

# Le Dictionnaire de l'Apprentissage Automatique d'**A'**alto

Alexander Jung<sup>1</sup>, Konstantina Olioumtsevits<sup>1</sup>, Juliette Gronier<sup>2</sup>,  
et Salvatore Rastelli<sup>1</sup>

<sup>1</sup>Aalto University    <sup>2</sup>ENS Lyon

July 14, 2025



please cite as: A. Jung, K. Olioumtsevits, J. Gronier and S.  
Rastelli *The Aalto Dictionary of Machine Learning*. Espoo,  
Finland: Aalto University, 2025.

## Remerciements

Ce dictionnaire de l'apprentissage automatique a évolué au fil du développement et l'enseignement de plusieurs cours, parmi lesquels CS-E3210 Machine Learning: Basic Principles, CS-C3240 Machine Learning, CS-E4800 Artificial Intelligence, CS-EJ3211 Machine Learning with Python, CS-EJ3311 Deep Learning with Python, CS-E4740 Federated Learning, and CS-E407507 Human-Centered Machine Learning. Ces cours ont été proposés à Aalto University <https://www.aalto.fi/en>, à des apprenants adultes via le Finnish Institute of Technology (FITech) <https://fitech.io/en/>, et à des étudiants et étudiantes internationaux dans le cadre de l'alliance universitaire européenne Unite! <https://www.aalto.fi/en/unite>.

Nous remercions les étudiants et étudiantes pour leurs retours de qualité qui ont contribué à façonner ce dictionnaire. En particulier, un grand merci à Mikko Seesto pour sa relecture minutieuse.

Cette traduction française s'appuie notamment sur le Glossaire de l'intelligence artificielle (IA) proposé par la CNIL

<https://www.cnil.fr/fr/intelligence-artificielle/glossaire-ia>, ainsi que sur les ressources du site FranceTerme, géré par le Ministère de la Culture <https://www.culture.fr/franceterme>, \*et le site de l'Office québécois de la langue française <https://www.oqlf.gouv.qc.ca>.

## Notations et symboles

### Ensembles et fonctions

$a \in \mathcal{A}$	L'objet $a$ est un élément de l'ensemble $\mathcal{A}$ .
$a := b$	On note $a$ comme abréviation de $b$ .
$ \mathcal{A} $	Le cardinal (i.e., le nombre d'éléments) d'un ensemble fini $\mathcal{A}$ .
$\mathcal{A} \subseteq \mathcal{B}$	$\mathcal{A}$ est un sous-ensemble de $\mathcal{B}$ .
$\mathcal{A} \subset \mathcal{B}$	$\mathcal{A}$ est un sous-ensemble strict de $\mathcal{B}$ (i.e., non égal à $\mathcal{B}$ ).
$\mathbb{N}$	Les entiers naturels $1, 2, \dots$
$\mathbb{R}$	Les nombres réels $x$ $[1]$ .
$\mathbb{R}_+$	Les réels positifs ou nuls $x \geq 0$ .
$\mathbb{R}_{++}$	Les réels strictement positifs $x > 0$ .
$\{0, 1\}$	L'ensemble composé des deux réels 0 et 1.
$[0, 1]$	L'intervalle fermé des nombres réels $x$ tels que $0 \leq x \leq 1$ .

$\operatorname{argmin}_{\mathbf{w}} f(\mathbf{w})$	<p>L'ensemble des points qui minimisent la fonction à valeurs réelles fonction <math>f(\mathbf{w})</math>.</p> <p>Voir aussi: fonction.</p>
$\mathbb{S}^{(n)}$	<p>L'ensemble des vecteurs de norme unitaire dans <math>\mathbb{R}^{n+1}</math>.</p> <p>Voir aussi: norme.</p>
$\exp(a)$	<p>La fonction exponentielle évaluée en un réel <math>a \in \mathbb{R}</math>.</p>
$\log(a)$	<p>Le logarithme d'un réel strictement positif <math>a \in \mathbb{R}_{++}</math>.</p>
$f(\cdot) : \mathcal{A} \rightarrow \mathcal{B} : a \mapsto f(a)$	<p>Une fonction (ou application) d'un ensemble <math>\mathcal{A}</math> dans un ensemble <math>\mathcal{B}</math>, qui associe à chaque entrée <math>a \in \mathcal{A}</math> une image bien définie <math>f(a) \in \mathcal{B}</math>. L'ensemble <math>\mathcal{A}</math> est le domaine de définition de la fonction <math>f</math> et l'ensemble <math>\mathcal{B}</math> est l'ensemble d'arrivée de <math>f</math>. L'apprentissage automatique vise à apprendre une fonction <math>h</math> qui prend en entrée les caractéristiques <math>\mathbf{x}</math> d'un point de données et renvoie une prédiction <math>h(\mathbf{x})</math> pour son étiquette étiquette <math>y</math>.</p> <p>Voir aussi: fonction, application, apprentissage automatique, hypothèse, caractéristique, point de données, prédiction, étiquette.</p>
$\operatorname{epi}(f)$	<p>L' épigraphe d'une fonction à valeurs réelles <math>f : \mathbb{R}^d \rightarrow \mathbb{R}</math>.</p> <p>Voir aussi: épigraphe, fonction.</p>

$\frac{\partial f(w_1, \dots, w_d)}{\partial w_j}$	<p>La dérivée partielle (si elle existe) d'une fonction à valeurs réelles <math>f : \mathbb{R}^d \rightarrow \mathbb{R}</math> par rapport à <math>w_j</math> [2, Ch. 9].</p> <p>Voir aussi: fonction.</p>
$\nabla f(\mathbf{w})$	<p>Le gradient d'une fonction à valeurs réelles dérivable <math>f : \mathbb{R}^d \rightarrow \mathbb{R}</math> est le vecteur <math>\nabla f(\mathbf{w}) = \left( \frac{\partial f}{\partial w_1}, \dots, \frac{\partial f}{\partial w_d} \right)^T \in \mathbb{R}^d</math> [2, Ch. 9].</p> <p>Voir aussi: gradient, dérivable, fonction.</p>

## Matrices et Vecteurs

$\mathbf{x} = (x_1, \dots, x_d)^T$	Un vecteur de taille $d$ , dont la $j$ -ième composante est $x_j$ .
$\mathbb{R}^d$	L'ensemble des vecteurs $\mathbf{x} = (x_1, \dots, x_d)^T$ constitués de $d$ composantes réelles $x_1, \dots, x_d \in \mathbb{R}$ .
$\mathbf{I}_{l \times d}$	Une matrice identité généralisée de $l$ lignes et $d$ colonnes. Les composantes de $\mathbf{I}_{l \times d} \in \mathbb{R}^{l \times d}$ valent 1 sur la diagonale principale et 0 ailleurs.
$\mathbf{I}_d, \mathbf{I}$	Une matrice identité carrée de taille $d \times d$ . Si la dimension est claire dans le contexte, on peut omettre l'indice.
$\ \mathbf{x}\ _2$	La norme euclidienne (ou $\ell_2$ ) du vecteur $\mathbf{x} = (x_1, \dots, x_d)^T \in \mathbb{R}^d$ définie par $\ \mathbf{x}\ _2 := \sqrt{\sum_{j=1}^d x_j^2}$ . Voir aussi: norme
$\ \mathbf{x}\ $	Une certaine norme du vecteur $\mathbf{x} \in \mathbb{R}^d$ [3]. Sauf indication contraire, on entend par là la norme euclidienne $\ \mathbf{x}\ _2$ . Voir aussi: norme
$\mathbf{x}^T$	La transposée d'une matrice ayant pour unique colonne le vecteur $\mathbf{x} \in \mathbb{R}^d$ .
$\mathbf{X}^T$	La transposée d'une matrice $\mathbf{X} \in \mathbb{R}^{m \times d}$ . Une matrice carrée à valeurs réelles $\mathbf{X} \in \mathbb{R}^{m \times m}$ est dite symétrique si $\mathbf{X} = \mathbf{X}^T$ .
$\mathbf{X}^{-1}$	La matrice inverse d'une matrice $\mathbf{X} \in \mathbb{R}^{d \times d}$ . Voir aussi: matrice inverse.

$\mathbf{0} = (0, \dots, 0)^T$	Le vecteur de $\mathbb{R}^d$ dont toutes les composantes valent 0.
$\mathbf{1} = (1, \dots, 1)^T$	Le vecteur de $\mathbb{R}^d$ dont toutes les composantes valent 1.
$(\mathbf{v}^T, \mathbf{w}^T)^T$	Le vecteur de longueur $d + d'$ obtenu en concaténant les $\mathbf{v} \in \mathbb{R}^d$ avec celles de $\mathbf{w} \in \mathbb{R}^{d'}$ .
$\text{span}\{\mathbf{B}\}$	Le sous-espace engendré par une matrice $\mathbf{B} \in \mathbb{R}^{a \times b}$ , c'est-à-dire l'ensemble de toutes les combinaisons linéaires des colonnes de $\mathbf{B}$ : $\text{span}\{\mathbf{B}\} = \{\mathbf{B}\mathbf{a} : \mathbf{a} \in \mathbb{R}^b\} \subseteq \mathbb{R}^a$ .
$\det(\mathbf{C})$	Le déterminant de la matrice $\mathbf{C}$ . Voir aussi: déterminant
$\mathbf{A} \otimes \mathbf{B}$	Le produit de Kronecker des matrices $\mathbf{A}$ et $\mathbf{B}$ [4].

## Théorie des probabilités

$\mathbf{x} \sim p(\mathbf{z})$  La variable aléatoire (VA)  $\mathbf{x}$  suit la loi de probabilité  $p(\mathbf{z})$  [5, 6].

Voir aussi : VA, loi de probabilité

---

$\mathbb{E}_p\{f(\mathbf{z})\}$  L'espérance d'une VA  $f(\mathbf{z})$  obtenue en appliquant une fonction déterministe  $f$  à une VA  $\mathbf{z}$  dont la loi de probabilité est  $\mathbb{P}(\mathbf{z})$ . Si la loi de probabilité est claire dans le contexte, on écrit simplement  $\mathbb{E}\{f(\mathbf{z})\}$ .

Voir aussi : espérance, VA, loi de probabilité

---

$\text{cov}(x, y)$  La covariance entre deux VA à valeurs réelles définies sur un même espace probabilisé.

Voir aussi : covariance, VA, espace probabilisé

---

$\mathbb{P}(\mathbf{x}, y)$  Une loi de probabilité (conjointe) d'une VA dont les réalisations sont des points de données avec des caractéristiques  $\mathbf{x}$  et une étiquette  $y$ .

Voir aussi : loi de probabilité, VA, point de données, étiquette

---

$\mathbb{P}(\mathbf{x}|y)$  Une loi de probabilité conditionnelle d'une VA  $\mathbf{x}$  étant donnée la valeur d'une autre VA  $y$  [7, Sec. 3.5].

Voir aussi : loi de probabilité, VA.



$\mathbb{P}(\mathbf{x}; \mathbf{w})$  Une loi de probabilité paramétrée d'une VA  $\mathbf{x}$ . La loi de probabilité dépend d'un vecteur de paramètres  $\mathbf{w}$ . Par exemple,  $\mathbb{P}(\mathbf{x}; \mathbf{w})$  pourrait être une loi normale multivariée avec un vecteur de paramètres  $\mathbf{w}$  donné par les composantes du vecteur de moyenne  $\mathbb{E}\{\mathbf{x}\}$  et la matrice de covariance  $\mathbb{E}\left\{(\mathbf{x} - \mathbb{E}\{\mathbf{x}\})(\mathbf{x} - \mathbb{E}\{\mathbf{x}\})^T\right\}$ .  
Voir aussi : loi de probabilité, VA, loi normale multivariée, moyenne, matrice de covariance

---

$\mathcal{N}(\mu, \sigma^2)$  La loi de probabilité d'une variable aléatoire normale centrée réduite (VA normale centrée réduite)  $x \in \mathbb{R}$  ayant comme moyenne (ou espérance)  $\mu = \mathbb{E}\{x\}$  et comme variance  $\sigma^2 = \mathbb{E}\{(x - \mu)^2\}$ .  
Voir aussi : VA normale centrée réduite, moyenne, espérance, variance, loi de probabilité

---

$\mathcal{N}(\boldsymbol{\mu}, \mathbf{C})$  La loi normale multivariée d'une VA normale centrée réduite vectorielle  $\mathbf{x} \in \mathbb{R}^d$  ayant comme moyenne (ou espérance)  $\boldsymbol{\mu} = \mathbb{E}\{\mathbf{x}\}$  et comme matrice de covariance  $\mathbf{C} = \mathbb{E}\{(\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^T\}$ .  
Voir aussi : loi normale multivariée, VA normale centrée réduite, moyenne, matrice de covariance

## Apprentissage automatique

$r$	<p>Un indice <math>r = 1, 2, \dots</math> qui énumère les points de données.</p> <p>Voir aussi : points de données.</p>
$m$	<p>Le nombre de points de données dans un jeu de données (c'est-à-dire la taille du jeu de données).</p> <p>Voir aussi : points de données, jeu de données.</p>
$\mathcal{D}$	<p>Un jeu de données <math>\mathcal{D} = \{\mathbf{z}^{(1)}, \dots, \mathbf{z}^{(m)}\}</math> est une liste de points de données individuels <math>\mathbf{z}^{(r)}</math>, pour <math>r = 1, \dots, m</math>.</p> <p>Voir aussi : points de données, jeu de données.</p>
$d$	<p>Le nombre de caractéristiques qui constituent un point de données.</p> <p>Voir aussi : caractéristiques, point de données.</p>
$x_j$	<p>La <math>j</math>-ième caractéristique d'un point de données. La première caractéristique est notée <math>x_1</math>, la deuxième <math>x_2</math>, et ainsi de suite.</p> <p>Voir aussi : caractéristiques, point de données.</p>
$\mathbf{x}$	<p>Le vecteur de caractéristiques <math>\mathbf{x} = (x_1, \dots, x_d)^T</math> d'un point de données, dont les composantes sont les différentes caractéristiques du point de données.</p> <p>Voir aussi : vecteur de caractéristiques, caractéristiques, point de données.</p>

$\mathcal{X}$	<p>L'espace des caractéristiques <math>\mathcal{X}</math> est l'ensemble de toutes les valeurs possibles que les caractéristiques <math>\mathbf{x}</math> d'un point de données peuvent prendre.</p> <p>Voir aussi : espace des caractéristiques, caractéristiques, point de données.</p>
$\mathbf{z}$	<p>Au lieu du symbole <math>\mathbf{x}</math>, on utilise parfois <math>\mathbf{z}</math> comme un autre symbole pour désigner un vecteur dont les composantes sont les différentes caractéristiques d'un point de données.</p> <p>On a besoin de deux symboles différents pour distinguer les caractéristiques brutes des caractéristiques apprises [8, Ch. 9].</p> <p>Voir aussi : caractéristiques, point de données.</p>
$\mathbf{x}^{(r)}$	<p>Le vecteur de caractéristiques du <math>r</math>-ième point de données dans un jeu de données.</p> <p>Voir aussi : caractéristiques, point de données, jeu de données.</p>
$x_j^{(r)}$	<p>La <math>j</math>-ième caractéristique du <math>r</math>-ième point de données dans un jeu de données.</p> <p>Voir aussi : caractéristiques, point de données, jeu de données.</p>
$y$	<p>L'étiquette (ou quantité d'intérêt) d'un point de données.</p> <p>Voir aussi : étiquette, point de données.</p>
$y^{(r)}$	<p>L'étiquette du <math>r</math>-ième point de données.</p> <p>Voir aussi : étiquette, point de données.</p>

$(\mathbf{x}^{(r)}, y^{(r)})$  Les caractéristiques et l'étiquette du  $r$ -ième point de données.

Voir aussi : caractéristiques, étiquette, point de données.

---

$\mathcal{Y}$  L'espace des étiquettes  $\mathcal{Y}$  d'une méthode d'apprentissage automatique comprend toutes les valeurs d'étiquette qu'un point de données peut porter. L'espace des étiquettes nominal peut être plus grand que l'ensemble des différentes valeurs d'étiquette présentes dans un jeu de données donné (par exemple, un ensemble d'entraînement (ou d'apprentissage)).

Les problèmes (ou méthodes) d'apprentissage automatique utilisant un espace des étiquettes numérique, comme  $\mathcal{Y} = \mathbb{R}$  ou  $\mathcal{Y} = \mathbb{R}^3$ , sont appelés problèmes (ou méthodes) de régression.

Les problèmes (ou méthodes) d'apprentissage automatique utilisant un espace des étiquettes discret, comme  $\mathcal{Y} = \{0, 1\}$  ou  $\mathcal{Y} = \{chat, chien, souris\}$ , sont appelés problèmes (ou méthodes) de classification.

Voir aussi : espace des étiquettes, étiquette, jeu de données, apprentissage automatique.

---

$\mathcal{B}$  Un mini-lot (ou sous-ensemble) de points de données choisis aléatoirement.

Voir aussi : lot, points de données.

$B$	<p>La taille (c'est-à-dire le nombre de points de données) d'un mini-lot.</p> <p>Voir aussi : lot, points de données.</p>
$h(\cdot)$	<p>Une fonction hypothèse qui lit les caractéristiques <math>\mathbf{x}</math> d'un point de données et produit une prédiction <math>\hat{y} = h(\mathbf{x})</math> pour son étiquette <math>y</math>.</p> <p>Voir aussi : hypothèse, caractéristiques, point de données, prédiction, étiquette.</p>
$\mathcal{Y}^{\mathcal{X}}$	<p>Étant donnés deux ensembles <math>\mathcal{X}</math> et <math>\mathcal{Y}</math>, on note <math>\mathcal{Y}^{\mathcal{X}}</math> l'ensemble de toutes les fonctions hypothèses possibles <math>h : \mathcal{X} \rightarrow \mathcal{Y}</math>.</p> <p>Voir aussi : espace des caractéristiques, espace des étiquettes, hypothèse.</p>
$\mathcal{H}$	<p>Un espace des hypothèses ou modèle utilisé par une méthode d'apprentissage automatique. L'espace des hypothèses est constitué des différentes hypothèses <math>h : \mathcal{X} \rightarrow \mathcal{Y}</math>, parmi lesquelles la méthode d'apprentissage automatique doit choisir.</p> <p>Voir aussi : espace des hypothèses, hypothèse, apprentissage automatique.</p>
$d_{\text{eff}}(\mathcal{H})$	<p>La dimension effective d'un espace des hypothèses <math>\mathcal{H}</math>.</p> <p>Voir aussi : dimension effective, espace des hypothèses.</p>

$B^2$	<p>Le biais au carré d'une fonction hypothèse <math>\hat{h}</math> apprise par une méthode d'apprentissage automatique. La méthode est entraînée sur des points de données modélisés comme des réalisations de VA. Puisque les données sont des réalisations de VA, la fonction hypothèse apprise <math>\hat{h}</math> est également une réalisation d'une VA.</p> <p>Voir aussi : biais, hypothèse, points de données, réalisation, VA.</p>
$V$	<p>La variance (des paramètres) de la fonction hypothèse apprise par une méthode d'apprentissage automatique. La méthode est entraînée sur des points de données modélisés comme des réalisations de VA. Puisque les données sont des réalisations de VA, la fonction hypothèse apprise <math>\hat{h}</math> est également une réalisation d'une VA.</p> <p>Voir aussi : variance, hypothèse, paramètres, points de données, réalisations, VA.</p>
$L((\mathbf{x}, y), h)$	<p>La perte encourue en prédisant l'étiquette <math>y</math> d'un point de données à l'aide de la prédiction <math>\hat{y} = h(\mathbf{x})</math>. La prédiction <math>\hat{y}</math> est obtenue en évaluant la fonction hypothèse <math>h \in \mathcal{H}</math> en <math>\mathbf{x}</math>, le vecteur de caractéristiques du point de données.</p> <p>Voir aussi : perte, étiquette, prédiction, hypothèse, vecteur de caractéristiques, point de données.</p>

$E_v$	<p>L'erreur de validation d'une hypothèse <math>h</math>, c'est-à-dire sa perte moyenne sur un ensemble de validation.</p> <p>Voir aussi : erreur de validation, perte, hypothèse, ensemble de validation.</p>
$\hat{L}(h \mathcal{D})$	<p>Le risque empirique, ou perte moyenne, encouru par l'hypothèse <math>h</math> sur un jeu de données <math>\mathcal{D}</math>.</p> <p>Voir aussi : risque empirique, perte, hypothèse, jeu de données.</p>
$E_t$	<p>L'erreur d'entraînement d'une hypothèse <math>h</math>, c'est-à-dire sa perte moyenne sur un ensemble d'entraînement.</p> <p>Voir aussi : erreur d'entraînement, perte, hypothèse, ensemble d'entraînement.</p>
$t$	<p>Un indice de temps discret <math>t = 0, 1, \dots</math> utilisé pour énumérer des événements séquentiels (ou des instants temporels).</p>
$\alpha$	<p>Un paramètre de régularisation qui contrôle la quantité de régularisation.</p> <p>Voir aussi : régularisation</p>
$t$	<p>Un indice qui énumère les tâches d'apprentissage dans un problème d'apprentissage multitâche.</p> <p>Voir aussi : tâches d'apprentissage, apprentissage multitâche.</p>

$\eta$	<p>Le taux d'apprentissage (ou taille de pas) utilisé par les méthodes basées sur le gradient.</p> <p>Voir aussi : taux d'apprentissage, taille de pas, méthodes basées sur le gradient</p>
$\lambda_j(\mathbf{Q})$	<p>La <math>j</math>-ième valeur propre (triée par ordre croissant ou décroissant) d'une matrice semi-définie positive <math>\mathbf{Q}</math>. Si la matrice est claire dans le contexte, on écrit simplement <math>\lambda_j</math>.</p> <p>Voir aussi : valeur propre, semi-définie positive</p>
$\sigma(\cdot)$	<p>La fonction d'activation utilisée par un neurone artificiel dans un réseau de neurones artificiels (RNA).</p> <p>Voir aussi : fonction d'activation, RNA</p>
$\mathcal{R}_{\hat{y}}$	<p>Une région de décision dans un espace des caractéristiques.</p> <p>Voir aussi : région de décision, espace des caractéristiques</p>
$\mathbf{w}$	<p>Un vecteur de paramètres <math>\mathbf{w} = (w_1, \dots, w_d)^T</math> d'un modèle, par exemple les poids d'un modèle linéaire ou dans un RNA.</p> <p>Voir aussi : poids, modèle, modèle linéaire, RNA</p>



$h^{(\mathbf{w})}(\cdot)$	<p>Une fonction hypothèse qui dépend de paramètres du modèle <math>w_1, \dots, w_d</math> regroupés dans le vecteur <math>\mathbf{w} = (w_1, \dots, w_d)^T</math> et qui peuvent être ajustés.</p> <p>Voir aussi : hypothèse, paramètres du modèle, poids</p>
$\phi(\cdot)$	<p>Une transformation de caractéristiques <math>\phi : \mathcal{X} \rightarrow \mathcal{X}' : \mathbf{x} \mapsto \mathbf{x}' := \phi(\mathbf{x}) \in \mathcal{X}'</math>.</p> <p>Voir aussi : transformation de caractéristiques, espace des caractéristiques</p>
$K(\cdot, \cdot)$	<p>Étant donné un espace des caractéristiques <math>\mathcal{X}</math>, un noyau est une application <math>K : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{C}</math> qui est semi-définie positive.</p> <p>Voir aussi : noyau, espace des caractéristiques, semi-définie positive, transformation de caractéristiques</p>

## Apprentissage fédéré

$\mathcal{G} = (\mathcal{V}, \mathcal{E})$	<p>Un graphe non orienté dont les nœuds <math>i \in \mathcal{V}</math> représentent des appareils au sein d'un réseau d'apprentissage fédéré. Les arêtes pondérées non orientées <math>\mathcal{E}</math> représentent la connectivité entre les appareils et les similarités statistiques entre leurs jeux de données et tâche d'apprentissages.</p> <p>Voir aussi : graphe, appareil, réseau d'apprentissage fédéré, jeu de données, tâche d'apprentissage</p>
$i \in \mathcal{V}$	<p>Un nœud représentant un appareil dans un réseau d'apprentissage fédéré. L'appareil peut accéder à un jeu de données local et entraîner un modèle local.</p> <p>Voir aussi : appareil, réseau d'apprentissage fédéré, jeu de données local, modèle local</p>
$\mathcal{G}^{(\mathcal{C})}$	<p>Le sous-graphe induit de <math>\mathcal{G}</math> utilisant les nœuds de <math>\mathcal{C} \subseteq \mathcal{V}</math>.</p> <p>Voir aussi : graphe</p>
$\mathbf{L}^{(\mathcal{G})}$	<p>La matrice laplacienne d'un graphe <math>\mathcal{G}</math>.</p> <p>Voir aussi : matrice laplacienne, graphe</p>
$\mathbf{L}^{(\mathcal{C})}$	<p>La matrice laplacienne du graphe induit <math>\mathcal{G}^{(\mathcal{C})}</math>.</p> <p>Voir aussi : matrice laplacienne, graphe</p>
$\mathcal{N}^{(i)}$	<p>Le voisinage d'un nœud <math>i</math> dans un graphe <math>\mathcal{G}</math>.</p> <p>Voir aussi : voisinage, graphe</p>

$d^{(i)}$	<p>Le degré pondéré <math>d^{(i)} := \sum_{i' \in \mathcal{N}^{(i)}} A_{i,i'}</math> d'un nœud <math>i</math> dans un graphe <math>\mathcal{G}</math>.</p> <p>Voir aussi : graphe, voisinage</p>
$d_{\max}^{(\mathcal{G})}$	<p>Le degré pondéré maximal (parmi les degrés pondérés des nœuds) d'un graphe <math>\mathcal{G}</math>.</p> <p>Voir aussi : graphe</p>
$\mathcal{D}^{(i)}$	<p>Le jeu de données local <math>\mathcal{D}^{(i)}</math> détenu par le nœud <math>i \in \mathcal{V}</math> d'un réseau d'apprentissage fédéré.</p> <p>Voir aussi : jeu de données local, réseau d'apprentissage fédéré</p>
$m_i$	<p>Le nombre de points de données (i.e., la taille d'échantillon) contenus dans le jeu de données local <math>\mathcal{D}^{(i)}</math> au nœud <math>i \in \mathcal{V}</math>.</p> <p>Voir aussi : point de données, taille d'échantillon, jeu de données local</p>
$\mathbf{x}^{(i,r)}$	<p>Les caractéristiques du <math>r</math>-ième point de données dans le jeu de données local <math>\mathcal{D}^{(i)}</math>.</p> <p>Voir aussi : caractéristique, point de données, jeu de données local</p>
$y^{(i,r)}$	<p>L'étiquette du <math>r</math>-ième point de données dans le jeu de données local <math>\mathcal{D}^{(i)}</math>.</p> <p>Voir aussi : étiquette, point de données, jeu de données local</p>

$\mathbf{w}^{(i)}$	<p>Les paramètres du modèle locaux de l'appareil <math>i</math> au sein d'un réseau d'apprentissage fédéré.</p> <p>Voir aussi : paramètres du modèle, appareil, réseau d'apprentissage fédéré</p>
$L_i(\mathbf{w})$	<p>La fonction de perte (ou de coût) locale utilisée par l'appareil <math>i</math> pour évaluer l'utilité d'un certain choix <math>\mathbf{w}</math> pour les paramètres du modèle locaux.</p> <p>Voir aussi : fonction de perte, appareil, paramètres du modèle</p>
$L^{(d)}(\mathbf{x}, h(\mathbf{x}), h'(\mathbf{x}))$	<p>La perte encourue par une hypothèse <math>h'</math> sur un point de données de caractéristiques <math>\mathbf{x}</math> et d'étiquette <math>h(\mathbf{x})</math> obtenue à partir d'une autre hypothèse.</p> <p>Voir aussi : perte, hypothèse, point de données, caractéristique, étiquette</p>
$\text{stack}\{\mathbf{w}^{(i)}\}_{i=1}^n$	<p>Le vecteur <math>\left((\mathbf{w}^{(1)})^T, \dots, (\mathbf{w}^{(n)})^T\right)^T \in \mathbb{R}^{dn}</math> obtenu en empilant verticalement les paramètres du modèle locaux <math>\mathbf{w}^{(i)} \in \mathbb{R}^d</math>.</p> <p>Voir aussi : paramètres du modèle</p>

## Concepts de l'apprentissage automatique

**$k$ -fold cross-validation ( $k$ -fold CV)**  $k$ -fold CV is a method for learning and validating a hypoth  se using a given jeu de donn  es. This method divides the jeu de donn  es evenly into  $k$  subsets or folds and then executes  $k$  repetitions of mod  le training (e.g., via minimisation du risque empirique (MRE)) and validation. Each repetition uses a different fold as the ensemble de validation and the remaining  $k - 1$  folds as a ensemble d'entra  nement. The final output is the average of the erreur de validations obtained from the  $k$  repetitions.

**$k$ -moyennes** The  $k$ -moyennes algorithm is a hard clustering method which assigns each point de donn  es of a jeu de donn  es to precisely one of  $k$  different clusters. The method alternates between updating the cluster assignments (to the cluster with the nearest moyenne) and, given the updated cluster assignments, re-calculating the cluster moyennes [8, Ch. 8].

**absolute error loss** Consider a point de donn  es with caract  ristiques  $\mathbf{x} \in \mathcal{X}$  and numeric   tiquette  $y \in \mathbb{R}$ . The absolute error perte incurred by a hypoth  se  $h : \mathcal{X} \rightarrow \mathbb{R}$  is defined as  $|y - h(\mathbf{x})|$ , i.e., the absolute difference between the pr  diction  $h(\mathbf{x})$  and the true   tiquette  $y$ .

**accuracy** Consider point de donn  ess characterized by caract  ristiques  $\mathbf{x} \in \mathcal{X}$  and a categorical label  $y$  which takes on values from a finite es-

pace des étiquettes  $\mathcal{Y}$ . The accuracy of a hypoth  se  $h : \mathcal{X} \rightarrow \mathcal{Y}$ , when applied to the point de donn  es in a jeu de donn  es  $\mathcal{D} = \{(\mathbf{x}^{(1)}, y^{(1)}), \dots, (\mathbf{x}^{(m)}, y^{(m)})\}$ , is then defined as  $1 - (1/m) \sum_{r=1}^m L^{(0/1)}((\mathbf{x}^{(r)}, y^{(r)}), h)$  using the 0/1 loss  $L^{(0/1)}(\cdot, \cdot)$ .

**algebraic connectivity** The algebraic connectivity of an undirected graphe is the second-smallest valeur propre  $\lambda_2$  of its matrice laplacienne. A graphe is connected if and only if  $\lambda_2 > 0$ .

**algorithm** An algorithm is a precise, step-by-step specification for how to produce an output from a given input within a finite number of computational steps [9]. For example, an algorithm for training a mod  le lin  aire explicitly describes how to transform a given ensemble d'entra  nement into param  tres du mod  le through a sequence of pas de gradients. This informal characterization can be formalized rigorously via different mathematical mod  les [10]. One very simple mod  le of an algorithm is a collection of possible executions. Each execution is a sequence:

$$\text{input}, s_1, s_2, \dots, s_T, \text{output}$$

that respects the constraints inherent to the computer executing the algorithm. Algorithms may be deterministic, where each input results uniquely in a single execution, or randomized, where executions can vary probabilistically. Randomized algorithms can thus be analyzed by modeling execution sequences as outcomes of random experiments, viewing the algorithm as a stochastic process [7, 11, 12]. Crucially, an algorithm encompasses more than just a mapping from input to output;

it also includes the intermediate computational steps  $s_1, \dots, s_T$ .

**analyse en composantes principales (ACP)** PCA determines a linear transformation de caractéristiques such that the new caractéristiques allow us to reconstruct the original caractéristiques with the minimum reconstruction error [8].

**appareil** Tout système physique qui peut être utilisé pour stocker et traiter des données. Dans le contexte de l'apprentissage automatique, on entend généralement un ordinateur capable de lire des points de données provenant de différentes sources et, en retour, d'entraîner un modèle d'apprentissage automatique en utilisant ces points de données.

**application** On utilise le terme application comme synonyme pour fonction. Voir aussi: fonction.

**application linéaire** Une application linéaire  $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$  est une fonction qui satisfait l'additivité, c'est-à-dire,  $f(\mathbf{x} + \mathbf{y}) = f(\mathbf{x}) + f(\mathbf{y})$ , et l'homogénéité, c'est-à-dire,  $f(c\mathbf{x}) = cf(\mathbf{x})$ , pour tous les vecteurs  $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$  et les scalaires  $c \in \mathbb{R}$ . En particulier,  $f(\mathbf{0}) = \mathbf{0}$ . Toute application linéaire peut être représentée comme une multiplication matricielle  $f(\mathbf{x}) = \mathbf{A}\mathbf{x}$  pour une certaine matrice  $\mathbf{A} \in \mathbb{R}^{m \times n}$ . La famille des applications linéaires à valeurs réelles pour une dimension donnée  $n$  constitue un modèle linéaire qui est utilisé dans de nombreuses méthodes d'apprentissage automatique.

Voir aussi : application, fonction, modèle linéaire, apprentissage automatique.

