

Machine Learning Introduction à la théorie et bonnes pratiques



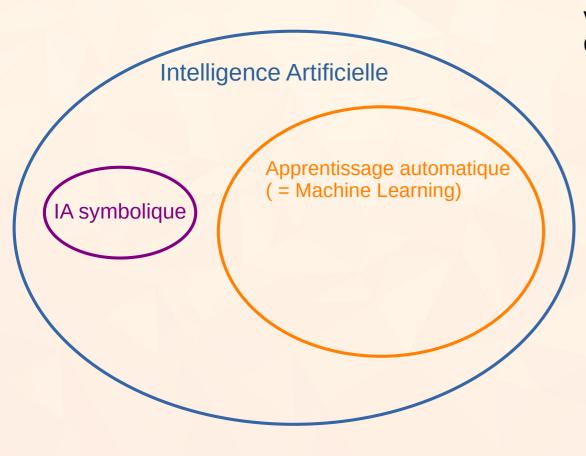
Sommaire



I) Définitions	3
Types et sous-catégories	
III) Données, loi de probabilité	5
IV) Non-supervisé	
V) Renforcement	
VÍ) Supervisé	
1) lois usuelles	
2) principe	
3) fonction de perte	
4) choix distribution	
5) fonction d'activation en sortie (Deep Learning)	
6) métrique	18
7) overfitting	
8) a) train/validation/test	
b) stratégie de validation	22
10) modèle	
11) résumé	
VII) Sources	

I - Définition

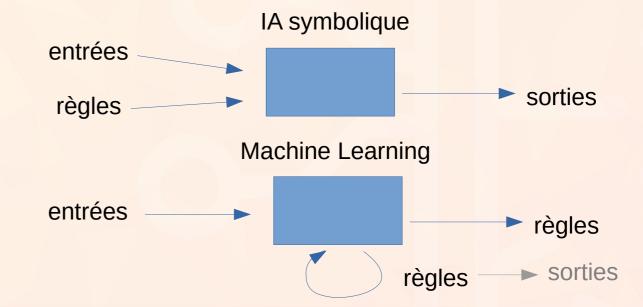




Intelligence artificielle: Science des théories et techniques visant à créer des machines/programmes qui font preuve d'intelligence (vague)

IA symbolique : résolution par manipulation logique de symboles

Machine Learning : résolution par confrontation et amélioration avec des données

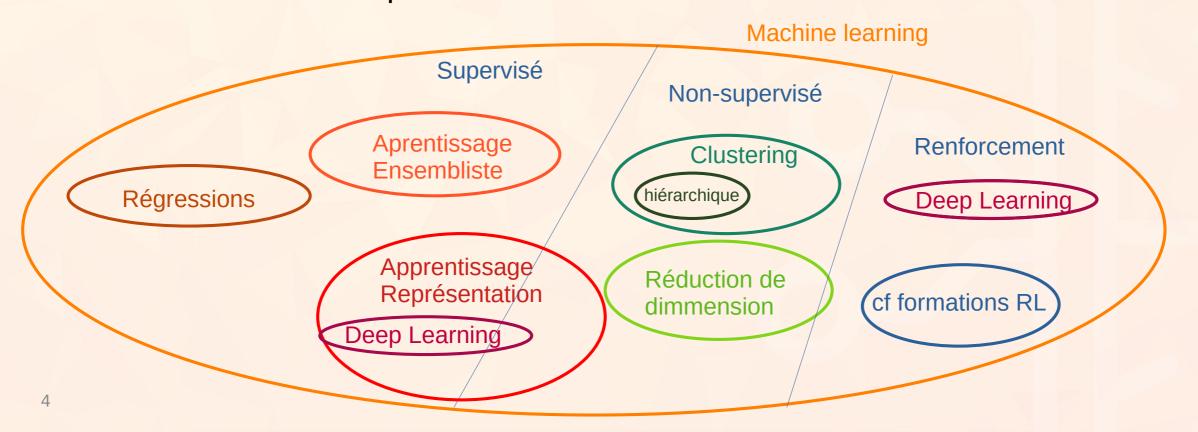


II - Types et sous-catégories



3 types de problèmes :

- Supervisé : entrées, labels
- Non-supervisé : entrées, pas de labels
- Renforcement : récompenses situationnelles



