

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

Alexander Jung and Konstantina Olioumtsevs

June 3, 2025



please cite as: A. Jung and K. Olioumtsevs, *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>.

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 $0 \leq x \leq 1$.

$\operatorname{argmin}_{\mathbf{w}} f(\mathbf{w})$	L'ensemble des point qui minimisent la fonction à valeurs réelles $f(\mathbf{w})$.
$\mathbb{S}^{(n)}$	L'ensemble des vecteurs de norme unitaire dans \mathbb{R}^{n+1} .
$\log a$	Le logarithme d'un réel strictement positif $a \in \mathbb{R}_{++}$.
$h(\cdot) : \mathcal{A} \rightarrow \mathcal{B} : a \mapsto h(a)$	<p>Une fonction (ou application) qui accepte tout élément $a \in \mathcal{A}$ d'un ensemble \mathcal{A} en entrée et fournit un élément bien défini $h(a) \in \mathcal{B}$ d'un ensemble \mathcal{B}. L'ensemble \mathcal{A} est le domaine de définition de la fonction h et l'ensemble \mathcal{B} est l'ensemble d'arrivée de h. L'apprentissage automatique vise à trouver (ou apprendre en la construisant) une fonction h (c'est-à-dire une hypothèse) qui prend en entrée les caractéristiques \mathbf{x} d'un point de données et renvoie une prédiction $h(\mathbf{x})$ pour son étiquette y.</p>
$\operatorname{epi}(f)$	L'épigraphe d'une fonction à valeurs réelles $f : \mathbb{R}^d \rightarrow \mathbb{R}$.
$\frac{\partial f(w_1, \dots, w_d)}{\partial w_j}$	La dérivée partielle (si elle existe) d'une fonction à valeurs réelles $f : \mathbb{R}^d \rightarrow \mathbb{R}$ par rapport à w_j [2, Ch. 9].
$\nabla f(\mathbf{w})$	<p>Le gradient d'une fonction dérivable à valeurs réelles $f : \mathbb{R}^d \rightarrow \mathbb{R}$ est le vecteur $\nabla f(\mathbf{w}) = \left(\frac{\partial f}{\partial w_1}, \dots, \frac{\partial f}{\partial w_d} \right)^T \in \mathbb{R}^d$ [2, Ch. 9].</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 ℓ_2) du vecteur $\mathbf{x} = (x_1, \dots, x_d)^T \in \mathbb{R}^d$ définie par $\ \mathbf{x}\ _2 := \sqrt{\sum_{j=1}^d x_j^2}$.
$\ \mathbf{x}\ $	Une certaine norme du vecteur $\mathbf{x} \in \mathbb{R}^d$ [3]. Sauf indication contraire, on entend par là la norme euclidienne $\ \mathbf{x}\ _2$.
\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 $\mathbf{X} = \mathbf{X}^T$.
$\mathbf{0} = (0, \dots, 0)^T$	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} .
$\mathbf{A} \otimes \mathbf{B}$	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].
$\mathbb{E}_p\{f(\mathbf{z})\}$	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 $p(\mathbf{z})$. Si la loi de probabilité est claire dans le contexte, on écrit simplement $\mathbb{E}\{f(\mathbf{z})\}$.
$p(\mathbf{x}, y)$	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 .
$p(\mathbf{x} y)$	Une loi de probabilité conditionnelle d'une VA \mathbf{x} étant donnée la valeur d'une autre VA y [7, Sec. 3.5].
$p(\mathbf{x}; \mathbf{w})$	Une loi de probabilité paramétrée d'une VA \mathbf{x} . La loi de probabilité dépend d'un vecteur de paramètres \mathbf{w} . Par exemple, $p(\mathbf{x}; \mathbf{w})$ 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 $\mathbb{E}\left\{(\mathbf{x} - \mathbb{E}\{\mathbf{x}\})(\mathbf{x} - \mathbb{E}\{\mathbf{x}\})^T\right\}$.

$\mathcal{N}(\mu, \sigma^2)$ 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 $\sigma^2 = \mathbb{E}\{(x - \mu)^2\}$.

$\mathcal{N}(\boldsymbol{\mu}, \mathbf{C})$ 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 $\mathbf{C} = \mathbb{E}\{(\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^T\}$.

Apprentissage automatique

r	Un indice $r = 1, 2, \dots$ qui énumère les points de données.
m	Le nombre de points de données dans un jeu de données (c'est-à-dire la taille du jeu de données).
\mathcal{D}	Un jeu de données $\mathcal{D} = \{\mathbf{z}^{(1)}, \dots, \mathbf{z}^{(m)}\}$ est une liste de points de données individuels $\mathbf{z}^{(r)}$, pour $r = 1, \dots, m$.
d	Le nombre de caractéristiques qui constituent un point de données.
x_j	La j -ième caractéristique d'un point de données. La première caractéristique est noté x_1 , le deuxième x_2 , et ainsi de suite.
\mathbf{x}	Le vecteur de caractéristiques $\mathbf{x} = (x_1, \dots, x_d)^T$ d'un point de données, dont les composantes sont les différentes caractéristiques du point de données.
\mathcal{X}	L'espace des caractéristiques \mathcal{X} est l'ensemble de toutes les valeurs possibles que les caractéristiques \mathbf{x} d'un point de données peuvent prendre.
\mathbf{z}	Au lieu du symbole \mathbf{x} , on utilise parfois \mathbf{z} comme un autre symbole pour désigner un vecteur dont les composantes sont les différentes caractéristiques d'un point de données. On a besoin de deux symboles différents pour distinguer les caractéristiques brutes des caractéristiques apprises [8, Ch. 9].