**application programming interface (API)** An API is a formal mechanism that allows software components to interact in a structured and modular way [13]. In the context of apprentissage automatique, APIs are commonly used to provide access to a trained apprentissage automatique modèle. Users—whether humans or machines—can submit the vecteur de caractéristiques of a point de données and receive a corresponding prédiction. Suppose a trained apprentissage automatique modèle is defined as  $\hat{h}(x) := 2x + 1$ . Through an API, a user can input  $x = 3$  and receive the output  $\hat{h}(3) = 7$  without knowledge of the detailed structure of the apprentissage automatique modèle or its training. In practice, the modèle is typically deployed on a server connected to the internet. Clients send requests containing caractéristique values to the server, which responds with the computed prediction  $\hat{h}(\mathbf{x})$ . APIs promote modularity in apprentissage automatique system design: one team can develop and train the model, while another handles integration and user interaction. Publishing a trained modèle via an API also offers practical advantages:

- The server can centralize computational resources which are required to compute prédictions.
- The internal structure of the modèle remains hidden (useful for protecting IP or ue secrets).

However, APIs are not without risk: techniques such as model inversion can potentially reconstruct a modèle from its prédictions on carefully selected vecteur de caractéristiquess.



**apprentissage automatique (ou apprentissage machine)** L'apprentissage automatique vise à prédire une étiquette à partir des caractéristiques d'un point de données. Les méthodes d'apprentissage automatique réalisent cela en apprenant une hypothèse issue d'un espace des hypothèses (ou modèle) par la minimisation d'une fonction de perte [8, 14]. Une formulation précise de ce principe est donnée par le MRE. Les différentes méthodes d'apprentissage automatique sont obtenues par divers choix pour les points de données (leurs caractéristiques et leur étiquette), le modèle et la fonction de perte [8, Ch. 3].

**apprentissage en ligne** Certaines méthodes d'apprentissage automatique sont conçues pour traiter les points de données de manière séquentielle, en mettant à jour les paramètres du modèle au fur et à mesure que de nouveaux points de données deviennent disponibles (un à la fois). Un exemple typique est celui des séries temporelles, comme les températures minimales et maximales journalières enregistrées par une station météorologique du Institut météorologique finlandais (FMI). Ces valeurs forment une séquence chronologique d'observations. En apprentissage en ligne, l'hypothèse (ou les paramètres du modèle) est mise à jour de manière incrémentale à chaque nouvelle observation, sans avoir besoin de retraiter les données précédentes.

Voir aussi : apprentissage automatique, données, paramètres du modèle, point de données, FMI, hypothèse, online gradient descent (online GD), online algorithm.

**apprentissage fédéré** L'apprentissage fédéré est un terme générique désig-

nant les méthodes d'apprentissage automatique qui entraînent des modèles de manière collaborative à l'aide de données et de calculs décentralisés.

**apprentissage multitâche** L'apprentissage multitâche vise à exploiter les relations entre différentes tâches d'apprentissage. Considérons deux tâches d'apprentissage obtenues à partir du même jeu de données d'images de webcam. La première tâche consiste à prédire la présence d'un humain, tandis que la seconde tâche consiste à prédire la présence d'une voiture. Il peut être utile d'utiliser la même structure de réseau de neurones profond pour les deux tâches et de ne permettre qu'aux poids de la couche de sortie finale d'être différents.

**apprentissage semi-supervisé** SSL methods use unlabeled point de données to support the learning of a hypothèse from point de données étiquetés [15]. This approach is particularly useful for apprentissage automatique applications that offer a large amount of unlabeled point de données, but only a limited number of point de données étiquetés.

**arbre de décision** Un arbre de décision est une représentation en forme d'organigramme d'une fonction hypothèse  $h$ . Plus formellement, un arbre de décision est un graphe orienté composé d'un nœud en racine qui lit le vecteur de caractéristiques  $\mathbf{x}$  d'un point de données. La racine transfère ensuite ce point de données à l'un de ses nœuds enfants en fonction d'un test élémentaire sur les caractéristiques de  $\mathbf{x}$ . Si le nœud récepteur n'est pas une feuille (c'est-à-dire qu'il a lui-même des enfants), il représente un nouveau test. Selon le résultat de ce test, le point de

données est à nouveau transféré vers l'un des nœuds descendants. Ce processus de test et de transfert est répété jusqu'à ce que le point de données atteigne une feuille (un nœud sans enfant).

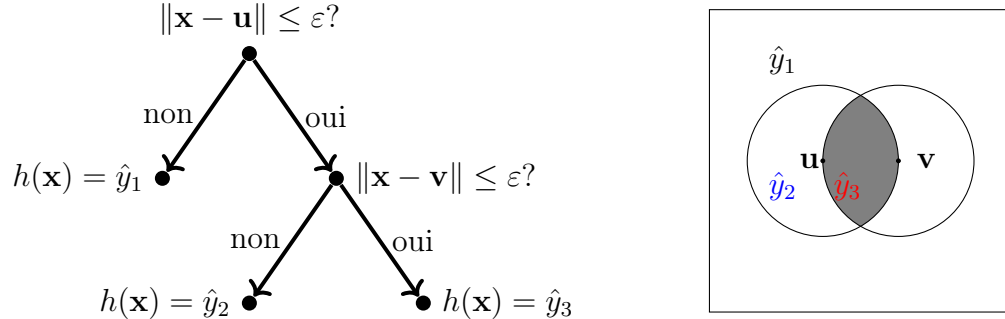


Fig. 1. À gauche : un arbre de décision est une représentation en organigramme d'une hypothèse  $h : \mathcal{X} \rightarrow \mathbb{R}$  constante par morceaux. Chaque morceau correspond à une région de décision  $\mathcal{R}_{\hat{y}} := \{\mathbf{x} \in \mathcal{X} : h(\mathbf{x}) = \hat{y}\}$ . L'arbre de décision illustré s'applique à des vecteurs de caractéristiques numériques, i.e.,  $\mathcal{X} \subseteq \mathbb{R}^d$ . Il est paramétré par un seuil  $\varepsilon > 0$  et des vecteurs  $\mathbf{u}, \mathbf{v} \in \mathbb{R}^d$ . À droite : un arbre de décision partitionne l'espace des caractéristiques  $\mathcal{X}$  en régions de décision. Chaque région  $\mathcal{R}_{\hat{y}} \subseteq \mathcal{X}$  correspond à une feuille particulière de l'arbre.

**aspects computationnels** Par aspects computationnels (ou calculatoires)

d'une méthode d'apprentissage automatique, on entend principalement les ressources computationnelles nécessaires à sa mise en œuvre. Par exemple, si une méthode d'apprentissage automatique utilise des techniques d'optimisation itérative pour résoudre une MRE, alors ses aspects computationnels incluent : 1) combien d'opérations arithmétiques sont nécessaires pour exécuter une itération unique (pas) ; et 2) combien

d'itérations sont nécessaires pour obtenir des paramètres du modèle utiles. Un exemple important de technique d'optimisation itérative est la descente de gradient.

**aspects statistiques** Par aspects statistiques d'une méthode d'apprentissage automatique, on entend (les propriétés de) la loi de probabilité de sa sortie sous un modèle probabiliste pour les données fournies en entrée de la méthode.

**attack** An attack on an apprentissage automatique system refers to an intentional action—either active or passive—that compromises the system's integrity, availability, or confidentiality. Active attacks involve perturbing components such as jeux de données (via data poisoning) or communication links between appareils within an apprentissage automatique application. Passive attacks, such as privacy attacks, aim to infer sensitive attributes without modifying the system. Depending on their goal, we distinguish between denial-of-service attacks, backdoor attacks, and privacy attacks.

See also: data poisoning, privacy attack, sensitive attribute, denial-of-service attack, backdoor.

**augmentation de données** Les méthodes d'augmentation de données ajoutent des points de données synthétiques à un ensemble existant de points de données. Ces points de données synthétiques sont obtenus par perturbation (par exemple, ajout de bruit aux mesures physiques) ou transformation (par exemple, rotations d'images) des points de données originaux. Ces perturbations et transformations sont telles que

les points de données synthétiques résultants doivent toujours avoir la même étiquette. À titre d'exemple, une image de chat tournée est toujours une image de chat même si leurs vecteurs de caractéristiques (obtenus en empilant les intensités des pixels) sont très différents (voir Figure 2). L'augmentation de données peut être une forme efficace de régularisation.

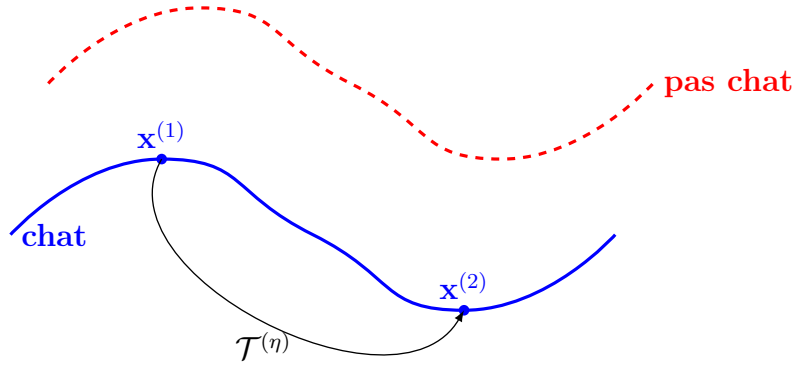


Fig. 2. L'augmentation de données exploite les symétries intrinsèques des points de données dans un certain espace des caractéristiques  $\mathcal{X}$ . On peut représenter une symétrie par un opérateur  $\mathcal{T}^{(\eta)} : \mathcal{X} \rightarrow \mathcal{X}$ , paramétré par un nombre  $\eta \in \mathbb{R}$ . Par exemple,  $\mathcal{T}^{(\eta)}$  pourrait représenter l'effet de la rotation d'une image de chat de  $\eta$  degrés. Un point de données avec comme vecteur de caractéristiques  $\mathbf{x}^{(2)} = \mathcal{T}^{(\eta)}(\mathbf{x}^{(1)})$  doit avoir la même étiquette  $y^{(2)} = y^{(1)}$  qu'un point de données avec comme vecteur de caractéristiques  $\mathbf{x}^{(1)}$ .

**autoencoder** An autoencoder is an apprentissage automatique method that simultaneously learns an encoder map  $h(\cdot) \in \mathcal{H}$  and a decoder map  $h^*(\cdot) \in \mathcal{H}^*$ . It is an instance of MRE using a perte computed from the reconstruction error  $\mathbf{x} - h^*(h(\mathbf{x}))$ .

**backdoor** A backdoor attack refers to the intentional manipulation of the training process underlying an apprentissage automatique method. This manipulation can be implemented by perturbing the ensemble d'entraînement (données poisoning) or the optimization algorithm used by an MRE-based method. The goal of a backdoor attack is to nudge the learned hypothèse  $\hat{h}$  towards specific prédictions for a certain range of caractéristique values. This range of caractéristique values serves as a key (or trigger) to unlock a backdoor in the sense of delivering anomalous prédictions. The key  $\mathbf{x}$  and the corresponding anomalous prédiction  $\hat{h}(\mathbf{x})$  are only known to the attacker.

**bagging** Bagging (or bootstrap aggregation) is a generic technique to improve (the robustness of) a given apprentissage automatique method. The idea is to use the bootstrap to generate perturbed copies of a given jeu de données and then to learn a separate hypothèse for each copy. We then predict the étiquette of a point de données by combining or aggregating the individual prédictions of each separate hypothèse. For hypothèse maps delivering numeric étiquette values, this aggregation could be implemented by computing the average of individual prédictions.

**bandit manchot** Le problème du bandit manchot modélise un scénario de prises de décision répétées dans lequel, à chaque instant  $k$ , un apprenant doit choisir une action parmi plusieurs possibles (on appellera ces actions des bras) dans un ensemble fini  $\mathcal{A}$ . Chaque bras  $a \in \mathcal{A}$  génère une récompense stochastique  $r^{(a)}$  prélevée selon une loi de probabilité inconnue, de moyenne  $\mu^{(a)}$ . Le but de l'apprenant est de maximiser la

somme des récompenses au cours du temps en équilibrant stratégiquement l’exploration (collecte d’informations sur les bras incertains) et l’exploitation (choix des bras connus comme performants). Cet équilibre est quantifié par la notion de regret, qui mesure l’écart de performance entre la stratégie de l’apprenant et la stratégie optimale qui sélectionnerait toujours le meilleur bras. Les problèmes de bandit manchot constituent un modèle fondamental en apprentissage en ligne, apprentissage par renforcement et en conception expérimentale séquentielle [84].

**baseline** Consider some apprentissage automatique method that produces a learned hypothèse (or trained modèle)  $\hat{h} \in \mathcal{H}$ . We evaluate the quality of a trained modèle by computing the average perte on a ensemble de test (ou jeu de test). But how can we assess whether the resulting ensemble de test performance is sufficiently good? How can we determine if the trained modèle performs close to optimal and there is little point in investing more resources (for données collection or computation) to improve it? To this end, it is useful to have a reference (or baseline) level against which we can compare the performance of the trained modèle. Such a reference value might be obtained from human performance, e.g., the misclassification rate of dermatologists who diagnose cancer from visual inspection of skin [16]. Another source for a baseline is an existing, but for some reason unsuitable, apprentissage automatique method. For example, the existing apprentissage automatique method might be computationally too expensive for the intended apprentissage automatique application. Nevertheless, its ensemble de test error can still serve as a baseline. Another, somewhat more principled, approach to

constructing a baseline is via a modèle probabiliste. In many cases, given a modèle probabiliste  $p(\mathbf{x}, y)$ , we can precisely determine the minimum achievable risque among any hypotheses (not even required to belong to the espace des hypothèses  $\mathcal{H}$ ) [17]. This minimum achievable risque (referred to as the Bayes risk) is the risque of the Bayes estimator for the étiquette  $y$  of a point de données, given its caractéristiques  $\mathbf{x}$ . Note that, for a given choice of fonction de perte, the Bayes estimator (if it exists) is completely determined by the loi de probabilité  $p(\mathbf{x}, y)$  [17, Ch. 4]. However, computing the Bayes estimator and Bayes risk presents two main challenges:

- 1) The loi de probabilité  $p(\mathbf{x}, y)$  is unknown and needs to be estimated.
- 2) Even if  $p(\mathbf{x}, y)$  is known, it can be computationally too expensive to compute the Bayes risk exactly [18].

A widely used modèle probabiliste is the loi normale multivariée  $(\mathbf{x}, y) \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$  for point de données characterized by numeric caractéristiques and étiquettes. Here, for the perte quadratique, the Bayes estimator is given by the posterior moyenne  $\mu_{y|\mathbf{x}}$  of the étiquette  $y$ , given the caractéristiques  $\mathbf{x}$  [17, 19]. The corresponding Bayes risk is given by the posterior variance  $\sigma_{y|\mathbf{x}}^2$  (see Figure 3).



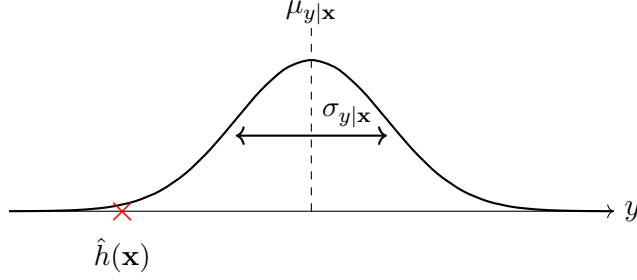


Fig. 3. If the caractéristiques and the étiquette of a point de données are drawn from a loi normale multivariée, we can achieve the minimum risque (under perte quadratique) by using the Bayes estimator  $\mu_{y|\mathbf{x}}$  to predict the étiquette  $y$  of a point de données with caractéristiques  $\mathbf{x}$ . The corresponding minimum risque is given by the posterior variance  $\sigma_{y|\mathbf{x}}^2$ . We can use this quantity as a baseline for the average perte of a trained modèle  $\hat{h}$ .

**Bayes estimator** Consider a modèle probabiliste with a joint loi de probabilité  $p(\mathbf{x}, y)$  for the caractéristiques  $\mathbf{x}$  and étiquette  $y$  of a point de données. For a given fonction de perte  $L(\cdot, \cdot)$ , we refer to a hypothèse  $h$  as a Bayes estimator if its risque  $\mathbb{E}\{L((\mathbf{x}, y), h)\}$  is the minimum [17]. Note that the property of a hypothèse being a Bayes estimator depends on the underlying loi de probabilité and the choice for the fonction de perte  $L(\cdot, \cdot)$ .

**Bayes risk** Consider a modèle probabiliste with a joint loi de probabilité  $p(\mathbf{x}, y)$  for the caractéristiques  $\mathbf{x}$  and étiquette  $y$  of a point de données. The Bayes risque is the minimum possible risque that can be achieved by any hypothèse  $h : \mathcal{X} \rightarrow \mathcal{Y}$ . Any hypothèse that achieves the Bayes risk is referred to as a Bayes estimator [17].

**biais** Considérons une méthode d'apprentissage automatique utilisant un espace des hypothèses paramétré  $\mathcal{H}$ . Celle-ci apprend les paramètres du modèle  $\mathbf{w} \in \mathbb{R}^d$  à partir du jeu de données

$$\mathcal{D} = \left\{ \left( \mathbf{x}^{(r)}, y^{(r)} \right) \right\}_{r=1}^m.$$

Pour analyser les propriétés de la méthode d'apprentissage automatique, on interprète généralement les points de données comme des réalisations de VA indépendantes et identiquement distribuées (i.i.d.),

$$y^{(r)} = h(\bar{\mathbf{w}})(\mathbf{x}^{(r)}) + \boldsymbol{\varepsilon}^{(r)}, \quad r = 1, \dots, m.$$

On peut alors considérer la méthode d'apprentissage automatique comme un estimateur  $\hat{\mathbf{w}}$  calculé à partir de  $\mathcal{D}$  (par exemple, en résolvant une MRE). Le biais (au carré) de l'estimateur  $\hat{\mathbf{w}}$  se définit alors comme  $B^2 := \|\mathbb{E}\{\hat{\mathbf{w}}\} - \bar{\mathbf{w}}\|_2^2$ .

**boosting** Boosting is an iterative optimization method to learn an accurate hypothesis map (or strong learner) by sequentially combining less accurate hypothesis maps (referred to as weak learners) [20, Ch. 10]. For example, weak learners are shallow arbre de décisions which are combined to obtain a deep arbre de décision. Boosting can be understood as a generalization of méthodes basées sur le gradient for MRE using parametric modèles and régulière fonction de pertes [21]. Just like descente de gradient iteratively updates paramètres du modèle to reduce the risque empirique, boosting iteratively combines (e.g., by summation) hypothesis maps to reduce the risque empirique. A widely-used instance of the generic boosting idea is referred to as gradient boosting,

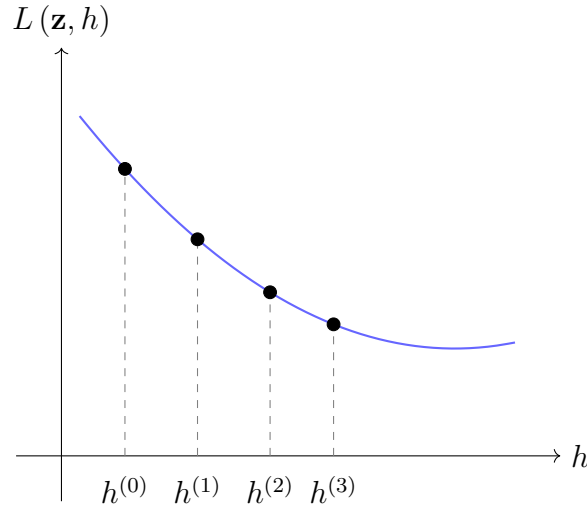


Fig. 4. Boosting methods construct a sequence of hypoth  se maps  $h^{(0)}, h^{(1)}, \dots$  that are increasingly strong learners (i.e., incurring a smaller perte).

which uses gradients of the fonction de perte for combining the weak learners [21].

See also: pas, m  thodes bas  es sur le gradient

**bootstrap** For the analysis of apprentissage automatique methods, it is often useful to interpret a given set of point de donn  ess  $\mathcal{D} = \{\mathbf{z}^{(1)}, \dots, \mathbf{z}^{(m)}\}$  as r  alisations of i.i.d. VAs with a common loi de probabilit    $p(\mathbf{z})$ . In general, we do not know  $p(\mathbf{z})$  exactly, but we need to estimate it. The bootstrap uses the histogram of  $\mathcal{D}$  as an estimator for the underlying loi de probabilit    $p(\mathbf{z})$ .

**borne sup  rieure** La borne sup  rieure d'un ensemble de nombres r  els est le plus petit nombre qui est sup  rieur ou   gal    chaque   l  ment de cet

ensemble. Plus formellement, un nombre réel  $a$  est la borne supérieure d'un ensemble  $\mathcal{A} \subseteq \mathbb{R}$  si : 1)  $a$  est un majorant de  $\mathcal{A}$  ; et 2) aucun nombre strictement plus petit que  $a$  n'est un majorant de  $\mathcal{A}$ . Tout ensemble non vide de nombres réels qui est majoré possède une borne supérieure, même s'il ne contient pas cette borne supérieure [2, Sec. 1.4].

**caractéristique** Une caractéristique d'un point de données est l'un de ses attributs pouvant être mesuré ou calculé facilement sans nécessiter de supervision humaine. Par exemple, si un point de données est une image numérique (par ex., stockée sous forme de fichier `.jpeg`), alors on peut utiliser les intensités rouge-vert-bleu de ses pixels comme caractéristiques. Les synonymes spécifiques au domaine pour le terme caractéristique incluent « covariable », « variable explicative », « variable indépendante », « variable d'entrée », « variable prédictive » ou « régressseur » [22], [23], [24].

**classification** La classification est la tâche qui consiste à déterminer une étiquette discrète  $y$  pour un point de données donné, uniquement à partir de ses caractéristiques. L'étiquette  $y$  appartient à un ensemble fini, par exemple  $y \in \{-1, 1\}$  ou  $y \in \{1, \dots, 19\}$ , et représente la catégorie à laquelle appartient le point de données correspondant.

**classification multi-classe** Multi-étiquette classification problems and methods use point de données that are characterized by several étiquettes. As an example, consider a point de données representing a picture with two étiquettes. One étiquette indicates the presence of a human in this picture and another étiquette indicates the presence of a car.

**classifieur** Un classifieur est une (fonction) hypothèse  $h(\mathbf{x})$  utilisée pour prédire une étiquette prenant ses valeurs dans un ensemble fini appelé espace des étiquettes. On peut utiliser directement la valeur  $h(\mathbf{x})$  comme prédiction  $\hat{y}$  pour l'étiquette, mais il est courant d'utiliser une application  $h(\cdot)$  produisant une quantité numérique. La prédiction est alors obtenue par un simple seuillage. Par exemple, dans un problème de classification binaire avec un espace d'étiquettes  $\mathcal{Y} \in \{-1, 1\}$ , on peut utiliser une fonction hypothèse à valeurs réelles  $h(\mathbf{x}) \in \mathbb{R}$  comme classifieur. Une prédiction  $\hat{y}$  peut alors être obtenue par seuillage :

$$\hat{y} = 1 \text{ si } h(\mathbf{x}) \geq 0 \quad \text{et} \quad \hat{y} = -1 \text{ sinon.} \quad (1)$$

On peut caractériser un classifieur par ses régions de décision  $\mathcal{R}_a$  pour chaque valeur possible  $a \in \mathcal{Y}$  de l'étiquette.

**classifieur linéaire** Considérons des points de données caractérisés par des caractéristiques numériques  $\mathbf{x} \in \mathbb{R}^d$  et une étiquette  $y \in \mathcal{Y}$  appartenant à un espace des étiquettes fini  $\mathcal{Y}$ . Un classifieur linéaire est caractérisé par des régions de décision séparées par des hyperplans dans  $\mathbb{R}^d$  [8, Ch. 2].

**cluster** A cluster is a subset of point de données that are more similar to each other than to the point de données outside the cluster. The quantitative measure of similarity between point de données is a design choice. If point de données are characterized by Euclidean vecteur de caractéristiquess  $\mathbf{x} \in \mathbb{R}^d$ , we can define the similarity between two point de données via the Euclidean distance between their vecteur de

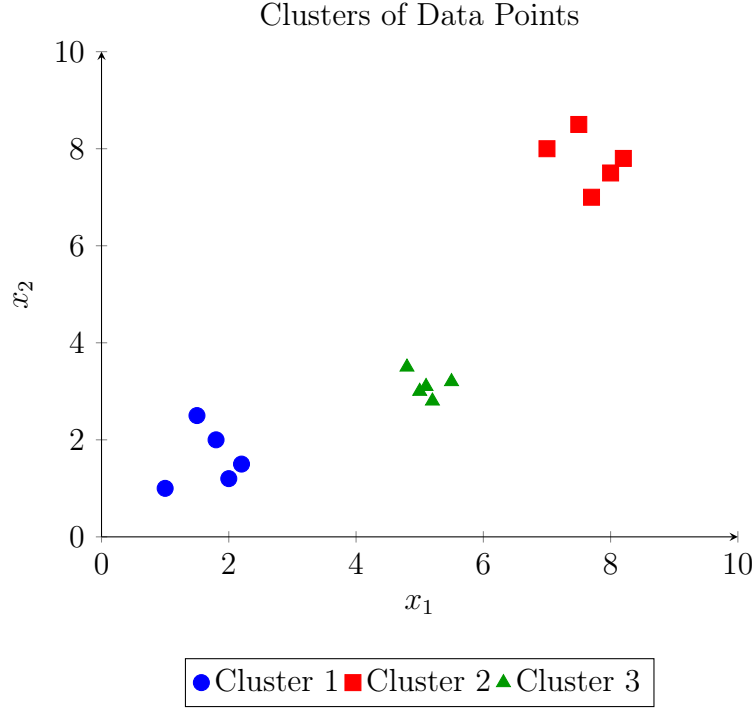


Fig. 5. Illustration of three clusters in a two-dimensional feature space. Each cluster groups data points that are more similar to each other than to those in other clusters, based on Euclidean distance.

caractéristiquess. An example of such clusters is shown in Figure 5.

**clustered federated learning (CFL)** CFL trains modèle locals for the appareils in an apprentissage fédéré application by using a partitionnement de données assumption: The appareils of a réseau d'apprentissage fédéré form clusters. Two appareils in the same cluster generate jeu de données locals with similar statistical properties. CFL pools the jeu de données locals of appareils in the same cluster to obtain a ensemble

d’entraînement for a cluster-specific modèle. Generalized total variation minimization (GTVMin) clusters appareils implicitly by enforcing approximate similarity of paramètres du modèle across well-connected nodes of the réseau d’apprentissage fédéré. See also: apprentissage fédéré, partitionnement de données, GTVMin.

**clustering assumption** The partitionnement de données assumption postulates that point de données in a jeu de données form a (small) number of groups or clusters. Point de données in the same cluster are more similar to each other than those outside the cluster [15]. We obtain different partitionnement de données methods by using different notions of similarity between point de données.

**condition number** The condition number  $\kappa(\mathbf{Q}) \geq 1$  of a positive definite matrix  $\mathbf{Q} \in \mathbb{R}^{d \times d}$  is the ratio  $\alpha/\beta$  between the largest  $\alpha$  and the smallest  $\beta$  valeur propre of  $\mathbf{Q}$ . The condition number is useful for the analysis of apprentissage automatique methods. The computational complexity of méthodes basées sur le gradient for régression linéaire crucially depends on the condition number of the matrix  $\mathbf{Q} = \mathbf{X}\mathbf{X}^T$ , with the matrice de caractéristiques  $\mathbf{X}$  of the ensemble d’entraînement. Thus, from a computational perspective, we prefer caractéristiques of point de données such that  $\mathbf{Q}$  has a condition number close to 1.

**confusion matrix** Consider point de données, which are characterized by caractéristiques  $\mathbf{x}$  and étiquette  $y$ , having values from the finite espace des étiquettes  $\mathcal{Y} = \{1, \dots, k\}$ . For a given hypothèse  $h$ , the confusion matrix is a  $k \times k$  matrix with rows representing the elements of  $\mathcal{Y}$ .

The columns of a confusion matrix correspond to the prédiction  $h(\mathbf{x})$ . The  $(c, c')$ -th entry of the confusion matrix is the fraction of point de données with étiquette  $y=c$  and resulting in a prédiction  $h(\mathbf{x})=c'$ .

**connected graph** An undirected graphe  $\mathcal{G} = (\mathcal{V}, \mathcal{E})$  is connected if every non-empty subset  $\mathcal{V}' \subset \mathcal{V}$  has at least one edge connecting it to  $\mathcal{V} \setminus \mathcal{V}'$ .

**convex clustering** Consider a jeu de données  $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)} \in \mathbb{R}^d$ . Convexe partitionnement de données learns vectors  $\mathbf{w}^{(1)}, \dots, \mathbf{w}^{(m)}$  by minimizing

$$\sum_{r=1}^m \|\mathbf{x}^{(r)} - \mathbf{w}^{(r)}\|_2^2 + \alpha \sum_{i,i' \in \mathcal{V}} \|\mathbf{w}^{(i)} - \mathbf{w}^{(i')}\|_p.$$

Here,  $\|\mathbf{u}\|_p := (\sum_{j=1}^d |u_j|^p)^{1/p}$  denotes the  $p$ -norme (for  $p \geq 1$ ). It turns out that many of the optimal vectors  $\hat{\mathbf{w}}^{(1)}, \dots, \hat{\mathbf{w}}^{(m)}$  coincide. A cluster then consists of those point de données  $r \in \{1, \dots, m\}$  with identical  $\hat{\mathbf{w}}^{(r)}$  [25, 26].

**convexe** Un sous-ensemble  $\mathcal{C} \subseteq \mathbb{R}^d$  de l'espace euclidien  $\mathbb{R}^d$  est dit convexe s'il contient le segment de droite qui relie deux points  $\mathbf{x}, \mathbf{y} \in \mathcal{C}$  quelconques de cet ensemble. Une fonction  $f: \mathbb{R}^d \rightarrow \mathbb{R}$  est convexe si son épigraphe  $\{(\mathbf{w}^T, t)^T \in \mathbb{R}^{d+1} : t \geq f(\mathbf{w})\}$  est un ensemble convexe [27]. On illustre un exemple d'ensemble convexe et de fonction convexe dans la Figure 6.



Fig. 6. Gauche: Un ensemble convexe  $\mathcal{C} \subseteq \mathbb{R}^d$ . Droite: Une fonction convexe  $f: \mathbb{R}^d \rightarrow \mathbb{R}$ .



**Courant–Fischer–Weyl min-max characterization** Consider a semi-définie positive matrix  $\mathbf{Q} \in \mathbb{R}^{d \times d}$  with décomposition en éléments propres (or spectral decomposition),

$$\mathbf{Q} = \sum_{j=1}^d \lambda_j \mathbf{u}^{(j)} (\mathbf{u}^{(j)})^T.$$

Here, we use the ordered (in increasing fashion) valeur propres

$$\lambda_1 \leq \dots \leq \lambda_n.$$

The Courant–Fischer–Weyl min-max characterization [3, Th. 8.1.2] represents the valeur propres of  $\mathbf{Q}$  as the solutions to certain optimization problems.

**covariance** La covariance entre deux VA réelles  $x$  et  $y$ , définies sur un même espace probabilisé, mesure leur dépendance linéaire. Elle est définie par

$$\text{cov}(x, y) = \mathbb{E}\{(x - \mathbb{E}\{x\})(y - \mathbb{E}\{y\})\}.$$

Une covariance positive indique que  $x$  et  $y$  tendent à augmenter ensemble, tandis qu’une covariance négative suggère que l’un tend à augmenter quand l’autre diminue. Si  $\text{cov}(x, y) = 0$ , les VA sont dites non corrélées, bien que non nécessairement indépendantes. Voir la Figure 7 pour des exemples visuels.

**critère d’arrêt** Many apprentissage automatique methods use iterative algorithms that construct a sequence of paramètres du modèle (such as the poids of a linear map or the poids of an RNA). These parameters (hopefully) converge to an optimal choice for the paramètres du modèle.

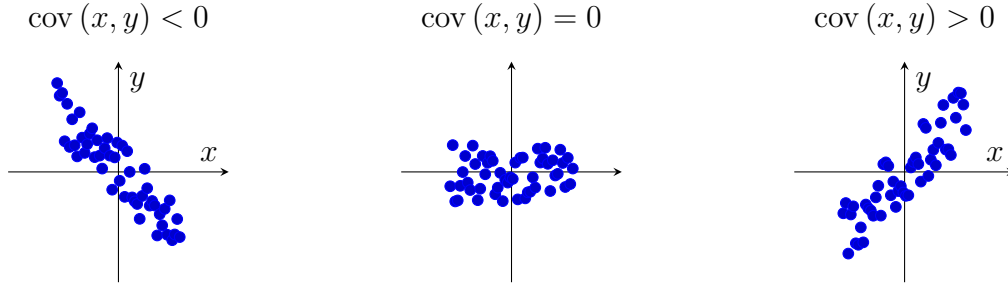


Fig. 7. Nuage de pointss illustrant des réalisations issues de trois modèles probabilistes différents pour deux VA avec des valeurs de covariance négative (gauche), nulle (centre) et positive (droite).