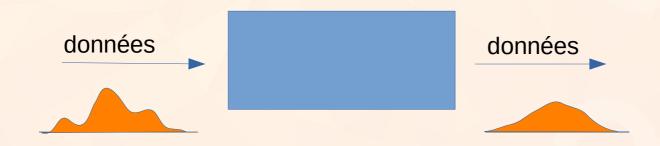
III - Données - Loi de Proba



Les données

Censées représenter la réalité :

- Modélisation incomplète
- Observabilité incomplète
- → données = variables aléatoires, lois de probabilités

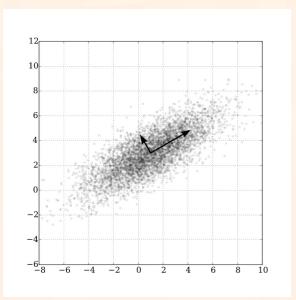


IV - Non-supervisé



Objectif : modifier la représentation des données 3 types :

- Représentation par réduction dimension : projection sans perte
- Représentation indépendante : éliminer les dépendances des axes
- Représentation parcimonieuse : données clairsemées (0 & 1)



PCA (Scikit-learn)

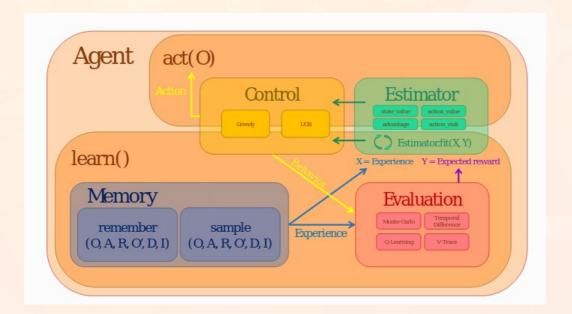
Method name	Parameters	Scalability	Usecase	Geometry (metric used)
K-Means	number of clusters	Very large n_samples, medium n_clusters with MiniBatch code	General-purpose, even cluster size, flat geometry, not too many clusters	Distances between points
Affinity propaga- tion	damping, sample preference	Not scalable with n_samples	Many clusters, uneven cluster size, non-flat geometry	Graph distance (e.g. nearest-neighbor graph)
Mean-shift	bandwidth	Not scalable with n_samples	Many clusters, uneven cluster size, non-flat geometry	Distances between points
Spectral cluster- ing	number of clusters	Medium n_samples, small n_clusters	Few clusters, even cluster size, non-flat geometry	Graph distance (e.g. nearest-neighbor graph)
Ward hierarchical clustering	number of clusters or distance thresh- old	Large n_samples and n_clusters	Many clusters, possibly connectivity constraints	Distances between points
Agglomerative clustering	number of clusters or distance thresh- old, linkage type, distance	Large n_samples and n_clusters	Many clusters, possibly connectivity constraints, non Euclidean distances	Any pairwise distance
DBSCAN	neighborhood size	Very large n_samples, medium n_clusters	Non-flat geometry, uneven cluster sizes	Distances between nearest points
OPTICS	minimum cluster membership	Very large n_samples, large n_clusters	Non-flat geometry, uneven cluster sizes, variable cluster density	Distances between points
Gaussian mix- tures	many	Not scalable	Flat geometry, good for density estimation	Mahalanobis distances to centers
Birch	branching factor, threshold, optional global clusterer.	Large n_clusters and n_samples	Large dataset, outlier removal, data reduction.	Euclidean distance be- tween points

Clustering (Scikit-learn)

V - Apprentissage Renforcement



- Implémentez votre environement (états, actions, récompenses, observations) avec Gym par exemple
- Implémentez votre agent avec une bibliothèque (LearnRL, Horizon, Dopamine)
- Tester les différentes possibilités!