$\mathbf{x}^{(r)}$	Le vecteur de caractéristiques du r -ième point de données dans un jeu de données.
$x_j^{(r)}$	La j -ième caractéristique du r -ième point de données dans un jeu de données.
y	L'étiquette (ou quantité d'intérêt) d'un point de données.
$y^{(r)}$	L'étiquette du r -ième point de données.
$(\mathbf{x}^{(r)}, y^{(r)})$	Les caractéristiques et l'étiquette du r -ième 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.

\mathcal{B}	Un mini-lot (ou sous-ensemble) de points de données choisis aléatoirement.
B	La taille (c'est-à-dire le nombre de points de données) d'un mini-lot.
$h(\cdot)$	Une fonction hypothèse qui lit les caractéristiques \mathbf{x} d'un point de données et produit une prédiction $\hat{y} = h(\mathbf{x})$ pour son étiquette y .
$\mathcal{Y}^{\mathcal{X}}$	Étant donnés deux ensembles \mathcal{X} et \mathcal{Y} , on note $\mathcal{Y}^{\mathcal{X}}$ l'ensemble de toutes les fonctions hypothèses possibles $h : \mathcal{X} \rightarrow \mathcal{Y}$.
\mathcal{H}	Un espace des hypothèses ou modèle utilisé par une méthode d'apprentissage automatique. L'espace des hypothèses est constitué des différentes hypothèses $h : \mathcal{X} \rightarrow \mathcal{Y}$, parmi lesquelles la méthode d'apprentissage automatique doit choisir.
$d_{\text{eff}}(\mathcal{H})$	La dimension effective d'un espace des hypothèses \mathcal{H} .
B^2	Le biais au carré d'une fonction hypothèse \hat{h} 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 \hat{h} est également une réalisation d'une VA.

V	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 \hat{h} est également une réalisation d'une VA.
$L((\mathbf{x}, y), h)$	La perte encourue en prédisant l'étiquette y d'un point de données à l'aide de la prédiction $\hat{y} = h(\mathbf{x})$. La prédiction \hat{y} est obtenue en évaluant la fonction hypothèse $h \in \mathcal{H}$ en \mathbf{x} , le vecteur de caractéristiques du point de données.
E_v	L'erreur de validation d'une hypothèse h , c'est-à-dire sa perte moyenne sur un ensemble de validation (ou jeu de validation).
$\hat{L}(h \mathcal{D})$	Le risque empirique, ou perte moyenne, encouru par l'hypothèse h sur un jeu de données \mathcal{D} .
E_t	L'erreur d'entraînement d'une hypothèse h , c'est-à-dire sa perte moyenne sur un ensemble d'entraînement.
t	Un indice de temps discret $t = 0, 1, \dots$ utilisé pour énumérer des événements séquentiels (ou des instants temporels).
t	Un indice qui énumère les tâches d'apprentissage dans un problème d'apprentissage multitâche.

η	Le taux d'apprentissage (ou taille de pas) utilisé par les méthodes basées sur le gradient.
α	Un paramètre de régularisation qui contrôle la quantité de régularisation.
$\lambda_j(\mathbf{Q})$	La j -ième valeur propre (triée par ordre croissant ou décroissant) d'une matrice semi-définie positive \mathbf{Q} . Si la matrice est claire dans le contexte, on écrit simplement λ_j .
$\sigma(\cdot)$	La fonction d'activation utilisée par un neurone artificiel dans un réseau de neurones artificiels (RNA).
$\mathcal{R}_{\hat{y}}$	Une région de décision dans un espace des caractéristiques.
\mathbf{w}	Un vecteur de paramètres $\mathbf{w} = (w_1, \dots, w_d)^T$ d'un modèle, par exemple les poids d'un modèle linéaire ou dans un RNA.
$h^{(\mathbf{w})}(\cdot)$	Une fonction hypothèse qui dépend de paramètres du modèle w_1, \dots, w_d regroupés dans le vecteur $\mathbf{w} = (w_1, \dots, w_d)^T$ et qui peuvent être ajustés.
$\phi(\cdot)$	Une feature map $\phi : \mathcal{X} \rightarrow \mathcal{X}' : \mathbf{x} \mapsto \mathbf{x}' := \phi(\mathbf{x}) \in \mathcal{X}'$.
$K(\cdot, \cdot)$	Étant donné un espace des caractéristiques \mathcal{X} , un kernel est une application $K : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{C}$ qui est sdv.

Apprentissage fédéré

$\mathcal{G} = (\mathcal{V}, \mathcal{E})$	Un graphe non orienté dont les nœuds $i \in \mathcal{V}$ représentent des dispositifs au sein d'un federated learning network (FL network). Les arêtes pondérées non orientées \mathcal{E} représentent la connectivité entre les dispositifs et les similarités statistiques entre leurs jeux de données et tâche d'apprentissages.
$i \in \mathcal{V}$	Un nœud représentant un dispositif dans un FL network. Le dispositif peut accéder à un jeu de données local et entraîner un local model.
$\mathcal{G}^{(\mathcal{C})}$	Le sous-graphe induit de \mathcal{G} utilisant les nœuds de $\mathcal{C} \subseteq \mathcal{V}$.
$\mathbf{L}^{(\mathcal{G})}$	La Laplacian matrix d'un graphe \mathcal{G} .
$\mathbf{L}^{(\mathcal{C})}$	La Laplacian matrix du graphe induit $\mathcal{G}^{(\mathcal{C})}$.
$\mathcal{N}^{(i)}$	Le voisinage d'un nœud i dans un graphe \mathcal{G} .
$d^{(i)}$	Le degré pondéré $d^{(i)} := \sum_{i' \in \mathcal{N}^{(i)}} A_{i,i'}$ d'un nœud i dans un graphe \mathcal{G} .
$d_{\max}^{(\mathcal{G})}$	Le degré pondéré maximal (parmi les degrés pondérés des nœuds) d'un graphe \mathcal{G} .
$\mathcal{D}^{(i)}$	Le jeu de données local $\mathcal{D}^{(i)}$ détenu par le nœud $i \in \mathcal{V}$ d'un FL network.

