

# **Apprentissage**

## **Automatique Introduction**

**A.Larhlimi**

Apprentissage Automatique -- Introduction -- 1

# **Aujourd'hui**

- Introduction générale:
  - Exemples, définitions, problématiques, approches, vocabulaire
- Deux approches élémentaires à connaître:
  - Modélisation bayésienne
  - k plus proches voisins (kNN)

# « Machine Learning »

- Un domaine scientifique hybride:
  - Statistique
  - Intelligence artificielle
  - « Computer science »
  - Traitement du signal

- Utilisant des techniques généralistes:
  - Optimisation numérique
  - Hardware
  - Gestion de base de données

Apprentissage Automatique -- Introduction -- 3

## **Pourquoi le « Machine Learning »?**

- Thème à la mode: Intelligence Artificielle, « deep learning », « big data »...
- Raison épistémologique

- On ne sait pas modéliser les problèmes complexes ... mais on dispose d'exemples en grand nombre représentant la variété des situations
- « Data driven » vs. « Model Based »
- Raison scientifique
  - L'apprentissage est une faculté essentielle du vivant
- Raison économique
  - La récolte de données est plus facile que le développement d'expertise

## **Domaines techniques utilisant du ML**

- ML comme outil de conception
  - Vision & Reconnaissance des formes
  - Traitement du langage
  - Traitement de la parole
  - Robotique
  - « Data Mining »
  - Recherche dans BDD
  - Recommandations
  - Marketing...
- ML comme outil explicatif
  - Neuroscience
  - Psychologie
  - Sciences cognitives

# Données = carburant du ML

CERN /  
Large Hadron Collider



Google :  
24 PetaOctets/jour

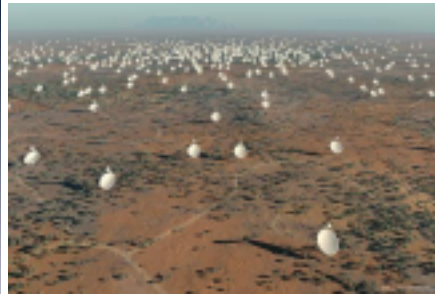
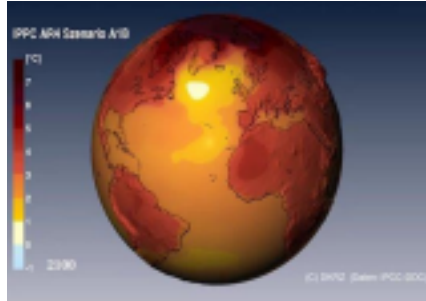
DKRZ (Climat)  
500 Po

Copernicus :  
> 1Po/an



Square Kilometer Array 1376 Po/an (en  
2024)





BIG DATA

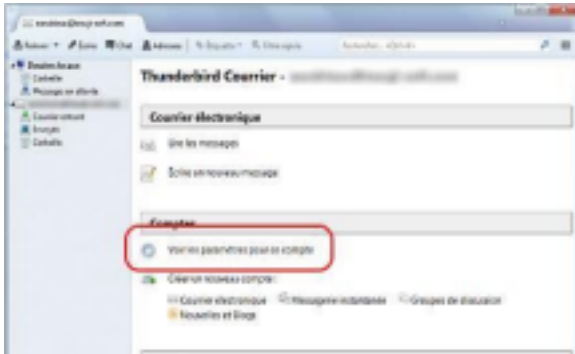
Apprentissage Automatique -- Introduction -- 6

# Apprentissage automatique :

## applications *Anti-Spam (Classifieur Bayesien)*







*1997 : DeepBlue bat Kasparov*

*2017: Alpha GO bat Ke Jie*

*2019: AlphaStar champion de StarCraft*

*Tri postal automatique (détection de chiffres manuscrits par réseaux de neurones)*



# Apprentissage automatique applications



Recommandation ciblée

(régression logistique)



Appareil photo avec détection



de visages (*boosting*)



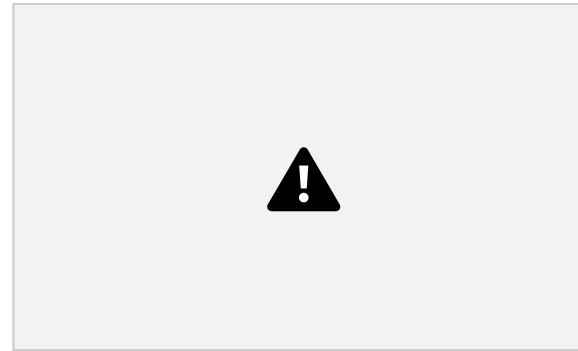
Apprentissage Automatique -- Introduction -- 8

# Apprentissage automatique

Chat Bots  
(*Réseaux de neurones*)

Diagnostic médical  
(*Réseaux de neurones*)





## Traduction multi-lingue (Réseaux de neurones)



# « Deep Learning » : le mot clé

**inévitable**



**Données**

**Algorithmes**

*apprentissage*



**Moyens de calcul**

**Logiciels**

*Une rupture scientifique et  
technologique en*



Apprentissage Automatique -- Introduction -- 11

# MACHINE LEARNING

# Problématique

# générale



# Dans ce cours

L'apprentissage automatique est:

- une démarche de **conception** d'une **fonction de prédiction**
- par une modélisation ou programmation **non explicite** à partir **d'exemples** (signaux, images, texte, mesures...)

# Formalisation

- Donnée à interpréter ( $x$ )
  - Mesures, texte, image, enregistrement, vidéo ou caractéristiques extraites de ...
- Prédiction ( $y$ )
  - Décision, choix, action, réponse, préférence, groupe,

commande, valeur...

- Echantillons ( $\diamond\diamond = \{ \diamond\diamond\blacksquare, \diamond\diamond\blacksquare \}$ )
  - Exemples de données et de (bonnes) prédictions
  - « Base d'apprentissage »:  $\diamond\diamond$


## Formalisation

**Donnée à**

interpréter Prédiction   $\mapsto$  

**Echantillons  
exemples**

**Prédicteur**

- Hypothèse forte: les échantillons contiennent toute l'information exploitable et utile
- Prédicteur = « interpolateur » à partir des données 

## Deux phases

### Prédiction

Donnée

Programme ou Inférence Résultat  
Modèle

### Apprentissage

données  
Programme ou  
Modèle

Estimation Base de

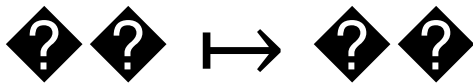
Apprentissage Automatique -- Introduction -- 16

**Deux phases**

**Prédiction**



# Apprentissage



**Caractéristiques du prédicteur:**  
**Paramètres,**  
**Poids,**  
**Prototypes...**

# Formalisation

Fonction paramétrique de prédiction

$$\hat{y} = f(\mathbf{x}; \theta)$$

Apprentissage = trouver le  $\theta$  qui optimise un critère  $J$

$$\theta = \arg \min_{\theta} J(\theta, \theta')$$



A partir d'une base d'apprentissage

$$\diamond ? \diamond ? = \{ \diamond ? \blacksquare, \diamond ? \blacksquare \}$$

## Exemple: Reconnaissance de chiffres

**manuscripts**



- Comment définir les éléments ?

??, ??, ??, ??

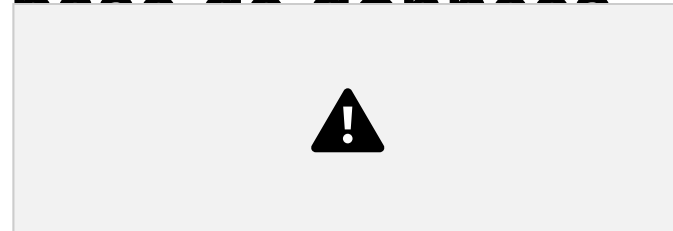
- Les fonctions d'apprentissage et de prédiction?