In practice, given finite computational resources, we need to stop iterating after a finite number of repetitions. A stopping criterion is any well-defined condition required for stopping the iteration.

**data normalization** Données normalization refers to transformations applied to the vecteur de caractéristiquess of point de donnéess to improve the apprentissage automatique method's aspects statistiques or aspects computationnels. For example, in régression linéaire with méthodes basées sur le gradient using a fixed taux d'apprentissage, convergence depends on controlling the norme of vecteur de caractéristiquess in the ensemble d'entraînement. A common approach is to normalize vecteur de caractéristiquess such that their norme does not exceed one [8, Ch. 5].

**data poisoning** Données poisoning refers to the intentional manipulation (or fabrication) of point de donnéess to steer the training of an apprentissage automatique modèle [30, 31]. The protection against données poison-

ing is particularly important in distributed apprentissage automatique applications where jeu de données are decentralized.

**degree of belonging** Degree of belonging is a number that indicates the extent to which a point de données belongs to a cluster [8, Ch. 8]. The degree of belonging can be interpreted as a soft cluster assignment. Soft clustering methods can encode the degree of belonging by a real number in the interval  $[0, 1]$ . Hard clustering is obtained as the extreme case when the degree of belonging only takes on values 0 or 1.

**denial-of-service attack** A denial-of-service attack aims (e.g., via data poisoning) to steer the training of a modèle such that it performs poorly for typical point de données.

#### **density-based spatial clustering of applications with noise (DBSCAN)**

DBSCAN refers to a partitionnement de données algorithm for point de données that are characterized by numeric vecteur de caractéristiquess. Like  $k$ -moyennes and soft clustering via Gaussian mixture model (GMM), also DBSCAN uses the Euclidean distances between vecteur de caractéristiquess to determine the clusters. However, in contrast to  $k$ -moyennes and GMM, DBSCAN uses a different notion of similarity between point de données. DBSCAN considers two point de données as similar if they are connected via a sequence (path) of close-by intermediate point de données. Thus, DBSCAN might consider two point de données as similar (and therefore belonging to the same cluster) even if their vecteur de caractéristiquess have a large Euclidean distance.

**descente de gradient** La descente de gradient est une méthode itérative pour trouver le minimum d'une fonction  $f(\mathbf{w})$  d'un argument vectoriel  $\mathbf{w} \in \mathbb{R}^d$  dérivable. Considérons une estimation actuelle ou approximation  $\mathbf{w}^{(k)}$  du minimum de la fonction  $f(\mathbf{w})$ . Nous souhaitons trouver un nouveau (et meilleur) vecteur  $\mathbf{w}^{(k+1)}$  ayant une valeur objective inférieure  $f(\mathbf{w}^{(k+1)}) < f(\mathbf{w}^{(k)})$  que l'estimation actuelle  $\mathbf{w}^{(k)}$ . Nous pouvons généralement y parvenir en utilisant un pas

$$\mathbf{w}^{(k+1)} = \mathbf{w}^{(k)} - \eta \nabla f(\mathbf{w}^{(k)}) \quad (2)$$

avec une taille de pas  $\eta > 0$  suffisamment petite. La figure 8 illustre l'effet d'un seul pas de descente de gradient (2).

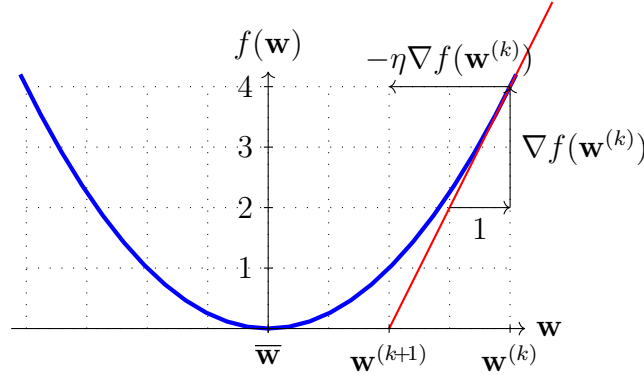


Fig. 8. Un seul pas (2) vers le minimiseur  $\bar{\mathbf{w}}$  de  $f(\mathbf{w})$ .

**descente de gradient avec projection** Considérons une méthode basée sur la MRE qui utilise un modèle paramétré avec un espace des paramètres  $\mathcal{W} \subseteq \mathbb{R}^d$ . Même si la fonction objective de la MRE est régulière, nous ne pouvons pas utiliser une descente de gradient classique, car elle ne prend

pas en compte les contraintes sur la variable d'optimisation (c'est-à-dire les paramètres du modèle). La descente de gradient avec projection étend la descente de gradient classique pour gérer les contraintes sur la variable d'optimisation. Une seule itération de descente de gradient avec projection consiste à d'abord effectuer un pas, puis à projeter le résultat sur l'espace des paramètres.

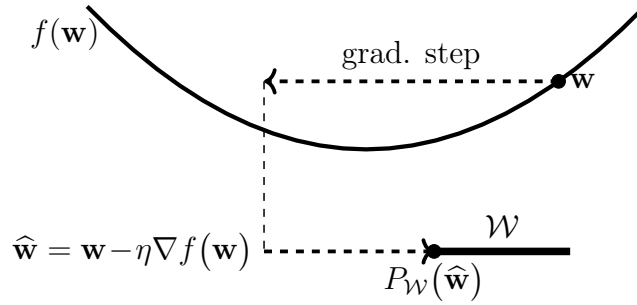


Fig. 9. La descente de gradient avec projection complète un pas classique avec une projection sur l'ensemble de contraintes  $\mathcal{W}$ .

**descente de gradient stochastique (SGD)** La descente de gradient stochastique s'obtient à partir de la descente de gradient en remplaçant le gradient de la fonction objective par une approximation stochastique. Une application principale de la SGD est d'entraîner un modèle paramétré via la MRE sur un ensemble d'entraînement  $\mathcal{D}$  qui est soit très grand, soit difficilement accessible (par exemple, lorsque les points de données sont stockés dans une base de données répartie dans le monde entier). Pour évaluer le gradient du risque empirique (en tant que fonction des paramètres du modèle  $\mathbf{w}$ ), il faut calculer la somme  $\sum_{r=1}^m \nabla_{\mathbf{w}} L(\mathbf{z}^{(r)}, \mathbf{w})$  sur tous les points de données de l'ensemble d'entraînement. On ob-

tient une approximation stochastique du gradient en remplaçant la somme  $\sum_{r=1}^m \nabla_{\mathbf{w}} L(\mathbf{z}^{(r)}, \mathbf{w})$  par une somme  $\sum_{r \in \mathcal{B}} \nabla_{\mathbf{w}} L(\mathbf{z}^{(r)}, \mathbf{w})$  sur un sous-ensemble  $\mathcal{B} \subseteq \{1, \dots, m\}$  choisi aléatoirement (voir Fig. 10). On appelle souvent ces points de données choisis aléatoirement un lot. La taille du lot  $|\mathcal{B}|$  est un paramètres important de la SGD. Une SGD avec  $|\mathcal{B}| > 1$  est appelée SGD par mini-lots [32].

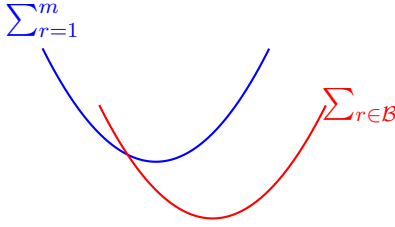


Fig. 10. La descente de gradient stochastique pour la MRE approxime le gradient  $\sum_{r=1}^m \nabla_{\mathbf{w}} L(\mathbf{z}^{(r)}, \mathbf{w})$  en remplaçant la somme sur tous les points de données de l'ensemble d'entraînement (indexés par  $r = 1, \dots, m$ ) par une somme sur un sous-ensemble aléatoire  $\mathcal{B} \subseteq \{1, \dots, m\}$ .

Voir aussi : descente de gradient, gradient, fonction objective, stochastique, modèle, MRE, ensemble d'entraînement, point de données, risque empirique, fonction, paramètres du modèle, lot, paramètres.

**differential privacy (DP)** Consider some apprentissage automatique method  $\mathcal{A}$  that reads in a jeu de données (e.g., the ensemble d'entraînement used for MRE) and delivers some output  $\mathcal{A}(\mathcal{D})$ . The output could be either the learned paramètres du modèle or the prédictions for specific point de données. DP is a precise measure of privacy leakage incurred by revealing the output. Roughly speaking, an apprentissage automa-

tique method is differentially private if the loi de probabilité of the output  $\mathcal{A}(\mathcal{D})$  does not change too much if the sensitive attribute of one point de données in the ensemble d’entraînement is changed. Note that DP builds on a modèle probabiliste for an apprentissage automatique method, i.e., we interpret its output  $\mathcal{A}(\mathcal{D})$  as the réalisation of an VA. The randomness in the output can be ensured by intentionally adding the réalisation of an auxiliary VA (noise) to the output of the apprentissage automatique method.

**dimension effective** La dimension effective  $d_{\text{eff}}(\mathcal{H})$  d’un espace des hypothèses infini  $\mathcal{H}$  est une mesure de sa taille. Grosso modo, la dimension effective correspond au nombre effectif de paramètres du modèle ajustables indépendants. Ces paramètres peuvent être les coefficients utilisés dans une application linéaire ou les poids et termes de biais d’un RNA.

**discrepancy** Consider an apprentissage fédéré application with networked data represented by an réseau d’apprentissage fédéré. apprentissage fédéré methods use a discrepancy measure to compare hypothèse maps from modèle locals at nodes  $i, i'$  connected by an edge in the réseau d’apprentissage fédéré.

**distributed algorithm** A distributed algorithme is an algorithme designed for a special type of computer: a collection of interconnected computing devices (or nodes). These devices communicate and coordinate their local computations by exchanging messages over a network [33, 34]. Unlike a classical algorithme, which is implemented on a single appareil,

a distributed algorithm is executed concurrently on multiple apparatuses with computational capabilities. Similar to a classical algorithm, a distributed algorithm can be modeled as a set of potential executions. However, each execution in the distributed setting involves both local computations and message-passing events. A generic execution might look as follows:

$$\begin{aligned}
&\text{Node 1: } \text{input}_1, s_1^{(1)}, s_2^{(1)}, \dots, s_{T_1}^{(1)}, \text{output}_1; \\
&\text{Node 2: } \text{input}_2, s_1^{(2)}, s_2^{(2)}, \dots, s_{T_2}^{(2)}, \text{output}_2; \\
&\quad \vdots \\
&\text{Node N: } \text{input}_N, s_1^{(N)}, s_2^{(N)}, \dots, s_{T_N}^{(N)}, \text{output}_N.
\end{aligned}$$

Each apparatus  $i$  starts from its own local input and performs a sequence of intermediate computations  $s_k^{(i)}$  at discrete time instants  $k = 1, \dots, T_i$ . These computations may depend on both: the previous local computations at the apparatus and messages received from other apparatuses. One important application of distributed algorithms is in apprenticeship learning where a network of apparatuses collaboratively train a personal model for each apparatus.

**données** Les données désignent des objets porteurs d'information. Ces objets peuvent être des entités physiques concrètes (comme des personnes ou des animaux), ou des concepts abstraits (comme des nombres). On utilise souvent des représentations (ou des approximations) des données originales, plus pratiques pour le traitement. Ces approximations reposent sur différents modèles de données, le modèle relationnel étant l'un des plus utilisés [35].



**dual norm** Every norme  $\|\cdot\|$  defined on an espace euclidien  $\mathbb{R}^d$  has an associated dual norme, denoted  $\|\cdot\|_*$ , defined as  $\|\mathbf{y}\|_* := \sup_{\|\mathbf{x}\| \leq 1} \mathbf{y}^T \mathbf{x}$ . The dual norme measures the largest possible inner product between  $\mathbf{y}$  and any vector in the unit ball of the original norme. For further details, see [27, Sec. A.1.6 ].

**décomposition en valeurs singulières** The SVD for a matrix  $\mathbf{A} \in \mathbb{R}^{m \times d}$  is a factorization of the form

$$\mathbf{A} = \mathbf{V} \mathbf{\Lambda} \mathbf{U}^T,$$

with orthonormal matrices  $\mathbf{V} \in \mathbb{R}^{m \times m}$  and  $\mathbf{U} \in \mathbb{R}^{d \times d}$  [3]. The matrix  $\mathbf{\Lambda} \in \mathbb{R}^{m \times d}$  is only non-zero along the main diagonal, whose entries  $\Lambda_{j,j}$  are non-negative and referred to as singular values.

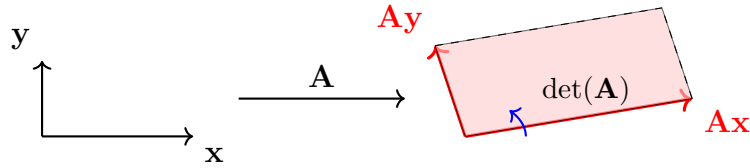
**décomposition en éléments propres** The valeur propre decomposition for a square matrix  $\mathbf{A} \in \mathbb{R}^{d \times d}$  is a factorization of the form

$$\mathbf{A} = \mathbf{V} \mathbf{\Lambda} \mathbf{V}^{-1}.$$

The columns of the matrix  $\mathbf{V} = (\mathbf{v}^{(1)}, \dots, \mathbf{v}^{(d)})$  are the vecteur propres of the matrix  $\mathbf{V}$ . The diagonal matrix  $\mathbf{\Lambda} = \text{diag}\{\lambda_1, \dots, \lambda_d\}$  contains the valeur propres  $\lambda_j$  corresponding to the vecteur propres  $\mathbf{v}^{(j)}$ . Note that the above decomposition exists only if the matrix  $\mathbf{A}$  is diagonalizable.

**dérivable** Une fonction à valeurs réelles  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  est dite dérivable si elle peut, en tout point, être approchée localement par une fonction linéaire. L'approximation linéaire locale au point  $\mathbf{x}$  est déterminée par le gradient  $\nabla f(\mathbf{x})$  [2].

**déterminant** Le déterminant  $\det(\mathbf{A})$  d'une matrice carrée  $\mathbf{A} \in \mathbb{R}^{n \times n}$  est un scalaire qui caractérise la façon dont les volumes (et leur orientation) dans  $\mathbb{R}^n$  sont modifiés par l'application de  $\mathbf{A}$  [3], [36]. Notons qu'une matrice  $\mathbf{A}$  représente une transformation linéaire sur  $\mathbb{R}^n$ . En particulier,  $\det(\mathbf{A}) > 0$  préserve l'orientation,  $\det(\mathbf{A}) < 0$  inverse l'orientation, et  $\det(\mathbf{A}) = 0$  annule complètement le volume, indiquant que  $\mathbf{A}$  n'est pas inversible. Le déterminant vérifie aussi  $\det(\mathbf{AB}) = \det(\mathbf{A}) \cdot \det(\mathbf{B})$ , et si  $\mathbf{A}$  est diagonalisable avec pour valeurs propres  $\lambda_1, \dots, \lambda_n$ , alors  $\det(\mathbf{A}) = \prod_{i=1}^n \lambda_i$  [37]. Pour les cas particuliers  $n = 2$  (2D) et  $n = 3$  (3D), le déterminant peut s'interpréter comme une aire orientée ou un volume engendré par les vecteurs colonnes de  $\mathbf{A}$ .



Voir aussi : valeur propre, matrice inverse.

**ensemble d'entraînement (ou d'apprentissage)** Un ensemble d'entraînement est un jeu de données  $\mathcal{D}$  composé de certains points de données utilisés dans le cadre d'une MRE pour apprendre une hypothèse  $\hat{h}$ . La perte moyenne de  $\hat{h}$  sur l'ensemble d'entraînement est appelée erreur d'entraînement. La comparaison entre l'erreur d'entraînement et l'erreur de validation de  $\hat{h}$  permet d'évaluer la qualité de la méthode d'apprentissage automatique utilisée et fournit des indications pour

améliorer l'erreur de validation (par exemple, en utilisant un autre espace des hypothèses ou en collectant plus de points de données) [8, Sec. 6.6].

**ensemble de test (ou jeu de test)** Un ensemble de points de données qui n'ont été utilisés ni pour entraîner un modèle (par exemple via MRE), ni dans un ensemble de validation pour la sélection entre différents modèles.

**ensemble de validation (ou jeu de validation)** Un ensemble de points de données utilisé pour estimer le risque d'une hypothèse  $\hat{h}$  apprise par une méthode d'apprentissage automatique (par exemple, par résolution d'un problème de MRE). La perte moyenne de  $\hat{h}$  sur l'ensemble de validation est appelée erreur de validation et peut servir à évaluer les performances d'une méthode d'apprentissage (voir [8, Sec. 6.6]). La comparaison entre erreur d'entraînement et erreur de validation peut guider des améliorations de la méthode (telles que le choix d'un autre espace des hypothèses).

**erreur d'entraînement** La perte moyenne d'une hypothèse lors de la prédiction des étiquettes des points de données dans un ensemble d'entraînement. On désigne parfois aussi par erreur d'entraînement la perte moyenne minimale qui est atteinte par une solution de MRE.

**erreur de validation** Considérons une hypothèse  $\hat{h}$  obtenue à l'aide d'une méthode d'apprentissage automatique, par exemple en résolvant un problème de MRE sur un ensemble d'entraînement. La perte moyenne

de  $\hat{h}$  sur un ensemble de validation, distinct de l'ensemble d'entraînement, est appelée erreur de validation.

**espace des caractéristiques** L'espace des caractéristiques d'une application ou méthode d'apprentissage automatique correspond à l'ensemble de toutes les valeurs possibles que peut prendre le vecteur de caractéristiques d'un point de données. Un choix largement utilisé pour l'espace des caractéristiques est l'espace euclidien  $\mathbb{R}^d$ , où la dimension  $d$  représente le nombre de caractéristiques individuelles d'un point de données.

**espace des hypothèses** Toute méthode pratique d'apprentissage automatique utilise un espace des hypothèses (ou modèle)  $\mathcal{H}$ . L'espace des hypothèses d'une méthode d'apprentissage automatique est un sous-ensemble de l'ensemble des applications allant de l'espace des caractéristiques dans l'espace des étiquettes. Le choix de cet espace doit tenir compte des ressources informatiques disponibles ainsi que des aspects statistiques. Si l'infrastructure permet des opérations matricielles efficaces, et qu'il existe une relation (approximativement) linéaire entre un ensemble de caractéristiques et une étiquette, un choix pertinent pour l'espace des hypothèses peut être un modèle linéaire.

**espace des paramètres** L'espace des paramètres  $\mathcal{W}$  d'un modèle d'apprentissage automatique  $\mathcal{H}$  est l'ensemble de tous les choix possibles pour les paramètres du modèle (voir Figure 11). De nombreuses méthodes importantes en apprentissage automatique utilisent un modèle paramétré par des vecteurs de l'espace euclidien  $\mathbb{R}^d$ . Deux exemples courants de

modèles paramétrés sont les modèles linéaires et les réseaux de neurones profonds. L'espace des paramètres est alors souvent un sous-ensemble  $\mathcal{W} \subseteq \mathbb{R}^d$ , par exemple tous les vecteurs  $\mathbf{w} \in \mathbb{R}^d$  dont la norme est inférieure à un.

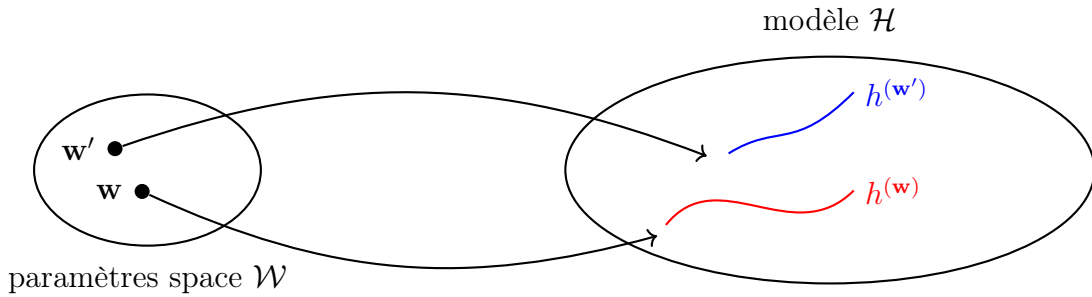


Fig. 11. L'espace des paramètres  $\mathcal{W}$  d'un modèle d'apprentissage automatique  $\mathcal{H}$  contient tous les choix possibles pour les paramètres du modèle. Chaque choix  $\mathbf{w}$  pour les paramètres du modèle sélectionne une hypothèse  $h(\mathbf{w}) \in \mathcal{H}$ .

Voir aussi: paramètres, modèle, paramètres du modèle.

**espace des étiquettes** Considérons une application d'apprentissage automatique impliquant des points de données caractérisés par des caractéristiques et des étiquettes. L'espace des étiquettes est constitué de toutes les valeurs possibles que l'étiquette d'un point de données peut prendre. Les méthodes de régression, visant à prédire des étiquettes numériques, utilisent souvent l'espace des étiquettes  $\mathcal{Y} = \mathbb{R}$ . Les méthodes de classification binaire utilisent un espace des étiquettes constitué de deux éléments différents, par exemple  $\mathcal{Y} = \{-1, 1\}$ ,  $\mathcal{Y} = \{0, 1\}$ , ou .

See also: apprentissage automatique, point de données, caractéristique,

étiquette, régression, classification.

**espace euclidien** L'espace euclidien  $\mathbb{R}^d$  de dimension  $d \in \mathbb{N}$  est constitué des vecteurs  $\mathbf{x} = (x_1, \dots, x_d)$ , avec  $d$  composantes réelles  $x_1, \dots, x_d \in \mathbb{R}$ . Un tel espace euclidien est muni d'une structure géométrique définie par le produit scalaire  $\mathbf{x}^T \mathbf{x}' = \sum_{j=1}^d x_j x'_j$  entre deux vecteurs  $\mathbf{x}, \mathbf{x}' \in \mathbb{R}^d$  quelconques [2].

**espace probabilisé** Un espace probabilisé est un modèle mathématique d'un processus physique (une expérience aléatoire) avec un résultat incertain. Formellement, un espace probabilisé  $\mathcal{P}$  est un triplet  $(\Omega, \mathcal{F}, P)$  où

- $\Omega$  est un espace échantillon contenant tous les résultats élémentaires possibles d'une expérience aléatoire ;
- $\mathcal{F}$  est une tribu (ou sigma-algèbre), une collection de sous-ensembles de  $\Omega$  (appelés événements) qui satisfait certaines propriétés de fermeture par opérations sur les ensembles ;
- $P$  est une mesure de probabilité, une fonction qui attribue une probabilité  $P(\mathcal{A}) \in [0, 1]$  à chaque événement  $\mathcal{A} \in \mathcal{F}$ . Cette fonction doit satisfaire  $P(\Omega) = 1$  et

$$P\left(\bigcup_{i=1}^{\infty} \mathcal{A}_i\right) = \sum_{i=1}^{\infty} P(\mathcal{A}_i)$$

pour toute suite dénombrable d'événements deux à deux disjoints  $\mathcal{A}_1, \mathcal{A}_2, \dots$  dans  $\mathcal{F}$ .

Les espaces probabilisés fournissent la base pour définir les VA et raisonner sur incertitude dans les applications d'apprentissage automatique [6, 19, 38].

**espérance** Soit un vecteur de caractéristiques numérique  $\mathbf{x} \in \mathbb{R}^d$  que l'on interprète comme la réalisation d'une VA suivant une loi de probabilité  $p(\mathbf{x})$ . On définit l'espérance de  $\mathbf{x}$  comme l'intégrale  $\mathbb{E}\{\mathbf{x}\} := \int \mathbf{x}p(\mathbf{x})$  [2, 6, 39]. Remarquons que l'espérance n'est définie que si cette intégrale existe, c'est-à-dire si la VA est intégrable.

**estimation error** Consider point de données, each with vecteur de caractéristiques  $\mathbf{x}$  and étiquette  $y$ . In some applications, we can model the relation between the vecteur de caractéristiques and the étiquette of a point de données as  $y = \bar{h}(\mathbf{x}) + \varepsilon$ . Here, we use some true underlying hypothèse  $\bar{h}$  and a noise term  $\varepsilon$  which summarizes any modeling or labeling errors. The estimation error incurred by an apprentissage automatique method that learns a hypothèse  $\hat{h}$ , e.g., using MRE, is defined as  $\hat{h}(\mathbf{x}) - \bar{h}(\mathbf{x})$ , for some vecteur de caractéristiques. For a parametric espace des hypothèses, which consists of hypothèse maps determined by paramètres du modèle  $\mathbf{w}$ , we can define the estimation error as  $\Delta\mathbf{w} = \hat{\mathbf{w}} - \bar{\mathbf{w}}$  [20, 40].

**expectation-maximization (EM)** Consider a modèle probabiliste  $\mathbb{P}(\mathbf{z}; \mathbf{w})$  for the point de données  $\mathcal{D}$  generated in some apprentissage automatique application. The maximum de vraisemblance estimator for the paramètres du modèle  $\mathbf{w}$  is obtained by maximizing  $\mathbb{P}(\mathcal{D}; \mathbf{w})$ . However, the resulting optimization problem might be computationally challenging. Espérance-maximization approximates the maximum de vraisemblance estimator by introducing a latent VA  $\mathbf{z}$  such that maximizing  $\mathbb{P}(\mathcal{D}, \mathbf{z}; \mathbf{w})$  would be easier [20, 41, 42]. Since we do not observe

$\mathbf{z}$ , we need to estimate it from the observed jeu de données  $\mathcal{D}$  using a conditional espérance. The resulting estimate  $\hat{\mathbf{z}}$  is then used to compute a new estimate  $\hat{\mathbf{w}}$  by solving  $\max_{\mathbf{w}} \mathbb{P}(\mathcal{D}, \hat{\mathbf{z}}; \mathbf{w})$ . The crux is that the conditional espérance  $\hat{\mathbf{z}}$  depends on the paramètres du modèle  $\hat{\mathbf{w}}$ , which we have updated based on  $\hat{\mathbf{z}}$ . Thus, we have to re-calculate  $\hat{\mathbf{z}}$ , which, in turn, results in a new choice  $\hat{\mathbf{w}}$  for the paramètres du modèle. In practice, we repeat the computation of the conditional espérance (i.e., the E-step) and the update of the paramètres du modèle (i.e., the M-step) until some critère d'arrêt is met.

**expert** En apprentissage automatique, l'objectif est d'apprendre une hypothèse  $h$  capable de prédire avec précision l'étiquette d'un point de données à partir de ses caractéristiques. L'erreur de prédiction est mesurée à l'aide d'une fonction de perte. Idéalement, on cherche à obtenir une hypothèse minimisant la perte sur tout point de données. On peut préciser cet objectif informel avec l'hypothèse d'indépendance et de distribution identique (hypothèse i.i.d.) et le Bayes risk, qui sert de baseline pour la perte moyenne d'une hypothèse. Une autre approche pour définir une baseline consiste à utiliser l'hypothèse  $h'$  apprise par une méthode d'apprentissage automatique existante. On appelle alors cette hypothèse  $h'$  un expert [43]. Les méthodes de minimisation de regret cherchent à apprendre une hypothèse dont la perte est comparable à celle du meilleur expert [43, 44].

**explainable empirical risk minimization (EERM)** Explainable MRE is an instance of SRM that adds a régularisation term to the aver-



age perte in the fonction objective of MRE. The régularisation term is chosen to favor hypothèse maps that are intrinsically explainable for a specific user. This user is characterized by their prédictions provided for the point de données in a ensemble d'entraînement [45].

**explainable machine learning (explainable ML)** Explainable apprentissage automatique methods aim at complementing each prédiction with an explication of how the prédiction has been obtained. The construction of an explicit explication might not be necessary if the apprentissage automatique method uses a sufficiently simple (or interpretable) modèle [46].

**explicabilité** On définit l'explicabilité (subjective) d'une méthode d'apprentissage automatique comme le niveau de simulabilité [49] des prédictions fournies par un système d'apprentissage automatique à un utilisateur humain. Des mesures quantitatives de l'explicabilité (subjective) d'un modèle entraîné peuvent être construites en comparant ses prédictions avec les prédictions fournies par un utilisateur sur un ensemble de test [45, 49]. Alternativement, on peut utiliser des modèles probabilistes pour les données et mesurer l'explicabilité d'un modèle d'apprentissage automatique entraîné via l'entropie différentielle conditionnelle de ses prédictions, étant donné les prédictions de l'utilisateur [50, 51].

**explication** Une approche pour rendre les méthodes d'apprentissage automatique plus transparentes consiste à fournir une explication en complément de la prédiction produite par une méthode d'apprentissage automatique. Les explications peuvent prendre différentes formes. Il

peut s’agir d’un texte en langage naturel ou d’une mesure quantitative indiquant l’importance des différentes caractéristiques d’un point de données [47]. On peut aussi utiliser des formes visuelles d’explication, comme des cartes d’intensité pour la classification d’images [48].

### **Explications locales interprétables et agnostiques au modèle (LIME)**

Considérons un modèle entraîné (ou une hypothèse apprise)  $\hat{h} \in \mathcal{H}$ , qui associe le vecteur de caractéristiques d’un point de données à la prédiction  $\hat{y} = \hat{h}$ . Les explications locales interprétables et agnostiques au modèle (LIME) sont une technique permettant d’expliquer le comportement de  $\hat{h}$  localement autour d’un point de données de vecteur de caractéristiques  $\mathbf{x}^{(0)}$  [52]. L’explication est donnée sous la forme d’une approximation locale  $g \in \mathcal{H}'$  de  $\hat{h}$  (voir Fig. ??). Cette approximation peut être obtenue par une instance de MRE avec un ensemble d’entraînement soigneusement conçu. En particulier, l’ensemble d’entraînement est composé de points de données ayant un vecteur de caractéristiques  $\mathbf{x}$  proche de  $\mathbf{x}^{(0)}$  et une (pseudo-)étiquette  $\hat{h}(\mathbf{x})$ . Remarquons que l’on peut utiliser un modèle  $\mathcal{H}'$  différent du modèle original  $\mathcal{H}$  pour l’approximation. Par exemple, on peut utiliser un arbre de décision pour approximer localement un réseau de neurones profond. Un autre choix très courant pour  $\mathcal{H}'$  est le modèle linéaire.

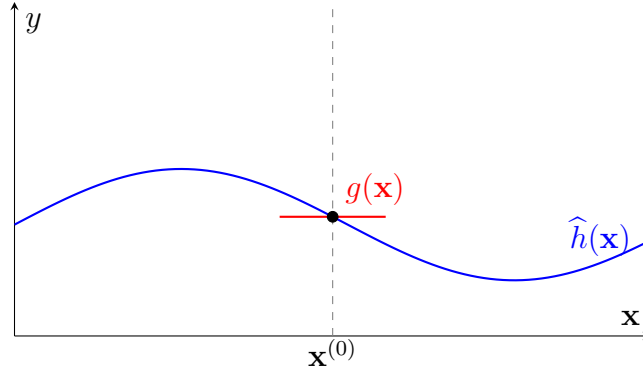


Fig. 12. Pour expliquer (comprendre) un modèle entraîné  $\hat{h} \in \mathcal{H}$ , autour d'une vecteur de caractéristiques donnée  $\mathbf{x}^{(0)}$ , on peut utiliser une approximation locale  $g \in \mathcal{H}'$ .

Voir aussi: modèle, explication, MRE, ensemble d'entraînement, étiquette, arbre de décision, réseau de neurones profond, modèle linéaire.

**feature learning** Consider an apprentissage automatique application with point de données characterized by raw caractéristiques  $\mathbf{x} \in \mathcal{X}$ . Caractéristique learning refers to the task of learning a map

$$\Phi : \mathcal{X} \rightarrow \mathcal{X}' : \mathbf{x} \mapsto \mathbf{x}',$$

that reads in raw caractéristiques  $\mathbf{x} \in \mathcal{X}$  of a point de données and delivers new caractéristiques  $\mathbf{x}' \in \mathcal{X}'$  from a new espace des caractéristiques  $\mathcal{X}'$ . Different caractéristique learning methods are obtained for different design choices of  $\mathcal{X}, \mathcal{X}'$ , for a espace des hypothèses  $\mathcal{H}$  of potential maps  $\Phi$ , and for a quantitative measure of the usefulness of a specific  $\Phi \in \mathcal{H}$ .

For example, analyse en composantes principales (ACP) uses  $\mathcal{X} := \mathbb{R}^d$ ,  $\mathcal{X}' := \mathbb{R}^{d'}$  with  $d' < d$ , and a espace des hypothèses

$$\mathcal{H} := \{ \Phi : \mathbb{R}^d \rightarrow \mathbb{R}^{d'} : \mathbf{x}' := \mathbf{F}\mathbf{x} \text{ with some } \mathbf{F} \in \mathbb{R}^{d' \times d} \}.$$

ACP measures the usefulness of a specific map  $\Phi(\mathbf{x}) = \mathbf{F}\mathbf{x}$  by the minimum linear reconstruction error incurred on a jeu de données,

$$\min_{\mathbf{G} \in \mathbb{R}^{d' \times d}} \sum_{r=1}^m \left\| \mathbf{G}\mathbf{F}\mathbf{x}^{(r)} - \mathbf{x}^{(r)} \right\|_2^2.$$

**FedAvg** FedAvg refers to a family of iterative apprentissage fédéré algorithmes. It uses a server-client setting and alternatives between client-wise modèles locaux re-training, followed by the aggregation of updated paramètres du modèle at the server [53].

