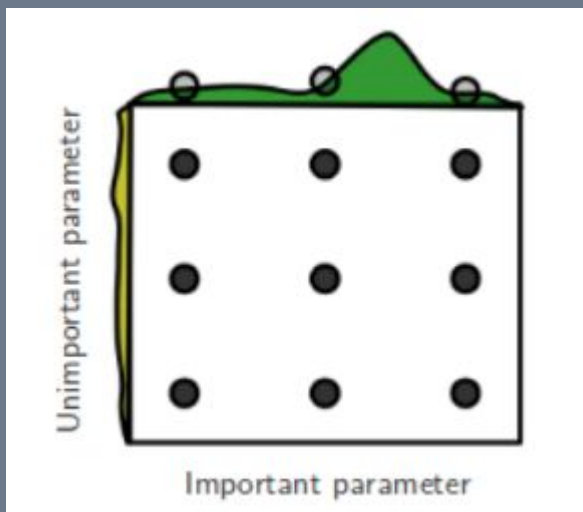
- Hyperparam 1 : 0.0001 -> 1000
- Hyperparam 2 : 2 -> 200
- Hyperparam 3 : 1 -> 4

Objectif : 2 combinaisons

=> Comment choisir les valeurs optimales ?

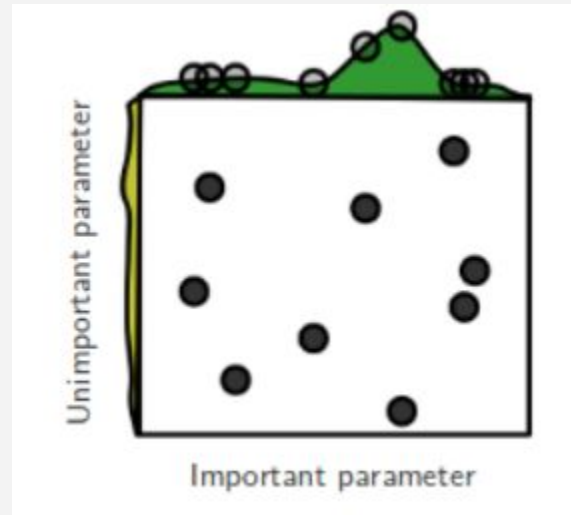
Grid search

- Hyperparam 1 : [0.1, 1, 10]
- Hyperparam 2 : [5, 75, 150]
- Hyperparam 3 : [1, 2, 4]



Random search

- Hyperparam 1 : norm(0.01, 10)
- Hyperparam 2 : randint(5, 200)
- Hyperparam 3 : randint(1, 4)



Tuning

Exemple :

Algorithme A

- Hyperparam 1 : 0.0001 -> 1000
- Hyperparam 2 : 2 -> 200
- Hyperparam 3 : 1 -> 4

Objectif : 2 combinaisons

Grid Search :

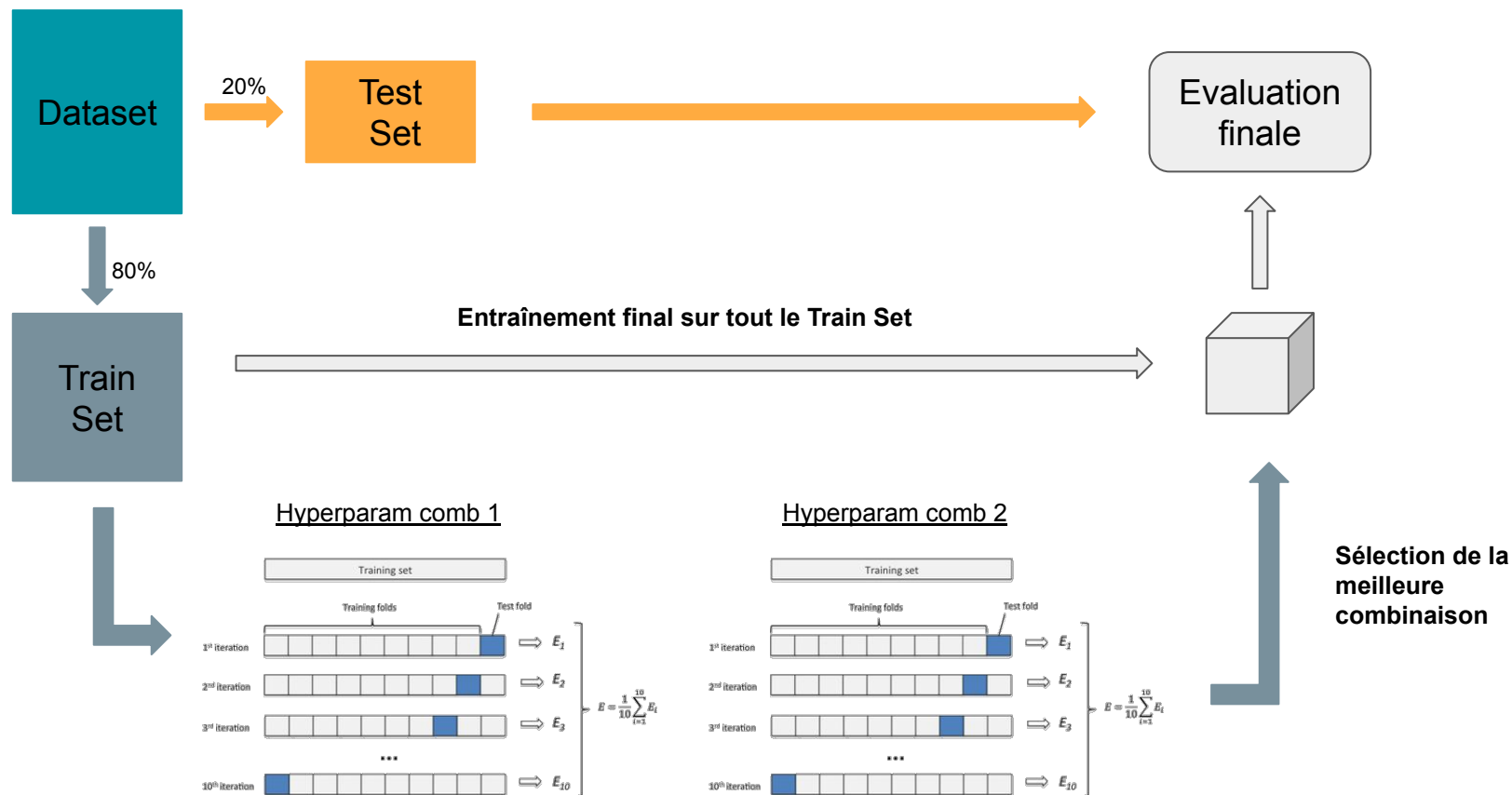
- H1 = 0.1, H2 = 5, H3 = 1
- H1 = 1, H2 = 150, H3 = 4

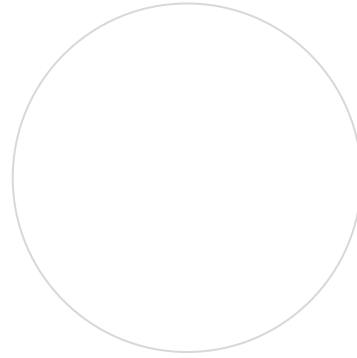
Random Search :

- H1 = 8.5, H2 = 112, H3 = 3
- H1 = 0.5, H2 = 65, H3 = 2

=> Le random Search permet d'explorer plus et est moins biaisé par les hypothèses faites par la personne mettant en place l'algorithme

Tuning + Cross Validation





Fin du chapitre 4