??  $\mapsto$  ??

??, ??  $\mapsto$  ??

# Etape 1: choix de la base de données

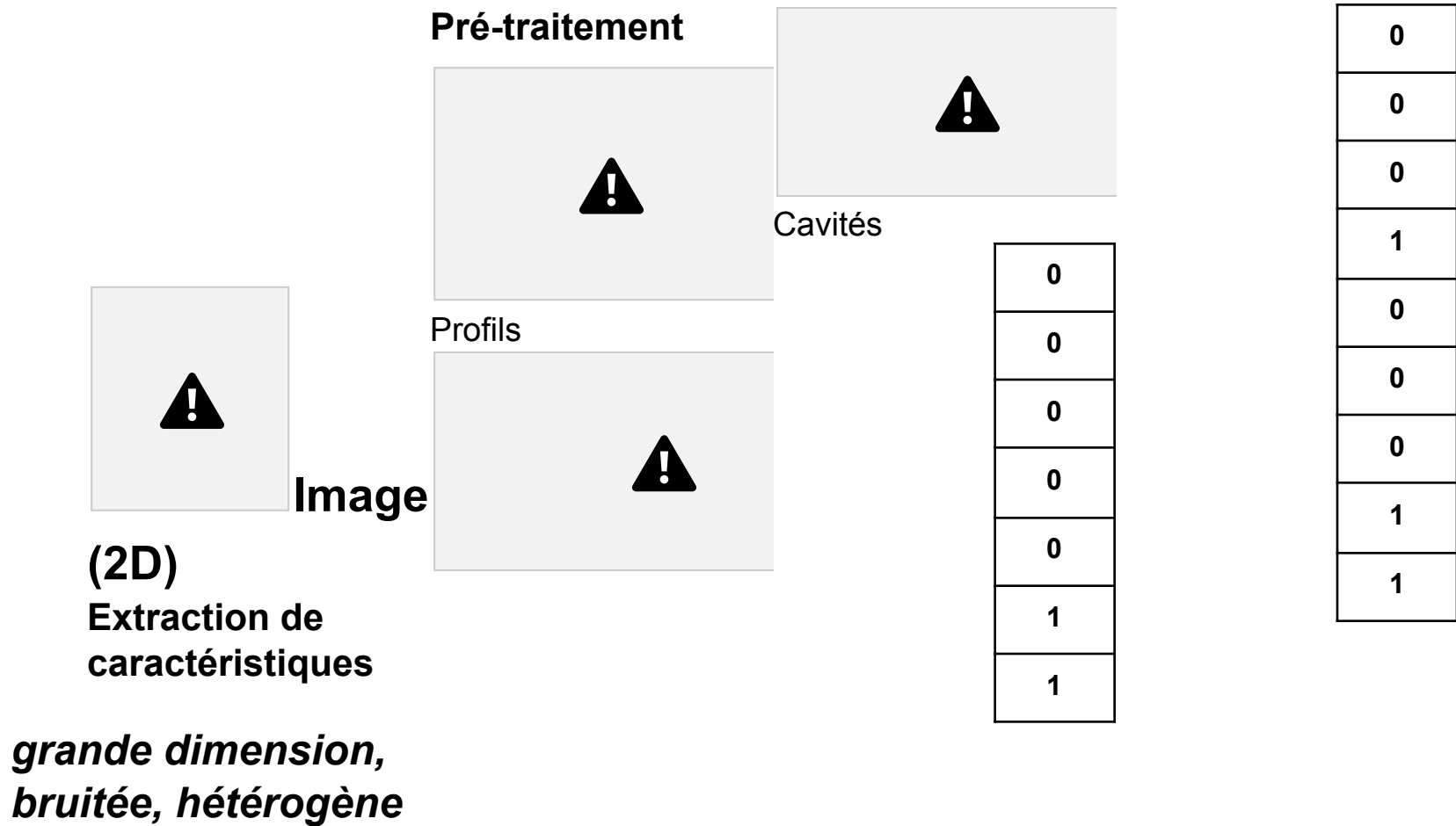
- Elle existe:
  - Scikit-learn:
  - MNIST:
  - SVHN:



- Il faut la construire:

- Recueil de données existantes •
- Expérimentations (photos, mesures...)

## **Etape 2: mise en forme des données**



***petite dimension, homogène, «  
propre »***

Occupation  
de zones

?? = Vecteur

Apprentissage Automatique -- Introduction -- 21

## **Etape 3: choix de l'approche**

- Quel type de fonction et de problème d'apprentissage?
  - Classification
  - On connaît les classes cibles □ Apprentissage supervisé
- Nature des données?
  - Vecteurs de taille fixe mais grands □ algorithmes avec bon contrôle de la régularisation

- Taille de la base de données?
  - Grande (> 10000 exemples) □ optimisation efficace
- Nature fonctionnelle des prédicteurs?
  - Arbres de décision, SVM, Réseaux de neurones...

## Types d'apprentissage

- **Apprentissage supervisé**
  - Les données d'apprentissage contiennent les objectifs de prédiction (annotations)

- **Apprentissage non supervisé**
  - Les données d'apprentissage sont brutes
- **Apprentissage semi-supervisé**
  - Les données d'apprentissage sont partiellement annotées
- **Apprentissage par transfert**
  - Les données d'apprentissage sont proches du problème visé
- **Apprentissage par renforcement**
  - Les prédictions sont issues d'une séquence d'actions et sont caractérisées par une mesure de qualité (« reward »)

## Types de prédictions

- **Classification**



- Binaire: spam / non spam
- Identification: « tata Monique »

- **Régression**

- Prédiction de température, de cours de bourse
- Localisation d'objet dans image
- Commande

- **Structure**

- Graphe des articulations d'une personne

- **Regroupement**

- Photos dans base de données personnelle

- **Texte**

- « C'est un chat qui saute sur une table. »

## Nature fonctionnelle du prédicteur

- Dépend de la forme des données (vecteurs, listes, réels/discret) et du type de prédiction
- Exemples
  - Plus proches voisins
  - Machines à vecteurs de supports (SVM)
  - Arbre de décision
  - Ensembles de classifieurs (forêts aléatoires, « boosting »...)
  - Réseaux de neurones
  - Règles/Programmation logique
  - Modèles probabilistes (Réseaux bayésiens, Chaînes ou champs de Markov...)

- Etc.

## **Etape 4: optimisation**

Apprentissage =

- définir un espace fonctionnel et un critère paramétrique (coût, énergie...)
- appliquer un optimiseur et régler ses paramètres •
- vérifier que l'apprentissage se passe bien • évaluation de la capacité de généralisation • convergence



# Optimisation

- **Optimisation convexe**
  - Ex. Minimisation séquentielle de problème quadratique
- **Optimisation stochastique**
  - Ex. Descente de gradient stochastique, Algorithmes génétiques

- **Optimisation sous contraintes**
  - Ex. Programmation linéaire
- **Optimisation combinatoire**
  - Ex. Algorithmes gloutons

## **Etape 5: évaluation**



# Métriques d'évaluation

- Dépend du type de prédiction
- Classification
  - Taux d'erreur moyen
  - Matrice de confusion
  - Précision/rappel
  - Courbe ROC
- Régression
  - Erreur quadratique
- Détection
  - Taux de recouvrement moyen

## **Résumé des étapes de conception**

1. Constituer des bases de données
2. Préparer les données: Analyser, visualiser, prétraiter, transformer, extraire
3. Concevoir le modèle (type de prédicteur, principe d'apprentissage)
4. Définir un critère et Optimiser (l'apprentissage)



proprement dit)

## 5. Evaluer

# EXTRACTION DE CARACTÉRISTIQUES

# Travailler avec des données

Deux activités complémentaires:

## Préparer les données

- Etape coûteuse mais indispensable
- Objectif: rendre possible l'apprentissage avec des données: •

Propres, homogènes, recalées, calibrées, organisées, facilement accessibles, renseignées...