See also: apprentissage fédéré, algorithme.

**FedGD** An apprentissage fédéré distributed algorithm that can be implemented as message passing across an réseau d'apprentissage fédéré.

See also: apprentissage fédéré, algorithme, pas, méthodes basées sur le gradient.

**FedProx** FedProx refers to an iterative apprentissage fédéré algorithme that alternates between separately training modèle locaux and combining the updated local paramètres du modèle. In contrast to FedAvg, which uses descente de gradient stochastique (SGD) to train modèle locaux, FedProx uses a opérateur proximal for the training [54].

**FedRelax** An apprentissage fédéré distributed algorithm.

See also: apprentissage fédéré, algorithme.

**FedSGD** An apprentissage fédéré distributed algorithm that can be implemented as message passing across an réseau d'apprentissage fédéré.  
See also: apprentissage fédéré, algorithmes, pas, méthodes basées sur le gradient, SGD.

**flow-based clustering** Flow-based partitionnement de données groups the nodes of an undirected graphe by applying  $k$ -moyennes partitionnement de données to node-wise vecteur de caractéristiquess. These vecteur de caractéristiquess are built from network flows between carefully selected sources and destination nodes [55].

**fonction** Une *fonction* est une règle mathématique qui associe à chaque élément  $u \in \mathcal{U}$  exactement un élément  $v \in \mathcal{V}$  [2]. On écrit cela  $f : \mathcal{U} \rightarrow \mathcal{V}$ , où  $\mathcal{U}$  est le domaine de définition et  $\mathcal{V}$  l'ensemble d'arrivée de  $f$ . Autrement dit, une fonction  $f$  définit une sortie unique  $f(u) \in \mathcal{V}$  pour chaque entrée  $u \in \mathcal{U}$ .

**fonction d'activation** On associe à chaque neurone artificiel dans un RNA une fonction d'activation  $\sigma(\cdot)$  qui prend en entrée une combinaison pondérée des entrées du neurone  $x_1, \dots, x_d$  et produit une sortie unique  $a = \sigma(w_1x_1 + \dots + w_dx_d)$ . Notons que chaque neurone est paramétré par les poids  $w_1, \dots, w_d$ .

**fonction de densité de probabilité** La fonction de densité de probabilité  $p(x)$  d'une VA réelle  $x \in \mathbb{R}$  est une représentation particulière de sa loi de probabilité. Si la fonction de densité de probabilité existe, elle peut être utilisée pour calculer la probabilité que  $x$  prenne une valeur dans un ensemble (mesurable)  $\mathcal{B} \subseteq \mathbb{R}$  avec  $\mathbb{P}(x \in \mathcal{B}) = \int_{\mathcal{B}} p(x')dx'$  [7, Ch. 3].

La fonction de densité de probabilité d'une VA vectorielle  $\mathbf{x} \in \mathbb{R}^d$  (si elle existe) permet de calculer la probabilité que  $\mathbf{x}$  appartienne à une région (mesurable)  $\mathcal{R}$  avec  $\mathbb{P}(\mathbf{x} \in \mathcal{R}) = \int_{\mathcal{R}} p(\mathbf{x}') dx'_1 \dots dx'_d$  [7, Ch. 3].

**fonction de perte (ou de coût)** Une fonction de perte (ou de coût) est une application

$$L : \mathcal{X} \times \mathcal{Y} \times \mathcal{H} \rightarrow \mathbb{R}_+ : ((\mathbf{x}, y), h) \mapsto L((\mathbf{x}, y), h).$$

Elle associe un réel positif ou nul (i.e., la perte)  $L((\mathbf{x}, y), h)$  à une paire composée d'un point de données, de caractéristiques  $\mathbf{x}$  et étiquette  $y$ , et d'une hypothèse  $h \in \mathcal{H}$ . La valeur  $L((\mathbf{x}, y), h)$  mesure l'écart entre l'étiquette réelle  $y$  et la prédiction  $h(\mathbf{x})$ . Des valeurs plus faibles (proches de zéro) de  $L((\mathbf{x}, y), h)$  indiquent un écart plus faible entre la prédiction  $h(\mathbf{x})$  et l'étiquette  $y$ . La figure 13 représente une fonction de perte pour un point de données donné, de caractéristiques  $\mathbf{x}$  et d'étiquette  $y$ , en fonction de l'hypothèse  $h \in \mathcal{H}$ .

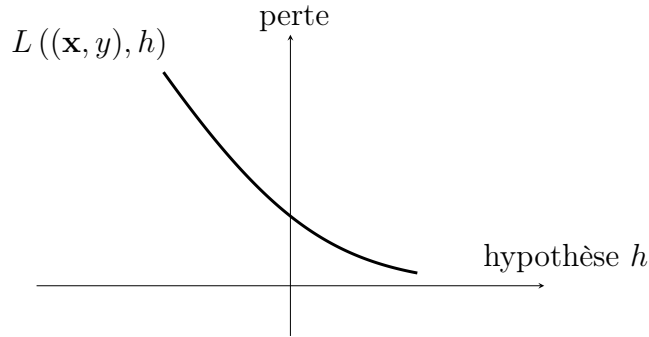


Fig. 13. Une fonction de perte  $L((\mathbf{x}, y), h)$  pour un point de données fixé, de vecteur de caractéristiques  $\mathbf{x}$  et d'étiquette  $y$ , et une hypothèse variable  $h$ . Les méthodes d'apprentissage automatique cherchent à trouver (ou apprendre) une hypothèse minimisant la perte.

**fonction objective** Une fonction objective est une application qui associe à chaque valeur d'une variable d'optimisation, comme les paramètres du modèle  $\mathbf{w}$  d'une hypothèse  $h^{(\mathbf{w})}$ , une valeur objective  $f(\mathbf{w})$ . La valeur objective  $f(\mathbf{w})$  peut être le risque ou le risque empirique d'une hypothèse  $h^{(\mathbf{w})}$ .

**frontière de décision** Considérons une fonction hypothèse  $h$  qui lit un vecteur de caractéristiques  $\mathbf{x} \in \mathbb{R}^d$  et renvoie une valeur à partir d'un ensemble fini  $\mathcal{Y}$ . La frontière de décision de  $h$  est l'ensemble des vecteurs  $\mathbf{x} \in \mathbb{R}^d$  qui se trouvent entre différentes régions de décision. Plus précisément, un vecteur  $\mathbf{x}$  appartient à la frontière de décision si et seulement si chaque voisinage  $\{\mathbf{x}' : \|\mathbf{x} - \mathbf{x}'\| \leq \varepsilon\}$ , pour tout  $\varepsilon > 0$ , contient au moins deux vecteurs avec des images différentes par la fonction.

**Gaussian mixture model (GMM)** A GMM is a particular type of modèle probabiliste for a numeric vector  $\mathbf{x}$  (e.g., the caractéristiques of a point de données). Within a GMM, the vector  $\mathbf{x}$  is drawn from a randomly selected loi normale multivariée  $p^{(c)} = \mathcal{N}(\boldsymbol{\mu}^{(c)}, \mathbf{C}^{(c)})$  with  $c = I$ . The index  $I \in \{1, \dots, k\}$  is an VA with probabilities  $\mathbb{P}(I = c) = p_c$ . Note that a GMM is parametrized by the probability  $p_c$ , the moyenne vector  $\boldsymbol{\mu}^{(c)}$ , and the matrice de covariance  $\boldsymbol{\Sigma}^{(c)}$  for each  $c = 1, \dots, k$ . GMMs are widely used for partitionnement de données, density estimation, and as a generative modèle.

**Gaussian process (GP)** A GP is a collection of VAs  $\{f(\mathbf{x})\}_{\mathbf{x} \in \mathcal{X}}$  indexed by input values  $\mathbf{x}$  from some input space  $\mathcal{X}$ , such that for any finite

subset  $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)} \in \mathcal{X}$ , the corresponding VAs  $f(\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)})$  have a joint multivariate Gaussian distribution:

$$(f(\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)})) \sim \mathcal{N}(\boldsymbol{\mu}, \mathbf{K}).$$

For a fixed input space  $\mathcal{X}$ , a GP is fully specified (or parametrized) by

- a moyenne function  $\mu(\mathbf{x}) = \mathbb{E}\{f(\mathbf{x})\}$ ,
- and a covariance function  $K(\mathbf{x}, \mathbf{x}') = \mathbb{E}\{(f(\mathbf{x}) - \mu(\mathbf{x}))(f(\mathbf{x}') - \mu(\mathbf{x}'))\}$ .

**Example.** We can interpret the temperature distribution across Finland (at a specific point in time) as the réalisation of a GP  $f(\mathbf{x})$ , where each input  $\mathbf{x} = (\text{lat}, \text{lon})$  denotes a geographic location. Temperature observations from FMI weather stations provide samples of  $f(\mathbf{x})$  at specific locations (see Fig. 14). A GP allows to predict the temperature nearby FMI weather stations and to quantify the uncertainty of these predictions.



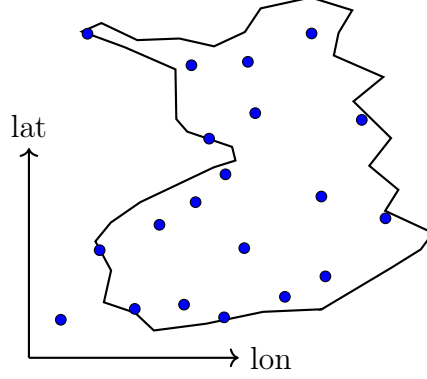


Fig. 14. We can interpret the temperature distribution over Finland as a réalisation of a GP indexed by geographic coordinates and sampled at FMI weather stations (indicated by blue dots).

**generalized total variation (GTV)** GTV is a measure of the variation of trained modèle locals  $h^{(i)}$  (or their paramètres du modèle  $\mathbf{w}^{(i)}$ ) assigned to the nodes  $i = 1, \dots, n$  of an undirected weighted graphe  $\mathcal{G}$  with edges  $\mathcal{E}$ . Given a measure  $d^{(h, h')}$  for the discrepancy between hypothèse maps  $h, h'$ , the GTV is

$$\sum_{\{i, i'\} \in \mathcal{E}} A_{i, i'} d^{(h^{(i)}, h^{(i')})}.$$

Here,  $A_{i, i'} > 0$  denotes the weight of the undirected edge  $\{i, i'\} \in \mathcal{E}$ .

**generalized total variation minimization (GTVMin)** GTV minimization is an instance of regularized empirical risk minimization (RERM) using the GTV of local paramètres du modèle as a regularizer [58].

**geometric median (GM)** The GM of a set of input vectors  $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}$  in  $\mathbb{R}^d$  is a point  $\mathbf{z} \in \mathbb{R}^d$  that minimizes the sum of distances to the

vectors [27] such that

$$\mathbf{z} \in \operatorname{argmin}_{\mathbf{y} \in \mathbb{R}^d} \sum_{r=1}^m \|\mathbf{y} - \mathbf{x}^{(r)}\|_2. \quad (3)$$

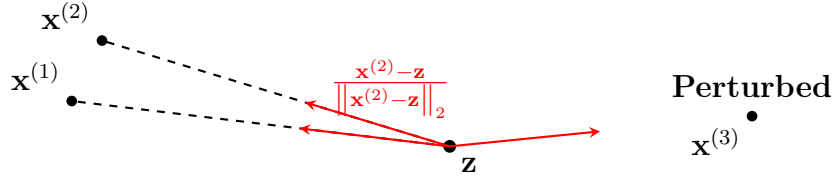


Fig. 15. Consider a solution  $\mathbf{z}$  of (3) that does not coincide with any of the input vectors. The optimality condition for (3) requires that the unit vectors from  $\mathbf{z}$  to the input vectors sum to zero.

Figure 15 illustrates a fundamental property of the GM: If  $\mathbf{z}$  does not coincide with any of the input vectors, then the unit vectors pointing from  $\mathbf{z}$  to each  $\mathbf{x}^{(r)}$  must sum to zero - this is the zero-subgradient (optimality) condition of (3). It turns out that the solution to (3) cannot be arbitrarily pulled away from trustworthy input vectors as long as they are the majority [59, Th. 2.2].

**gradient** Pour une fonction à valeurs réelles  $f : \mathbb{R}^d \rightarrow \mathbb{R} : \mathbf{w} \mapsto f(\mathbf{w})$ , s'il existe un vecteur  $\mathbf{g}$  tel que  $\lim_{\mathbf{w} \rightarrow \mathbf{w}'} \frac{f(\mathbf{w}) - (f(\mathbf{w}') + \mathbf{g}^T(\mathbf{w} - \mathbf{w}'))}{\|\mathbf{w} - \mathbf{w}'\|} = 0$ , alors on le nomme le gradient de  $f$  en  $\mathbf{w}'$ . S'il existe, le gradient est unique et est noté  $\nabla f(\mathbf{w}')$  ou  $\nabla f(\mathbf{w})|_{\mathbf{w}'}$  [2].

**grand modèle de langage (GML)** Large language modèles is an umbrella term for apprentissage automatique methods that process and generate

human-like text. These methods typically use réseau de neurones profonds with billions (or even trillions) of paramètres. A widely used choice for the network architecture is referred to as Transformers [60]. The training of large language modèles is often based on the task of predicting a few words that are intentionally removed from a large text corpus. Thus, we can construct point de données étiquetés simply by selecting some words of a text as étiquettes and the remaining words as caractéristiques of point de données. This construction requires very little human supervision and allows for generating sufficiently large ensemble d'entraînements for large language modèles.

**graph clustering** Graphe partitionnement de données aims at partitionnement de données point de données that are represented as the nodes of a graphe  $\mathcal{G}$ . The edges of  $\mathcal{G}$  represent pairwise similarities between point de données. Sometimes we can quantify the extend of these similarities by an poids d'arête [55, 61].

**graphe** Un graphe  $\mathcal{G} = (\mathcal{V}, \mathcal{E})$  est une paire qui consiste en un ensemble de nœuds  $\mathcal{V}$  et un ensemble d'arêtes  $\mathcal{E}$ . Dans sa forme la plus générale, un graphe est spécifié par une application qui associe à chaque arête  $e \in \mathcal{E}$  une paire de nœuds [62]. Une famille importante de graphes est celle des graphes simples non orientés. Un graphe simple non orienté est obtenu en identifiant chaque arête  $e \in \mathcal{E}$  à deux nœuds différents  $\{i, i'\}$ . Les graphes pondérés précisent également des poids numériques  $A_e$  pour chaque arête  $e \in \mathcal{E}$ .

**généralisation** Beaucoup de systèmes d'apprentissage automatique (et d'intelligence

artificielle (IA)) actuels reposent sur la MRE : au noyau, ils entraînent un modèle (c'est-à-dire apprennent une hypothèse  $\hat{h} \in \mathcal{H}$ ) en minimisant la perte moyenne (ou risque empirique) sur des points de données  $\mathbf{z}^{(1)}, \dots, \mathbf{z}^{(m)}$ , qui constituent un ensemble d'entraînement  $\mathcal{D}^{(\text{train})}$ . La généralisation désigne la capacité d'une méthode d'apprentissage automatique à bien fonctionner en dehors de cet ensemble d'entraînement. Toute théorie mathématique sur la généralisation nécessite une notion mathématique du « en dehors de l'ensemble d'entraînement ». Par exemple, la théorie statistique de l'apprentissage utilise un modèle probabiliste tel que l'hypothèse i.i.d. pour la génération des données : les points de données du ensemble d'entraînement sont des réalisations i.i.d. d'une loi de probabilité sous-jacente  $p(\mathbf{z})$ . Un tel modèle probabiliste permet d'explorer ce qui est en dehors de l'ensemble d'entraînement en tirant d'autres réalisations i.i.d. selon  $p(\mathbf{z})$ . De plus, grâce à l'hypothèse i.i.d., on peut définir le risque d'un modèle entraîné  $\hat{h} \in \mathcal{H}$  comme la perte espérée  $\bar{L}(\hat{h})$ . On peut également utiliser des bornes de concentration ou des résultats de convergence pour des suites de VA i.i.d. afin de borner l'écart entre le risque empirique  $\hat{L}(\hat{h}|\mathcal{D}^{(\text{train})})$  d'un modèle entraîné et son risque [56]. Il est aussi possible d'étudier la généralisation sans utiliser de modèle probabiliste. Par exemple, on peut considérer des perturbations (déterministes) des points de données de l'ensemble d'entraînement pour étudier ce qui est en dehors. En général, on souhaite que le modèle entraîné soit robuste, c'est-à-dire que ses prédictions ne changent pas trop pour de petites perturbations des points de données. Par exemple, pour un modèle entraîné à détecter

un objet dans une photo prise avec un smartphone, le résultat de la détection ne devrait pas changer si l'on masque un petit nombre de pixels choisis aléatoirement dans l'image [57].

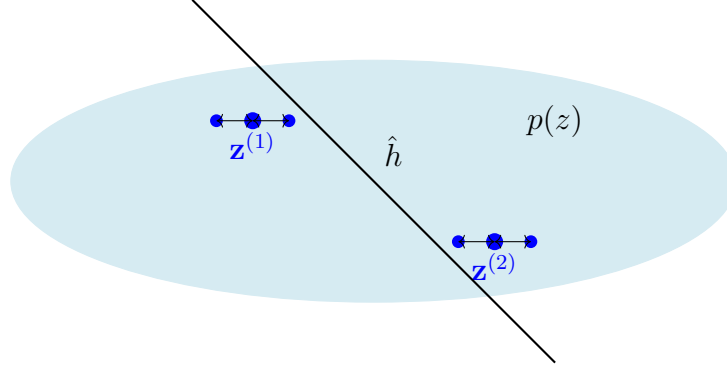


Fig. 16. Deux points de données  $\mathbf{z}^{(1)}, \mathbf{z}^{(2)}$  utilisés comme ensemble d'entraînement pour apprendre une hypothèse  $\hat{h}$  via MRE. On peut évaluer  $\hat{h}$  en dehors de l'ensemble d'entraînement soit avec l'hypothèse i.i.d. avec une loi de probabilité sous-jacente  $p(\mathbf{z})$ , soit en perturbant les points de données.

**hard clustering** Hard partitionnement de données refers to the task of partitioning a given set of point de données into (a few) non-overlapping clusters. The most widely used hard partitionnement de données method is  $k$ -moyennes.

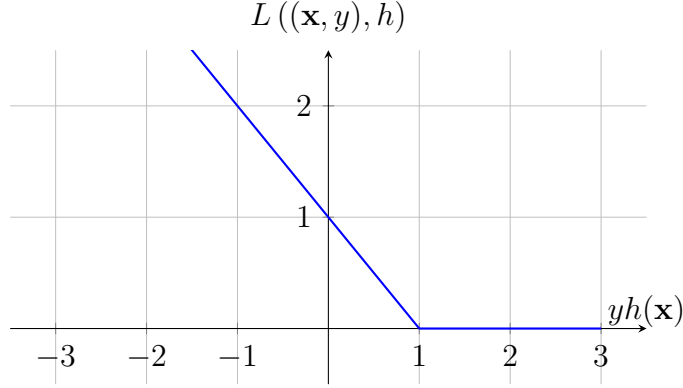
**high-dimensional regime** The high-dimensional regime of MRE is characterized by the dimension effective of the modèle being larger than the taille d'échantillon, i.e., the number of (labeled) point de données in the ensemble d'entraînement. For example, régression linéaire methods operate in the high-dimensional regime whenever the number  $d$

of caractéristiques used to characterize point de données exceeds the number of point de données in the ensemble d'entraînement. Another example of apprentissage automatique methods that operate in the high-dimensional regime is large RNAs, which have far more tunable poids (and bias terms) than the total number of point de données in the ensemble d'entraînement. High-dimensional statistics is a recent main thread of probabilité theory that studies the behavior of apprentissage automatique methods in the high-dimensional regime [63, 64].

**Hilbert space** A Hilbert space is a complete inner product space [65]. That is, it is a vector space equipped with an inner product between pairs of vectors, and it satisfies the additional requirement of completeness: every Cauchy sequence of vectors converges to a limit within the space. A canonical example of a Hilbert space is the espace euclidien  $\mathbb{R}^d$ , for some dimension  $d$ , consisting of vectors  $\mathbf{u} = (u_1, \dots, u_d)^T$  and the standard inner product  $\mathbf{u}^T \mathbf{v}$ .

**hinge loss** Consider a point de données characterized by a vecteur de caractéristiques  $\mathbf{x} \in \mathbb{R}^d$  and a binary étiquette  $y \in \{-1, 1\}$ . The hinge perte incurred by a real-valued hypothèse map  $h(\mathbf{x})$  is defined as

$$L((\mathbf{x}, y), h) := \max\{0, 1 - yh(\mathbf{x})\}. \quad (4)$$



A regularized variant of the hinge perte is used by the support vector machine (SVM) [66].

**histogram** Consider a jeu de données  $\mathcal{D}$  that consists of  $m$  point de données  $\mathbf{z}^{(1)}, \dots, \mathbf{z}^{(m)}$ , each of them belonging to some cell  $[-U, U] \times \dots \times [-U, U] \subseteq \mathbb{R}^d$  with side length  $U$ . We partition this cell evenly into smaller elementary cells with side length  $\Delta$ . The histogram of  $\mathcal{D}$  assigns each elementary cell to the corresponding fraction of point de données in  $\mathcal{D}$  that fall into this elementary cell. A visual example of such a histogram is provided in Figure 17.

**horizontal federated learning (HFL)** HFL uses jeu de données locaux constituted by different point de données but uses the same caractéristiques to characterize them [67]. For example, weather forecasting uses a network of spatially distributed weather (observation) stations. Each weather station measures the same quantities, such as daily temperature, air pressure, and precipitation. However, different weather stations measure the characteristics or caractéristiques of different spatiotemporal

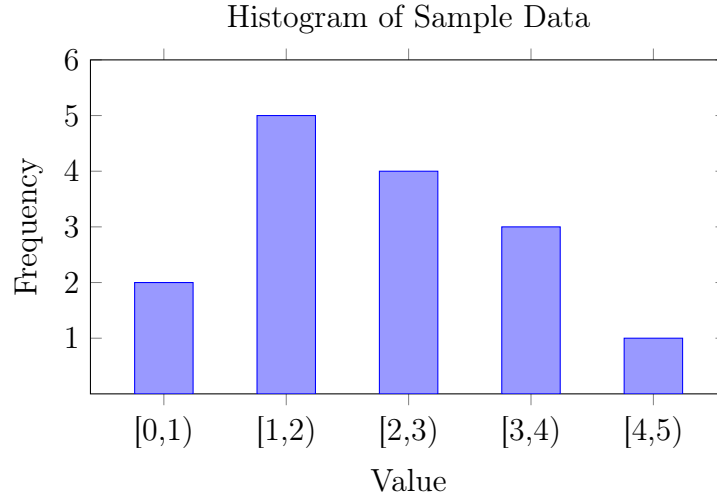


Fig. 17. A histogram representing the frequency of data points falling within discrete value ranges (bins). Each bar height shows the count of samples in the corresponding interval.

regions. Each spatiotemporal region represents an individual point de données, each characterized by the same caractéristiques (e.g., daily temperature or air pressure).

See also: apprentissage fédéré, vertical federated learning (vertical FL), clustered federated learning (CFL).

**Huber loss** The Huber perte unifies the perte quadratique and the absolute error loss.

**Huber regression** Huber régression refers to MRE-based methods that use the Huber loss as a measure of the prédiction error. Two important special cases of Huber régression are least absolute deviation regression and régression linéaire. Tuning the threshold parameter of the Huber



loss allows the user to trade the robustness of the absolute error loss against the computational benefits of the régulière perte quadratique.

**hypothèse** Une hypothèse désigne une application (ou fonction)  $h : \mathcal{X} \rightarrow \mathcal{Y}$  allant de l'espace des caractéristiques  $\mathcal{X}$  vers l'espace des étiquettes  $\mathcal{Y}$ . Étant donné un point de données avec des caractéristiques  $\mathbf{x}$ , on utilise une fonction hypothèse  $h$  pour estimer (ou approximer) son étiquette  $y$  à l'aide de la prédiction  $\hat{y} = h(\mathbf{x})$ . L'apprentissage automatique consiste à apprendre (ou trouver) une hypothèse  $h$  telle que  $y \approx h(\mathbf{x})$  pour tout point de données (de caractéristiques  $\mathbf{x}$  et étiquette  $y$ ).

**hypothèse d'indépendance et de distribution identique (hypothèse i.i.d.)**

L'hypothèse i.i.d. interprète les points de données d'un jeu de données comme des réalisations de VA i.i.d..

**incertitude** L'incertitude désigne le degré de confiance — ou de manque de confiance — associé à une quantité comme une prédiction de modèle, une estimation de paramètre ou une observation de point de données. En apprentissage automatique, l'incertitude provient de diverses sources, comme des données bruitées, un nombre limité d'échantillons d'entraînement, ou une ambiguïté dans les hypothèses du modèle. La théorie des probabilités fournit un cadre rigoureux pour représenter et quantifier cette incertitude.

**indépendantes et identiquement distribuées (i.i.d.)** Il peut être utile d'interpréter des points de données  $\mathbf{z}^{(1)}, \dots, \mathbf{z}^{(m)}$  comme des réalisations de VA i.i.d. suivant une loi de probabilité commune. Si ces VA sont

à valeurs continues, leur fonction de densité de probabilité conjointe est  $p(\mathbf{z}^{(1)}, \dots, \mathbf{z}^{(m)}) = \prod_{r=1}^m p(\mathbf{z}^{(r)})$ , où  $p(\mathbf{z})$  est la fonction de densité de probabilité marginale commune des VA sous-jacentes (i.e., dont les points de données sont les réalisations).

**Institut météorologique finlandais (FMI)** Le FMI est une agence gouvernementale responsable de la collecte et du rapport sur les données météorologiques en Finlande.

**intelligence artificielle (IA)** L'intelligence artificielle (IA) fait référence à des systèmes qui se comportent de manière rationnelle au sens de maximiser une récompense à long terme. L'approche de l'apprentissage automatique en matière d'IA consiste à entraîner un modèle pour prédire des actions optimales. Ces prédictions sont calculées à partir d'observations sur l'état de l'environnement. Le choix de la fonction de perte distingue les applications d'IA des applications d'apprentissage automatique plus basiques. Les systèmes d'IA ont rarement accès à un ensemble d'entraînement étiqueté qui permettrait de mesurer la perte moyenne pour tout choix possible des paramètres du modèle. À la place, les systèmes d'IA utilisent des signaux de récompense observés pour obtenir une estimation (ponctuelle) de la perte engendrée par le choix actuel des paramètres du modèle.

**intelligence artificielle digne de confiance (IA digne de confiance)**

Outre les aspects computationnels et aspects statistiques, un troisième aspect fondamental du développement des méthodes d'apprentissage automatique est leur fiabilité [68]. L'Union européenne a proposé sept

exigences clés pour une IA digne de confiance (généralement basée sur des méthodes d'apprentissage automatique) [69] :

- 1) Facteur humain et contrôle humain ;
- 2) Robustesse technique et sécurité ;
- 3) Respect de la vie privée et gouvernance des données ;
- 4) Transparence ;
- 5) Diversité, non-discrimination et équité ;
- 6) Bien-être sociétal et environnemental ;
- 7) Responsabilisation.

**interprétabilité** Une méthode d'apprentissage automatique est interprétable pour un utilisateur humain si celui-ci peut comprendre le processus de décision de la méthode. Une approche pour définir précisément l'interprétabilité repose sur le concept de simulabilité, c'est-à-dire la capacité d'un humain à simuler mentalement le comportement du modèle [49, 51, 70–72]. L'idée est la suivante : si un utilisateur humain comprend une méthode d'apprentissage automatique, alors il devrait être capable d'anticiper ses prédictions sur un ensemble de test. Nous illustrons un tel ensemble de test dans la Fig. 18 qui montre également deux hypothèses apprises,  $\hat{h}$  et  $\hat{h}'$ . La méthode d'apprentissage automatique produisant l'hypothèse  $\hat{h}$  est interprétable pour un utilisateur humain familier avec le concept de application linéaire. Puisque  $\hat{h}$  correspond à une application linéaire, l'utilisateur peut anticiper les prédictions de  $\hat{h}$  sur l'ensemble de test. En revanche, la méthode

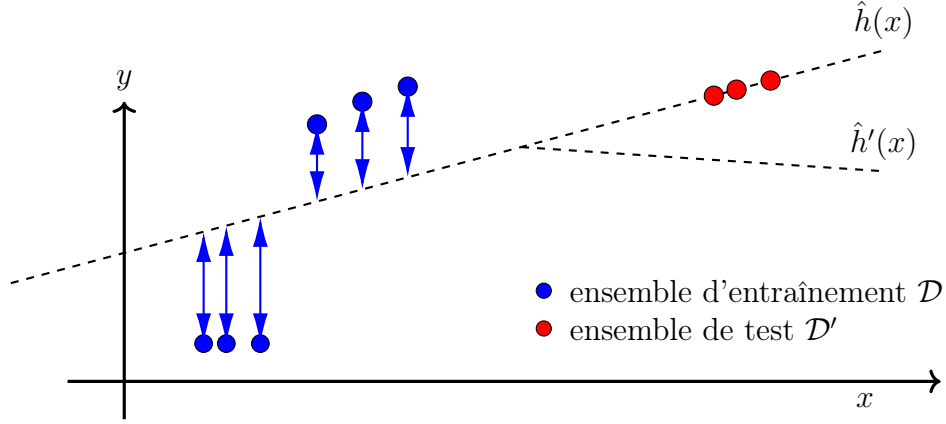


Fig. 18. Nous pouvons évaluer l’interprétabilité des modèles d’apprentissage automatique entraînés  $\hat{h}$  et  $\hat{h}'$  en comparant leurs prédictions aux pseudo-étiquettes générées par un utilisateur humain pour  $\mathcal{D}'$ .

d’apprentissage automatique fournissant  $\hat{h}'$  n’est pas interprétable, car son comportement ne correspond plus aux attentes de l’utilisateur. La notion d’interprétabilité est étroitement liée à celle d’explicabilité, car toutes deux visent à rendre les méthodes d’apprentissage automatique plus compréhensibles pour les humains. Comme illustré dans la Figure 18, l’interprétabilité d’une méthode d’apprentissage automatique  $\hat{h}$  exige que l’utilisateur humain puisse anticiper ses prédictions sur un ensemble de test arbitraire. Cela contraste avec l’explicabilité, où l’utilisateur est aidé par des explications externes — comme des cartes de saillance ou des exemples de référence issus du ensemble d’entraînement — pour comprendre les prédictions de  $\hat{h}$  sur un jeu de test spécifique  $\mathcal{D}'$ .

Voir aussi : explicabilité, intelligence artificielle digne de confiance (IA digne de confiance), régularisation, LIME.

**jeu de données** Un jeu de données désigne une collection de points de données. Ces points de données portent des informations sur une certaine quantité d'intérêt (ou étiquette) dans une application de l'apprentissage automatique. Les méthodes d'apprentissage automatique utilisent des jeux de données pour l'entraînement du modèle (par exemple via la MRE) et la validation du modèle.

Il est important de noter que notre notion de jeu de données est très flexible, car elle autorise des types de points de données très variés. En effet, les points de données peuvent être des objets physiques concrets (comme des humains ou des animaux) ou des objets abstraits (comme des nombres).

À titre d'exemple, la Figure ?? illustre un jeu de données utilisant des vaches comme points de données.



Fig. 19. Un troupeau de vaches dans les Alpes

Bien souvent, un ingénieur en apprentissage automatique n'a pas d'accès direct à un jeu de données. En effet, accéder au jeu de données de la Figure ?? impliquerait de visiter le troupeau de vaches dans les Alpes. À la place, il faut utiliser une approximation (ou représentation) du jeu de données plus pratique à manipuler.

Divers modèles mathématiques ont été développés pour représenter ou approximer les jeux de données [73], [74], [75], [76].

L'un des modèles de données les plus utilisés est le modèle relationnel, qui organise les données sous forme de tableau (ou relation) [35], [73].

Un tableau est composé de lignes et de colonnes :

- Chaque ligne du tableau représente un seul point de données.
- Chaque colonne du tableau correspond à un attribut spécifique du point de données. Les méthodes d'apprentissage automatique peuvent utiliser ces attributs comme caractéristiques ou étiquettes du point de données.

Par exemple, la Table 1 montre une représentation du jeu de données de la Figure ???. Dans le modèle relationnel, l'ordre des lignes est sans importance, et chaque attribut (colonne) doit être défini précisément par un domaine spécifiant l'ensemble des valeurs possibles.

Dans les applications de l'apprentissage automatique, ces domaines d'attributs deviennent l'espace des caractéristiques et l'espace des étiquettes.

Nom	Poids	Âge	Taille	Température de l'estomac
Zenzi	100	4	100	25
Berta	140	3	130	23
Resi	120	4	120	31

Table 1: Une relation (ou table) représentant le jeu de données de la Figure ??.

Bien que le modèle relationnel soit utile pour de nombreuses applications

en apprentissage automatique, il peut s'avérer insuffisant vis-à-vis des exigences en matière de IA digne de confiance.