m_i	Le nombre de points de données (i.e., la sample size) contenus dans le jeu de données local $\mathcal{D}^{(i)}$ au nœud $i \in \mathcal{V}$.
$\mathbf{x}^{(i,r)}$	Les caractéristiques du r -ième point de données dans le jeu de données local $\mathcal{D}^{(i)}$.
$y^{(i,r)}$	L'étiquette du r -ième point de données dans le jeu de données local $\mathcal{D}^{(i)}$.
$\mathbf{w}^{(i)}$	Les paramètres du modèle locaux de l'appareil i au sein d'un FL network.
$L_i(\mathbf{w})$	La fonction de perte (ou de coût) locale utilisée par le dispositif i pour évaluer l'utilité d'un certain choix \mathbf{w} pour les paramètres du modèle locaux.
$L^{(d)}(\mathbf{x}, h(\mathbf{x}), h'(\mathbf{x}))$	La perte encourue par une hypothèse h' sur un point de données de caractéristiques \mathbf{x} et d'étiquette $h(\mathbf{x})$ obtenue à partir d'une autre hypothèse.
$\text{stack}\{\mathbf{w}^{(i)}\}_{i=1}^n$	Le vecteur $\left((\mathbf{w}^{(1)})^T, \dots, (\mathbf{w}^{(n)})^T \right)^T \in \mathbb{R}^{dn}$ obtenu en empilant verticalement les paramètres du modèle locaux $\mathbf{w}^{(i)} \in \mathbb{R}^d$.

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 ees 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 λ_2 of its Laplacian matrix. 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 gradient steps. 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 feature map such that the new caractéristiques allow us to reconstruct the original caractéristiques with the minimum reconstruction error [8].

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 trade 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 caractéristiques.

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 fédéré FL is an umbrella term for apprentissage automatique methods that train modèles in a collaborative fashion using decentralized données and computation.

apprentissage multitâche Multitask learning aims at leveraging relations between different tâche d'apprentissages. Consider two tâche d'apprentissages obtained from the same jeu de données of webcam snapshots. The first task is to predict the presence of a human, while the second task is to predict the presence of a car. It might be useful to use the same deep net structure for both tasks and only allow the poids of the final output layer to be different.

apprentissage semi-supervisé SSL methods use unlabeled point de données to support the learning of a hypothèse from labeled datapoints [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 labeled datapoints.

arbre de décision A decision tree is a flow-chart-like representation of a hypothèse map h . More formally, a decision tree is a directed graphe containing a root node that reads in the vecteur de caractéristiques \mathbf{x} of a point de données. The root node then forwards the point de données to one of its children nodes based on some elementary test on the caractéristiques \mathbf{x} . If the receiving child node is not a leaf node, i.e., it has itself children nodes, it represents another test. Based on the test result, the point de données is forwarded to one of its descendants. This testing and forwarding of the point de données is continued until the point de données ends up in a leaf node (having no children nodes).

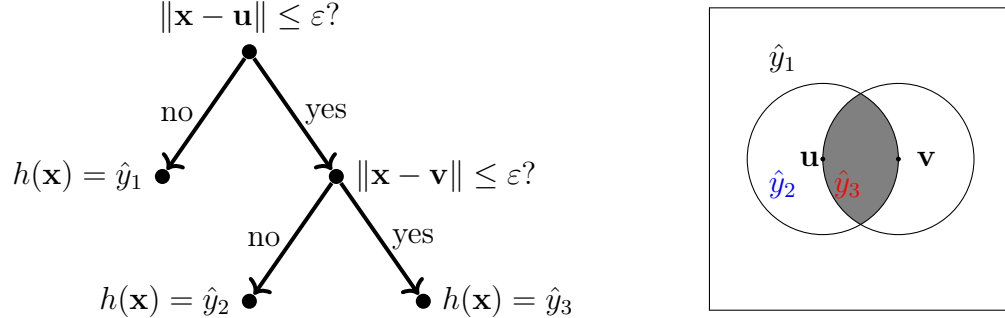


Figure 1: Left: A decision tree is a flow-chart-like representation of a piecewise constant hypoth se $h : \mathcal{X} \rightarrow \mathbb{R}$. Each piece is a r gion de d cision $\mathcal{R}_{\hat{y}} := \{\mathbf{x} \in \mathcal{X} : h(\mathbf{x}) = \hat{y}\}$. The depicted decision tree can be applied to numeric vecteur de caract ristiques, i.e., $\mathcal{X} \subseteq \mathbb{R}^d$. It is parametrized by the threshold $\varepsilon > 0$ and the vectors $\mathbf{u}, \mathbf{v} \in \mathbb{R}^d$. Right: A decision tree partitions the espace des caract ristiques \mathcal{X} into r gion de d cisions. Each r gion de d cision $\mathcal{R}_{\hat{y}} \subseteq \mathcal{X}$ corresponds to a specific leaf node in the decision tree.