- « Data engineering » (un nouveau métier!)

### Transformer les données

- Objectif: Extraire l'information des données, leurs caractéristiques (« features »), construire leur forme

## **Extraction de caractéristiques**

- « Feature extraction » en anglais
- Données brutes pas exploitables directement:

- Bruitées
- Grandes dimensions (image, enregistrement)
- Information utile noyée
- Etape critique de la « reconnaissance des formes » • Caractéristiques trop simples: pas assez d'information, confusion • Caractéristiques trop riches: complexité, bruit, grande variabilité □ Compromis difficile à régler entre expressivité, invariance, robustesse, taille, coût de calcul...
- Deux cas de figure:
  - On sait ce qui est important et pourquoi (expertise « métier »)
    - modélisation
  - On ne sait pas décrire ce qui est important
    - on l'apprend!

$$\{x_L, \dots, x_N\} \in \mathcal{X}$$

# Chaîne de prédiction générique = $F_y$

$\mathbf{x})($

capteur

Prétraitement  
extraction de  
caractéristiques

Algorithme de  
prédiction

mesures, images, texte...

$N$   
 $\in RI$

sources

$\{ \}_K$

$=_{1,...SSS}$

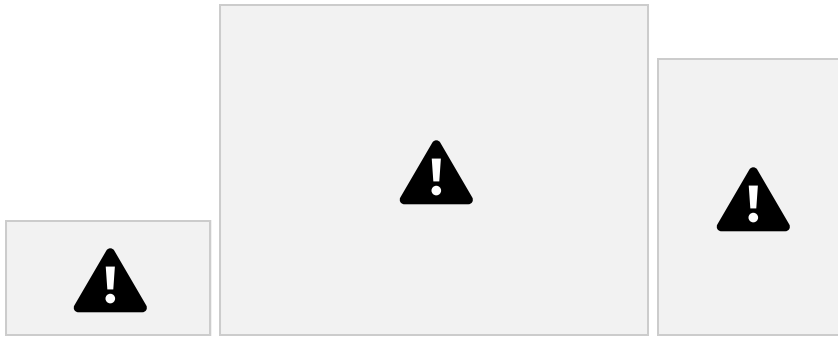
valeur/décision/action



caractéristiques

$d \mathbf{x} \in R$

*C'est un 3*



$\{ \neg \} a a A y_L, \dots = \in \_1$

## 1 – Construire une représentation

**(forme)** =  $Fy \mathbf{x}$ )(

capteur  
Prétraitement

extraction de  
caractéristiques

Algorithme de  
prédiction

mesures, images

$N$   
 $\in RI$

prédiction/décision/action

sources

$\{ \}_{K}$   
 $=_{1, \dots, S} S$

caractéristiques

$d \mathbf{x} \in R$

***Extraction d'information « métier »***

capteur

## **2 - Prédire**

Prétraitement  
extraction de  
caractéristiques



$$\{ \neg \} a a A y_L, \dots = \in_1$$

$$= F y \mathbf{x})($$

Algorithme  
de  
prédiction

mesures, images

$$N$$

$$\in RI$$

sources

$$\{ \} _K$$

$$= _1, \dots s s s$$

prédiction/décision/action

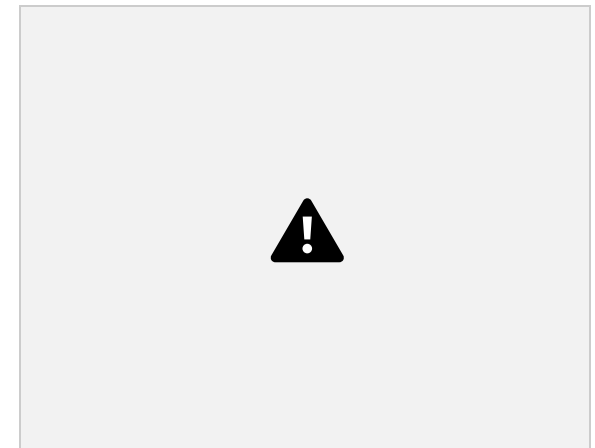
caractéristiques

$$^d \mathbf{x} \in R$$

***Etape « machine learning »***

# Exemples de caractéristiques en image

- Deux grandes classes: forme ou texture
- Forme
  - Dépend d'une étape de séparation du fond (segmentation, saillance)
- Caractéristiques structurelles ou
- Texture
  - Globales et/ou locales
  - Plus difficiles à associer à un objet précis





Apprentissage Automatique -- Introduction -- 37

$\{ \neg \} a a A y_L, \dots = \in_1$

**Prédire & extraire en même temps =**

$Fy \mathbf{x})($

capteur  
Prétraitement

extraction de  
caractéristiques

Algorithme de  
prédiction

mesures, images

$N$   
 $\in \mathcal{R}^I$

sources

$\{\}_{K}$

$=_{1,...,S} S$

prédiction/décision/action

caractéristiques

$\mathbf{x} \in \mathcal{R}^d$

***Etape « machine learning »  
moderne***

**« Deep features »**



On peut apprendre les caractéristiques image  
Réseaux convolutifs (cours DL)

# **UN CONCEPT CENTRAL: LA GÉNÉRALISATION**

**Exemple : régression polynomiale**

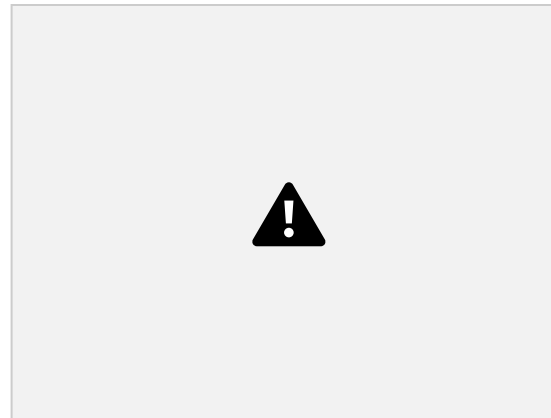
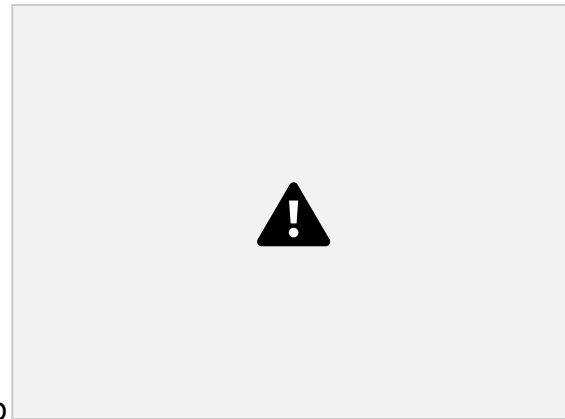
carré entre les points vrais  
et le polynôme estimé.

- La courbe verte est la véritable fonction à estimer (non polynomiale)

- Les données sont uniformément échantillonnées en  $\diamond \diamond$  mais bruitées en  $\diamond \diamond$ .

- L'erreur de régression est mesurée par la distance au

from Bishop



$${}_{110}{}_{22}xwxwxwxwx Fy \mathbf{w}, \varphi \varphi +++++ = \varphi_{MM} ) ( ) ( \dots$$

$$) ($$

- Prédiction utilise des **fonctions de base** encodant les données source (« features »):  $\diamond \diamond \diamond \diamond = \diamond \diamond \diamond \diamond$
- **Apprentissage** = Maximum de Vraisemblance:



- Que vaut cet apprentissage?