Des approches modernes, telles que les fiches descriptives des jeux de données, proposent une documentation plus complète, incluant des détails sur le processus de collecte des données, l'usage prévu et d'autres informations contextuelles [77].

**jeu de données local** Le concept de jeu de données local se situe entre les notions de point de données et de jeu de données. Un jeu de données local est constitué de plusieurs points de données, chacun étant caractérisé par des caractéristiques et une étiquettes. Contrairement à un jeu de données unique, utilisé dans les méthodes classiques d'apprentissage automatique, un jeu de données local peut être relié à d'autres jeux de données locaux par différentes formes de similarité. Ces similarités peuvent provenir de modèles probabilistes ou de l'infrastructure de communication, et sont représentées par les arêtes d'un réseau d'apprentissage fédéré.

**Kullback-Leibler divergence (KL divergence)** The KL divergence is a quantitative measure of how much one loi de probabilité is different from another loi de probabilité [78].

**least absolute deviation regression** Least absolute deviation regression is an instance of MRE using the absolute error loss. It is a special case of Huber regression.

**least absolute shrinkage and selection operator (Lasso)** The Lasso is an instance of SRM. It learns the poids  $\mathbf{w}$  of a linear map  $h(\mathbf{x}) = \mathbf{w}^T \mathbf{x}$

based on a ensemble d'entraînement. Lasso is obtained from régression linéaire by adding the scaled  $\ell_1$ -norme  $\alpha \|\mathbf{w}\|_1$  to the average perte quadratique incurred on the ensemble d'entraînement.

**loi de probabilité** Pour analyser les méthodes d'apprentissage automatique, il peut être utile d'interpréter les points de données comme des réalisations i.i.d. d'une VA. Les attributs de ces points de données sont alors régis par la loi (ou distribution) de probabilité de cette VA. La loi de probabilité d'une VA binaire  $y \in \{0, 1\}$  est entièrement déterminée par les probabilités  $\mathbb{P}(y = 0)$  et  $\mathbb{P}(y = 1) = 1 - \mathbb{P}(y = 0)$ . La loi de probabilité d'une VA à valeurs réelles  $x \in \mathbb{R}$  peut être spécifiée par une fonction de densité de probabilité  $p(x)$  telle que  $\mathbb{P}(x \in [a, b]) \approx p(a)|b - a|$ . Dans le cas le plus général, une loi de probabilité est définie par une mesure de probabilité [6, 19].

**loi des grands nombres** La loi des grands nombres désigne la convergence de la moyenne d'un nombre croissant (et grand) de VA i.i.d. vers la moyenne de leur loi de probabilité commune. Il existe plusieurs versions de la loi des grands nombres selon les notions de convergence utilisées [79].

**loi normale multivariée** La loi normale multivariée  $\mathcal{N}(\mathbf{m}, \mathbf{C})$  est un modèle probabiliste important pour les vecteurs de caractéristiques numériques. C'est une famille de lois de probabilité pour une VA vectorielle  $\mathbf{x} \in \mathbb{R}^d$  [7], [19], [80]. Chaque membre (i.e. une loi de probabilité) de cette famille est spécifié par sa moyenne  $\mathbf{m}$  et sa matrice de covariance  $\mathbf{C}$ . Si la matrice de covariance est inversible, la loi de probabilité de  $\mathbf{x}$  peut



s'écrire :

$$p(\mathbf{x}) \propto \exp \left( - (1/2) (\mathbf{x} - \mathbf{m})^T \mathbf{C}^{-1} (\mathbf{x} - \mathbf{m}) \right).$$

**lot** Dans le contexte de la SGD, un lot désigne un sous-ensemble choisi aléatoirement dans l'ensemble d'entraînement complet. On utilise les points de données de ce sous-ensemble pour estimer le gradient de l'erreur d'entraînement et, par la suite, mettre à jour les paramètres du modèle.

**matrice de caractéristiques** Considérons un jeu de données  $\mathcal{D}$  avec  $m$  points de données de vecteurs de caractéristiques  $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)} \in \mathbb{R}^d$ . Il est pratique de rassembler les vecteurs de caractéristiques individuels dans une matrice de caractéristiques  $\mathbf{X} := (\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)})^T$  de taille  $m \times d$ .

**matrice de covariance** La matrice de covariance d'une VA  $\mathbf{x} \in \mathbb{R}^d$  est définie comme  $\mathbb{E} \left\{ (\mathbf{x} - \mathbb{E}\{\mathbf{x}\}) (\mathbf{x} - \mathbb{E}\{\mathbf{x}\})^T \right\}$ .

**matrice inverse** On définit la matrice inverse  $\mathbf{A}^{-1}$  d'une matrice carrée  $\mathbf{A} \in \mathbb{R}^{n \times n}$  de rang maximal, c'est-à-dire dont les colonnes sont linéairement indépendantes. Dans ce cas, on dit que  $\mathbf{A}$  est inversible, et son inverse satisfait :

$$\mathbf{A} \mathbf{A}^{-1} = \mathbf{A}^{-1} \mathbf{A} = \mathbf{I}.$$

Une matrice carrée est inversible si et seulement si son déterminant est non nul. Les matrices inverses sont fondamentales pour la résolution de systèmes d'équations linéaires et dans la solution explicite de la

régression linéaire [36], [81]. Le concept de matrice inverse peut être étendu aux matrices non carrées ou de rang non maximal. On peut définir une « inverse à gauche »  $\mathbf{B}$  telle que  $\mathbf{BA} = \mathbf{I}$ , ou une « inverse à droite »  $\mathbf{C}$  telle que  $\mathbf{AC} = \mathbf{I}$ . Pour les matrices rectangulaires ou singulières, la pseudo-inverse de Moore–Penrose, notée  $\mathbf{A}^+$ , fournit une généralisation unifiée de la matrice inverse [3].

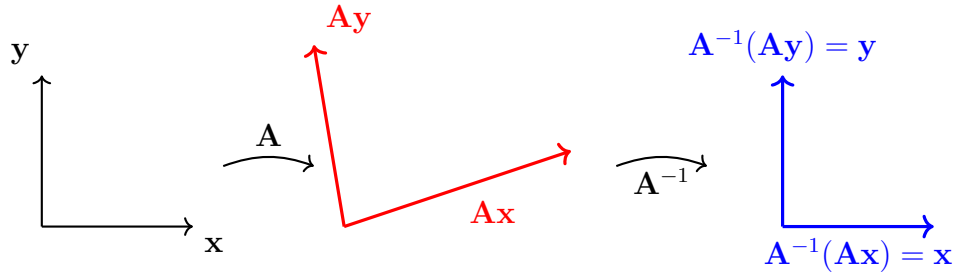


Fig. 20. Une matrice  $\mathbf{A}$  représente une transformation linéaire de  $\mathbb{R}^2$ . La matrice inverse  $\mathbf{A}^{-1}$  représente la transformation inverse.

Voir aussi : déterminant, régression linéaire, pseudo-inverse.

**matrice laplacienne** La structure d'un graphe  $\mathcal{G}$ , avec pour nœuds  $i = 1, \dots, n$ , peut être analysée à l'aide des propriétés de matrices spéciales associées à  $\mathcal{G}$ . L'une de ces matrices est la matrice laplacienne de  $\mathcal{G}$  :  $\mathbf{L}^{(\mathcal{G})} \in \mathbb{R}^{n \times n}$ , définie pour un graphe  $\mathcal{G}$  non orienté et pondéré [61, 82]. Elle est définie terme à terme par (voir Figure 21)

$$L_{i,i'}^{(\mathcal{G})} := \begin{cases} -A_{i,i'} & \text{pour } i \neq i', \{i, i'\} \in \mathcal{E}, \\ \sum_{i'' \neq i} A_{i,i''} & \text{pour } i = i', \\ 0 & \text{sinon.} \end{cases} \quad (5)$$

Ici,  $A_{i,i'}$  désigne le poids d'arête d'une arête  $\{i, i'\} \in \mathcal{E}$ .

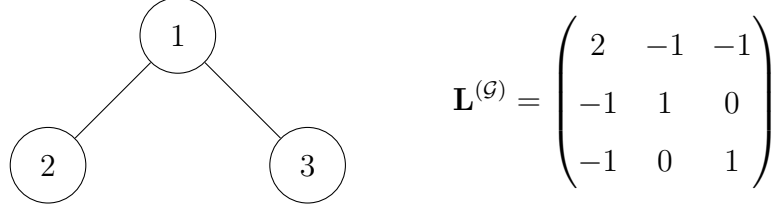


Fig. 21. À gauche : Un graphe non orienté  $\mathcal{G}$  avec trois nœuds  $i = 1, 2, 3$ . À droite : La matrice laplacienne  $\mathbf{L}^{(\mathcal{G})} \in \mathbb{R}^{3 \times 3}$  de  $\mathcal{G}$ .

**maximum** Le maximum d'un ensemble  $\mathcal{A} \subseteq \mathbb{R}$  de nombres réels est le plus grand élément de cet ensemble, si un tel élément existe. Un ensemble  $\mathcal{A}$  a un maximum s'il est majoré et atteint son borne supérieure [2, Sec. 1.4].

**maximum de vraisemblance** Considérons des points de données  $\mathcal{D} = \{\mathbf{z}^{(1)}, \dots, \mathbf{z}^{(m)}\}$  que l'on interprète comme les réalisations de VA i.i.d. avec une loi de probabilité commune  $\mathbb{P}(\mathbf{z}; \mathbf{w})$  qui dépend des paramètres du modèles  $\mathbf{w} \in \mathcal{W} \subseteq \mathbb{R}^n$ . Les méthodes de maximum de vraisemblance apprennent les paramètres du modèles  $\mathbf{w}$  en maximisant la probabilité  $\mathbb{P}(\mathcal{D}; \mathbf{w}) = \prod_{r=1}^m \mathbb{P}(\mathbf{z}^{(r)}; \mathbf{w})$  du jeu de données observé. Ainsi, l'estimateur du maximum de vraisemblance est une solution au problème d'optimisation  $\max_{\mathbf{w} \in \mathcal{W}} \mathbb{P}(\mathcal{D}; \mathbf{w})$ .

**mean squared estimation error (MSEE)** Consider an apprentissage automatique method that learns paramètres du modèle  $\hat{\mathbf{w}}$  based on some jeu de données  $\mathcal{D}$ . If we interpret the point de données in  $\mathcal{D}$  as i.i.d. réalisations of an VA  $\mathbf{z}$ , we define the estimation error  $\Delta \mathbf{w} := \hat{\mathbf{w}} - \overline{\mathbf{w}}$ .

Here,  $\bar{\mathbf{w}}$  denotes the true parameters of the model of the probability of  $\mathbf{z}$ . The mean squared estimation error is defined as the expectation  $\mathbb{E}\{\|\Delta\mathbf{w}\|^2\}$  of the squared Euclidean norm of the estimation error [17, 40].

**minimisation du risque empirique (MRE)** La minimisation du risque empirique (MRE) (risque empirique) est le problème d'optimisation qui consiste à trouver une hypothèse (dans un modèle) qui minimise la perte moyenne (ou risque empirique) sur un jeu de données  $\mathcal{D}$  donné (c'est-à-dire, l'ensemble d'entraînement). De nombreuses méthodes d'apprentissage automatique sont obtenues à partir du risque empirique via des choix de conception spécifiques pour le jeu de données, le modèle et la perte [8, Ch. 3].

Voir aussi: minimum, risque empirique, hypothèse, modèle, perte, jeu de données, ensemble d'entraînement, apprentissage automatique.

**minimum** Étant donné un ensemble de nombres réels, le minimum est le plus petit de ces nombres. Notons que pour certains ensembles, comme l'ensemble des nombres réels négatifs, le minimum n'existe pas.

**missing data** Consider a jeu de données constituted by point de données collected via some physical appareil. Due to imperfections and failures, some of the caractéristique or étiquette values of point de données might be corrupted or simply missing. Données imputation aims at estimating these missing values [83]. We can interpret données imputation as an apprentissage automatique problem where the étiquette of a point de données is the value of the corrupted caractéristique.

**model inversion** TBD.

**model selection** In apprentissage automatique, modèle selection refers to the process of choosing between different candidate modèles. In its most basic form, modèle selection amounts to: 1) training each candidate modèle; 2) computing the erreur de validation for each trained modèle; and 3) choosing the modèle with the smallest erreur de validation [8, Ch. 6].

**modèle** Dans le contexte de l'apprentissage automatique, le terme « modèle » désigne typiquement l'espace des hypothèses sous-jacent à une méthode d'apprentissage automatique [8], [56]. Cependant, ce terme est également utilisé dans d'autres domaines avec des significations différentes. Par exemple, un modèle probabiliste désigne un ensemble paramétré de lois de probabilité.

**modèle linéaire** Considérons des points de données, chacun étant caractérisé par un vecteur de caractéristiques numérique  $\mathbf{x} \in \mathbb{R}^d$ . Un modèle linéaire est un espace des hypothèses constitué de toutes les applications linéaires,

$$\mathcal{H}^{(d)} := \{h(\mathbf{x}) = \mathbf{w}^T \mathbf{x} : \mathbf{w} \in \mathbb{R}^d\}. \quad (6)$$

Notons que (6) définit une famille entière d'espace des hypothèses, paramétrée par le nombre  $d$  de caractéristiques qui sont combinées linéairement pour former la prédiction  $h(\mathbf{x})$ . Le choix de  $d$  est guidé par les aspects computationnels (par exemple, réduire  $d$  signifie moins de calcul), les aspects statistiques (par exemple, augmenter  $d$  peut réduire l'erreur de prédiction) et l'interprétabilité. Un modèle linéaire utilisant

peu de caractéristiques soigneusement sélectionnées a tendance à être considéré comme plus interprétable [46, 52].

**modèle local** Considérons une collection de jeux de données locaux qui sont assignés aux nœuds d'un réseau d'apprentissage fédéré. Un modèle local  $\mathcal{H}^{(i)}$  est un espace des hypothèses assigné à un nœud  $i \in \mathcal{V}$ . Différents nœuds peuvent se voir assigner des espace des hypothèsess différents, c'est-à-dire qu'en général  $\mathcal{H}^{(i)} \neq \mathcal{H}^{(i')}$  pour des nœuds différents  $i, i' \in \mathcal{V}$ .

**modèle probabiliste** Un modèle probabiliste interprète les points de données comme des réalisations de VA selon une loi de probabilité conjointe. Cette loi de probabilité conjointe implique généralement des paramètres qui doivent être choisis manuellement ou appris via des méthodes d'inférence statistique telles que l'estimation par maximum de vraisemblance [17].

**moyenne** La moyenne d'une VA  $\mathbf{x}$ , à valeurs dans un espace euclidien  $\mathbb{R}^d$ , est son espérance  $\mathbb{E}\{\mathbf{x}\}$ . Elle est définie comme l'intégrale de Lebesgue de  $\mathbf{x}$  par rapport à la loi de probabilité sous-jacente  $P$ ,

$$\mathbb{E}\{\mathbf{x}\} = \int_{\mathbb{R}^d} \mathbf{x} dP(\mathbf{x}),$$

voir par exemple [6] ou [2]. Nous utilisons également ce terme pour désigner la moyenne d'une séquence finie  $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)} \in \mathbb{R}^d$ . Cependant, ces deux définitions sont essentiellement équivalentes. En effet, on peut utiliser la séquence  $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)} \in \mathbb{R}^d$  pour construire une VA discrète  $\tilde{\mathbf{x}} = \mathbf{x}^{(I)}$  où l'indice  $I$  est choisi uniformément au hasard dans l'ensemble  $\{1, \dots, m\}$ . La moyenne de  $\tilde{\mathbf{x}}$  est précisément la moyenne empirique  $\frac{1}{m} \sum_{r=1}^m \mathbf{x}^{(r)}$ .

**mutual information (MI)** The MI  $I(\mathbf{x}; y)$  between two VAs  $\mathbf{x}, y$  defined on the same espace probabilisé is given by [78]

$$I(\mathbf{x}; y) := \mathbb{E} \left\{ \log \frac{p(\mathbf{x}, y)}{p(\mathbf{x})p(y)} \right\}.$$

It is a measure of how well we can estimate  $y$  based solely on  $\mathbf{x}$ . A large value of  $I(\mathbf{x}; y)$  indicates that  $y$  can be well predicted solely from  $\mathbf{x}$ . This prédiction could be obtained by a hypothèse learned by an MRE-based apprentissage automatique method.

**méthode à noyau** Une méthode à noyau est une méthode d'apprentissage automatique qui utilise un noyau  $K$  pour transformer le vecteur de caractéristiques initial (brut)  $\mathbf{x}$  d'un point de données en un nouveau (transformé) vecteur de caractéristiques  $\mathbf{z} = K(\mathbf{x}, \cdot)$  [66, 85]. La motivation derrière cette transformation est que, grâce à un noyau approprié, les points de données possèdent une géométrie « plus favorable » dans l'espace des caractéristiques transformé. Par exemple, dans un problème de classification binaire, l'utilisation des vecteurs de caractéristiques transformés  $\mathbf{z}$  peut permettre d'appliquer des modèles linéaires, même si les points de données ne sont pas linéairement séparables dans l'espace des caractéristiques initial (voir Figure 22).

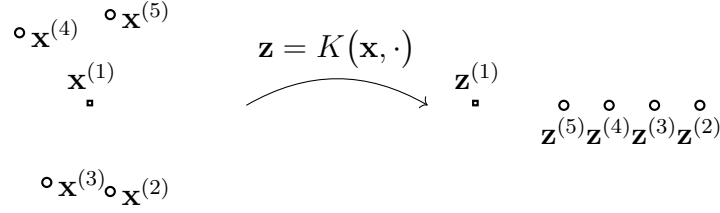


Fig. 22. Cinq points de données caractérisés par des vecteurs de caractéristiques  $\mathbf{x}^{(r)}$  et étiquettes  $y^{(r)} \in \{\circ, \square\}$ , pour  $r = 1, \dots, 5$ . Avec ces vecteurs de caractéristiques, il n'est pas possible de séparer les deux classes par une ligne droite (représentant la frontière de décision d'un classifieur linéaire). En revanche, le vecteurs de caractéristiques transformé  $\mathbf{z}^{(r)} = K(\mathbf{x}^{(r)}, \cdot)$  permet de séparer les points de données à l'aide d'un classifieur linéaire.

**méthodes basées sur le gradient** Les méthodes basées sur le gradient sont des techniques itératives pour trouver le minimum (ou le maximum) d'une fonction objective des paramètres du modèle dérivable. Ces méthodes construisent une suite d'approximations d'un choix optimal des paramètres du modèle qui aboutit à une valeur minimum (ou maximum) de la fonction objective. Comme leur nom l'indique, les méthodes basées sur le gradient utilisent les gradients de la fonction objective évalués lors des itérations précédentes pour construire de nouveaux paramètres du modèle (espérons-le) améliorés. Un exemple important d'une méthode basée sur le gradient est la descente de gradient.

**nearest neighbor (NN)** NN methods learn a hypothèse  $h : \mathcal{X} \rightarrow \mathcal{Y}$  whose function value  $h(\mathbf{x})$  is solely determined by the nearest voisins within



a given jeu de données. Different methods use different metrics for determining the nearest voisins. If point de données are characterized by numeric vecteur de caractéristiquess, we can use their Euclidean distances as the metric.

**networked data** Networked données consists of jeu de données locals that are related by some notion of pairwise similarity. We can represent networked données using a graphe whose nodes carry jeu de données locals and edges encode pairwise similarities. One example of networked données arises in apprentissage fédéré applications where jeu de données locals are generated by spatially distributed appareils.

**networked exponential families (nExpFam)** A collection of exponential families, each of them assigned to a node of an réseau d'apprentissage fédéré. The paramètres du modèle are coupled via the network structure by requiring them to have a small GTV [86].

**networked federated learning (NFL)** Networked apprentissage fédéré refers to methods that learn personalized modèles in a distributed fashion. These methods learn from jeu de données locals that are related by an intrinsic network structure.

**networked model** A networked modèle over an réseau d'apprentissage fédéré  $\mathcal{G} = (\mathcal{V}, \mathcal{E})$  assigns a modèle local (i.e., a espace des hypothèses) to each node  $i \in \mathcal{V}$  of the réseau d'apprentissage fédéré  $\mathcal{G}$ .

**node degree** The degree  $d^{(i)}$  of a node  $i \in \mathcal{V}$  in an undirected graphe is the number of its voisins, i.e.,  $d^{(i)} := |\mathcal{N}^{(i)}|$ .

**non régulière (ou non lisse)** On qualifie une fonction de non régulière si elle n'est pas régulière [87].

Voir aussi : fonction, régulière.

**norme** Une norme est une fonction qui associe à chaque élément (vecteur) d'un espace vectoriel un réel positif ou nul. Cette fonction doit être homogène, définie positive, et satisfaire l'inégalité triangulaire [37].

**noyau** Considérons des points de données caractérisés par un vecteur de caractéristiques  $\mathbf{x} \in \mathcal{X}$  avec un espace des caractéristiques générique  $\mathcal{X}$ . Un noyau (à valeurs réelles)  $K : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$  associe à chaque paire de vecteurs de caractéristiques  $\mathbf{x}, \mathbf{x}' \in \mathcal{X}$  un nombre réel  $K(\mathbf{x}, \mathbf{x}')$ . La valeur  $K(\mathbf{x}, \mathbf{x}')$  est souvent interprétée comme une mesure de similarité entre  $\mathbf{x}$  et  $\mathbf{x}'$ . Les méthode à noyaux utilisent un noyau pour transformer le vecteur de caractéristiques  $\mathbf{x}$  en un nouveau vecteur de caractéristiques  $\mathbf{z} = K(\mathbf{x}, \cdot)$ . Ce nouveau vecteur de caractéristiques appartient à un espace des caractéristiques linéaire  $\mathcal{X}'$ , qui est (en général) différent de l'espace des caractéristiques original  $\mathcal{X}$ . L'espace des caractéristiques  $\mathcal{X}'$  possède une structure mathématique spécifique : c'est un espace de Hilbert à noyau reproduisant [66, 85].

**nuage de points** Une technique de visualisation qui représente des points de données par des marqueurs dans un plan bidimensionnel. La Fig. 23 montre un exemple de nuage de points.

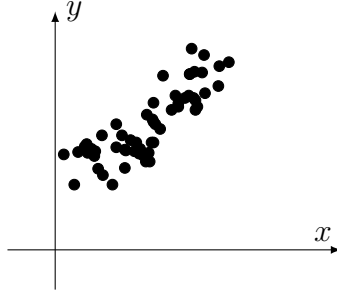


Fig. 23. Un nuage de points avec des marqueurs cercles, où les points de données représentent les conditions météorologiques quotidiennes en Finlande. Chaque point de données est caractérisé par sa température minimale diurne  $x$  comme caractéristique et sa température maximale diurne  $y$  comme étiquette. Les températures ont été mesurées à la station météo FMI Helsinki Kaisaniemi durant la période du 01.09.2024 au 28.10.2024.

Un nuage de points permet une inspection visuelle des points de données naturellement représentés par des vecteurs de caractéristiques dans des espaces de grande dimension.

Voir aussi : réduction de dimension.

**online algorithm** An online algorithm processes input données incrementally, receiving point de données sequentially and making decisions or producing outputs (or decisions) immediately without having access to the entire input in advance [43], [44]. Unlike an offline algorithm, which has the entire input available from the start, an online algorithm must handle incertitude about future inputs and cannot revise past decisions. Similar to an offline algorithm, we also represent an online algorithm formally as a collection of possible executions. However, the

execution sequence for an online algorithm has a distinct structure:

$$\text{in}_1, s_1, \text{out}_1, \text{in}_2, s_2, \text{out}_2, \dots, \text{in}_T, s_T, \text{out}_T.$$

Each execution begins from an initial state (i.e.,  $\text{in}_1$ ) and proceeds through alternating computational steps, outputs (or decisions), and inputs. Specifically, at step  $k$ , the algorithm performs a computational step  $s_k$ , generates an output  $\text{out}_k$ , and then subsequently receives the next input (point de données)  $\text{in}_{k+1}$ . A notable example of an online algorithm in apprentissage automatique is online GD, which incrementally updates paramètres du modèle as new point de données arrive.

See also: apprentissage en ligne, online GD.

**online gradient descent (online GD)** Consider an apprentissage automatique method that learns paramètres du modèle  $\mathbf{w}$  from some espace des paramètres  $\mathcal{W} \subseteq \mathbb{R}^d$ . The learning process uses point de données  $\mathbf{z}^{(t)}$  that arrive at consecutive time-instants  $t = 1, 2, \dots$ . Let us interpret the point de données  $\mathbf{z}^{(t)}$  as i.i.d. copies of an VA  $\mathbf{z}$ . The risque  $\mathbb{E}\{L(\mathbf{z}, \mathbf{w})\}$  of a hypothèse  $h^{(\mathbf{w})}$  can then (under mild conditions) be obtained as the limit  $\lim_{T \rightarrow \infty} (1/T) \sum_{t=1}^T L(\mathbf{z}^{(t)}, \mathbf{w})$ . We might use this limit as the fonction objective for learning the paramètres du modèle  $\mathbf{w}$ . Unfortunately, this limit can only be evaluated if we wait infinitely long in order to collect all point de données. Some apprentissage automatique applications require methods that learn online: as soon as a new point de données  $\mathbf{z}^{(t)}$  arrives at time  $t$ , we update the current paramètres du modèle  $\mathbf{w}^{(t)}$ . Note that the new point de données  $\mathbf{z}^{(t)}$  contributes

the component  $L(\mathbf{z}^{(t)}, \mathbf{w})$  to the risque. As its name suggests, online descent de gradient updates  $\mathbf{w}^{(t)}$  via a (projected) pas

$$\mathbf{w}^{(t+1)} := P_{\mathcal{W}}(\mathbf{w}^{(t)} - \eta_t \nabla_{\mathbf{w}} L(\mathbf{z}^{(t)}, \mathbf{w})). \quad (7)$$

Note that (7) is a pas for the current component  $L(\mathbf{z}^{(t)}, \cdot)$  of the risque. The update (7) ignores all the previous components  $L(\mathbf{z}^{(t')}, \cdot)$ , for  $t' < t$ . It might therefore happen that, compared to  $\mathbf{w}^{(t)}$ , the updated paramètres du modèle  $\mathbf{w}^{(t+1)}$  increase the retrospective average perte  $\sum_{t'=1}^{t-1} L(\mathbf{z}^{(t')}, \cdot)$ . However, for a suitably chosen taux d'apprentissage  $\eta_t$ , online descent de gradient can be shown to be optimal in practically relevant settings. By optimal, we mean that the paramètres du modèle  $\mathbf{w}^{(T+1)}$  delivered by online descent de gradient after observing  $T$  point de données  $\mathbf{z}^{(1)}, \dots, \mathbf{z}^{(T)}$  are at least as good as those delivered by any other learning method [44, 88].

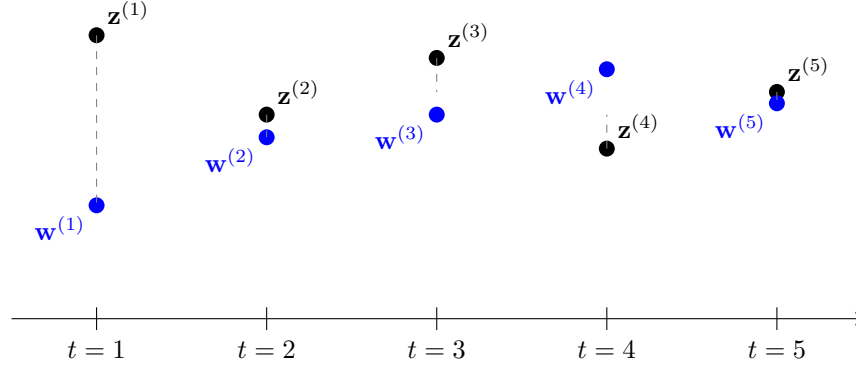


Fig. 24. An instance of online descent de gradient that updates the paramètres du modèle  $\mathbf{w}^{(t)}$  using the point de données  $\mathbf{z}^{(t)} = x^{(t)}$  arriving at time  $t$ . This instance uses the perte quadratique  $L(\mathbf{z}^{(t)}, w) = (x^{(t)} - w)^2$ .

**optimism in the face of uncertainty** apprentissage automatique methods learn parameters of the model  $\mathbf{w}$  according to some performance criterion  $\bar{f}(\mathbf{w})$ . However, they usually cannot access  $\bar{f}(\mathbf{w})$  directly but rely on an estimate (or approximation)  $f(\mathbf{w})$  of  $\bar{f}(\mathbf{w})$ . As a case in point, MRE-based methods use the average loss on a given set of data (i.e., the training set) as an estimate for the risk of a hypothesis. Using a probabilistic model, one can construct a confidence interval  $[l^{(\mathbf{w})}, u^{(\mathbf{w})}]$  for each choice  $\mathbf{w}$  for the parameters of the model. One simple construction is  $l^{(\mathbf{w})} := f(\mathbf{w}) - \sigma/2$ ,  $u^{(\mathbf{w})} := f(\mathbf{w}) + \sigma/2$ , with  $\sigma$  being a measure of the (expected) deviation of  $f(\mathbf{w})$  from  $\bar{f}(\mathbf{w})$ . We can also use other constructions for this interval as long as they ensure that  $\bar{f}(\mathbf{w}) \in [l^{(\mathbf{w})}, u^{(\mathbf{w})}]$  with a sufficiently high probability. An optimist chooses the parameters of the model according to the most favourable - yet still plausible - value  $\tilde{f}(\mathbf{w}) := l^{(\mathbf{w})}$  of the performance criterion. Two examples of machine learning methods that use such an optimistic construction of an objective function are SRM [56, Ch. 11] and upper confidence bound (UCB) methods for sequential decision making [84, Sec. 2.2].

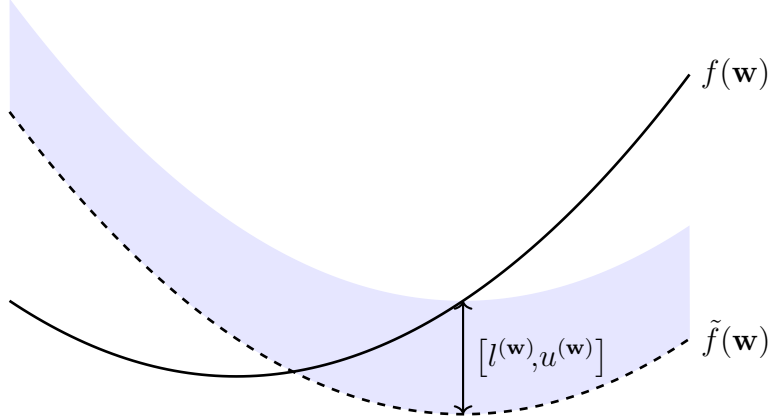


Fig. 25. apprentissage automatique methods learn paramètres du modèle  $\mathbf{w}$  by using some estimate of  $f(\mathbf{w})$  for the ultimate performance criterion  $\tilde{f}(\mathbf{w})$ . Using a modèle probabiliste, one can use  $f(\mathbf{w})$  to construct confidence intervals  $[l^{(\mathbf{w})}, u^{(\mathbf{w})}]$  which contain  $\tilde{f}(\mathbf{w})$  with high probability. The best plausible performance measure for a specific choice  $\mathbf{w}$  of paramètres du modèle is  $\tilde{f}(\mathbf{w}) := l^{(\mathbf{w})}$ .

**opérateur proximal** Étant donné une fonction convexe  $f(\mathbf{w}')$ , on définit son opérateur proximal comme suit [92, 97] :

$$\mathbf{prox}_{f(\cdot), \rho}(\mathbf{w}) := \operatorname{argmin}_{\mathbf{w}' \in \mathbb{R}^d} \left[ f(\mathbf{w}') + \frac{\rho}{2} \|\mathbf{w} - \mathbf{w}'\|_2^2 \right] \quad \text{avec } \rho > 0.$$

Comme illustré à la Figure 26, évaluer l'opérateur proximal revient à minimiser une version pénalisée de  $f(\mathbf{w}')$ . Le terme de pénalité est la distance euclidienne quadratique pondérée à un vecteur donné  $\mathbf{w}$ . L'opérateur proximal peut être interprété comme une généralisation du pas, défini pour une fonction régulière et convexe  $f(\mathbf{w}')$ . En effet, effectuer un pas de gradient avec une taille de pas  $\eta$  à partir du vecteur

actuel  $\mathbf{w}$  revient à appliquer l'opérateur proximal à la fonction linéarisée  $\tilde{f}(\mathbf{w}') = (\nabla f(\mathbf{w}))^T (\mathbf{w}' - \mathbf{w})$ , avec  $\rho = 1/\eta$ .

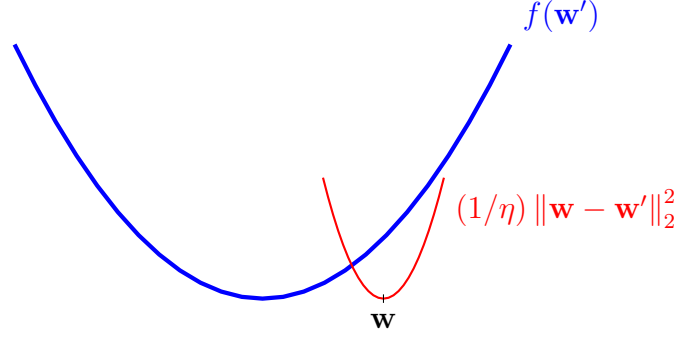


Fig. 26. Un pas généralisé met à jour un vecteur  $\mathbf{w}$  en minimisant une version pénalisée de la fonction  $f(\cdot)$ . Le terme de pénalité correspond à la distance euclidienne quadratique pondérée entre la variable d'optimisation  $\mathbf{w}'$  et le vecteur donné  $\mathbf{w}$ .