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'entra nement (donn es poisoning) or the optimization algorithm used by an MRE-based method. The goal of a backdoor attack is to nudge the learned hypoth se \hat{h} towards specific pr dictions 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.

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 probabilistic model. In many cases, given a probabilistic model $p(\mathbf{x}, y)$, 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 probabilistic model 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 squared error loss, 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 2).

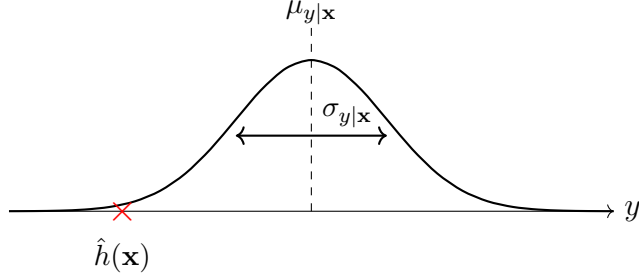


Figure 2: If the caractéristiques and the étiquette of a point de données are drawn from a loi normale multivariée, we can achieve the minimum risque (under squared error loss) by using the Bayes estimator $\mu_{y|\mathbf{x}}$ to predict the étiquette y of a point de données with caractéristiques \mathbf{x} . 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 probabilistic model 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 probabilistic model 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.),

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

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 hypothèse map (or strong learner) by sequentially combining less accurate hypothèse 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 lisse 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) hypothèse maps to reduce the risque empirique. A widely-used instance

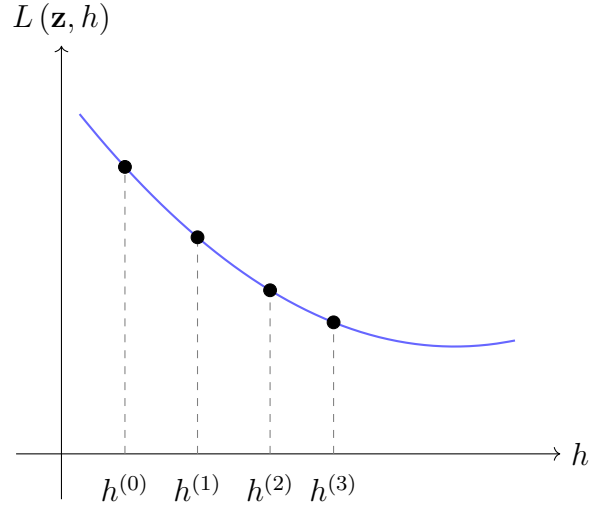


Figure 3: Boosting methods construct a sequence of hypothesis maps $h^{(0)}, h^{(1)}, \dots$ that are increasingly strong learners (i.e., incurring a smaller perte).

of the generic boosting idea is referred to as gradient boosting, which uses gradients of the fonction de perte for combining the weak learners [21].

See also: gradient step, méthodes basées sur le gradient

bootstrap For the analysis of apprentissage automatique methods, it is often useful to interpret a given set of point de données $\mathcal{D} = \{\mathbf{z}^{(1)}, \dots, \mathbf{z}^{(m)}\}$ as réalisations of i.i.d. VAs with a common loi de probabilité $p(\mathbf{z})$. 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 probabilité $p(\mathbf{z})$.

borne supérieure The supremum of a set of real numbers is the smallest number that is greater than or equal to every element in the set. More formally, a real number a is the supremum of a set $\mathcal{A} \subseteq \mathbb{R}$ if: 1) a is an upper bound of \mathcal{A} ; and 2) no number smaller than a is an upper bound of \mathcal{A} . Every non-empty set of real numbers that is bounded above has a supremum, even if it does not contain its supremum as an element [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 « régresseur » [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.

classifier A classifier is a hypothèse (map) $h(\mathbf{x})$ used to predict a étiquette taking values from a finite espace des étiquettes. We might use the function value $h(\mathbf{x})$ itself as a prédiction \hat{y} for the étiquette. However, it is customary to use a map $h(\cdot)$ that delivers a numeric quantity. The prédiction is then obtained by a simple thresholding step. For example, in a binary classification problem with $\mathcal{Y} \in \{-1, 1\}$, we might use a real-valued hypothèse map $h(\mathbf{x}) \in \mathbb{R}$ as a classifier. A prédiction \hat{y} can then be obtained via thresholding,

$$\hat{y} = 1 \text{ for } h(\mathbf{x}) \geq 0 \text{ and } \hat{y} = -1 \text{ otherwise.} \quad (1)$$

We can characterize a classifier by its région de décisions \mathcal{R}_a , for every possible étiquette value $a \in \mathcal{Y}$.

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 caractéristiquess. An example of such clusters is shown in Figure 4.

clustered federated learning (CFL) CFL trains local models for the dispositifs in an apprentissage fédéré application by using a partitionnement