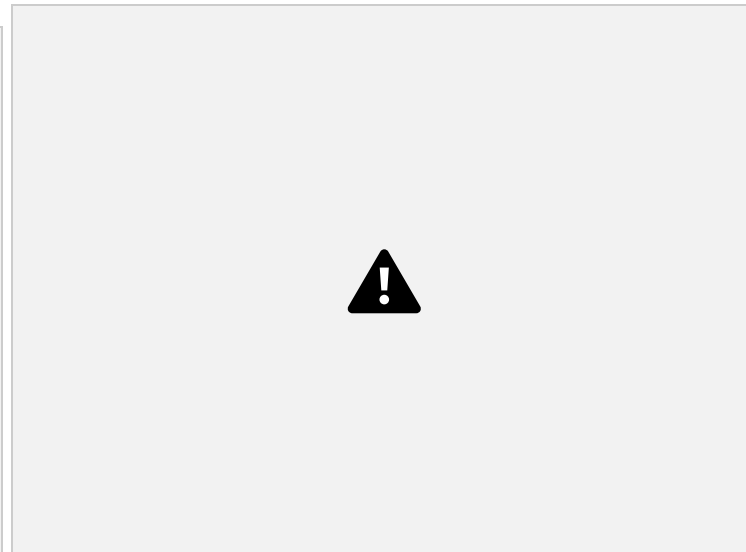
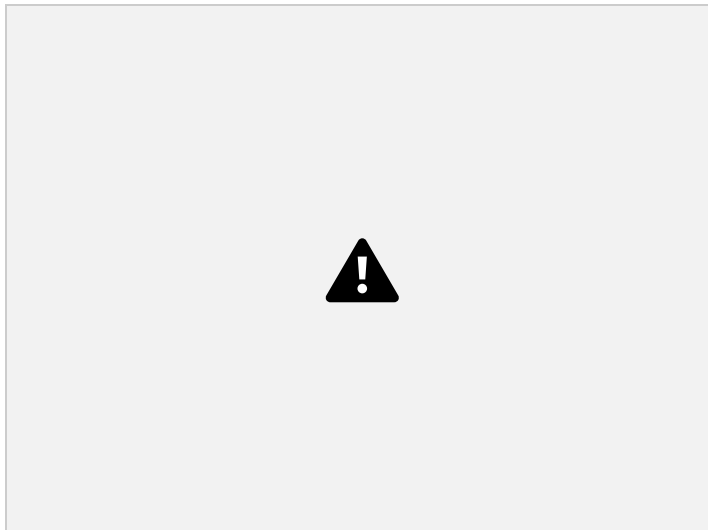
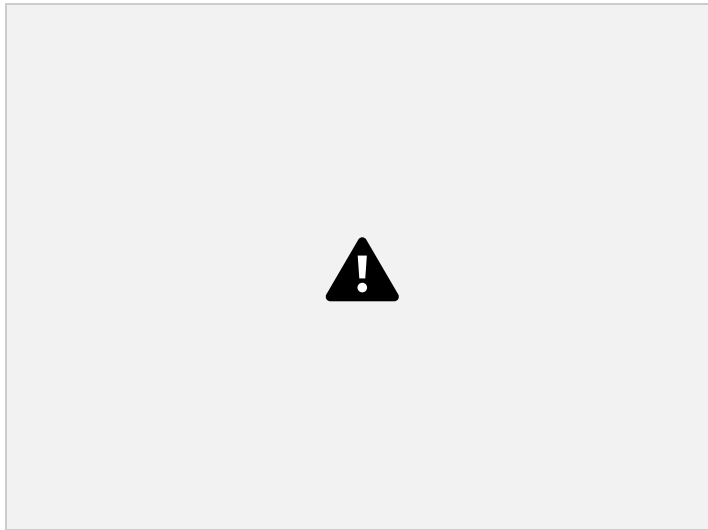


# Régression polynomiale: $\diamond \diamond \diamond \diamond =$



from Bishop





# “Training vs. Test”

Erreur de régression est calculée sur des données  $\mathcal{D} = \{ \mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N \}$ :

$$E_{\text{RMS}} = \sqrt{\frac{1}{N} \sum_{i=1}^N (y_i - \hat{y}_i)^2}$$

Mais quelles données?

- Données d'apprentissage (Training):
  - c'est un moyen de modélisation

- Données opérationnelles (Test):
  - c'est la situation réelle
  - Celles pour lesquelles on veut une bonne prédiction

# Comportement des erreurs

$$E_{\text{ww}} = \frac{1}{N} \sum_{n=1}^N (y_n - \hat{y}_n)^2$$

$i$

$i i_1$



Prédicteur  
pas assez précis  
Bon régime

Sur apprentissage = « Overfitting »

# Erreur de généralisation

- La mesure de bon fonctionnement est l'erreur sur des données nouvelles  
généraliser  $\neq$  mémoriser (par cœur)
- Problème: les données nouvelles sont par nature inconnues! (sinon, elles seraient utilisées)
- Il est nécessaire de faire des hypothèses sur leur nature et sur le modèle de prédiction.

Deux phénomènes à contrôler (éviter)

- **Simplisme:** modélisation trop grossière pour rendre compte de la variété des données
  - Erreur d'apprentissage et de test importantes
- **Sur-apprentissage (« Overfitting »):** modèle trop complexe se spécialisant sur les données d'apprentissage

- Ecart entre erreur d'apprentissage et erreur de test

## Construire son chantier d'apprentissage

- Préparer les données
  - Simplifier/compléter/formater/homogénéiser/calibrer...
- Diviser en deux ensembles:
  - **Apprentissage** (“Train”) pour optimiser les paramètres du modèle.
  - **Test** pour estimer la qualité de l'apprentissage dans son contexte d'utilisation, i.e. **l'erreur de généralisation**.
- L'ensemble de test n'est jamais utilisé pour l'apprentissage (optimisation), seulement pour son

évaluation.

Apprentissage Automatique -- Introduction -- 47

# DEUX APPROCHES ÉLÉMENTAIRES

Modélisation bayésienne

Plus proches voisins



# Théorie Bayésienne de la décision

- On considère les données  $\mathbf{x}$ ,  $\mathbf{y}$  comme des variables aléatoires.
- On les modélise par des lois de probabilités:
  - $p(\mathbf{x})$ ,  $p(\mathbf{y})$  : lois a priori (ou marginales)
  - $p(\mathbf{x}, \mathbf{y})$  : loi jointe
  - $p(\mathbf{x} | \mathbf{y})$  : vraisemblance conditionnelle
  - $p(\mathbf{y} | \mathbf{x})$  : loi a posteriori
- Classification:  $\mathbf{y} \in 1, 2 \dots \mathbf{y}$  est une étiquette • On

cherche à prédire une unique étiquette  $y^*$  à partir de  $x$

$$x \mapsto y^*$$

- Théorie de la décision démontre que le meilleur choix est:

$$y^* = \arg \max$$

$$y \mid x$$

## Théorie Bayésienne de la décision

- Deux questions:
  - Comment calculer  $y \mid x$  = apprentissage
  - Comment trouver le max = prédiction
- « Astuce »: utiliser la loi de Bayes

$$x_1, x_2, \dots, x_n \mid y = \frac{p(x_1, x_2, \dots, x_n, y)}{p(y)}$$

- On connaît en général la fréquence d'occurrence des classes  $y$
- On sait plus facilement calculer la **vraisemblance**:  $p(y \mid x_1, x_2, \dots, x_n)$   
 $\mid p(y)$  • « Si je sais dans quelle classe je suis, je sais décrire le comportement/distribution de mes données »
- Le max sur  $y$  ne dépend que de  $p(x_1, x_2, \dots, x_n \mid y)$

$$\text{et } y^* = \arg \max_y p(y \mid x_1, x_2, \dots, x_n)$$

$$p(x_1, x_2, \dots, x_n \mid y) \mid p(y)$$

## Approche Bayésienne multivariée

- Calcul de la loi conditionnelle: Modèle multivarié
- Par ex. modèle gaussien décrivant  $\begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix} + \begin{bmatrix} \sigma_{11} & \sigma_{12} \\ \sigma_{21} & \sigma_{22} \end{bmatrix} \begin{bmatrix} \epsilon_1 \\ \epsilon_2 \end{bmatrix} \dots$

$$\begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \in \mathbb{R}^2: P \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}; \begin{bmatrix} \sigma_{11} & \sigma_{12} \\ \sigma_{21} & \sigma_{22} \end{bmatrix} = 1$$

$$(2\pi)^{-2/2} \exp - \frac{1}{2} (\begin{bmatrix} x_1 \\ x_2 \end{bmatrix} - \begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix})' \begin{bmatrix} \sigma_{11} & \sigma_{12} \\ \sigma_{21} & \sigma_{22} \end{bmatrix}^{-1} (\begin{bmatrix} x_1 \\ x_2 \end{bmatrix} - \begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix})$$