**outlier** Many apprentissage automatique methods are motivated by the hypothèse i.i.d., which interprets point de données as réalisations of i.i.d. VAs with a common loi de probabilité. The hypothèse i.i.d. is useful for applications where the statistical properties of the données generation process are stationary (or time-invariant) [89]. However, in some applications the données consists of a majority of regular point de données that conform with an hypothèse i.i.d. as well as a small number of point de données that have fundamentally different statistical properties compared to the regular point de données. We refer to a point de données that substantially deviates from the statistical properties of most point de données as an outlier. Different methods for outlier



detection use different measures for this deviation. Statistical learning theory studies fundamental limits on the ability to mitigate outliers reliably [90, 91].

**paramètres** Les paramètres d'un modèle en apprentissage automatique sont des quantités ajustables (c'est-à-dire apprenables ou modifiables) qui permettent de choisir parmi différentes fonctions hypothèse. Par exemple, le modèle linéaire  $\mathcal{H} := \{h^{(\mathbf{w})} : h^{(\mathbf{w})}(x) = w_1x + w_2\}$  correspond à l'ensemble des fonctions hypothèse  $h^{(\mathbf{w})}(x) = w_1x + w_2$  avec un choix particulier des paramètres  $\mathbf{w} = (w_1, w_2)^T \in \mathbb{R}^2$ . Un autre exemple de paramètres est le poids attribué à une connexion entre deux neurones dans un RNA.

**paramètres du modèle** Les paramètres d'un modèle sont des quantités utilisées pour sélectionner une fonction hypothèse spécifique à partir d'un modèle. On peut considérer une liste de paramètres de modèle comme un identifiant unique d'une fonction hypothèse, de la même manière qu'un numéro de sécurité sociale identifie une personne en France.

**partitionnement de données** Clustering methods decompose a given set of point de données into a few subsets, which are referred to as clusters. Each cluster consists of point de données that are more similar to each other than to point de données outside the cluster. Different clustering methods use different measures for the similarity between point de données and different forms of cluster representations. The clustering method  $k$ -moyennes uses the average caractéristique vector

(cluster moyenne) of a cluster as its representative. A popular soft clustering method based on GMM represents a cluster by a loi normale multivariée.

**pas de gradient (pas)** Étant donnée une fonction dérivable à valeurs réelles  $f(\cdot) : \mathbb{R}^d \rightarrow \mathbb{R}$  et un vecteur  $\mathbf{w} \in \mathbb{R}^d$ , le pas de gradient met à jour  $\mathbf{w}$  en ajoutant le gradient négatif mis à l'échelle  $\nabla f(\mathbf{w})$  pour obtenir le nouveau vecteur (voir Figure 27)

$$\hat{\mathbf{w}} := \mathbf{w} - \eta \nabla f(\mathbf{w}). \quad (8)$$

Mathématiquement, le pas est un opérateur (typiquement non-linéaire)  $\mathcal{T}^{(f,\eta)}$  paramétré par la fonction  $f$  et la taille de pas  $\eta$ .

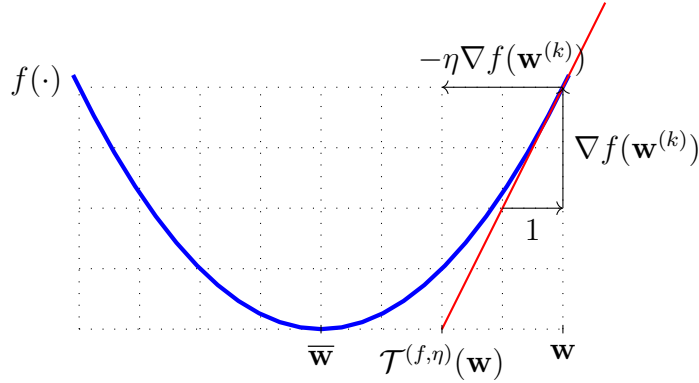


Fig. 27. Le pas classique (8) transforme un vecteur donné  $\mathbf{w}$  en le vecteur mis à jour  $\mathbf{w}'$ . Il définit un opérateur  $\mathcal{T}^{(f,\eta)}(\cdot) : \mathbb{R}^d \rightarrow \mathbb{R}^d : \mathbf{w} \mapsto \hat{\mathbf{w}}$ .

Notez que le pas (8) optimise localement - dans un voisinage dont la taille est déterminée par la taille de pas  $\eta$  - une approximation linéaire de la fonction  $f(\cdot)$ . Une généralisation naturelle de (8) est d'optimiser

localement la fonction elle-même - au lieu de son approximation linéaire - telle que

$$\hat{\mathbf{w}} = \underset{\mathbf{w}' \in \mathbb{R}^d}{\operatorname{argmin}} f(\mathbf{w}') + (1/\eta) \|\mathbf{w} - \mathbf{w}'\|_2^2. \quad (9)$$

Nous utilisons intentionnellement le même symbole  $\eta$  pour le paramètre dans (9) que celui utilisé pour la taille de pas dans (8). Plus le  $\eta$  choisi dans (9) est grand, plus la mise à jour avancera vers la réduction de la valeur de la fonction  $f(\hat{\mathbf{w}})$ . Notez que, tout comme le pas (8), la mise à jour (9) définit aussi un opérateur (typiquement non-linéaire) paramétré par la fonction  $f(\cdot)$  et le paramètre  $\eta$ . Pour une fonction convexe  $f(\cdot)$ , cet opérateur est connu sous le nom de opérateur proximal de  $f(\cdot)$  [92].

**perte (ou coût)** En apprentissage automatique, on utilise une fonction de perte  $L(\mathbf{z}, h)$  pour mesurer l'erreur commise lorsqu'une hypothèse est appliquée à un point de données. Par léger abus de langage, on utilise le terme *perte* à la fois pour désigner la fonction de perte  $L$  elle-même et la valeur spécifique  $L(\mathbf{z}, h)$  associée à une donnée  $\mathbf{z}$  et une hypothèse  $h$ .

**perte logistique** Considérons un point de données caractérisé par des caractéristiques  $\mathbf{x}$  et une étiquette binaire  $y \in \{-1, 1\}$ . On utilise une hypothèse à valeurs réelles  $h$  pour prédire l'étiquette  $y$  à partir des caractéristiques  $\mathbf{x}$ . La perte logistique associée à cette prédiction est définie comme

$$L((\mathbf{x}, y), h) := \log(1 + \exp(-yh(\mathbf{x}))). \quad (10)$$

Il est important de noter que l'expression (10) de la perte logistique

s'applique uniquement dans le cas d'un espace des étiquettes  $\mathcal{Y} = \{-1, 1\}$  et lorsque l'on utilise la règle de seuillage définie en (1).

**perte quadratique** La perte quadratique mesure l'erreur de prédiction d'une hypothèse  $h$  lorsqu'elle prédit une étiquette numérique  $y \in \mathbb{R}$  à partir des caractéristiques  $\mathbf{x}$  d'un point de données. Elle est définie par

$$L((\mathbf{x}, y), h) := (y - \underbrace{h(\mathbf{x})}_{=\hat{y}})^2.$$

**poids** Considérons un espace des hypothèses paramétré  $\mathcal{H}$ . On utilise le terme poids pour désigner des paramètres du modèle numériques utilisés pour pondérer les caractéristiques ou leurs transformations afin de calculer  $h^{(\mathbf{w})} \in \mathcal{H}$ . Un modèle linéaire utilise des poids  $\mathbf{w} = (w_1, \dots, w_d)^T$  pour calculer la combinaison linéaire  $h^{(\mathbf{w})}(\mathbf{x}) = \mathbf{w}^T \mathbf{x}$ . Les poids sont également utilisés dans les RNA pour former des combinaisons linéaires de caractéristiques ou des sorties de neurones dans les couches cachées.

**poids d'arête** Chaque arête  $\{i, i'\}$  d'un réseau d'apprentissage fédéré est associée à un poids d'arête non négatif  $A_{i,i'} \geq 0$ . Un poids d'arête nul  $A_{i,i'} = 0$  indique l'absence d'une arête entre les nœuds  $i, i' \in \mathcal{V}$ .

**point de données** Un point de données correspond à tout objet qui transmet de l'information [78]. Les points de données peuvent être des étudiants, des signaux radio, des arbres, des forêts, des images, des VA, des nombres réels ou des protéines. On caractérise les points de données à l'aide de deux types d'attributs. Le premier type d'attributs est appelé caractéristique. Les caractéristiques sont des attributs d'un point de données qui peuvent être mesurés ou calculés automatiquement. L'autre

type d'attributs est appelé étiquette. L'étiquette d'un point de données représente un fait (ou une quantité d'intérêt) de plus haut niveau. Contrairement aux caractéristiques, déterminer l'étiquette d'un point de données nécessite généralement des experts humains (experts du domaine). De manière générale, l'apprentissage automatique vise à prédire l'étiquette d'un point de données uniquement à partir de ses caractéristiques.

**point de données étiqueté** Un point de données dont l'étiquette est connue ou a été déterminée par un certain moyen, pouvant nécessiter une intervention humaine.

**predictor** A predictor is a real-valued hypothesis map. Given a point de données with caractéristiques  $\mathbf{x}$ , the value  $h(\mathbf{x}) \in \mathbb{R}$  is used as a prédiction for the true numeric étiquette  $y \in \mathbb{R}$  of the point de données.

**principe de minimisation des données** La réglementation européenne sur la protection des données inclut un principe de minimisation des données. Ce principe impose au responsable du traitement de limiter la collecte des informations personnelles à ce qui est directement pertinent et nécessaire pour atteindre un objectif spécifié. Les données doivent être conservées uniquement aussi longtemps que nécessaire pour remplir cet objectif [28, Article 5(1)(c)], [29].

**privacy attack** A privacy attack on an apprentissage automatique system aims to infer sensitive attributes of individuals by exploiting partial access to a trained apprentissage automatique modèle. One form of a privacy attack is model inversion.

See also: attack, sensitive attribute, model inversion, IA digne de confiance, Règlement général sur la protection des données (RGPD).

**privacy funnel** The privacy funnel is a method for learning privacy-friendly caractéristiques of point de données [93].

**privacy leakage** Consider an apprentissage automatique application that processes a jeu de données  $\mathcal{D}$  and delivers some output, such as the prédictions obtained for new point de données. Privacy leakage arises if the output carries information about a private (or sensitive) caractéristique of a point de données (which might be a human) of  $\mathcal{D}$ . Based on a modèle probabiliste for the données generation, we can measure the privacy leakage via the MI between the output and the sensitive caractéristique. Another quantitative measure of privacy leakage is DP. The relations between different measures of privacy leakage have been studied in the literature (see [94]).

**privacy protection** Consider some apprentissage automatique method  $\mathcal{A}$  that reads in a jeu de données  $\mathcal{D}$  and delivers some output  $\mathcal{A}(\mathcal{D})$ . The output could be the learned paramètres du modèle  $\hat{\mathbf{w}}$  or the prédiction  $\hat{h}(\mathbf{x})$  obtained for a specific point de données with caractéristiques  $\mathbf{x}$ . Many important apprentissage automatique applications involve point de données representing humans. Each point de données is characterized by caractéristiques  $\mathbf{x}$ , potentially a étiquette  $y$ , and a sensitive attribute  $s$  (e.g., a recent medical diagnosis). Roughly speaking, privacy protection means that it should be impossible to infer, from the output  $\mathcal{A}(\mathcal{D})$ , any of the sensitive attributes of point de données

in  $\mathcal{D}$ . Mathematically, privacy protection requires non-invertibility of the map  $\mathcal{A}(\mathcal{D})$ . In general, just making  $\mathcal{A}(\mathcal{D})$  non-invertible is typically insufficient for privacy protection. We need to make  $\mathcal{A}(\mathcal{D})$  sufficiently non-invertible.

**probabilistic principal component analysis (PPCA)** Probabilistic ACP extends basic ACP by using a modèle probabiliste for point de données. The modèle probabiliste of probabilistic ACP reduces the task of dimensionality reduction to an estimation problem that can be solved using EM methods.

**probabilité** On associe une valeur de probabilité, typiquement choisie dans l'intervalle  $[0, 1]$ , à chaque événement pouvant se produire dans une expérience aléatoire [6, 7, 39, 95].

**projection** Considérons un sous-ensemble  $\mathcal{W} \subseteq \mathbb{R}^d$  de l'espace euclidien de dimension  $d$ . On définit la projection  $P_{\mathcal{W}}(\mathbf{w})$  d'un vecteur  $\mathbf{w} \in \mathbb{R}^d$  sur  $\mathcal{W}$  comme

$$P_{\mathcal{W}}(\mathbf{w}) = \operatorname{argmin}_{\mathbf{w}' \in \mathcal{W}} \|\mathbf{w} - \mathbf{w}'\|_2. \quad (11)$$

Autrement dit,  $P_{\mathcal{W}}(\mathbf{w})$  est le vecteur dans  $\mathcal{W}$  qui est le plus proche de  $\mathbf{w}$ . La projection est bien définie uniquement pour les sous-ensembles  $\mathcal{W}$  pour lesquels le minimum ci-dessus existe [27].

**proximable** A convexe function for which the opérateur proximal can be computed efficiently is sometimes referred to as proximable or simple [96].

**prédiction** Une prédiction est une estimation ou une approximation d’une certaine quantité d’intérêt. L’apprentissage automatique se concentre sur l’apprentissage ou la recherche d’une fonction hypothèse qui prend en entrée les caractéristiques  $\mathbf{x}$  d’un point de données et fournit une prédiction  $\hat{y} := h(\mathbf{x})$  pour son étiquette  $y$ .

**pseudo-inverse** La pseudo-inverse de Moore–Penrose  $\mathbf{A}^+$  d’une matrice  $\mathbf{A} \in \mathbb{R}^{m \times d}$  généralise la notion de matrice inverse [3]. La pseudo-inverse apparaît naturellement dans le cadre de la régression Ridge appliquée à un jeu de données avec des étiquettes arbitraires  $\mathbf{y}$  et une matrice de caractéristiques  $\mathbf{X} = \mathbf{A}$  [20, Ch. 3]. Les paramètres du modèle appris par la régression Ridge sont donnés par

$$\hat{\mathbf{w}}^{(\alpha)} = (\mathbf{A}^T \mathbf{A} + \alpha \mathbf{I})^{-1} \mathbf{A}^T \mathbf{y}, \quad \alpha > 0.$$

On peut alors définir la pseudo-inverse  $\mathbf{A}^+ \in \mathbb{R}^{d \times m}$  avec la limite [98, Ch. 3]

$$\lim_{\alpha \rightarrow 0^+} \hat{\mathbf{w}}^{(\alpha)} = \mathbf{A}^+ \mathbf{y}.$$

Voir aussi : matrice inverse, régression Ridge, jeu de données, étiquette, matrice de caractéristiques, paramètres du modèle, régression Ridge.

**quadratic function** A function  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  of the form

$$f(\mathbf{w}) = \mathbf{w}^T \mathbf{Q} \mathbf{w} + \mathbf{q}^T \mathbf{w} + a,$$

with some matrix  $\mathbf{Q} \in \mathbb{R}^{d \times d}$ , vector  $\mathbf{q} \in \mathbb{R}^d$ , and scalar  $a \in \mathbb{R}$ .



**random forest** A random forest is a set of different arbre de décisions. Each of these arbre de décisions is obtained by fitting a perturbed copy of the original jeu de données.

**rectified linear unit (ReLU)** The ReLU is a popular choice for the fonction d'activation of a neuron within an RNA. It is defined as  $\sigma(z) = \max\{0, z\}$ , with  $z$  being the weighted input of the artificial neuron.

**regret** Le regret d'une hypothèse  $h$  par rapport à une autre hypothèse  $h'$  (considérée comme baseline) est défini comme la différence entre la perte engendrée par  $h$  et celle engendrée par  $h'$  [43]. L'hypothèse de référence  $h'$  est aussi appelée un expert.

**regularized empirical risk minimization (RERM)** Basic MRE learns a hypothèse (or trains a modèle)  $h \in \mathcal{H}$  based solely on the risque empirique  $\hat{L}(h|\mathcal{D})$  incurred on a ensemble d'entraînement  $\mathcal{D}$ . To make MRE less prone to surapprentissage, we can implement régularisation by including a (scaled) regularizer  $\mathcal{R}\{h\}$  in the learning objective. This leads to regularized empirical risk minimization (RERM),

$$\hat{h} \in \underset{h \in \mathcal{H}}{\operatorname{argmin}} \hat{L}(h|\mathcal{D}) + \alpha \mathcal{R}\{h\}. \quad (12)$$

The parameter  $\alpha \geq 0$  controls the régularisation strength. For  $\alpha = 0$ , we recover standard MRE without régularisation. As  $\alpha$  increases, the learned hypothèse is increasingly biased toward small values of  $\mathcal{R}\{h\}$ . The component  $\alpha \mathcal{R}\{h\}$  in the fonction objective of (12) can be intuitively understood as a surrogate for the increased average perte that may occur when predicting étiquettes for point de données outside

the ensemble d'entraînement. This intuition can be made precise in various ways. For example, consider a modèle linéaire trained using perte quadratique and the regularizer  $\mathcal{R}\{h\} = \|\mathbf{w}\|_2^2$ . In this setting,  $\alpha\mathcal{R}\{h\}$  corresponds to the expected increase in perte caused by adding VA normale centrée réduites to the vecteur de caractéristiquess in the ensemble d'entraînement [8, Ch. 3]. A principled construction for the regularizer  $\mathcal{R}\{h\}$  arises from approximate upper bounds on the generalization error. The resulting RERM instance is known as SRM [99, Sec. 7.2].

**regularized loss minimization (RLM)** See RERM.

**regularizer** A regularizer assigns each hypothèse  $h$  from a espace des hypothèses  $\mathcal{H}$  a quantitative measure  $\mathcal{R}\{h\}$  for how much its prédiction error on a ensemble d'entraînement might differ from its prédiction errors on point de donnéess outside the ensemble d'entraînement. Régression Ridge uses the regularizer  $\mathcal{R}\{h\} := \|\mathbf{w}\|_2^2$  for linear hypothèse maps  $h^{(\mathbf{w})}(\mathbf{x}) := \mathbf{w}^T \mathbf{x}$  [8, Ch. 3]. Lasso uses the regularizer  $\mathcal{R}\{h\} := \|\mathbf{w}\|_1$  for linear hypothèse maps  $h^{(\mathbf{w})}(\mathbf{x}) := \mathbf{w}^T \mathbf{x}$  [8, Ch. 3].

**Rényi divergence** The Rényi divergence measures the (dis)similarity between two loi de probabilités [100].

**risque** Considérons une hypothèse  $h$  utilisée pour prédire l'étiquette  $y$  d'un point de données basée sur ses caractéristiques  $\mathbf{x}$ . Nous mesurons la qualité d'une prédiction particulière en utilisant une fonction de perte  $L((\mathbf{x}, y), h)$ . Si nous interprétons les points de données comme les réalisations de VA i.i.d., alors  $L((\mathbf{x}, y), h)$  devient la réalisation d'une

VA. La hypothèse i.i.d. nous permet de définir le risque d'une hypothèse comme l'espérance de la perte  $\mathbb{E}\{L((\mathbf{x}, y), h)\}$ . Notons que le risque de  $h$  dépend à la fois du choix spécifique de la fonction de perte et de la loi de probabilité des points de données.

**risque empirique** Le risque empirique  $\hat{L}(h|\mathcal{D})$  d'une hypothèse sur un jeu de données  $\mathcal{D}$  correspond à la perte moyenne encourue par  $h$  lorsqu'elle est appliquée aux différents points de données de  $\mathcal{D}$ .

**robustness** Robustness is a key requirement for IA digne de confiance. It refers to the property of an apprentissage automatique system to maintain acceptable performance even when subjected to different forms of perturbations. These perturbations can be to the caractéristiques of a point de données in order to manipulate the prédiction delivered by a trained apprentissage automatique modèle. Robustness also includes the stability of MRE-based methods against perturbations of the ensemble d'entraînement. Such perturbations can occur within data poisoning attacks.

See also: IA digne de confiance, apprentissage automatique, caractéristique, point de données, prédiction, modèle, stability, MRE, ensemble d'entraînement, data poisoning, attack.

**Règlement général sur la protection des données (RGPD)** Le RGPD a été promulgué par l'Union européenne (UE) et est entré en vigueur le 25 mai 2018 [28]. Il garantit la protection de la vie privée et des droits liés aux données des individus au sein de l'UE. Le RGPD a des implications importantes sur la manière dont les données sont collectées,

stockées et utilisées dans les applications d'apprentissage automatique.

Parmi ses dispositions principales, on trouve :

- Principe de minimisation des données : les systèmes d'apprentissage automatique ne doivent utiliser que la quantité de données personnelles strictement nécessaire à leur finalité.
- Transparence et explicabilité : les systèmes d'apprentissage automatique doivent permettre aux utilisateurs de comprendre comment sont prises les décisions les concernant.
- Droits des personnes concernées : les utilisateurs doivent pouvoir accéder à leurs données personnelles, les rectifier, les supprimer, et s'opposer aux décisions automatisées ainsi qu'au profilage.
- Responsabilité : les organisations doivent garantir une sécurité robuste des données et prouver leur conformité au RGPD par la documentation et des audits réguliers.

**réalisation** Considérons une VA  $x$  qui associe à chaque élément (c'est-à-dire un résultat ou événement élémentaire)  $\omega \in \mathcal{P}$  d'un espace probabilisé  $\mathcal{P}$  un élément  $a$  d'un espace mesurable  $\mathcal{N}$  [2, 6, 39]. Une réalisation de  $x$  est tout élément  $a' \in \mathcal{N}$  pour lequel il existe un élément  $\omega' \in \mathcal{P}$  tel que  $x(\omega') = a'$ .

**récompense** Une récompense désigne une quantité observée (ou mesurée) qui permet d'estimer la perte subie par la prédiction (ou décision) d'une hypothèse  $h(\mathbf{x})$ . Par exemple, dans une application d'apprentissage automatique pour véhicules autonomes,  $h(\mathbf{x})$  pourrait représenter la

direction actuelle du volant d'un véhicule. On peut construire une récompense à partir des mesures d'un capteur de collision indiquant si le véhicule se dirige vers un obstacle. Une faible récompense est donnée à la direction  $h(\mathbf{x})$  si le véhicule avance dangereusement vers un obstacle.

**réduction de dimension** Les méthodes de réduction de dimension associent des caractéristiques brutes (généralement nombreuses) à un ensemble (relativement petit) de nouvelles caractéristiques. Ces méthodes peuvent être utilisées pour visualiser des points de données en apprenant deux caractéristiques pouvant servir de coordonnées pour une représentation dans un nuage de points.

**région de décision** Considérons une fonction hypothèse qui renvoie des valeurs d'un ensemble fini  $\mathcal{Y}$ . Pour chaque valeur (catégorie) d'étiquette  $a \in \mathcal{Y}$ , l'hypothèse  $h$  détermine un sous-ensemble de valeurs de caractéristiques  $\mathbf{x} \in \mathcal{X}$  telles que  $h(\mathbf{x}) = a$ . On appelle ce sous-ensemble une région de décision de l'hypothèse  $h$ .

**régression** Les problèmes de régression se concentrent sur la prédiction d'une étiquette numérique uniquement à partir des caractéristiques d'un point de données [8, Ch. 2].

**régression linéaire** La régression linéaire vise à apprendre une fonction hypothèse linéaire pour prédire une étiquette numérique à partir des caractéristiques numériques d'un points de données. La qualité d'une fonction hypothèse linéaire est mesurée par la moyenne de la perte quadratique subie sur un ensemble de points de données étiquetés, que nous appelons l'ensemble d'entraînement.

**régression logistique** La régression logistique apprend une fonction hypothèse (ou classifieur) linéaire  $h(\mathbf{x}) = \mathbf{w}^T \mathbf{x}$  pour prédire une étiquette binaire  $y$  à partir du vecteur de caractéristiques numérique  $\mathbf{x}$  d'un point de données. La qualité d'une telle fonction hypothèse linéaire est mesurée à l'aide de la perte logistique moyenne sur un ensemble de points de données étiquetés (c'est-à-dire l'ensemble d'entraînement).

**régression polynomiale** La régression polynomiale vise à apprendre une fonction hypothèse polynomiale pour prédire une étiquette numérique à partir des caractéristiques numériques d'un point de données. Pour des points de données caractérisés par une seule caractéristique numérique, la régression polynomiale utilise l'espace des hypothèses  $\mathcal{H}_d^{(\text{poly})} := \{h(x) = \sum_{j=0}^{d-1} x^j w_j\}$ . La qualité d'une fonction hypothèse polynomiale est mesurée via la perte quadratique moyenne encourue sur un ensemble de points de données étiquetés (appelé ensemble d'entraînement).

**régression Ridge** La régression Ridge apprend les poids  $\mathbf{w}$  d'une fonction hypothèse linéaire  $h^{(\mathbf{w})}(\mathbf{x}) = \mathbf{w}^T \mathbf{x}$ . La qualité d'un choix particulier des paramètres du modèle  $\mathbf{w}$  est mesurée par la somme de deux composantes. La première composante est la moyenne de la perte quadratique subie par  $h^{(\mathbf{w})}$  sur un ensemble de points de données étiquetés (i.e., l'ensemble d'entraînement). La deuxième composante est la norme euclidienne au carré, mise à l'échelle,  $\alpha \|\mathbf{w}\|_2^2$  avec un paramètre de régularisation  $\alpha > 0$ . Ajouter  $\alpha \|\mathbf{w}\|_2^2$  à la moyenne de la perte quadratique équivaut à remplacer chaque points de données initial par la réalisation d'une infinité sw VA i.i.d. centrées autour de ces points de données.

**régularisation** Un défi majeur des applications modernes d'apprentissage automatique est qu'elles utilisent souvent de grands modèles, avec une dimension effective de l'ordre du milliard. Entraîner un modèle de grande dimension à l'aide de méthodes de MRE basiques conduit souvent au surapprentissage : l'hypothèse apprise a de bonnes performances sur l'ensemble d'entraînement mais insuffisantes en dehors de celui-ci. La régularisation désigne des modifications apportées à une instance donnée de MRE afin d'éviter le surapprentissage, c'est-à-dire pour garantir que l'hypothèse apprise fonctionne presque aussi bien en dehors de l'ensemble d'entraînement. Il existe trois manières de mettre en œuvre la régularisation :

- 1) Élaguer le modèle : on réduit le modèle original  $\mathcal{H}$  pour obtenir un modèle plus petit  $\mathcal{H}'$ . Dans le cas d'un modèle paramétrique, cette réduction peut se faire via des contraintes sur les paramètres du modèle (par exemple  $w_1 \in [0.4, 0.6]$  pour le poids de la caractéristique  $x_1$  dans la régression linéaire).
- 2) Pénaliser la perte : on modifie la fonction objective de la MRE en ajoutant un terme de pénalité à l'erreur d'entraînement. Ce terme estime combien la perte (ou le risque) attendue est plus grande que la perte moyenne sur l'ensemble d'entraînement.
- 3) Augmentation de données : on peut agrandir l'ensemble d'entraînement  $\mathcal{D}$  en ajoutant des copies perturbées des points de données originaux de  $\mathcal{D}$ . Une telle perturbation consiste par exemple à ajouter la réalisation d'une VA au vecteur de caractéristiques d'un point de données.

La figure 28 illustre ces trois approches de régularisation. Ces approches sont étroitement liées et parfois entièrement équivalentes : la augmentation de données qui utilise des VA normales centrées réduites pour perturber les vecteurs de caractéristiques de l'ensemble d'entraînement dans le cas de la régression linéaire a le même effet que l'ajout du terme de pénalité  $\lambda \|\mathbf{w}\|_2^2$  à l'erreur d'entraînement (ce qui correspond à la régression Ridge). Le choix de la méthode de régularisation peut dépendre des ressources de calcul disponibles. Par exemple, il peut être bien plus facile de mettre en œuvre une augmentation de données que de réaliser un élagage de modèle.

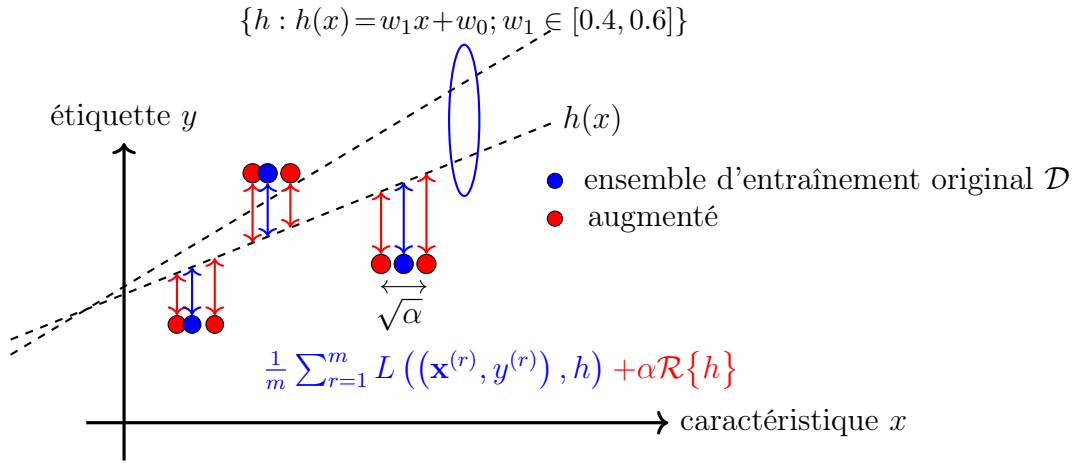


Fig. 28. Trois approches pour la régularisation: 1) augmentation de données; 2) pénalisation de la perte; et 3) élagage du modèle (via des contraintes sur les paramètres du modèle).

**régulière (ou lisse)** Une fonction à valeurs réelles  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  est dite régulière (ou lisse) si elle est dérivable et si son gradient  $\nabla f(\mathbf{w})$  est



continu en tout point  $\mathbf{w} \in \mathbb{R}^d$  (on parle aussi de fonction de classe  $\mathcal{C}^1$ ) [87, 101]. Une fonction régulière  $f$  est dite dérivable de gradient  $\beta$ -lipschitzien (ou  $\beta$ -smooth) si son gradient  $\nabla f(\mathbf{w})$  vérifie :

$$\|\nabla f(\mathbf{w}) - \nabla f(\mathbf{w}')\| \leq \beta \|\mathbf{w} - \mathbf{w}'\|, \text{ pour tout } \mathbf{w}, \mathbf{w}' \in \mathbb{R}^d.$$

La constante  $\beta$  mesure le degré de régularité de la fonction  $f$  : plus  $\beta$  est petit, plus  $f$  est lisse. Les problèmes d'optimisation comportant une fonction objective régulière peuvent être résolus efficacement par des méthodes basées sur le gradient. En effet, les méthodes basées sur le gradient approximent la fonction objective localement autour d'un point courant  $\mathbf{w}$  en utilisant son gradient. Cette approximation est pertinente lorsque le gradient ne varie pas trop rapidement. Cette affirmation intuitive peut être rendue rigoureuse en étudiant l'effet d'un seul pas avec une taille de pas  $\eta = 1/\beta$  (voir Figure 29).

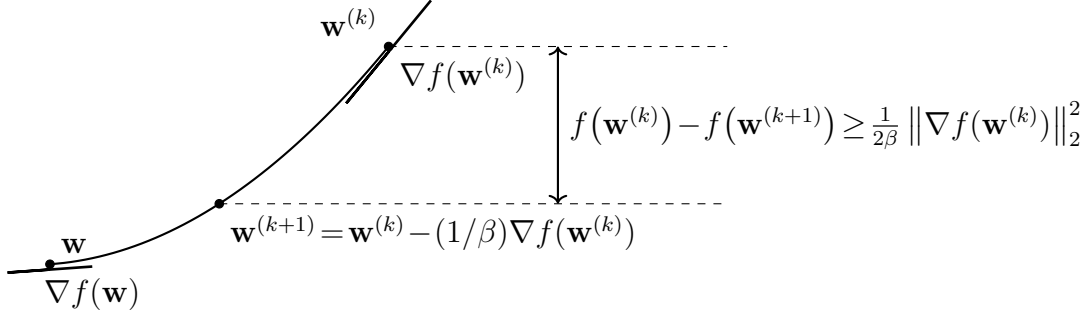


Fig. 29. Considérons une fonction objective  $f(\mathbf{w})$  qui est  $\beta$ -smooth. Effectuer un pas avec une taille de pas  $\eta = 1/\beta$  diminue la fonction objective d'au moins  $\frac{1}{2\beta} \|\nabla f(\mathbf{w}^{(k)})\|_2^2$  [87, 101, 102]. Notez que la taille de pas  $\eta = 1/\beta$  devient plus grande lorsque  $\beta$  diminue. Ainsi, pour des fonctions objectives plus lisses (c'est-à-dire avec un plus petit  $\beta$ ), on peut effectuer des pas plus grands.