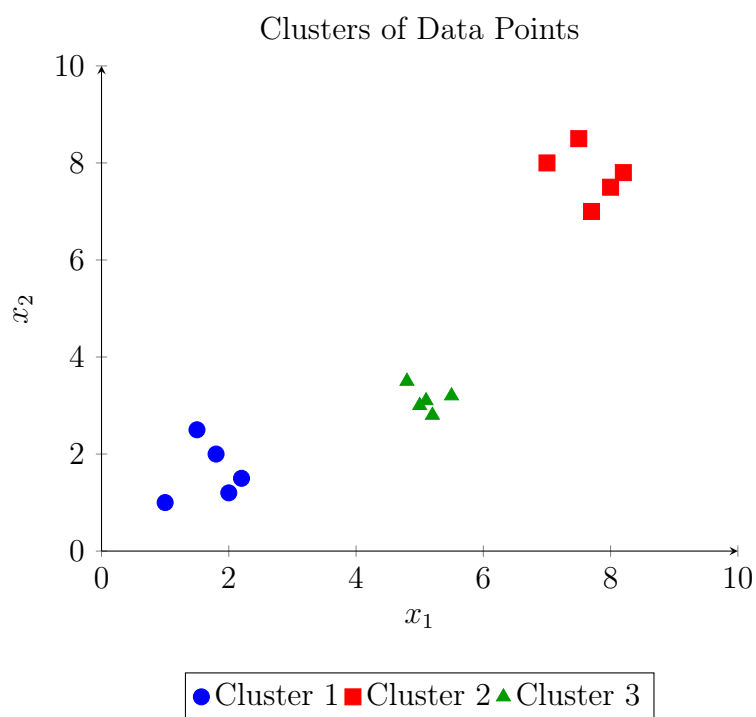


Figure 4: Illustration of three clusters in a two-dimensional feature space. Each cluster groups data points that are more similar to each other than to those in other clusters, based on Euclidean distance.

de données assumption: The dispositifs of a FL network form clusters. Two dispositifs in the same cluster generate jeu de données locals with similar statistical properties. CFL pools the jeu de données locals of dispositifs in the same cluster to obtain a ensemble d’entraînement for a cluster-specific modèle. Generalized total variation minimization (GTVMin) clusters dispositifs implicitly by enforcing approximate similarity of paramètres du modèle across well-connected nodes of the FL network. 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.

computational aspects By computational aspects of an apprentissage automatique method, we mainly refer to the computational resources required for its implementation. For example, if an apprentissage automatique method uses iterative optimization techniques to solve MRE, then its computational aspects include: 1) how many arithmetic operations are needed to implement a single iteration (gradient step); and 2) how many iterations are needed to obtain useful paramètres du modèle. One important example of an iterative optimization technique is descente de gradient.

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 α/β between the largest α and the smallest β 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 linear regression crucially depends on the condition number of the matrix $\mathbf{Q} = \mathbf{X}\mathbf{X}^T$, with the feature matrix \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

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

Here, $\|\mathbf{u}\|_p := \left(\sum_{j=1}^d |u_j|^p \right)^{1/p}$ 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 donn  ess $r \in \{1, \dots, m\}$ with identical $\widehat{\mathbf{w}}^{(r)}$ [25, 26].

convexe A subset $\mathcal{C} \subseteq \mathbb{R}^d$ of the Euclidean space \mathbb{R}^d is referred to as convex if it contains the line segment between any two points $\mathbf{x}, \mathbf{y} \in \mathcal{C}$ in that set. A function $f: \mathbb{R}^d \rightarrow \mathbb{R}$ is convex if its epigraph $\{(\mathbf{w}^T, t)^T \in \mathbb{R}^{d+1} : t \geq f(\mathbf{w})\}$ is a convex set [27]. We illustrate one example of a convex set and a convex function in Figure 5.



Figure 5: Left: A convex set $\mathcal{C} \subseteq \mathbb{R}^d$. Right: A convex function $f: \mathbb{R}^d \rightarrow \mathbb{R}$.

Courant–Fischer–Weyl min-max characterization Consider a sdp matrix $\mathbf{Q} \in \mathbb{R}^{d \times d}$ with d  composition en   l  ments 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.

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. 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 augmentation Données augmentation methods add synthetic point de données to an existing set of point de données. These synthetic point de données are obtained by perturbations (e.g., adding noise to physical measurements) or transformations (e.g., rotations of images) of the original point de données. These perturbations and transformations are such that the resulting synthetic point de données should still have the same étiquette. As a case in point, a rotated cat image is still a cat image even if their vecteur de caractéristique (obtained by stacking pixel color intensities) are very different (see Figure 6). Données augmentation can be an efficient form of régularisation.

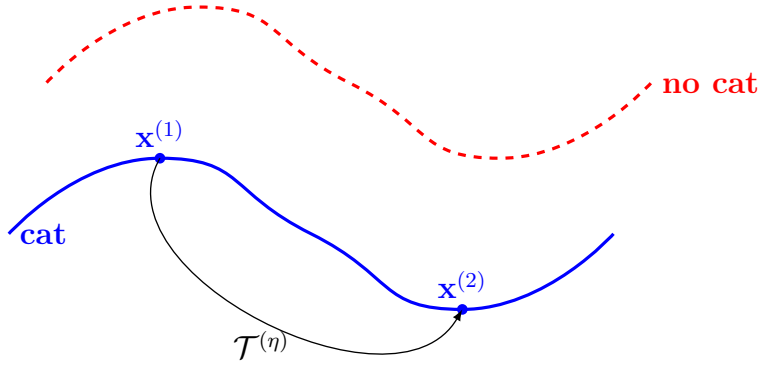


Figure 6: Données augmentation exploits intrinsic symmetries of point de données in some espace des caractéristiques \mathcal{X} . We can represent a symmetry by an operator $\mathcal{T}^{(\eta)} : \mathcal{X} \rightarrow \mathcal{X}$, parametrized by some number $\eta \in \mathbb{R}$. For example, $\mathcal{T}^{(\eta)}$ might represent the effect of rotating a cat image by η degrees. A point de données with vecteur de caractéristiques $\mathbf{x}^{(2)} = \mathcal{T}^{(\eta)}(\mathbf{x}^{(1)})$ must have the same étiquette $y^{(2)} = y^{(1)}$ as a point de données with vecteur de caractéristiques $\mathbf{x}^{(1)}$.