- Permet de décrire les corrélations entre dimensions.
- Mais demande de connaître la forme des distributions + limitation à petites dimensions.
- Si modélisation gaussienne et deux classes, la prédiction se réduit à calculer une fonction de degré 2

$$\begin{bmatrix} 1 \\ 1 \end{bmatrix} = \Sigma \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$

$$\begin{bmatrix} 1 \\ 1 \end{bmatrix} = \Sigma \begin{bmatrix} R & R \\ R & R \end{bmatrix} \begin{bmatrix} 1 \\ 1 \end{bmatrix} = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$$

$$\begin{bmatrix} 1 \\ 1 \end{bmatrix} = \Sigma \begin{bmatrix} 90 & 01 \\ 1 & 10 \\ 09 & 54 \\ 45 & -1 \end{bmatrix} \begin{bmatrix} 1 \\ 1 \end{bmatrix}$$

# Approche gaussienne multivariée

$$P(\mathbf{x}; \boldsymbol{\mu}, \boldsymbol{\Sigma})$$

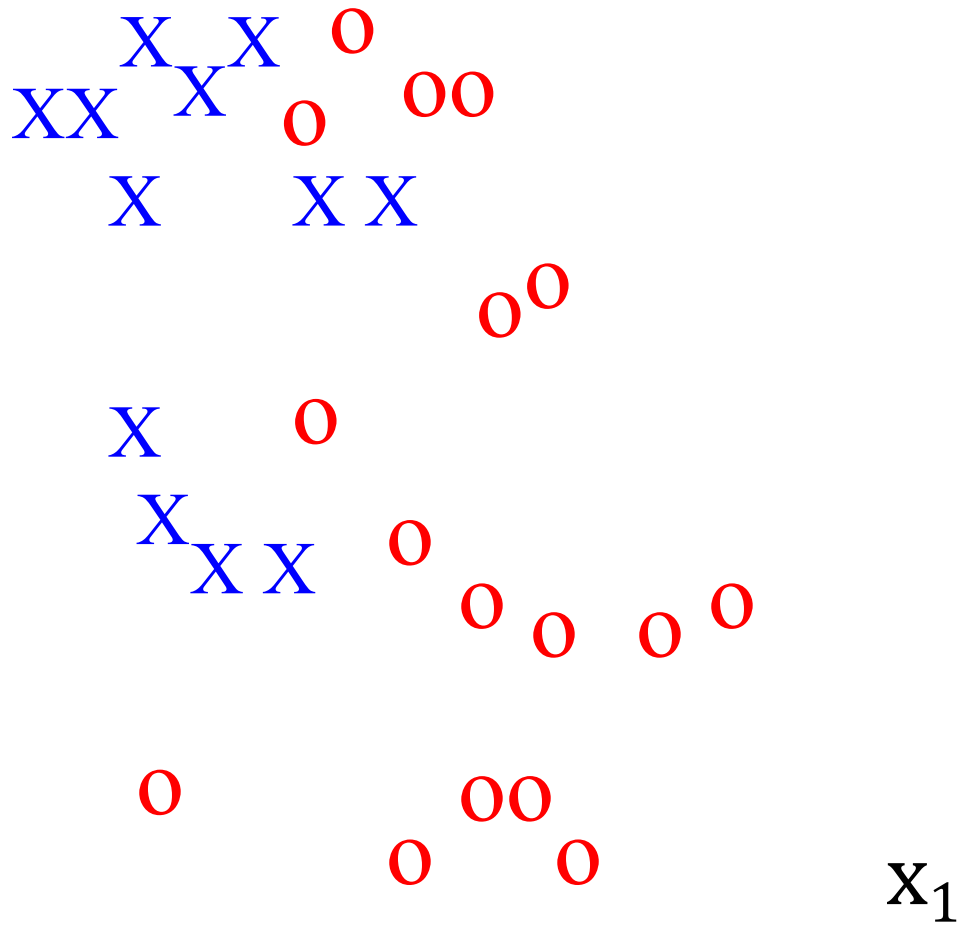
$$\mathbf{x}_2$$

$$\mathbf{x}$$

$$\mathbf{x}$$

$$\mathbf{x}$$

$$P(\mathbf{x}; \boldsymbol{\mu}, \boldsymbol{\Sigma})$$



*Séparatrice = Forme quadratique*

$$\begin{aligned} & (\diamond\diamond - \diamond\diamond_1)' \Sigma_1^{-1} (\diamond\diamond - \diamond\diamond_1) - \diamond\diamond - \diamond\diamond \Big| \Sigma_2^{-1} \diamond\diamond - \diamond\diamond \Big| \geq \\ & \diamond\diamond\diamond\diamond\diamond\diamond\diamond \end{aligned}$$

# Approche Bayésienne Naïve

- Calcul de la loi conditionnelle: hypothèse d'indépendance.

[illegible]

- On calcule la vraisemblance globale dimension par dimension □ Problème 1D, modèles plus faciles à estimer (gaussien, binomial, histogrammes, mélange de gaussiennes...)
- Permet de traiter des problèmes de plus grande dimension

- En pratique, on calcule plutôt la log-vraisemblance pour des questions de stabilité numérique

$$\log \frac{P(\mathbf{x}|\mathbf{y})}{P(\mathbf{y})} = \log \frac{P(\mathbf{x}, \mathbf{y})}{P(\mathbf{y})}$$

$$\mathbf{x}^* = \arg \max$$

$$\log P(\mathbf{x}|\mathbf{y}) = \log P(\mathbf{x}, \mathbf{y}) - \log P(\mathbf{y})$$