**réseau d'apprentissage fédéré** Un réseau d'apprentissage fédéré est un graphe non orienté pondéré dont les nœuds représentent des générateurs de données visant à entraîner un modèle local (ou personnalisé). Chaque nœud dans un réseau d'apprentissage fédéré représente un appareil capable de collecter un jeu de données local et, à son tour, d'entraîner un modèle local. Les méthodes d'apprentissage fédéré apprennent une hypothèse locale  $h^{(i)}$ , pour chaque nœud  $i \in \mathcal{V}$ , telle qu'elle engendre une faible perte sur les jeux de données locaux.

**réseau de neurones artificiels (RNA)** Un RNA est une représentation graphique (circulation de signaux) d'une fonction qui associe les caractéristiques d'un point de données en entrée à une prédiction de l'étiquette correspondante en sortie. L'unité fondamentale d'un RNA est le neurone

artificiel, qui applique une fonction d'activation à ses entrées pondérées. Les sorties de ces neurones servent d'entrées à d'autres neurones, formant des couches interconnectées.

**réseau de neurones profond** Un réseau de neurones profond est un RNA avec un nombre (relativement) élevé de couches cachées. L'apprentissage profond est un terme générique désignant les méthodes d'apprentissage automatique qui utilisent un réseau de neurones profond comme modèle [103].

**sample** A finite sequence (or list) of point de données  $\mathbf{z}^{(1)}, \dots, \mathbf{z}^{(m)}$  that is obtained or interpreted as the réalisation of  $m$  i.i.d. VAs with a common loi de probabilité  $p(\mathbf{z})$ . The length  $m$  of the sequence is referred to as the taille d'échantillon.

**sample covariance matrix** The sample matrice de covariance  $\hat{\Sigma} \in \mathbb{R}^{d \times d}$  for a given set of vecteur de caractéristiquess  $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)} \in \mathbb{R}^d$  is defined as

$$\hat{\Sigma} = (1/m) \sum_{r=1}^m (\mathbf{x}^{(r)} - \hat{\mathbf{m}})(\mathbf{x}^{(r)} - \hat{\mathbf{m}})^T.$$

Here, we use the sample mean  $\hat{\mathbf{m}}$ .

**sample mean** The sample moyenne  $\mathbf{m} \in \mathbb{R}^d$  for a given jeu de données, with vecteur de caractéristiquess  $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)} \in \mathbb{R}^d$ , is defined as

$$\mathbf{m} = (1/m) \sum_{r=1}^m \mathbf{x}^{(r)}.$$

**semi-définie positive** Une matrice symétrique (à valeurs réelles)  $\mathbf{Q} = \mathbf{Q}^T \in \mathbb{R}^{d \times d}$  est dite semi-définie positive si  $\mathbf{x}^T \mathbf{Q} \mathbf{x} \geq 0$  pour tout vecteur  $\mathbf{x} \in \mathbb{R}^d$ .

La propriété d'être semi-définie positive peut être étendue des matrices aux applications noyau symétriques (à valeurs réelles)  $K : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$  (avec  $K(\mathbf{x}, \mathbf{x}') = K(\mathbf{x}', \mathbf{x})$ ) de la manière suivante : pour tout ensemble fini de vecteurs de caractéristiques  $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}$ , la matrice résultante  $\mathbf{Q} \in \mathbb{R}^{m \times m}$  avec pour coefficients  $Q_{r,r'} = K(\mathbf{x}^{(r)}, \mathbf{x}^{(r')})$  est semi-définie positive [85].

**sensitive attribute** apprentissage automatique revolves around learning a hypothèse map that allows us to predict the étiquette of a point de données from its caractéristiques. In some applications, we must ensure that the output delivered by an apprentissage automatique system does not allow us to infer sensitive attributes of a point de données. Which part of a point de données is considered a sensitive attribute is a design choice that varies across different application domains.

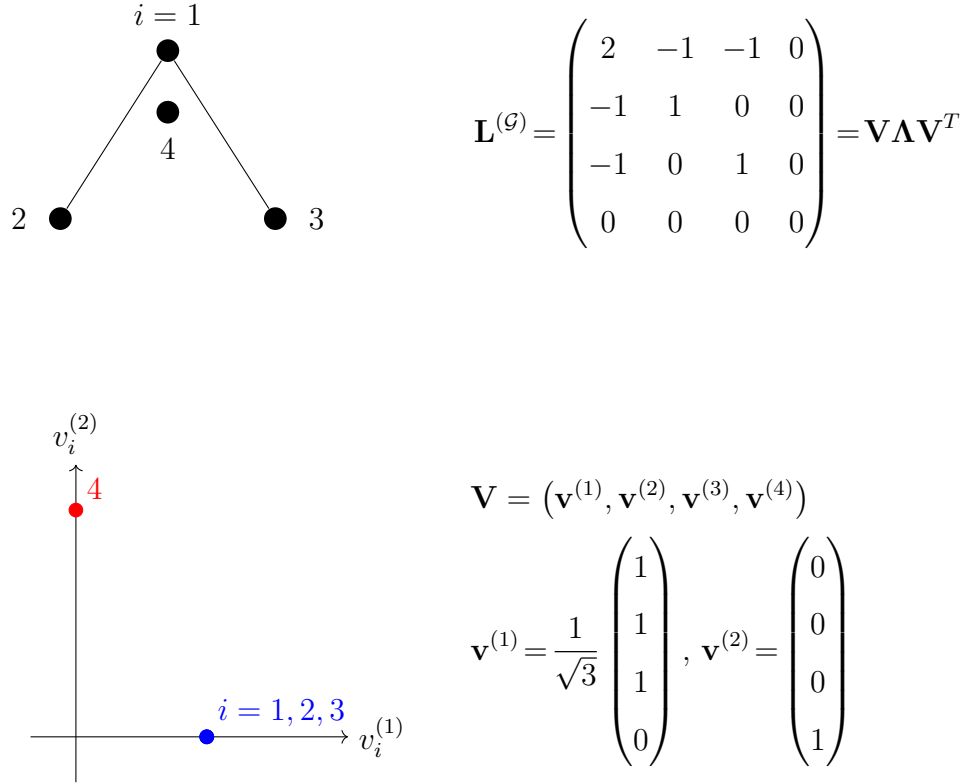
**similarity graph** Some apprentissage automatique applications generate point de données that are related by a domain-specific notion of similarity. These similarities can be represented conveniently using a similarity graphe  $\mathcal{G} = (\mathcal{V} := \{1, \dots, m\}, \mathcal{E})$ . The node  $r \in \mathcal{V}$  represents the  $r$ -th point de données. Two nodes are connected by an undirected edge if the corresponding point de données are similar.

**soft clustering** Soft partitionnement de données refers to the task of partitioning a given set of point de données into (a few) overlapping clusters. Each point de données is assigned to several different clusters with varying degrees of belonging. Soft partitionnement de données methods determine the degree of belonging (or soft cluster assignment)

for each point de données and each cluster. A principled approach to soft partitionnement de données is by interpreting point de données as i.i.d. réalisations of a GMM. We then obtain a natural choice for the degree of belonging as the conditional probabilité of a point de données belonging to a specific mixture component.

**sous-apprentissage** Considérons une méthode d'apprentissage automatique qui utilise la MRE pour apprendre une hypothèse minimisant le risque empirique sur un ensemble d'entraînement donné. On dit que cette méthode est en situation de sous-apprentissage si elle n'est pas capable d'apprendre une hypothèse avec un risque empirique suffisamment faible sur l'ensemble d'entraînement. En général, une méthode en situation de sous-apprentissage ne parviendra pas non plus à apprendre une hypothèse avec un risque faible.

**spectral clustering** Spectral partitionnement de données is a particular instance of graph clustering, i.e., it clusters point de données represented as the nodes  $i = 1, \dots, n$  of a graphe  $\mathcal{G}$ . Spectral partitionnement de données uses the vecteur propres of the matrice laplacienne  $\mathbf{L}^{(\mathcal{G})}$  to construct vecteur de caractéristiquess  $\mathbf{x}^{(i)} \in \mathbb{R}^d$  for each node (i.e., for each point de données)  $i = 1, \dots, n$ . We can feed these vecteur de caractéristiquess into espace euclidien-based partitionnement de données methods, such as  $k$ -moyennes or soft clustering via GMM. Roughly speaking, the vecteur de caractéristiquess of nodes belonging to a well-connected subset (or cluster) of nodes in  $\mathcal{G}$  are located nearby in the espace euclidien  $\mathbb{R}^d$  (see Figure 30).



**spectrogram** A spectrogram represents the time-frequency distribution of the energy of a time signal  $x(t)$ . Intuitively, it quantifies the amount of signal energy present within a specific time segment  $[t_1, t_2] \subseteq \mathbb{R}$  and frequency interval  $[f_1, f_2] \subseteq \mathbb{R}$ . Formally, the spectrogram of a signal is defined as the squared magnitude of its short-time Fourier transform (STFT) [104]. Figure 31 depicts a time signal along with its spectrogram.

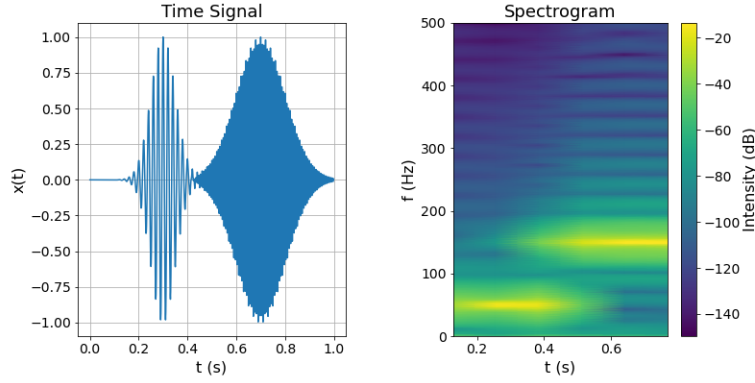


Fig. 31. Left: A time signal consisting of two modulated Gaussian pulses. Right: An intensity plot of the spectrogram.

The intensity plot of its spectrogram can serve as an image of a signal. A simple recipe for audio signal classification is to feed this signal image into réseau de neurones profonds originally developed for image classification and object detection [105]. It is worth noting that, beyond the spectrogram, several alternative representations exist for the time-frequency distribution of signal energy [106, 107].

**stability** Stability is a desirable property of a apprentissage automatique

method  $\mathcal{A}$  that maps a jeu de données  $\mathcal{D}$  (e.g., a ensemble d'entraînement) to an output  $\mathcal{A}(\mathcal{D})$ , such as learned paramètres du modèle or the prediction for a specific point de données. Intuitively,  $\mathcal{A}$  is stable if small changes in the input jeu de données  $\mathcal{D}$  lead to small changes in the output  $\mathcal{A}(\mathcal{D})$ . Several formal notions of stability exist that enable bounds on the généralisation error or risque of the method; see [56, Ch. 13]. To build intuition, consider the three datasets depicted in Fig. 32, each of which is equally likely under the same données-generating loi de probabilité. Since the optimal paramètres du modèle are determined by this underlying loi de probabilité, an accurate apprentissage automatique method  $\mathcal{A}$  should return the same (or very similar) output  $\mathcal{A}(\mathcal{D})$  for all three jeu de données. In other words, any useful  $\mathcal{A}$  must be robust to variability in sample réalisations from the same loi de probabilité, i.e., it must be stable.

**stochastic block model (SBM)** The stochastic block modèle is a probabilistic generative modèle for an undirected graphe  $\mathcal{G} = (\mathcal{V}, \mathcal{E})$  with a given set of nodes  $\mathcal{V}$  [108]. In its most basic variant, the stochastic block modèle generates a graphe by first randomly assigning each node  $i \in \mathcal{V}$  to a cluster index  $c_i \in \{1, \dots, k\}$ . A pair of different nodes in the graphe is connected by an edge with probabilité  $p_{i,i'}$  that depends solely on the étiquettes  $c_i, c_{i'}$ . The presence of edges between different pairs of nodes is statistically independent.

**stochastique** Une méthode est dite si elle comporte une composante aléatoire ou si elle est régie par des lois probabilistes. Les méthodes



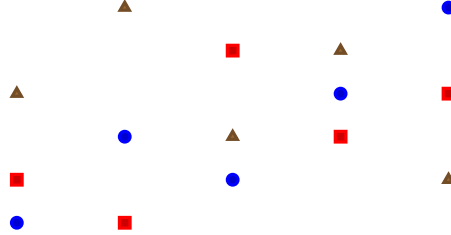


Fig. 32. Three jeu de données  $\mathcal{D}^{(*)}$ ,  $\mathcal{D}^{(\square)}$ , and  $\mathcal{D}^{(\triangle)}$ , each sampled independently from the same données-generating loi de probabilité. A stable apprentissage automatique method should return similar outputs when trained on any of these jeu de données.

d'apprentissage automatique utilisent l'aléatoire pour réduire la complexité computationnelle (voir, par exemple, SGD) ou pour modéliser l'incertitude dans les modèles probabilistes.

Voir aussi : incertitude, modèle probabiliste, SGD.

**strongly convex** A continuously dérivable real-valued function  $f(\mathbf{x})$  is strongly convexe with coefficient  $\sigma$  if  $f(\mathbf{y}) \geq f(\mathbf{x}) + \nabla f(\mathbf{x})^T (\mathbf{y} - \mathbf{x}) + (\sigma/2) \|\mathbf{y} - \mathbf{x}\|_2^2$  [87], [102, Sec. B.1.1].

**structural risk minimization (SRM)** Structural risk minimization (SRM) is an instance of RERM, which the modèle  $\mathcal{H}$  can be expressed as a countable union of sub-models:  $\mathcal{H} = \bigcup_{n=1}^{\infty} \mathcal{H}^{(n)}$ . Each sub-model  $\mathcal{H}^{(n)}$  permits the derivation of an approximate upper bound on the generalization error incurred when applying MRE to train  $\mathcal{H}^{(n)}$ . These

individual bounds—one for each sub-model—are then combined to form a regularizer used in the RERM objective. These approximate upper bounds (one for each  $\mathcal{H}^{(n)}$ ) are then combined to construct a regularizer for RERM [56, Sec. 7.2].

**subgradient** For a real-valued function  $f : \mathbb{R}^d \rightarrow \mathbb{R} : \mathbf{w} \mapsto f(\mathbf{w})$ , a vector  $\mathbf{a}$  such that  $f(\mathbf{w}) \geq f(\mathbf{w}') + (\mathbf{w} - \mathbf{w}')^T \mathbf{a}$  is referred to as a subgradient of  $f$  at  $\mathbf{w}'$  [109, 110].

**subgradient descent** Subgradient descent is a généralisation of descente de gradient that does not require differentiability of the function to be minimized. This généralisation is obtained by replacing the concept of a gradient with that of a subgradient. Similar to gradients, also subgradients allow us to construct local approximations of an fonction objective. The fonction objective might be the risque empirique  $\widehat{L}(h^{(\mathbf{w})}|\mathcal{D})$  viewed as a function of the paramètres du modèle  $\mathbf{w}$  that select a hypothèse  $h^{(\mathbf{w})} \in \mathcal{H}$ .

**support vector machine (SVM)** The SVM is a binary classification method that learns a linear hypothèse map. Thus, like régression linéaire and régression logistique, it is also an instance of MRE for the modèle linéaire. However, the SVM uses a different fonction de perte from the one used in those methods. As illustrated in Figure 33, it aims to maximally separate point de données from the two different classes in the espace des caractéristiques (i.e., maximum margin principle). Maximizing this separation is equivalent to minimizing a regularized variant of the hinge loss (4) [41, 66, 111].

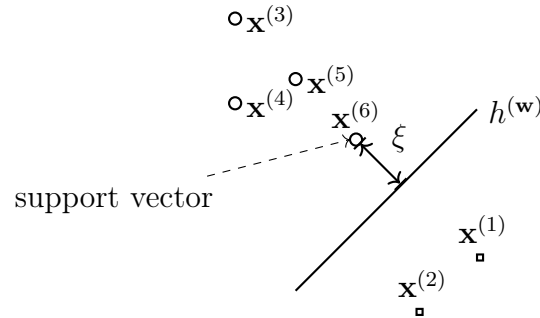


Fig. 33. The SVM learns a hypothèse (or classifieur)  $h^{(\mathbf{w})}$  with minimal average soft-margin hinge loss. Minimizing this perte is equivalent to maximizing the margin  $\xi$  between the frontière de décision de  $h^{(\mathbf{w})}$  and each class of the ensemble d'entraînement.

The above basic variant of SVM is only useful if the point de données from different categories can be (approximately) linearly separated. For an apprentissage automatique application where the categories are not derived from a noyau.

**surapprentissage** Considérons une méthode d'apprentissage automatique qui utilise la MRE pour apprendre une hypothèse avec le risque empirique minimal sur un ensemble d'entraînement donné. Une telle méthode fait du surapprentissage sur l'ensemble d'entraînement si elle apprend une hypothèse avec un petit risque empirique sur l'ensemble d'entraînement mais une perte significativement plus grande en dehors de cet ensemble.

**taille d'échantillon** Le nombre de points de données individuels contenus

dans un jeu de données.

**taille de pas** Voir taux d'apprentissage.

**taux d'apprentissage** Considérons une méthode itérative d'apprentissage automatique pour trouver ou apprendre une hypothèse utile  $h \in \mathcal{H}$ . Une telle méthode itérative répète des étapes computationnelles (de mise à jour) similaires qui ajustent ou modifient l'hypothèse actuelle afin d'obtenir une hypothèse améliorée. Un exemple bien connu de cette méthode itérative est la descente de gradient et ses variantes, SGD et descente de gradient avec projection. Un paramètre clé d'une méthode itérative est le taux d'apprentissage. Le taux d'apprentissage contrôle l'ampleur selon laquelle l'hypothèse courante peut être modifiée durant une seule itération. Un exemple bien connu de tel paramètre est la taille de pas utilisée lors d'une descente de gradient [8, Ch. 5].

**total variation** See GTV.

**transformation de caractéristiques** Une transformation de caractéristiques est une application qui transforme les caractéristiques originales d'un point de données en de nouvelles caractéristiques. Les nouvelles caractéristiques obtenues peuvent être préférables aux caractéristiques d'origine pour plusieurs raisons. Par exemple, l'agencement des points de données peut devenir plus simple (ou plus linéaire) dans le nouvel espace des caractéristiques, permettant ainsi l'utilisation de modèles linéaires dans ce nouvel espace. Cette idée est un moteur central du développement des méthodes à noyau [85]. Par ailleurs, les couches cachées d'un réseau de neurones profond peuvent être interprétées

comme une transformation de caractéristiques entraînable, suivie d'un modèle linéaire sous forme de couche de sortie. Une autre raison d'apprendre une transformation de caractéristiques peut être de réduire le surapprentissage et d'assurer une meilleure interprétabilité en apprenant un petit nombre de caractéristiques pertinentes [52]. Le cas particulier d'une transformation de caractéristiques produisant deux caractéristiques numériques est particulièrement utile pour la visualisation des données. En effet, on peut représenter les points de données dans un nuage de points en utilisant ces deux caractéristiques comme coordonnées.

**transparence** La transparence est une exigence fondamentale pour une IA digne de confiance [112]. Dans le contexte des méthodes d'apprentissage automatique, le terme est souvent utilisé de manière interchangeable avec explicabilité [50, 113]. Cependant, dans le cadre plus large des systèmes d'IA, la transparence va au-delà de l'explicabilité et inclut de fournir des informations sur les limitations, la fiabilité et l'utilisation prévue du système. Dans les systèmes de diagnostic médical, la transparence exige de révéler le niveau de confiance associé aux prédictions produites par un modèle entraîné. Dans l'évaluation du crédit, les décisions prises par des systèmes d'IA doivent être accompagnées d'explications sur les facteurs contributifs, tels que le revenu ou l'historique de crédit. Ces explications permettent aux humains (par exemple, un demandeur de prêt) de comprendre et de contester les décisions automatisées. Certaines méthodes d'apprentissage automatique offrent intrinsèquement une certaine transparence. Par exemple, la régression logistique fournit

une mesure quantitative de la fiabilité d’une classification à travers la valeur  $|h(\mathbf{x})|$ . Les arbres de décision en sont un autre exemple, car ils permettent d’utiliser des règles de décision lisibles par l’humain [46]. La transparence implique aussi de signaler clairement lorsqu’un utilisateur interagit avec un système d’IA. Par exemple, un chatbot alimenté par l’IA doit informer l’utilisateur qu’il interagit avec un système automatisé et non un humain. Enfin, la transparence suppose une documentation complète précisant l’objectif et les choix de conception du système d’IA. Des outils comme les fiches techniques de modèle [77] ou les cartes descriptives de systèmes d’IA [114] aident les praticiens à comprendre les cas d’usage prévus ainsi que les limitations du système [115].

**tâche d’apprentissage** Considérons un jeu de données  $\mathcal{D}$  constitué de plusieurs points de données, chacun étant caractérisé par des caractéristiques  $\mathbf{x}$ . Par exemple, le jeu de données  $\mathcal{D}$  peut être constitué des images d’une base de données particulière. Parfois, il peut être utile de représenter un jeu de données  $\mathcal{D}$ , ainsi que le choix des caractéristiques, par une loi de probabilité  $p(\mathbf{x})$ . Une tâche d’apprentissage associée à  $\mathcal{D}$  consiste en un choix spécifique pour l’étiquette d’un point de données et l’espace des étiquettes correspondant. Étant donné un choix de fonction de perte et de modèle, une tâche d’apprentissage donne lieu à une instance de MRE. Ainsi, on pourrait aussi définir une tâche d’apprentissage via une instance de MRE, c’est-à-dire via une fonction objective. Remarquons que, pour un même jeu de données, on obtient différentes tâches d’apprentissage en utilisant différents choix de caractéristiques et d’étiquette d’un point de données. Ces

tâches d'apprentissage sont liées, puisqu'elles sont basées sur le même jeu de données, et les résoudre conjointement (via des méthodes de apprentissage multitâche) est en général préférable à des résolutions distinctes [116], [117], [118].

**upper confidence bound (UCB)** Consider a apprentissage automatique application that requires selecting, at each time step  $k$ , an action  $a_k$  from a finite set of alternatives  $\mathcal{A}$ . The utility of selecting action  $a_k$  is quantified by a numeric récompense signal  $r^{(a_k)}$ . A widely used modèle probabiliste for this type of sequential decision-making problem is the stochastic multi-armed bandit setting [84]. In this model, the reward  $r^{(a)}$  is viewed as the réalisation of a VA with unknown moyenne  $\mu^{(a)}$ . Ideally, we would always choose the action with the largest expected récompense  $\mu^{(a)}$ , but these means are unknown and must be estimated from observed données. Simply choosing the action with the largest estimate  $\hat{\mu}^{(a)}$  can lead to suboptimal outcomes due to estimation uncertainty. The UCB strategy addresses this by selecting actions not only based on their estimated means but also by incorporating a term that reflects the uncertainty in these estimates—favouring actions with high potential reward and high uncertainty. Theoretical guarantees for the performance of UCB strategies, including logarithmic regret bounds, are established in [84].

**valeur propre** On qualifie de valeur propre d'une matrice carrée  $\mathbf{A} \in \mathbb{R}^{d \times d}$  le nombre  $\lambda \in \mathbb{R}$  s'il existe un vecteur non nul  $\mathbf{x} \in \mathbb{R}^d \setminus \{\mathbf{0}\}$  tels que  $\mathbf{Ax} = \lambda\mathbf{x}$ .

**validation** Considérons une hypothèse  $\hat{h}$  apprise à l'aide d'une méthode d'apprentissage automatique, par exemple en résolvant la MRE sur un ensemble d'entraînement  $\mathcal{D}$ . La validation désigne la pratique consistant à évaluer la perte encourue par l'hypothèse  $\hat{h}$  sur un ensemble de points de données qui ne sont pas contenus dans le ensemble d'entraînement  $\mathcal{D}$ .

**Vapnik–Chervonenkis dimension (VC dimension)** The VC dimension of an infinite space des hypothèses is a widely-used measure for its size. We refer to the literature (see [56]) for a precise definition of VC dimension as well as a discussion of its basic properties and use in apprentissage automatique.

**variable aléatoire (VA)** Une VA est une fonction qui associe chaque événement élémentaire d'un espace probabilisé  $\mathcal{P}$  à une valeur dans un espace d'arrivée [19], [6]. L'espace probabilisé est composé d'événements élémentaires et est muni d'une mesure de probabilité qui attribue des probabilités aux sous-ensembles de  $\mathcal{P}$ . Les différents types de VA comprennent :

- les VA binaires, qui associent chaque événement élémentaire à un élément d'un ensemble binaire (par exemple,  $\{-1, 1\}$  ou  $\{\text{chat}, \text{pas chat}\}$ );
- les VA à valeurs réelles, qui prennent des valeurs dans  $\mathbb{R}$ ;
- les VA vectorielles, qui associent chaque événement élémentaire à un vecteur de l'espace euclidien  $\mathbb{R}^d$ .

La théorie des probabilités utilise le concept d'espaces mesurables pour



définir rigoureusement et étudier les propriétés de (grandes) collections de VA [6].

**variable aléatoire normale centrée réduite** Une VA normale centrée réduite est une VA réelle  $x$  dont la fonction de densité de probabilité est donnée par [7], [19], [79]

$$p(x) = \frac{1}{\sqrt{2\pi}} \exp^{-x^2/2}.$$

Étant donnée une VA normale centrée réduite  $x$ , on peut construire une VA normale  $x'$  ayant pour moyenne  $\mu$  et variance  $\sigma^2$  via  $x' := \sigma(x + \mu)$ . La loi de probabilité d'une VA normale est appelée loi normale (ou loi gaussienne), notée  $\mathcal{N}(\mu, \sigma)$ .

Un vecteur aléatoire gaussien  $\mathbf{x} \in \mathbb{R}^d$  ayant pour matrice de covariance  $\mathbf{C}$  et pour moyenne  $\boldsymbol{\mu}$  peut être construit via  $\mathbf{x} := \mathbf{A}(\mathbf{z} + \boldsymbol{\mu})$ , où  $\mathbf{A}$  est une matrice telle que  $\mathbf{A}\mathbf{A}^T = \mathbf{C}$ , et  $\mathbf{z} := (z_1, \dots, z_d)^T$  est un vecteur dont les composantes sont des VA normales centrées réduites i.i.d.  $z_1, \dots, z_d$ .

Les vecteurs aléatoires gaussiens constituent un cas particulier des processus gaussiens, qui sont des transformations linéaires de suites infinies de VA normales centrées réduites [119].

Les VA normales sont largement utilisées comme modèles probabilistes pour l'analyse statistique en apprentissage automatique. Leur importance provient en partie du théorème central limite, qui stipule que la moyenne d'un nombre croissant de VA indépendantes (pas nécessairement normales) converge vers une VA normale [38].

Voir aussi : loi de probabilité, espace probabilisé.

**variance** La variance d'une VA réelle  $x$  est définie comme l'espérance  $\mathbb{E}\{(x - \mathbb{E}\{x\})^2\}$  de la différence au carré entre  $x$  et son espérance  $\mathbb{E}\{x\}$ . On étend cette définition aux VA vectorielles  $\mathbf{x}$  avec  $\mathbb{E}\{\|\mathbf{x} - \mathbb{E}\{\mathbf{x}\}\|_2^2\}$ .

**vecteur de caractéristiques** Un vecteur de caractéristiques est un vecteur  $\mathbf{x} = (x_1, \dots, x_d)^T$  dont les composantes sont des caractéristiques individuelles  $x_1, \dots, x_d$ . De nombreuses méthodes d'apprentissage automatique utilisent des vecteurs de caractéristiques appartenant à un espace euclidien de dimension finie  $\mathbb{R}^d$ . Cependant, pour certaines méthodes d'apprentissage automatique, il peut être plus pratique de travailler avec des vecteurs de caractéristiques appartenant à un espace vectoriel de dimension infinie (par exemple, voir la méthode à noyau).

**vecteur propre** An eigenvector of a matrix  $\mathbf{A} \in \mathbb{R}^{d \times d}$  is a non-zero vector  $\mathbf{x} \in \mathbb{R}^d \setminus \{\mathbf{0}\}$  such that  $\mathbf{A}\mathbf{x} = \lambda\mathbf{x}$  with some valeur propre  $\lambda$ .

**vertical federated learning (vertical FL)** Vertical apprentissage fédéré uses jeu de données locals that are constituted by the same point de données but characterizing them with different caractéristiques [120]. For example, different healthcare providers might all contain information about the same population of patients. However, different healthcare providers collect different measurements (e.g., blood values, electrocardiography, lung X-ray) for the same patients.

**voisinage** Le voisinage d'un nœud  $i \in \mathcal{V}$  est le sous-ensemble de nœuds constitué des voisins de  $i$ .

**voisins** Les voisins d'un nœud  $i \in \mathcal{V}$  dans un réseau d'apprentissage fédéré

sont les nœuds  $i' \in \mathcal{V} \setminus \{i\}$  qui sont connectés (via une arête) au nœud  $i$ .

**zero-gradient condition** Consider the unconstrained optimization problem  $\min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w})$  with a régulière and convexe fonction objective  $f(\mathbf{w})$ . A necessary and sufficient condition for a vector  $\hat{\mathbf{w}} \in \mathbb{R}^d$  to solve this problem is that the gradient  $\nabla f(\hat{\mathbf{w}})$  is the zero vector,

$$\nabla f(\hat{\mathbf{w}}) = \mathbf{0} \Leftrightarrow f(\hat{\mathbf{w}}) = \min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w}).$$

**0/1 loss** The 0/1 perte  $L^{(0/1)}((\mathbf{x}, y), h)$  measures the quality of a classifieur  $h(\mathbf{x})$  that delivers a prédiction  $\hat{y}$  (e.g., via thresholding (1)) for the étiquette  $y$  of a point de données with caractéristiques  $\mathbf{x}$ . It is equal to 0 if the prédiction is correct, i.e.,  $L^{(0/1)}((\mathbf{x}, y), h) = 0$  when  $\hat{y} = y$ . It is equal to 1 if the prédiction is wrong, i.e.,  $L^{(0/1)}((\mathbf{x}, y), h) = 1$  when  $\hat{y} \neq y$ .

**épigraphe** L'épigraphe d'une fonction à valeurs réelles  $f : \mathbb{R}^n \rightarrow \mathbb{R} \cup \{+\infty\}$  est l'ensemble des points situés sur sa courbe ou au dessus :

$$\text{epi}(f) = \{(\mathbf{x}, t) \in \mathbb{R}^n \times \mathbb{R} \mid f(\mathbf{x}) \leq t\}.$$

Une fonction est convexe si et seulement si son épigraphe est un ensemble convexe [27], [109].

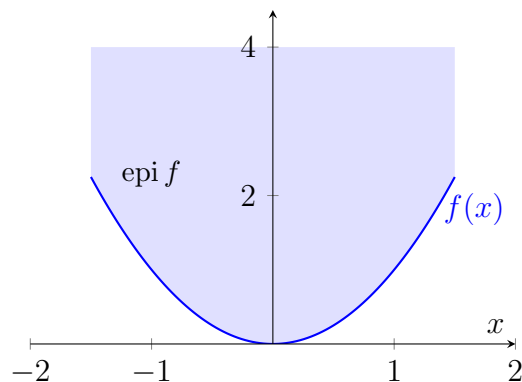


Fig. 34. Épigraphe de la fonction  $f(x) = x^2$  (i.e., la zone colorée).

Voir aussi : fonction, convexe.

**étiquette** Une étiquette est un fait ou une quantité d'intérêt de plus haut niveau associée à un point de données. Par exemple, si le point de données est une image, l'étiquette peut indiquer si l'image contient un chat ou non. Les synonymes de « étiquette », couramment utilisés dans certains domaines, incluent « variable réponse », « variable de sortie » et « cible » [22], [23], [24].