VI - Supervisé



Données : entrées & labels

Entrées

$$X = \begin{array}{c} X_{1,1} X_{1,2} - X_{1,d} \\ X_{2,1} X_{2,2} - X_{2,d} \\ & | & | & | & | \\ X_{n,1} X_{n,2} - X_{n,d} \end{array} \qquad Y = \begin{array}{c} Y_{1,1} Y_{1,2} - Y_{1,r} \\ Y_{2,1} Y_{2,2} - Y_{2,r} \\ & | & | & | \\ Y_{n,1} Y_{n,2} - Y_{n,r} \end{array}$$

Supervisé - Lois usuelles



Entrées

$$X = \begin{array}{c} X_{1,1} X_{1,2} - X_{1,d} \\ X_{2,1} X_{2,2} - X_{2,d} \\ & | & | \\ X_{n,1} X_{n,2} - X_{n,d} \end{array}$$

Labels

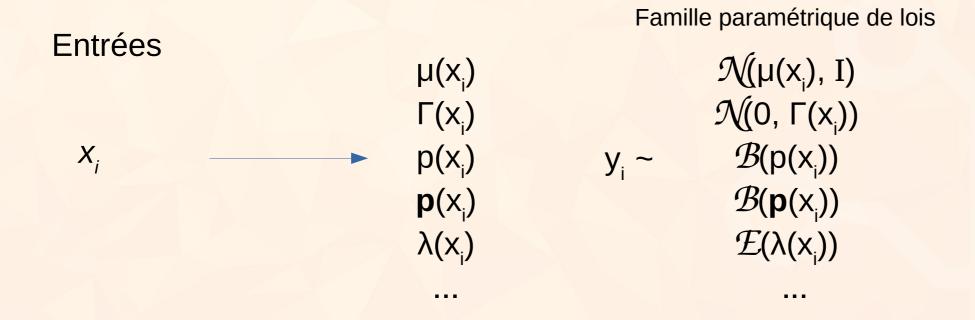


Famille paramétrique de lois

$$\mathcal{N}(\mu, I)$$
 $\mathcal{N}(0, \Gamma)$
 $\mathcal{B}(p)$
 $\mathcal{B}(p_1, ..., p_{k-1})$
 $\mathcal{E}(\lambda)$

Supervisé - Lois usuelles





Supervisé - Principe



$$y_i \sim \mathcal{N}(\mu(x_i), I)$$

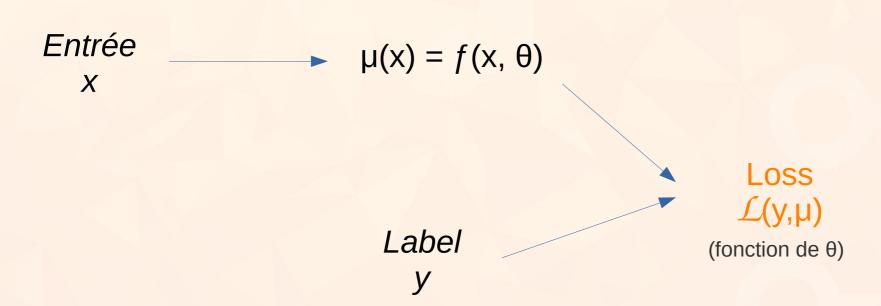
Entrée
$$\mu(x) = f(x, \theta)$$

f := model $\theta := paramètres (poids)$

Supervisé - Principe



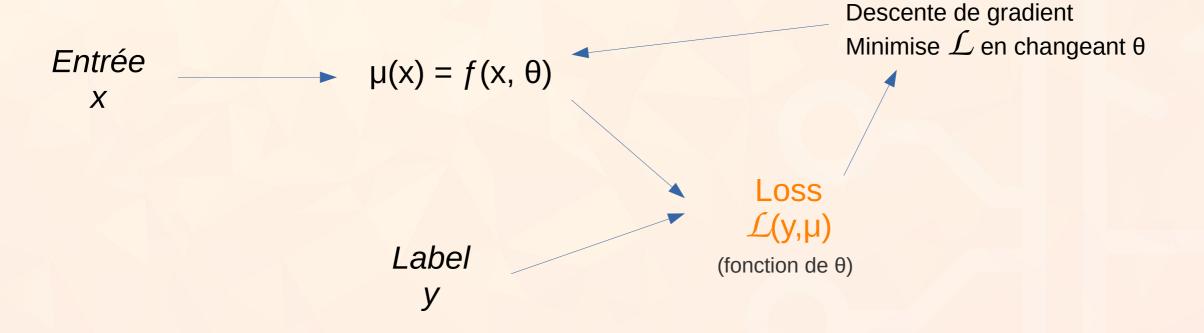
$$y_i \sim \mathcal{N}(\mu(x_i), I)$$



Supervisé - Principe



$$y_i \sim \mathcal{N}(\mu(x_i), I)$$



Supervisé - Fonction de perte



Fonction de perte (de cout, d'erreur) :