## Approche bayésienne naïve

$x_2$   
X  
X  
X X

X X X  
X X 0 0 0



X

XX

00

X

0

$$X_X$$

0

X

# O

# O

O

$$\mathbf{X}_1$$

00

# O

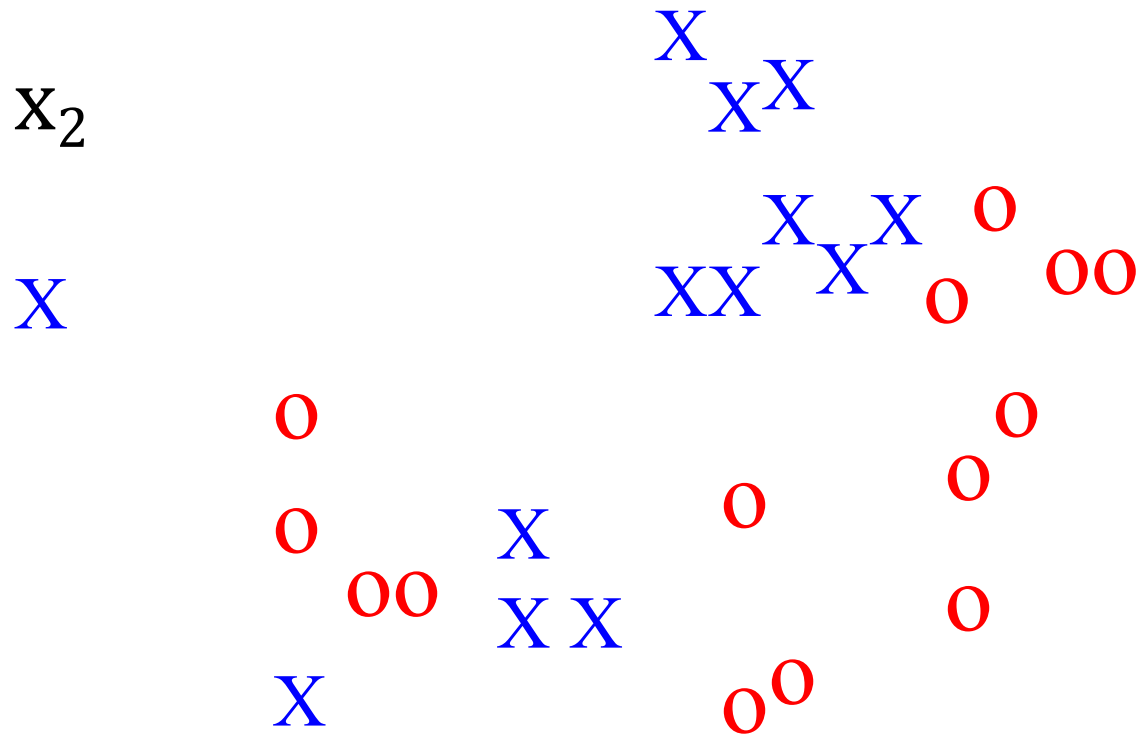
O

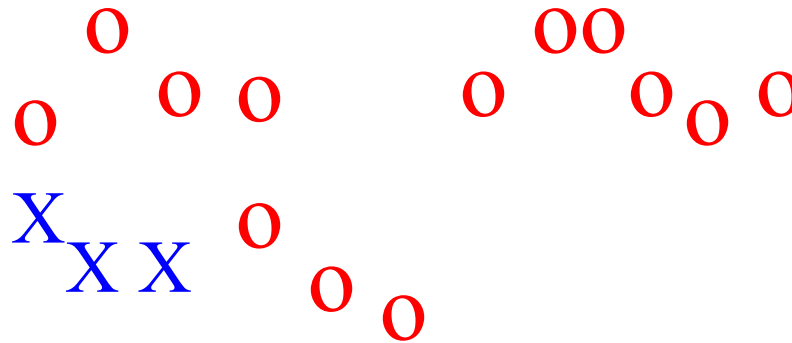
# O

O

X

# Approche bayésienne naïve vs. multivariée





$x^1$   $x$

# Approche bayésienne naïve vs. multivariée

$x_2$

$x_x$   
 $x$

$oo$

$x$

$oo$

$o$

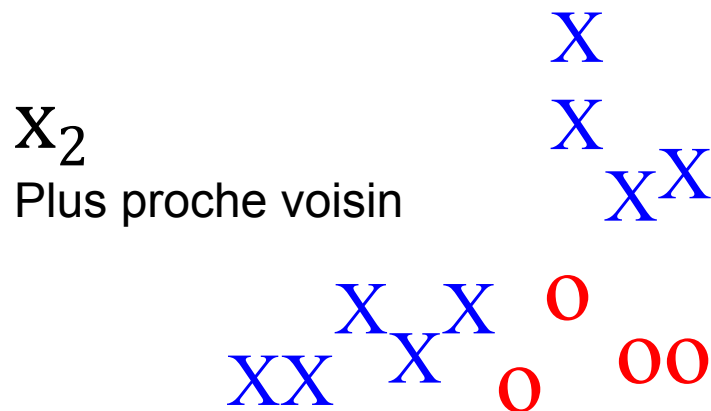
$x^1_x$   
 $x$

## Approche bayésienne: résumé

- Théorie probabiliste de la décision  $\square$  calcul de la loi a posteriori
- Expression de la loi a posteriori:
  - Hypothèse d'indépendance conditionnelle.
  - Modèle gaussien multivarié
- Apprentissage
  - Estimation de lois paramétriques simples
- Prédiction
  - Calcul de log-vraisemblance et max sur hypothèses

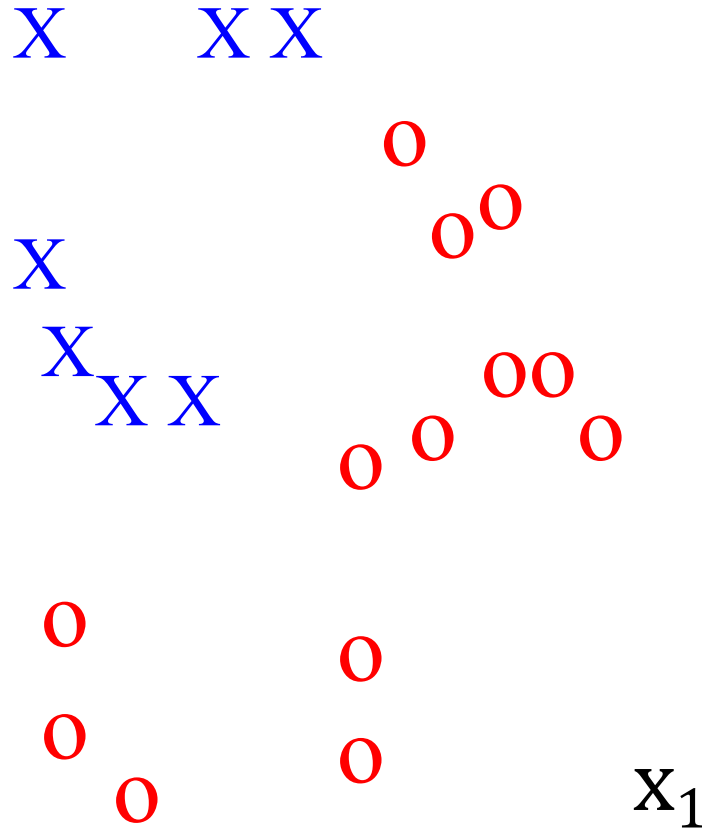
- Quand l'utiliser? (limitations)
  - Petits problèmes bien modélisés (gaussien multivarié) •
 Caractéristiques non corrélées (bayésien naïf, mais ça peut aussi marcher si c'est corrélé)

## Classification ppv



Prédiction = Classe **X**  
À classer

+



# Classification ppv

$X_2$

À classer

X

X  
X X

Prédiction = Classe o

X X X X  
X X X  
X

O O O O

X X

voisin

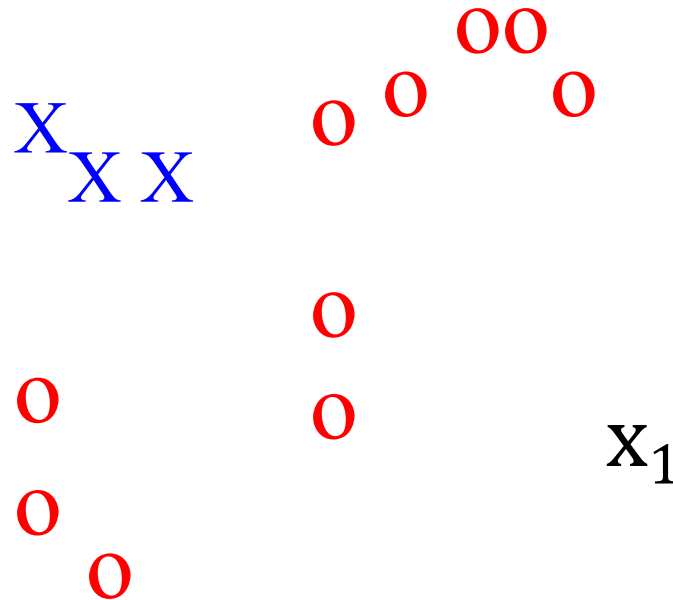
X

+

O O O

Plus proche





## Plus proche(s) voisin(s)

- Principe:
  - Deux échantillons **proches** dans l'espace de représentation ont les mêmes prédictions