# Index

- 0/1 loss, 131
- $k$ -fold cross-validation ( $k$ -fold CV),  
21
- $k$ -moyennes, 21
- absolute error loss, 21
- accuracy, 21
- algebraic connectivity, 22
- algorithme, 22
- analyse en composantes principales  
(ACP), 23
- appareil, 23
- application linéaire, 23
- application programming interface  
(API), 24
- apprentissage automatique, 25
- apprentissage en ligne, 25
- apprentissage fédéré, 25
- apprentissage multitâche, 26
- apprentissage semi-supervisé, 26
- arbre de décision, 26
- aspects computationnels, 27
- aspects statistiques, 28
- attack, 28
- augmentation de données, 28
- autoencoder, 29
- backdoor, 30
- bagging, 30
- bandit manchot, 30
- baseline, 31
- Bayes estimator, 33
- Bayes risk, 33
- biais, 34
- boosting, 34
- bootstrap, 35
- borne supérieure, 35
- caractéristique, 36
- classification, 36
- classification multi-classe, 36
- classifieur, 37
- classifieur linéaire, 37
- cluster, 37
- clustered federated learning (CFL),  
38
- clustering assumption, 39
- condition number, 39
- confusion matrix, 39

connected graph, 40  
 convex clustering, 40  
 Courant–Fischer–Weyl min-max  
     characterization, 41  
 covariance, 41  
 critère d’arrêt, 41  
  
 data normalization, 42  
 data poisoning, 42  
 degree of belonging, 43  
 denial-of-service attack, 43  
 density-based spatial clustering of  
     applications with noise  
     (DBSCAN), 43  
 descente de gradient, 44  
 descente de gradient stochastique  
     (SGD), 45  
 differential privacy (DP), 46  
 dimension effective, 47  
 discrepancy, 47  
 distributed algorithm, 47  
 données, 48  
 décomposition en valeurs  
     singulières, 49  
 décomposition en éléments propres,  
     49  
  
 dérivable, 49  
 déterminant, 50  
  
 ensemble d’entraînement (ou  
     d’apprentissage), 50  
 ensemble de test (ou jeu de test),  
     51  
 ensemble de validation (ou jeu de  
     validation), 51  
 erreur d’entraînement, 51  
 erreur de validation, 51  
 espace des caractéristiques, 52  
 espace des paramètres, 52  
 espace des étiquettes, 53  
 espace euclidien, 54  
 espace probabilisé, 54  
 espérance, 55  
 estimation error, 55  
 expectation-maximization (EM),  
     55  
 expert, 56  
 explainable empirical risk  
     minimization (EERM), 56  
 explainable machine learning  
     (explainable ML), 57  
 explicabilité, 57

- explication, 57
- Explications locales interprétables
  - et agnostiques au modèle (LIME), 58
- feature learning, 59
- FedAvg, 60
- FedGD, 60
- FedProx, 60
- FedRelax, 60
- FedSGD, 61
- flow-based clustering, 61
- fonction, 61
- fonction d'activation, 61
- fonction de densité de probabilité, 61
- fonction de perte (ou de coût), 62
- fonction objective, 63
- frontière de décision, 63
- Gaussian mixture model (GMM), 63
- Gaussian Process (GP), 63
- generalized total variation (GTV), 65
- generalized total variation minimization (GTVMin), 65
- geometric median (GM), 65
- gradient, 66
- grand modèle de langage (GML), 66
- graph clustering, 67
- graphe, 67
- hard clustering, 69
- high-dimensional regime, 69
- Hilbert space, 70
- hinge loss, 70
- histogram, 71
- horizontal federated learning (HFL), 71
- Huber loss, 72
- Huber regression, 72
- hypothèse, 73
- hypothèse d'indépendance et de distribution identique (hypothèse i.i.d.), 73
- IA digne de confiance, 74
- incertitude, 73
- indépendantes et identiquement distribuées (i.i.d.), 73

Institut météorologique finlandais (FMI), 74	mean squared estimation error (MSEE), 83
intelligence artificielle (IA), 74	minimum, 84
interprétabilité, 75	missing data, 84
jeu de données, 77	model selection, 85
jeu de données local, 79	modèle, 85
Kullback-Leibler divergence (KL divergence), 79	modèle linéaire, 85
least absolute deviation regression, 79	modèle local, 86
least absolute shrinkage and selection operator (Lasso), 79	modèle probabiliste, 86
loi de probabilité, 80	moyenne, 86
loi des grands nombres, 80	mutual information (MI), 87
loi normale multivariée, 80	méthode à noyau, 87
lot, 81	méthodes basées sur le gradient, 88
map, 23	nearest neighbor (NN), 88
matrice de caractéristiques, 81	networked data, 89
matrice de covariance, 81	networked exponential families (nExpFam), 89
matrice inverse, 81	networked federated learning (NFL), 89
matrice laplacienne, 82	networked model, 89
maximum, 83	node degree, 89
maximum de vraisemblance, 83	non régulière (ou non lisse), 90
	norme, 90
	noyau, 90
	nuage de points, 90



online algorithm, 91	probabilistic principal component
online gradient descent (online	analysis (PPCA), 103
GD), 92	probabilité, 103
optimism in the face of uncertainty,	projection, 103
94	proximable, 103
opérateur proximal, 95	prédiction, 104
outlier, 96	pseudo-inverse, 104
paramètres, 97	quadratic function, 104
paramètres du modèle, 97	Rényi divergence, 106
partitionnement de données, 97	random forest, 105
pas de gradient (pas), 98	rectified linear unit (ReLU), 105
perte (ou coût), 99	regret, 105
perte logistique, 99	regularized empirical risk
perte quadratique, 100	minimization (RERM),
poids, 100	105
poids d'arête, 100	regularized loss minimization
point de données, 100	(RLM), 106
point de données étiqueté, 101	regularizer, 106
predictor, 101	risque, 106
principe de minimisation des	risque empirique, 107
données, 101	robustness, 107
privacy attack, 101	Règlement général sur la
privacy funnel, 102	protection des données
privacy leakage, 102	(RGPD), 107
privacy protection, 102	réalisation, 108

- réduction de dimension, 109
- région de décision, 109
- régression, 109
- régression linéaire, 109
- régression logistique, 110
- régression polynomiale, 110
- régression Ridge, 110
- régularisation, 111
- régulière (ou lisse), 112
- réseau d'apprentissage fédéré, 114
- réseau de neurones artificiels  
(RNA), 114
- réseau de neurones profond, 115
- sample, 115
- sample covariance matrix, 115
- sample mean, 115
- semi-définie positive, 115
- sensitive attribute, 116
- similarity graph, 116
- soft clustering, 116
- sous-apprentissage, 117
- spectral clustering, 117
- spectrogram, 119
- stability, 119
- stochastic block model (SBM), 120
- stochastique, 120
- strongly convex, 121
- structural risk minimization  
(SRM), 121
- subgradient, 122
- subgradient descent, 122
- support vector machine (SVM),  
122
- surapprentissage, 123
- taille d'échantillon, 123
- taille de pas, 124
- taux d'apprentissage, 124
- total variation, 124
- transformation de caractéristiques,  
124
- transparence, 125
- tâche d'apprentissage, 126
- upper confidence bound (UCB),  
127
- valeur propre, 127
- validation, 128
- Vapnik–Chervonenkis dimension  
(VC dimension), 128
- variable aléatoire (VA), 128

variable aléatoire normale centrée	FL), 130
réduite, 129	voisinage, 130
variance, 130	voisins, 130
vecteur de caractéristiques, 130	zero-gradient condition, 131
vecteur propre, 130	épigraphe, 131
vertical federated learning (vertical	étiquette, 132

## References

- [1] W. Rudin, *Real and Complex Analysis*, 3rd ed. New York, NY, USA: McGraw-Hill, 1987.
- [2] W. Rudin, *Principles of Mathematical Analysis*, 3rd ed. New York, NY, USA: McGraw-Hill, 1976.
- [3] G. H. Golub and C. F. Van Loan, *Matrix Computations*, 4th ed. Baltimore, MD, USA: The Johns Hopkins Univ. Press, 2013.
- [4] G. H. Golub and C. F. Van Loan, “An analysis of the total least squares problem,” *SIAM J. Numer. Anal.*, vol. 17, no. 6, pp. 883–893, Dec. 1980, doi: 10.1137/0717073.
- [5] A. Klenke, *Probability Theory: A Comprehensive Course*, 3rd ed. Cham, Switzerland: Springer Nature, 2020.
- [6] P. Billingsley, *Probability and Measure*, 3rd ed. New York, NY, USA: Wiley, 1995.
- [7] D. P. Bertsekas and J. N. Tsitsiklis, *Introduction to Probability*, 2nd ed. Belmont, MA, USA: Athena Scientific, 2008.
- [8] A. Jung, *Machine Learning: The Basics*. Singapore, Singapore: Springer Nature, 2022.
- [9] T. H. Cormen, C. E. Leiserson, R. L. Rivest, and C. Stein, *Introduction to Algorithms*. Cambridge, MA, USA: MIT Press, 2022. [Online]. Available: <http://ebookcentral.proquest.com/lib/aalto-ebooks/detail.action?docID=6925615>

- [10] M. Sipser, *Introduction to the Theory of Computation*, 3rd ed. Andover, U.K.: Cengage Learning, 2013.
- [11] R. Motwani and P. Raghavan, *Randomized Algorithms*. Cambridge, U.K.: Cambridge Univ. Press, 1995.
- [12] R. G. Gallager, *Stochastic Processes: Theory for Applications*. New York, NY, USA: Cambridge Univ. Press, 2013.
- [13] L. Richardson and M. Amundsen, *RESTful Web APIs*. Sebastopol, CA, USA: O'Reilly Media, 2013.
- [14] T. Hastie, R. Tibshirani, and M. Wainwright, *Statistical Learning with Sparsity: The Lasso and Generalizations*. Boca Raton, FL, USA: CRC Press, 2015.
- [15] O. Chapelle, B. Schölkopf, and A. Zien, Eds. *Semi-Supervised Learning*. Cambridge, MA, USA: MIT Press, 2006.
- [16] M. P. Salinas et al., "A systematic review and meta-analysis of artificial intelligence versus clinicians for skin cancer diagnosis," *npj Digit. Med.*, vol. 7, no. 1, May 2024, Art. no. 125, doi: 10.1038/s41746-024-01103-x.
- [17] E. L. Lehmann and G. Casella, *Theory of Point Estimation*, 2nd ed. New York, NY, USA: Springer-Verlag, 1998.
- [18] G. F. Cooper, "The computational complexity of probabilistic inference using bayesian belief networks," *Artif. Intell.*, vol. 42, no. 2–3, pp. 393–405, Mar. 1990, doi: 10.1016/0004-3702(90)90060-D.

- [19] R. M. Gray, *Probability, Random Processes, and Ergodic Properties*, 2nd ed. New York, NY, USA: Springer Science+Business Media, 2009.
- [20] T. Hastie, R. Tibshirani, and J. Friedman, *The Elements of Statistical Learning: Data Mining, Inference, and Prediction*, 2nd ed. New York, NY, USA: Springer Science+Business Media, 2009.
- [21] J. H. Friedman, “Greedy function approximation: A gradient boosting machine,” *Ann. Statist.*, vol. 29, no. 5, pp. 1189–1232, Oct. 2001, doi: 10.1214/aos/1013203451.
- [22] D. N. Gujarati and D. C. Porter, *Basic Econometrics*, 5th ed. New York, NY, USA: McGraw-Hill/Irwin, 2009.
- [23] Y. Dodge, Ed. *The Oxford Dictionary of Statistical Terms*. New York, NY, USA: Oxford Univ. Press, 2003.
- [24] B. S. Everitt, *The Cambridge Dictionary of Statistics*, 2nd ed. Cambridge, U.K.: Cambridge Univ. Press, 2002.
- [25] D. Sun, K.-C. Toh, and Y. Yuan, “Convex clustering: Model, theoretical guarantee and efficient algorithm,” *J. Mach. Learn. Res.*, vol. 22, no. 9, pp. 1–32, Jan. 2021. [Online]. Available: <http://jmlr.org/papers/v22/18-694.html>
- [26] K. Pelckmans, J. De Brabanter, J. A. K. Suykens, and B. De Moor, “Convex clustering shrinkage,” presented at the PASCAL Workshop Statist. Optim. Clustering Workshop, 2005.

- [27] S. Boyd and L. Vandenberghe, *Convex Optimization*. Cambridge, U.K.: Cambridge Univ. Press, 2004.
- [28] European Union, “Regulation (EU) 2016/679 of the European Parliament and of the Council of 27 April 2016 on the protection of natural persons with regard to the processing of personal data and on the free movement of such data, and repealing Directive 95/46/EC (General Data Protection Regulation) (Text with EEA relevance),” L 119/1, May 4, 2016. [Online]. Available: <https://eur-lex.europa.eu/eli/reg/2016/679/oj>
- [29] European Union, “Regulation (EU) 2018/1725 of the European Parliament and of the Council of 23 October 2018 on the protection of natural persons with regard to the processing of personal data by the Union institutions, bodies, offices and agencies and on the free movement of such data, and repealing Regulation (EC) No 45/2001 and Decision No 1247/2002/EC (Text with EEA relevance),” L 295/39, Nov. 21, 2018. [Online]. Available: <https://eur-lex.europa.eu/eli/reg/2018/1725/oj>
- [30] X. Liu, H. Li, G. Xu, Z. Chen, X. Huang, and R. Lu, “Privacy-enhanced federated learning against poisoning adversaries,” *IEEE Trans. Inf. Forensics Security*, vol. 16, pp. 4574–4588, 2021, doi: 10.1109/TIFS.2021.3108434.
- [31] J. Zhang, B. Chen, X. Cheng, H. T. T. Binh, and S. Yu, “PoisonGAN: Generative poisoning attacks against federated learning in edge computing systems,” *IEEE Internet Things J.*, vol. 8, no. 5, pp. 3310–3322, Mar. 2021, doi: 10.1109/JIOT.2020.3023126.

- [32] L. Bottou, “On-line learning and stochastic approximations,” in *On-Line Learning in Neural Networks*, D. Saad, Ed. New York, NY, USA: Cambridge Univ. Press, 1999, ch. 2, pp. 9–42.
- [33] G. Tel, *Introduction to Distributed Algorithms*, 2nd ed. Cambridge, U.K.: Cambridge Univ. Press, 2000.
- [34] D. P. Bertsekas and J. N. Tsitsiklis, *Parallel and Distributed Computation: Numerical Methods*. Belmont, MA, USA: Athena Scientific, 2015.
- [35] E. F. Codd, “A relational model of data for large shared data banks,” *Commun. ACM*, vol. 13, no. 6, pp. 377–387, Jun. 1970, doi: 10.1145/362384.362685.
- [36] G. Strang, *Computational Science and Engineering*. Wellesley, MA, USA: Wellesley-Cambridge Press, 2007.
- [37] R. A. Horn and C. R. Johnson, *Matrix Analysis*, 2nd ed. New York, NY, USA: Cambridge Univ. Press, 2013.
- [38] S. Ross, *A First Course in Probability*, 9th ed. Boston, MA, USA: Pearson Education, 2014.
- [39] P. R. Halmos, *Measure Theory*. New York, NY, USA: Springer-Verlag, 1974.
- [40] S. M. Kay, *Fundamentals of Statistical Signal Processing: Estimation Theory*. Upper Saddle River, NJ, USA: Prentice Hall PTR, 1993.
- [41] C. M. Bishop, *Pattern Recognition and Machine Learning*. New York, NY, USA: Springer Science+Business Media, 2006.



- [42] M. J. Wainwright and M. I. Jordan, “Graphical models, exponential families, and variational inference,” *Found. Trends Mach. Learn.*, vol. 1, no. 1–2, pp. 1–305, Nov. 2008, doi: 10.1561/22000000001.
- [43] N. Cesa-Bianchi and G. Lugosi, *Prediction, Learning, and Games*. New York, NY, USA: Cambridge Univ. Press, 2006.
- [44] E. Hazan, “Introduction to online convex optimization,” *Found. Trends Optim.*, vol. 2, no. 3–4, pp. 157–325, Aug. 2016, doi: 10.1561/24000000013.
- [45] L. Zhang, G. Karakasidis, A. Odnoblyudova, L. Dogruel, Y. Tian, and A. Jung, “Explainable empirical risk minimization,” *Neural Comput. Appl.*, vol. 36, no. 8, pp. 3983–3996, Mar. 2024, doi: 10.1007/s00521-023-09269-3.
- [46] C. Rudin, “Stop explaining black box machine learning models for high-stakes decisions and use interpretable models instead,” *Nature Mach. Intell.*, vol. 1, no. 5, pp. 206–215, May 2019, doi: 10.1038/s42256-019-0048-x.
- [47] C. Molnar, *Interpretable Machine Learning: A Guide for Making Black Box Models Explainable*, 3rd ed., 2025. [Online]. Available: <https://christophm.github.io/interpretable-ml-book/>
- [48] R. R. Selvaraju, M. Cogswell, A. Das, R. Vedantam, D. Parikh, and D. Batra, “Grad-CAM: Visual explanations from deep networks via gradient-based localization,” in *2017 IEEE Int. Conf. Comput. Vis.*, pp. 618–626, doi: 10.1109/ICCV.2017.74.

- [49] J. Colin, T. Fel, R. Cadène, and T. Serre, “What I cannot predict, I do not understand: A human-centered evaluation framework for explainability methods,” in *Adv. Neural Inf. Process. Syst.*, S. Koyejo, S. Mohamed, A. Agarwal, D. Belgrave, K. Cho, and A. Oh, Eds. vol. 35, 2022, pp. 2832–2845. [Online]. Available: [https://proceedings.neurips.cc/paper\\_files/paper/2022/hash/13113e938f2957891c0c5e8df811dd01-Abstract-Conference.html](https://proceedings.neurips.cc/paper_files/paper/2022/hash/13113e938f2957891c0c5e8df811dd01-Abstract-Conference.html)
- [50] A. Jung and P. H. J. Nardelli, “An information-theoretic approach to personalized explainable machine learning,” *IEEE Signal Process. Lett.*, vol. 27, pp. 825–829, 2020, doi: 10.1109/LSP.2020.2993176.
- [51] J. Chen, L. Song, M. J. Wainwright, and M. I. Jordan, “Learning to explain: An information-theoretic perspective on model interpretation,” in *Proc. 35th Int. Conf. Mach. Learn.*, J. Dy and A. Krause, Eds. vol. 80, 2018, pp. 883–892. [Online]. Available: <https://proceedings.mlr.press/v80/chen18j.html>
- [52] M. T. Ribeiro, S. Singh, and C. Guestrin, “Why should i trust you?: Explaining the predictions of any classifier,” in *Proc. 22nd ACM SIGKDD Int. Conf. Knowl. Discovery Data Mining*, Aug. 2016, pp. 1135–1144, doi: 10.1145/2939672.2939778.
- [53] H. B. McMahan, E. Moore, D. Ramage, S. Hampson, and B. A. y Arcas, “Communication-efficient learning of deep networks from decentralized data,” in *Proc. 20th Int. Conf. Artif. Intell. Statist.*, A. Singh and J. Zhu, Eds. vol. 54, 2017, pp. 1273–1282. [Online]. Available: <https://proceedings.mlr.press/v54/mcmahan17a.html>

- [54] T. Li, A. K. Sahu, M. Zaheer, M. Sanjabi, A. Talwalkar, and V. Smith, “Federated optimization in heterogeneous networks,” in *Proc. Mach. Learn. Syst.*, I. Dhillon, D. Papailiopoulos, and V. Sze, Eds. vol. 2, 2020. [Online]. Available: [https://proceedings.mlsys.org/paper\\_files/paper/2020/hash/1f5fe83998a09396ebe6477d9475ba0c-Abstract.html](https://proceedings.mlsys.org/paper_files/paper/2020/hash/1f5fe83998a09396ebe6477d9475ba0c-Abstract.html)
- [55] Y. SarcheshmehPour, Y. Tian, L. Zhang, and A. Jung, “Flow-based clustering and spectral clustering: A comparison,” in *2021 55th Asilomar Conf. Signals, Syst., Comput.*, M. B. Matthews, Ed. pp. 1292–1296, doi: 10.1109/IEEECONF53345.2021.9723162.
- [56] S. Shalev-Shwartz and S. Ben-David, *Understanding Machine Learning: From Theory to Algorithms*. New York, NY, USA: Cambridge Univ. Press, 2014.
- [57] J. Su, D. V. Vargas, and K. Sakurai, “One pixel attack for fooling deep neural networks,” *IEEE Trans. Evol. Comput.*, vol. 23, no. 5, pp. 828–841, Oct. 2019, doi: 10.1109/TEVC.2019.2890858.
- [58] Y. SarcheshmehPour, Y. Tian, L. Zhang, and A. Jung, “Clustered federated learning via generalized total variation minimization,” *IEEE Trans. Signal Process.*, vol. 71, pp. 4240–4256, 2023, doi: 10.1109/TSP.2023.3322848.
- [59] H. P. Lopuhaä and P. J. Rousseeuw, “Breakdown points of affine equivariant estimators of multivariate location and covariance matrices,” *Ann. Statist.*, vol. 19, no. 1, pp. 229–248, Mar. 1991, doi: 10.1214/aos/1176347978.

- [60] A. Vaswani et al., “Attention is all you need,” in *Adv. Neural Inf. Process. Syst.*, I. Guyon, U. von Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, Eds. vol. 30, 2017, pp. 5998–6008. [Online]. Available: [https://papers.nips.cc/paper\\_files/paper/2017/hash/3f5ee243547dee91fbd053c1c4a845aa-Abstract.html](https://papers.nips.cc/paper_files/paper/2017/hash/3f5ee243547dee91fbd053c1c4a845aa-Abstract.html)
- [61] U. von Luxburg, “A tutorial on spectral clustering,” *Statist. Comput.*, vol. 17, no. 4, pp. 395–416, Dec. 2007, doi: 10.1007/s11222-007-9033-z.
- [62] R. T. Rockafellar, *Network Flows and Monotropic Optimization*. Belmont, MA, USA: Athena Scientific, 1998.
- [63] M. J. Wainwright, *High-Dimensional Statistics: A Non-Asymptotic Viewpoint*. Cambridge, U.K.: Cambridge Univ. Press, 2019.
- [64] P. Bühlmann and S. van de Geer, *Statistics for High-Dimensional Data: Methods, Theory and Applications*. Berlin, Germany: Springer-Verlag, 2011.
- [65] N. Young, *An Introduction to Hilbert Space*. New York, NY, USA: Cambridge Univ. Press, 1988.
- [66] C. H. Lampert, “Kernel methods in computer vision,” *Found. Trends Comput. Graph. Vis.*, vol. 4, no. 3, pp. 193–285, Sep. 2009, doi: 10.1561/06000000027.
- [67] Q. Yang, Y. Liu, Y. Cheng, Y. Kang, T. Chen, and H. Yu, “Horizontal federated learning,” in *Federated Learning*. Cham, Switzerland: Springer Nature, 2020, ch. 4, pp. 49–67.

- [68] D. Pfau and A. Jung, “Engineering trustworthy AI: A developer guide for empirical risk minimization,” Nov. 2024. [Online]. Available: <https://arxiv.org/abs/2410.19361>
- [69] High-Level Expert Group on Artificial Intelligence, “The assessment list for trustworthy artificial intelligence (ALTAI): For self assessment,” European Commission, Jul. 17, 2020. [Online]. Available: <https://digital-strategy.ec.europa.eu/en/library/assessment-list-trustworthy-artificial-intelligence-altai-self-assessment>
- [70] F. Doshi-Velez and B. Kim, “Towards a rigorous science of interpretable machine learning,” Mar. 2017. [Online]. Available: <https://arxiv.org/abs/1702.08608>
- [71] P. Hase and M. Bansal, “Evaluating explainable AI: Which algorithmic explanations help users predict model behavior?,” in *Proc. 58th Annu. Meeting Assoc. Comput. Linguistics*, D. Jurafsky, J. Chai, N. Schluter, and J. Tetreault, Eds. Jul. 2020, pp. 5540–5552. [Online]. Available: <https://aclanthology.org/2020.acl-main.491>
- [72] Z. C. Lipton, “The mythos of model interpretability: In machine learning, the concept of interpretability is both important and slippery,” *Queue*, vol. 16, no. 3, pp. 31–57, Jun. 2018, doi: 10.1145/3236386.3241340.
- [73] A. Silberschatz, H. F. Korth, and S. Sudarshan, *Database System Concepts*, 7th ed. New York, NY, USA: McGraw-Hill Education, 2019. [Online]. Available: <https://db-book.com/>

- [74] S. Abiteboul, R. Hull, and V. Vianu, *Foundations of Databases*. Reading, MA, USA: Addison-Wesley, 1995.
- [75] S. Hoberman, *Data Modeling Made Simple: A Practical Guide for Business and IT Professionals*, 2nd ed. Basking Ridge, NJ, USA: Technics Publications, 2009.
- [76] R. Ramakrishnan and J. Gehrke, *Database Management Systems*, 3rd ed. New York, NY, USA: McGraw-Hill, 2002.
- [77] T. Gebru et al., “Datasheets for datasets,” *Commun. ACM*, vol. 64, no. 12, pp. 86–92, Nov. 2021, doi: 10.1145/3458723.
- [78] T. M. Cover and J. A. Thomas, *Elements of Information Theory*, 2nd ed. Hoboken, NJ, USA: Wiley, 2006.
- [79] A. Papoulis and S. U. Pillai, *Probability, Random Variables, and Stochastic Processes*, 4th ed. New York, NY, USA: McGraw-Hill Higher Education, 2002.
- [80] A. Lapidoth, *A Foundation in Digital Communication*. Cambridge, U.K.: Cambridge Univ. Press, 2009.
- [81] R. A. Horn and C. R. Johnson, *Topics in Matrix Analysis*. Cambridge, UK: Cambridge Univ. Press, 1991.
- [82] A. Y. Ng, M. I. Jordan, and Y. Weiss, “On spectral clustering: Analysis and an algorithm,” in *Adv. Neural Inf. Process. Syst.*, T. Dietterich, S. Becker, and Z. Ghahramani, Eds. vol. 14, 2001, pp.

- 849–856. [Online]. Available: [https://papers.nips.cc/paper\\_files/paper/2001/hash/801272ee79cfde7fa5960571fee36b9b-Abstract.html](https://papers.nips.cc/paper_files/paper/2001/hash/801272ee79cfde7fa5960571fee36b9b-Abstract.html)
- [83] K. Abayomi, A. Gelman, and M. Levy, “Diagnostics for multivariate imputations,” *J. Roy. Statist. Soc.: Ser. C (Appl. Statist.)*, vol. 57, no. 3, pp. 273–291, Jun. 2008, doi: 10.1111/j.1467-9876.2007.00613.x.
  - [84] S. Bubeck and N. Cesa-Bianchi, “Regret analysis of stochastic and non-stochastic multi-armed bandit problems,” *Found. Trends Mach. Learn.*, vol. 5, no. 1, pp. 1–122, Dec. 2012, doi: 10.1561/22000000024.
  - [85] B. Schölkopf and A. J. Smola, *Learning with Kernels: Support Vector Machines, Regularization, Optimization, and Beyond*. Cambridge, MA, USA: MIT Press, 2002.
  - [86] A. Jung, “Networked exponential families for big data over networks,” *IEEE Access*, vol. 8, pp. 202 897–202 909, Nov. 2020, doi: 10.1109/ACCESS.2020.3033817.
  - [87] Y. Nesterov, *Introductory Lectures on Convex Optimization: A Basic Course*. Boston, MA, USA: Kluwer Academic, 2004.
  - [88] A. Rakhlin, O. Shamir, and K. Sridharan, “Making gradient descent optimal for strongly convex stochastic optimization,” in *Proc. 29th Int. Conf. Mach. Learn.*, J. Langford and J. Pineau, Eds. 2012, pp. 449–456. [Online]. Available: <https://icml.cc/Conferences/2012/papers/261.pdf>
  - [89] P. J. Brockwell and R. A. Davis, *Time Series: Theory and Methods*, 2nd ed. New York, NY, USA: Springer-Verlag, 1991.

- [90] M. Kearns and M. Li, “Learning in the presence of malicious errors,” *SIAM J. Comput.*, vol. 22, no. 4, pp. 807–837, Aug. 1993, doi: 10.1137/0222052.
- [91] G. Lugosi and S. Mendelson, “Robust multivariate mean estimation: The optimality of trimmed mean,” *Ann. Statist.*, vol. 49, no. 1, pp. 393–410, Feb. 2021, doi: 10.1214/20-AOS1961.
- [92] N. Parikh and S. Boyd, “Proximal algorithms,” *Found. Trends Optim.*, vol. 1, no. 3, pp. 127–239, Jan. 2014, doi: 10.1561/24000000003.
- [93] A. Makhdoumi, S. Salamatian, N. Fawaz, and M. Médard, “From the information bottleneck to the privacy funnel,” in *2014 IEEE Inf. Theory Workshop*, pp. 501–505, doi: 10.1109/ITW.2014.6970882.
- [94] A. Ünsal and M. Önen, “Information-theoretic approaches to differential privacy,” *ACM Comput. Surv.*, vol. 56, no. 3, Oct. 2023, Art. no. 76, doi: 10.1145/3604904.
- [95] O. Kallenberg, *Foundations of Modern Probability*. New York, NY, USA: Springer-Verlag, 1997.
- [96] L. Condat, “A primal–dual splitting method for convex optimization involving lipschitzian, proximable and linear composite terms,” *J. Optim. Theory Appl.*, vol. 158, no. 2, pp. 460–479, Aug. 2013, doi: 10.1007/s10957-012-0245-9.
- [97] H. H. Bauschke and P. L. Combettes, *Convex Analysis and Monotone Operator Theory in Hilbert Spaces*, 2nd ed. New York, NY, USA: Springer Science+Business Media, 2017.



- [98] A. Ben-Israel and T. N. E. Greville, *Generalized Inverses: Theory and Applications*, 2nd ed. New York, NY, USA: Springer-Verlag, 2003.
- [99] S. Shalev-Shwartz and A. Tewari, “Stochastic methods for  $\ell_1$  regularized loss minimization,” in *Proc. 26th Annu. Int. Conf. Mach. Learn.*, L. Bottou and M. Littman, Eds. Jun. 2009, pp. 929–936.
- [100] I. Csiszar, “Generalized cutoff rates and Renyi’s information measures,” *IEEE Trans. Inf. Theory*, vol. 41, no. 1, pp. 26–34, Jan. 1995, doi: 10.1109/18.370121.
- [101] S. Bubeck, “Convex optimization: Algorithms and complexity,” *Found. Trends Mach. Learn.*, vol. 8, no. 3–4, pp. 231–357, Nov. 2015, doi: 10.1561/22000000050.
- [102] D. P. Bertsekas, *Convex Optimization Algorithms*. Belmont, MA, USA: Athena Scientific, 2015.
- [103] I. Goodfellow, Y. Bengio, and A. Courville, *Deep Learning*. Cambridge, MA, USA: MIT Press, 2016.
- [104] L. Cohen, *Time-Frequency Analysis*. Upper Saddle River, NJ, USA: Prentice Hall PTR, 1995.
- [105] J. Li, L. Han, X. Li, J. Zhu, B. Yuan, and Z. Gou, “An evaluation of deep neural network models for music classification using spectrograms,” *Multimedia Tools Appl.*, vol. 81, no. 4, pp. 4621–4647, Feb. 2022, doi: 10.1007/s11042-020-10465-9.

- [106] B. Boashash, Ed. *Time Frequency Signal Analysis and Processing: A Comprehensive Reference*. Oxford, U.K.: Elsevier, 2003.
- [107] S. Mallat, *A Wavelet Tour of Signal Processing: The Sparse Way*, 3rd ed. Burlington, MA, USA: Academic, 2009.
- [108] E. Abbe, “Community detection and stochastic block models: Recent developments,” *J. Mach. Learn. Res.*, vol. 18, no. 177, pp. 1–86, Apr. 2018. [Online]. Available: <http://jmlr.org/papers/v18/16-480.html>
- [109] D. P. Bertsekas, A. Nedic, and A. E. Ozdaglar, *Convex Analysis and Optimization*. Belmont, MA, USA: Athena Scientific, 2003.
- [110] D. P. Bertsekas, *Nonlinear Programming*, 2nd ed. Belmont, MA, USA: Athena Scientific, 1999.
- [111] N. Cristianini and J. Shawe-Taylor, *An Introduction to Support Vector Machines and Other Kernel-based Learning Methods*. New York, NY, USA: Cambridge Univ. Press, 2000.
- [112] High-Level Expert Group on Artificial Intelligence, “Ethics guidelines for trustworthy AI,” European Commission, Apr. 8, 2019. [Online]. Available: <https://digital-strategy.ec.europa.eu/en/library/ethics-guidelines-trustworthy-ai>
- [113] C. Gallese, “The AI act proposal: A new right to technical interpretability?,” *SSRN Electron. J.*, Feb. 2023. [Online]. Available: <https://ssrn.com/abstract=4398206>

- [114] M. Mitchell et al., “Model cards for model reporting,” in *Proc. Conf. Fairness, Accountability, Transparency*, 2019, pp. 220–229, doi: 10.1145/3287560.3287596.
- [115] K. Shahriari and M. Shahriari, “IEEE standard review — Ethically aligned design: A vision for prioritizing human wellbeing with artificial intelligence and autonomous systems,” in *2017 IEEE Canada Int. Humanitarian Technol. Conf.*, pp. 197–201, doi: 10.1109/IHTC.2017.8058187.
- [116] R. Caruana, “Multitask learning,” *Mach. Learn.*, vol. 28, pp. 41–75, Jul. 1997, doi: 10.1023/A:1007379606734.
- [117] A. Jung, G. Hannak, and N. Goertz, “Graphical lasso based model selection for time series,” *IEEE Signal Process. Lett.*, vol. 22, no. 10, pp. 1781–1785, Oct. 2015, doi: 10.1109/LSP.2015.2425434.
- [118] A. Jung, “Learning the conditional independence structure of stationary time series: A multitask learning approach,” *IEEE Trans. Signal Process.*, vol. 63, no. 21, Nov. 2015, doi: 10.1109/TSP.2015.2460219.
- [119] C. E. Rasmussen and C. K. I. Williams, *Gaussian Processes for Machine Learning*. Cambridge, MA, USA: MIT Press, 2006.
- [120] Q. Yang, Y. Liu, Y. Cheng, Y. Kang, T. Chen, and H. Yu, “Vertical federated learning,” in *Federated Learning*. Cham, Switzerland: Springer Nature, 2020, ch. 5, pp. 69–81.