$$\mathcal{L}(y, \mu(x, \theta))$$
 petit $\longrightarrow y \sim \mathcal{N}(\mu(x, \theta), I)$

$$p(y|x,\theta) = \frac{1}{\sqrt{2\pi}} \exp{-\frac{||y - \mu(x,\theta)||^2}{2}}$$

1ère idée : estimateur du maximum de vraissemblance

$$\begin{aligned} \theta_{MV} &= \underset{\theta}{\operatorname{argmax}} \ p(Y|X,\theta) \\ &= \underset{\theta}{\operatorname{argmax}} \ \log p(Y|X,\theta) \\ &= \underset{\theta}{\operatorname{argmax}} \ \sum_{i} \log p(y_{i}|x_{i},\theta) \\ &= \underset{\theta}{\operatorname{argmax}} \ \frac{1}{2} \sum_{i} (-\log 2\pi - (y_{i} - \mu(x_{i},\theta))^{2}) \\ &= \underset{\theta}{\operatorname{argmax}} \ \sum_{i} (y_{i} - \mu(x_{i},\theta))^{2} \end{aligned} \qquad \begin{aligned} \mathcal{L}(Y,\mu(X,\theta)) &= \sum (y - \mu(x,\theta))^{2} \\ &= \sum (y - f(x,\theta))^{2} \end{aligned}$$

$$= \underset{\theta}{\operatorname{argmax}} \ \frac{1}{2} \sum_{i} (-\log 2\pi - (y_{i} - \mu(x_{i},\theta))^{2}) \\ &= \underset{\theta}{\operatorname{argmin}} \ \sum_{i} (y_{i} - \mu(x_{i},\theta))^{2} \end{aligned} \qquad \begin{aligned} \mathcal{L}(Y,\mu(X,\theta)) &= \sum (y - \mu(x,\theta))^{2} \\ &= \sum (y - f(x,\theta))^{2} \end{aligned}$$

$$= \underset{\theta}{\operatorname{Autre exemple pour }} \ y_{i} \sim \mathcal{B}(p(x_{i})) : \\ \mathcal{L}(Y,\mu(X,\theta)) &= -\sum y \log f(x,\theta) + (1-y) \log (1-f(x,\theta)) \end{aligned}$$

Supervisé - Fonction de perte



2ème idée : loi de Bayes et maximum a posteriori

$$\begin{split} p(A|B) &= \frac{p(B|A)P(A)}{P(B)} \\ \theta_{MAP} &= \operatorname*{argmax}_{\theta} p(\theta|Y,X) \\ &= \operatorname*{argmax}_{\theta} \log p(\theta|Y,X) \\ &= \operatorname*{argmax}_{\theta} \log p(Y|\theta,X) + \log p(\theta) \\ \theta &\sim \mathcal{N}\left(0,\frac{1}{\lambda}I\right) \\ \theta_{MAP} &= \operatorname*{argmax}_{\theta} \log p(Y|\theta,X) + \lambda \|\theta\|^2 \\ \mathcal{L}_{MAP}(\theta) &= \mathcal{L}_{MV}(\theta) + \lambda \|\theta\|^2 \\ \mathcal{L}_{MAP}(\theta) &= \sum (y - f(x,\theta))^2 + \lambda \|\theta\|^2 \end{split}$$

Supervisé - Choix distribution



Quelle distribution choisir?

Quelques idées (souvent bonnes!):

- Régression (prédire fonction dans IR) :
 𝒩(μ, I) → MSE + version multivariée
 𝒩(0, Γ) si on s'interesse aux valeurs proches de 0
 Ɛ(λ) si évènement temporel sans mémoire
- Classification (prédire la classe) $\mathcal{B}(p) \rightarrow 2$ labels (ou 1 label) $\mathcal{B}(p_1, ..., p_{k-1}) \rightarrow k$ labels
- Mélange ...