data minimization principle European données protection regulation includes a données minimization principle. This principle requires a données controller to limit the collection of personal information to what is directly relevant and necessary to accomplish a specified purpose. The données should be retained only for as long as necessary to fulfill that purpose [28, Article 5(1)(c)], [29].

data normalization Données normalization refers to transformations applied to the vecteur de caractéristiques of point de données to improve the apprentissage automatique method's statistical aspects or computational aspects. For example, in linear regression 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 données to steer the training of an apprentissage automatique modèle [30, 31]. The protection against données poisoning is particularly important in distributed apprentissage automatique applications where jeu de données are decentralized.

deep net A deep net is an RNA with a (relatively) large number of hidden layers. Deep learning is an umbrella term for apprentissage automatique methods that use a deep net as their modèle [32].

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 algorithme for point

de données that are characterized by numeric vecteur de caractéristiques. Like k -moyennes and soft clustering via Gaussian mixture model (GMM), also DBSCAN uses the Euclidean distances between vecteur de caractéristiques 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 caractéristiques have a large Euclidean distance.

descente de gradient Gradient descent is an iterative method for finding the minimum of a dérivable function $f(\mathbf{w})$ of a vector-valued argument $\mathbf{w} \in \mathbb{R}^d$. Consider a current guess or approximation $\mathbf{w}^{(k)}$ for the minimum of the function $f(\mathbf{w})$. We would like to find a new (better) vector $\mathbf{w}^{(k+1)}$ that has a smaller objective value $f(\mathbf{w}^{(k+1)}) < f(\mathbf{w}^{(k)})$ than the current guess $\mathbf{w}^{(k)}$. We can achieve this typically by using a gradient step

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

with a sufficiently small taille de pas $\eta > 0$. Figure 7 illustrates the effect of a single gradient descent step (2).

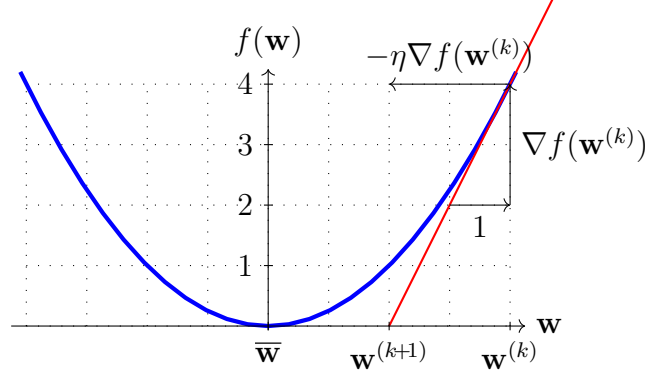


Figure 7: A single gradient step (2) towards the minimizer $\bar{\mathbf{w}}$ of $f(\mathbf{w})$.

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 automatique 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 probabilistic model 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 hy-

pothèses 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 FL network. apprentissage fédéré methods use a discrepancy measure to compare hypothèse maps from local models at nodes i, i' connected by an edge in the FL network.

dispositif Any physical system that can be used to store and process données. In the context of apprentissage automatique, we typically mean a computer that is able to read in point de données from different sources and, in turn, to train an apprentissage automatique modèle using these point de données.

distributed algorithm A distributed algorithm is an algorithm 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 algorithm, which is implemented on a single dispositif, a distributed algorithm is executed concurrently on multiple dispositifs 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 dispositif 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 dispositif and messages received from other dispositifs. One important application of distributed algorithms is in apprentissage fédéré where a network of dispositifs collaboratively train a personal modèle for each dispositif.

données Data refers to objects that carry information. These objects can be either concrete physical objects (such as persons or animals) or abstract concepts (such as numbers). We often use representations (or approximations) of the original data that are more convenient for data processing. These approximations are based on different data models, with the relational data model being one of the most widely used [35].

dual norm Every norme $\|\cdot\|$ defined on an Euclidean space \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 λ_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].

edge weight Each edge $\{i, i'\}$ of an FL network is assigned a non-negative edge weight $A_{i,i'} \geq 0$. A zero edge weight $A_{i,i'} = 0$ indicates the absence of an edge between nodes $i, i' \in \mathcal{V}$.

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) A set of point de données that have been used neither to train a modèle (e.g., via MRE) nor in a ensemble de validation to choose between different modèles.

ensemble de validation (ou jeu de validation) A set of point de données used to estimate the risque of a hypothèse \hat{h} that has been learned by some apprentissage automatique method (e.g., solving MRE). The average perte of \hat{h} on the validation set is referred to as the erreur de validation and can be used to diagnose an apprentissage automatique method (see [8, Sec. 6.6]). The comparison between erreur d'entraînement and erreur de validation can inform directions for improvement of the apprentissage automatique method (such as using a different espace des hypothèses).

erreur d'entraînement The average perte of a hypothèse when predicting the étiquettes of the point de données in a ensemble d'entraînement. We sometimes refer by training error also to minimal average perte which is achieved by a solution of MRE.

erreur de validation Consider a hypothèse \hat{h} which is obtained by some

apprentissage automatique method, e.g., using MRE on a ensemble d'entraînement. The average perte of \hat{h} on a ensemble de validation, which is different from the ensemble d'entraînement, is referred to as the validation error.