- Pour prédire, il suffit de trouver l'exemple annoté **le plus** proche, et d'associer son annotation (étiquette, valeur...)
- Que veut dire « proche »?
  - Nécessite la définition d'une métrique ou mesure de similarité  $d(x, x')$
  - Plusieurs métriques possibles: distance euclidienne (L2), city-block (L1), Minkowski, Mahalanobis...
  - On peut aussi « apprendre » la métrique ou mesure de similarité
- Que veut dire « le plus proche »?
  - Base d'échantillons annotés  $\mathcal{L} = \{x_1, x_2, x_3, x_4, \dots, x_n\}$
  - Recherche de l'échantillon le plus proche:  $x^* = \arg \min_{x \in \mathcal{L}} d(x, x')$
  - Assigne comme prédiction l'annotation du plus proche:  $y^* = y_{x^*}$

Valeur estimée

$$y = R \mathbf{x}$$

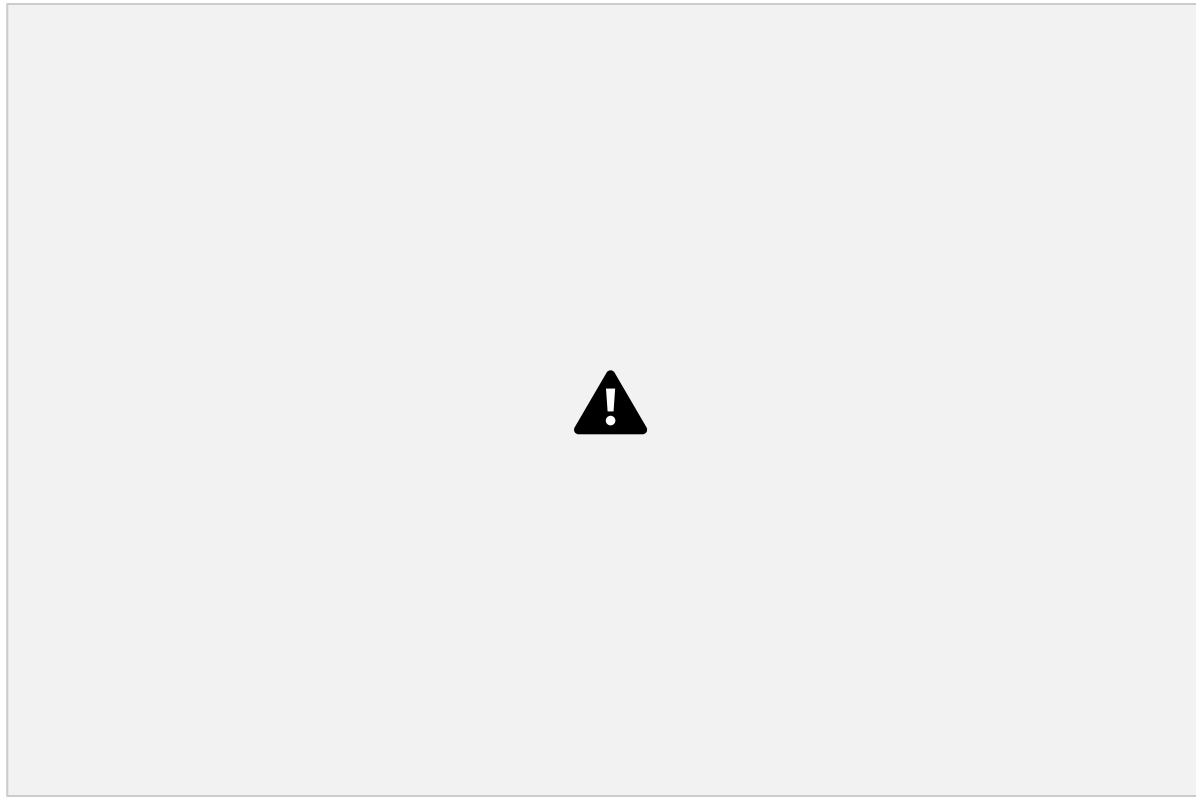
**x**

À estimer Plus proche voisin

# Fonction de classification

Données bruitées □ Régions isolées □ mauvaise régularité des prédictions





Chaque échantillon  
définit une région homogène de l'espace de représentation

## **k-plus proches voisins (« k-NN »)**

- Principe: décision à partir de plusieurs exemples de la base de

données d'apprentissage

- On ordonne les échantillons d'apprentissage en fonction de leur distance à la donnée à classer:

$$\diamond\diamond \blacksquare_1 \leq \diamond\diamond \blacksquare_2 \leq \dots \leq \diamond\diamond \blacksquare_{\diamond\diamond}$$

- On choisit les  $\diamond\diamond$  plus proches
- On prédit en choisissant la classe recueillant le plus de votes

$\diamond\diamond$

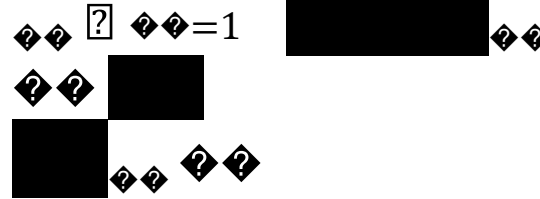
$$\diamond\diamond^* = \arg \min$$

$$\begin{aligned} & \diamond\diamond \boxed{?} \\ & \diamond\diamond \blacksquare_{\diamond\diamond} \\ & = 1 \\ & \diamond\diamond \end{aligned}$$

Où  $\diamond\diamond$  est la fonction de Kronecker (elle vaut 1 si égal, 0 sinon) • Si pas de max (ambiguïté sur la prédiction) on ne décide pas! • On peut aussi pondérer les votes:

$\diamond\diamond$

$$\diamond \diamond^* = \arg \min$$



# Fonction de classification 5 ppv Données

bruitées □ Régions isolées □ mauvaise régularité des prédictions









Chaque échantillon définit une région homogène de l'espace de représentation

Bornes statistiques asymptotiques ( $n \rightarrow \infty$ )

$$e_{\text{Bayes}} \leq e_{k\text{-NN}} \leq \frac{2e_{\text{Bayes}}}{k-1}$$

Où  $e_{\text{Bayes}}$  est l'erreur théorique optimale (Bayes),  $k$  est le nombre de classes et  $e_{k\text{-NN}}$  est l'erreur des k-ppv.

« L'erreur du k-NN est au plus deux fois moins bonne que l'erreur minimale théorique. »

## Coût de la prédiction du k-ppv

- Calcul de la prédiction dépend pour chaque exemple  $x_i$  d'un calcul + tri par rapport aux  $k$  exemples de la base:

$$x_i \leq x_{(1)} \leq x_{(2)} \leq \dots \leq x_{(k)} \leq x_{(k+1)}$$

- Pour N et d grands, coût important de la recherche exhaustive  $O(N \cdot d)$ . Il existe:

- Des algorithmes efficaces de recherche pour problèmes de tailles moyennes (KDtree)

J. Friedman, J. L. Bentley, and R. A. Finkel, "An algorithm for finding best matches in

logarithmic expected time,” *ACM Transaction on Mathematical Software*, vol. 3, no. 3, pp. 209–226, 1977.

- Des algorithmes d’approximation pour les grandes bases ( $>10^6$ ).

Jegou, H., Douze, M., & Schmid, C. (2011). Product quantization for nearest neighbor search. *IEEE transactions on pattern analysis and machine intelligence*, 33(1), 117-128. • Autre manière: pré-calculer les surfaces de séparation entre classes. La complexité de prédiction est alors liée à la complexité de la surface et/ou de son approximation. On verra comment d’autres approches permettent de l’estimer directement.