Deep Learning



Quelles fonctions d'activation choisir à la fin d'un NN?

- → Quelle domaine pour la cible ?
-]-∞, +∞[: sortie linéaire
-]0, +∞[: prédire le **log** + sortie linéaire
-]0, 1[: sortie sigmoïdale
-]0, 1[k de somme 1 : sortie softmax
- Autre : se ramener linéairement à un cas précédent
 - → 1 sortie par prédiction

Softmax

$$\sigma(\mathbf{z})_j = rac{\mathrm{e}^{z_j}}{\sum_{k=1}^K \mathrm{e}^{z_k}}$$
 pour tout $j \in \{1,\dots,K\}$.



Sigmoide



Supervisé - Métrique

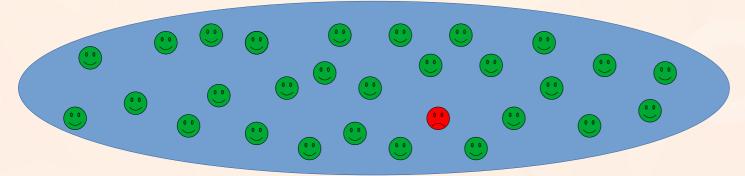


Choisir sa métrique

- → Plein de possibilités!
- Lien direct avec l'objectif informel
- Interprétable « physiquement » ou intuitivement
- Facile à calculer
- Attention aux pièges!

Exemple:

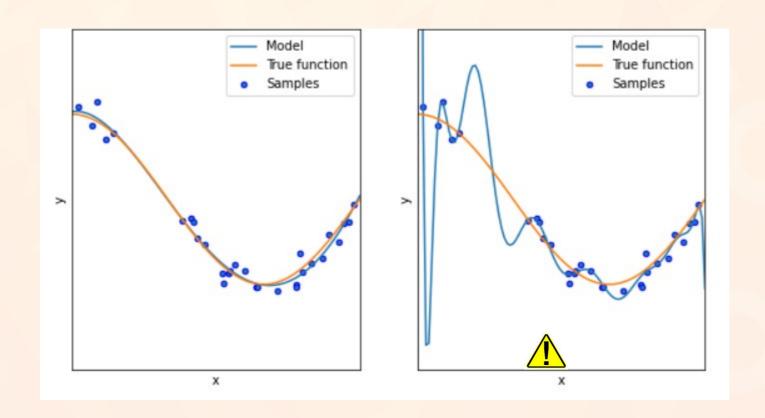
Accuracy en classification binaire = (TP+TN)/N



Supervisé - overfitting



Un mot sur l'overfitting ...



Supervisé - Train/Validation/Test



Le modèle

- Transformations
- Structure
- Régularisation
- Hyper-paramètres
- Paramètres entraînables



Supervisé - Train/Validation/Test



Hyper-paramètres

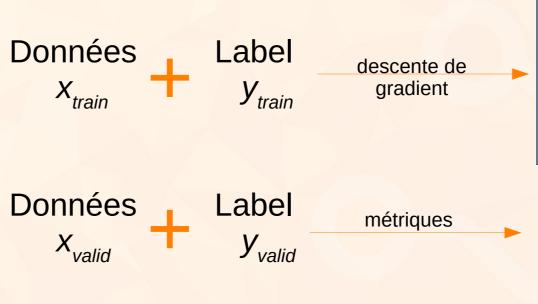
et caractéristiques

Paramètres

Le modèle

- Transformations
- Structure
- Régularisation
- Hyper-paramètres
- Paramètres entraînables





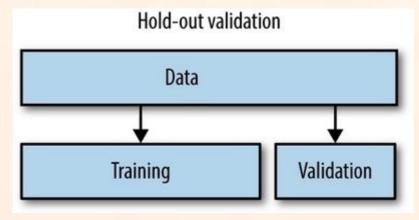


Supervisé - Stratégie Validation



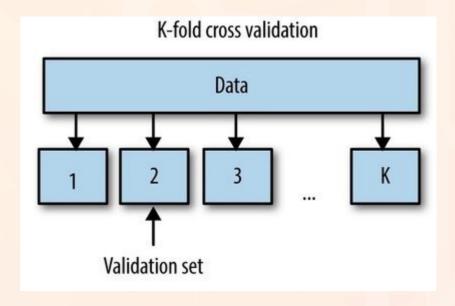
Choix de la stratégie de validation

- Hold-out validation :
 - → ensemble validation fixeperformance = performance sur validation