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'Euclidean space \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 statistical aspects. 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 étiquettes Considérons une application d'apprentissage automatique dans laquelle les points de données sont caractérisés par des caractéristiques et des étiquettes. L'espace des étiquettes correspond à

l'ensemble de toutes les valeurs possibles que peut prendre l'étiquette d'un point de données. Les méthodes de régression, qui visent 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 $\mathcal{Y} = \{\text{"image de chat"}, \text{"pas une image de chat"}\}$.

espace probabilisé A probability space is a mathematical model of a physical process (a random experiment) with an uncertain outcome. Formally, a probability space \mathcal{P} is a triplet (Ω, \mathcal{F}, P) where

- Ω is a sample space containing all possible elementary outcomes of a random experiment;
- \mathcal{F} is a sigma-algebra, a collection of subsets of Ω (called events) that satisfies certain closure properties under set operations;
- P is a probability measure, a function that assigns a probability $P(\mathcal{A}) \in [0, 1]$ to each event $\mathcal{A} \in \mathcal{F}$. The function must satisfy $P(\Omega) = 1$ and $P(\bigcup_{i=1}^{\infty} \mathcal{A}_i) = \sum_{i=1}^{\infty} P(\mathcal{A}_i)$ for any countable sequence of pairwise disjoint events $\mathcal{A}_1, \mathcal{A}_2, \dots$ in \mathcal{F} .

Probability spaces provide the foundation for defining VAs and to reason about uncertainty in apprentissage automatique applications [6, 19, 36].

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, 37]. 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 ε 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, 38].

Euclidean space The Euclidean space \mathbb{R}^d of dimension $d \in \mathbb{N}$ consists of vectors $\mathbf{x} = (x_1, \dots, x_d)$, with d real-valued entries $x_1, \dots, x_d \in \mathbb{R}$. Such an Euclidean space is equipped with a geometric structure defined by the inner product $\mathbf{x}^T \mathbf{x}' = \sum_{j=1}^d x_j x'_j$ between any two vectors $\mathbf{x}, \mathbf{x}' \in \mathbb{R}^d$ [2].

expectation-maximization (EM) Consider a probabilistic model $p(\mathbf{z}; \mathbf{w})$ for the point de données \mathcal{D} generated in some apprentissage automatique application. The maximum likelihood estimator for the paramètres du modèle \mathbf{w} is obtained by maximizing $p(\mathcal{D}; \mathbf{w})$. However, the resulting optimization problem might be computationally challenging. Espérance-maximization approximates the maximum likelihood estimator by introducing a latent VA \mathbf{z} such that maximizing $p(\mathcal{D}, \mathbf{z}; \mathbf{w})$ would be easier [20, 39, 40]. 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 $\max_{\mathbf{w}} p(\mathcal{D}, \hat{\mathbf{z}}; \mathbf{w})$. 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 apprentissage automatique aims to learn a hypothèse h that accurately predicts the étiquette of a point de données based on its caractéristiques. We measure the prédiction error using some fonction de perte. Ideally, we want to find a hypothèse that incurs minimal perte on any point de données. We can make this informal goal precise via the independent and identically distributed assumption (i.i.d. assumption) and by using the Bayes risk as the baseline for the (average) perte of a hypothèse. An alternative approach to obtaining a baseline is to use the hypothèse h' learned by an existing apprentissage automatique method. We refer to this hypothèse h' as an expert [41]. Regret minimization methods learn a hypothèse that incurs a perte comparable to the best expert [41, 42].

explainability We define the (subjective) explainability of an apprentissage automatique method as the level of simulatability [43] of the prédictions delivered by an apprentissage automatique system to a human

user. Quantitative measures for the (subjective) explainability of a trained modèle can be constructed by comparing its prédictions with the prédictions provided by a user on a ensemble de test [43, 44]. Alternatively, we can use probabilistic models for données and measure the explainability of a trained apprentissage automatique modèle via the conditional (differential) entropy of its prédictions, given the user prédictions [45, 46].

explainable empirical risk minimization (EERM) Explainable MRE is an instance of SRM that adds a régularisation term to the average perte in the objective function 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 [44].

explainable machine learning (explainable ML) Explainable apprentissage automatique methods aim at complementing each prédiction with an explanation of how the prédiction has been obtained. The construction of an explicit explanation might not be necessary if the apprentissage automatique method uses a sufficiently simple (or interpretable) modèle [47].

explanation One approach to make apprentissage automatique methods transparent is to provide an explanation along with the prédiction delivered by an apprentissage automatique method. Explanations can take on many different forms. An explanation could be some natural text or some quantitative measure for the importance of individual

caractéristiques of a point de données [48]. We can also use visual forms of explanations, such as intensity plots for image classification [49].

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 Φ , 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,

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

feature map Caractéristique map refers to a map that transforms the original caractéristiques of a point de données into new caractéristiques. The so-obtained new caractéristiques might be preferable over the original caractéristiques for several reasons. For example, the arrangement of

point de données might become simpler (or more linear) in the new espace des caractéristiques, allowing the use of modèle linéaires in the new caractéristiques. This idea is a main driver for the development of kernel methods [50]. Moreover, the hidden layers of a deep net can be interpreted as a trainable caractéristique map followed by a modèle linéaire in the form of the output layer. Another reason for learning a caractéristique map could be that learning a small number of new caractéristiques helps to avoid surapprentissage and ensures interpretability [51]. The special case of a caractéristique map delivering two numeric caractéristiques is particularly useful for données visualization. Indeed, we can depict point de données in a scatterplot by using two caractéristiques as the coordinates of a point de données.

feature matrix Consider a jeu de données \mathcal{D} with m point de données with vecteur de caractéristique $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)} \in \mathbb{R}^d$. It is convenient to collect the individual vecteur de caractéristique into a caractéristique matrix $\mathbf{X} := (\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)})^T$ of size $m \times d$.