## La malédiction des grandes dimensions

- Lorsque la dimension  $d$  de l’espace de représentation augmente, les points sont tous aussi proches ou aussi loin.
- On peut montrer, pour une distribution quelconque de  $d$  points tirés de manière indépendante dans  $[0,1]^d$ , que:

$$\lim_{n \rightarrow \infty} \frac{\sum_{i=1}^n x_i - \sum_{i=1}^n y_i}{n} = 0$$

- Ce n'est plus vrai si les distributions sont structurées...heureusement!
- On peut interpréter les techniques de Machine Learning comme des moyens de repérer les bonnes corrélations entre données.
- Conséquence pour les approches « plus proches voisins »: •
  - Ca ne marche que pour les faibles dimensions
  - Ou il faut **réduire** les dimensions de représentation avant de calculer les distances □ apprentissage non supervisé

- Avantages

- Schéma flexible, facile à mettre en œuvre, dépendant de la définition d'une similarité entre données.
- Bonnes propriétés statistiques ( $n \rightarrow \infty$ )

- Mais...

- Temps de calcul prohibitif pour grandes bases
  - Algorithmes efficaces de recherche optimaux ou sous-optimaux
- Régularité dépend des données, pas de l'apprentissage
  - Le k-PPV (« kNN ») pour lisser et réduire le bruit
- Malédiction des grandes dimensions (« Curse of dimensionality »)
  - Réduire la dimension de représentation



## « Plus proches voisins »: résumé

- Hypothèse de régularité = Si observations proches, même comportement
- Deux questions:
  - Que veut dire « proche »?
  - Comment trouver les plus proches?
- Apprentissage
  - Aucun
- Prédiction
  - Tri des distances aux échantillons + vote

- Quand l'utiliser? (limitations)
  - Efficace sur petits problèmes (dimensions & nombre d'exemples)
  - Pb du « curse of dimensionality » + temps de calcul
  - Disposer d'une mesure de similarité adaptée aux données

## A retenir

- « Programmer à partir des données »
  - Deux phases: apprentissage et prédiction
  - Plusieurs variétés de prédicteurs et d'apprentissage
- Démarche générique:
  - Constitution d'une base d'apprentissage
  - Analyse préliminaire des données + préparation
  - Conception du modèle

- Optimisation
- Evaluation
- Objectif principal: minimiser l'erreur de généralisation •  
Train vs. Test
- Deux approches élémentaires:
  - Modélisation bayésienne
  - Plus proches voisins

## Références

- K. Fukunaga, Introduction to Statistical Pattern Recognition (Second Edition), Academic Press, New York, 1990.
- P.A. Devijver and J. Kittler, Pattern Recognition, a Statistical Approach, Prentice Hall, Englewood Cliffs, 1982)
- R.O. Duda and P.E. Hart, Pattern classification and scene analysis, John Wiley & Sons, New York, 1973.
- L. Breiman, J.H. Friedman, R.A. Olshen, and C.J. Stone, Classification and regression trees,

Wadsworth, 1984.

- S. Haykin, Neural Networks, a Comprehensive Foundation. (Macmillan, New York, NY., 1994) • L. Devroye, L. Györfi and G. Lugosi, A Probabilistic Theory of Pattern Recognition, (Springer Verlag 1996)
- V. N. Vapnik, The nature of statistical learning theory (Springer-Verlag, 1995) • C. Bishop, Pattern Recognition and Machine Learning, (<https://www.microsoft.com/en-us/research/uploads/prod/2006/01/Bishop-Pattern-Recognition-and-Machine-Learning-2006.pdf>). • Jerome H. Friedman, Robert Tibshirani et Trevor Hastie, The Elements of Statistical Learning: Data Mining, Inference, and Prediction (<https://web.stanford.edu/~hastie/ElemStatLearn/>). • Ian Goodfellow and Yoshua Bengio and Aaron Courville, Deep Learning, An MIT Press book (<http://www.deeplearningbook.org>)
- Kevin Murphy, Machine Learning: a Probabilistic Perspective, (MIT Press, 2013) • Hal Daumé III, A Course in Machine Learning (<http://ciml.info/>)

## Bases de données

- UCI Repository: <http://www.ics.uci.edu/~mlearn/MLRepository.html> • UCI KDD Archive: <http://kdd.ics.uci.edu/summary.data.application.html> •

Statlib: <http://lib.stat.cmu.edu/>

- Delve: <http://www.cs.utoronto.ca/~delve/>

- Kaggle: <https://www.kaggle.com/>

- Benchmarks (Vision):

- ImageNet: <http://image-net.org/>

- MS COCO: <http://cocodataset.org/>

- MNIST et plus:

- [http://rodrigob.github.io/are\\_we\\_there\\_yet/build/classification\\_datasets\\_resu\\_lts.html](http://rodrigob.github.io/are_we_there_yet/build/classification_datasets_resu_lts.html)

- CV on line: <https://computervisiononline.com/datasets>

- Kitti: <http://www.cvlibs.net/datasets/kitti/>

- Waymo: <https://waymo.com/open>

- Journal of Machine Learning Research [www.jmlr.org](http://www.jmlr.org) •
- Machine Learning
- Neural Computation
- Neural Networks
- IEEE Transactions on Neural Networks
- IEEE Transactions on Pattern Analysis and Machine Intelligence •
- Annals of Statistics
- Journal of the American Statistical Association
- ...

## Conférences

- International Conference on Machine Learning (ICML) •
- European Conference on Machine Learning (ECML)
- Neural Information Processing Systems (NIPS)
- International Conference on Learning Representations (ICLR) •
- Uncertainty in Artificial Intelligence (UAI)
- International Joint Conference on Artificial Intelligence (IJCAI) •
- International Conference on Neural Networks (ICNN) • Conference of the
- American Association for Artificial Intelligence (AAAI) • IEEE Conference
- on Computer Vision and Pattern Recognition (CVPR) • European

Conference on Computer Vision (ECCV)

- International Conference on Computer Vision (ICCV)
- IEEE International Conference on Data Mining (ICDM) •

...

74

## Cours & tutoriaux

- Des MOOC (Français et Anglais)
- Des tutoriaux associés aux conférences (orientés recherche)
- Des cours en français:
  - <https://gricad-gitlab.univ-grenoble-alpes.fr/talks/fidle>

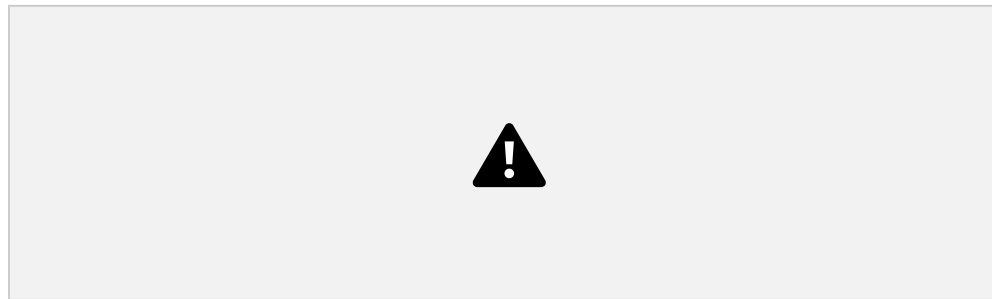


- [https://www.college-de-france.fr/site/stephane-mallat/\\_course.htm](https://www.college-de-france.fr/site/stephane-mallat/_course.htm)
- Des « cheat sheets »
  - <https://stanford.edu/~shervine/teaching/>

## Logiciels

- Environnement génériques: Matlab, ScikitLearn

- Environnements Deep Learning: Tensor Flow, Pytorch, mxnet...
- Beaucoup de codes sur GitHub



Apprentissage