- Stratified K-fold
 - → classes équilibrées sur chaque fold

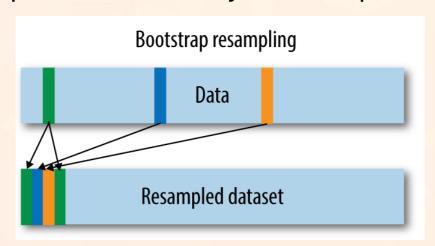
- Validation croisée K-fold
 - → K entraînements sur K-1 folds le dernier est l'ensemble de validation performance = moyenne des performances



Supervisé - Stratégie Validation

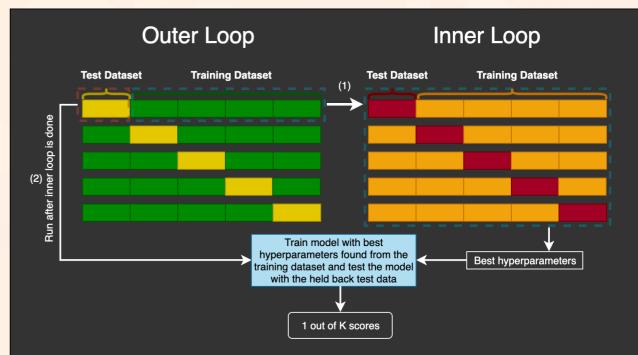


- K-fold + itérations randomisées
 - → P fois la stratégie K-folds on mêlange/divise les données aléatoirement performance = moyenne des performances



Nested K-fold

 → K-fold + K-fold sur chaque itération en faisant une recherche d'hyper-paramètres
 performance = moyenne des performances



Supervisé - Modèle



Mise au point du modèle

- → étape la plus longue
- → équilibre entre underfitting et overfitting

Astuces:

- → chercher des papiers de recherche sur des sujets similaires
- → commencer par un modèle pauvre qui bat le hasard ou des stratégies triviales
- ightarrow s'intéresser au preprocessing pour améliorer la prédiction sur D_{train}
- → essayer ensuite d'overfitter puis diminuer progressivement la taille
- → essayer des régularisations (loss, dropout, data augment, early stop, ...)
- \rightarrow se fixer des objectifs sur les métriques sur D_{test}

Supervisé - Résumé



Résumé (to-do list):

- Analyser données et objectif → type de ML, distributions, loss, FA de sortie
- Métriques prenant en compte les spécificités du problème et l'objectif
- Séparation des données et système de validation
- Mise au point du modèle
 - · Preprocessing
 - · Modèle simple
 - Overfitting
 - · Diminution taille
 - · Régularisation
- Test sur D_{test}

Supervisé - Résumé



Résumé (to-do list) :

- Analyser données et objectif → type de ML, distributions, loss, FA de sortie
- Métriques prenant en compte les spécificités du problème et l'objectif
- Séparation des données et système de validation
- Mise au point du modèle
 - · Preprocessing
 - · Modèle simple
 - Overfitting
 - · Diminution taille
 - · Régularisation
- Test sur D_{test}

Super, ça marche !

Sources

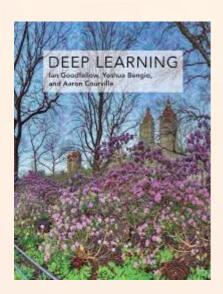


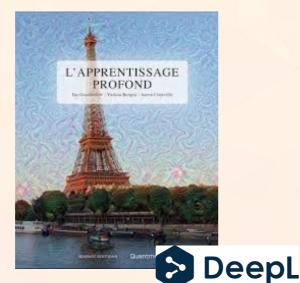
Sources:

- Deep Learning, Ian Goodfellow
- L'Apprentissage Profond avec Python, François Chollet
- Types de cross-validation :

 Types de cross-validation :

https://towardsdatascience.com/understanding-8-types-of-cross-validation-80c935a4976d







Vive le machine learning!



Merci de votre attention!

Valentin Goldité