FedAvg FedAvg refers to a family of iterative apprentissage fédéré algorithmes. It uses a server-client setting and alternatives between client-wise local models re-training, followed by the aggregation of updated paramètres du modèle at the server [52].

See also: apprentissage fédéré, algorithme.

federated learning network (FL network) A FL network is an undirected weighted graphe whose nodes represent données generators that aim to train a local (or personalized) modèle. Each node in a FL network

represents some dispositif capable of collecting a jeu de données local and, in turn, train a local model. Apprentissage fédéré methods learn a local hypothèse $h^{(i)}$, for each node $i \in \mathcal{V}$, such that it incurs small perte on the jeu de données locaux.

FedGD An apprentissage fédéré distributed algorithm that can be implemented as message passing across an FL network.

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

FedProx FedProx refers to an iterative apprentissage fédéré algorithme that alternates between separately training local models and combining the updated local paramètres du modèle. In contrast to FedAvg, which uses stochastic gradient descent (SGD) to train local models, FedProx uses a proximal operator for the training [53].

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 FL network.

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

Finnish Meteorological Institute (FMI) The FMI is a government agency responsible for gathering and reporting weather données in Finland.

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 [54].

fonction d’activation Each artificial neuron within an RNA is assigned an activation function $\sigma(\cdot)$ that maps a weighted combination of the neuron inputs x_1, \dots, x_d to a single output value $a = \sigma(w_1x_1 + \dots + w_dx_d)$. Note that each neuron is parametrized by the poids w_1, \dots, w_d .

fonction de perte (ou de coût) A perte function is a map

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

It assigns a non-negative real number (i.e., the perte) $L((\mathbf{x}, y), h)$ to a pair that consists of a point de données, with caractéristiques \mathbf{x} and étiquette y , and a hypothèse $h \in \mathcal{H}$. The value $L((\mathbf{x}, y), h)$ quantifies the discrepancy between the true étiquette y and the prédiction $h(\mathbf{x})$. Lower (closer to zero) values $L((\mathbf{x}, y), h)$ indicate a smaller discrepancy between prédiction $h(\mathbf{x})$ and étiquette y . Figure 8 depicts a perte function for a given point de données, with caractéristiques \mathbf{x} and étiquette y , as a function of the hypothèse $h \in \mathcal{H}$.

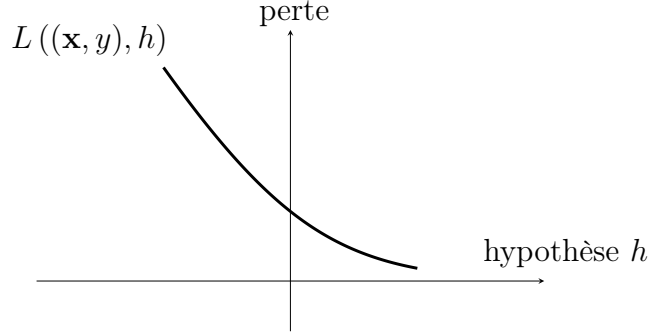


Figure 8: Some perte function $L((\mathbf{x}, y), h)$ for a fixed point de données, with vecteur de caractéristiques \mathbf{x} and étiquette y , and a varying hypothèse h . apprentissage automatique methods try to find (or learn) a hypothèse that incurs minimal perte.

frontière de décision Consider a hypothèse map h that reads in a caractéristique vector $\mathbf{x} \in \mathbb{R}^d$ and delivers a value from a finite set \mathcal{Y} . The decision boundary of h is the set of vectors $\mathbf{x} \in \mathbb{R}^d$ that lie between different région de décisions. More precisely, a vector \mathbf{x} belongs to the decision boundary if and only if each voisinage $\{\mathbf{x}' : \|\mathbf{x} - \mathbf{x}'\| \leq \varepsilon\}$, for any $\varepsilon > 0$, contains at least two vectors with different function values.

Gaussian mixture model (GMM) A GMM is a particular type of probabilistic model 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 $p(I = c) = p_c$. 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 $f(\mathbf{x}^{(1)}), \dots, f(\mathbf{x}^{(m)})$ have a joint multivariate Gaussian distribution:

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

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 Finnish Meteorological Institute (FMI) weather stations provide samples of $f(\mathbf{x})$ at specific locations (see Fig. 9). A GP allows to predict the temperature nearby FMI weather stations and to quantify the uncertainty of these predictions.

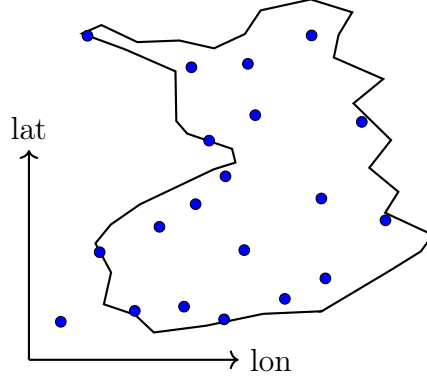


Figure 9: We can interpret the temperature distribution over Finland as a réalisation of a GP indexed by geographic coordinates and sampled at FMI weather stations (indicated by blue dots).

generalization Many current apprentissage automatique (and intelligence artificielle (IA)) systems are based on MRE: At their core, they train a modèle (i.e., learn a hypothèse $\hat{h} \in \mathcal{H}$) by minimizing the average perte (or risque empirique) on some point de données $\mathbf{z}^{(1)}, \dots, \mathbf{z}^{(m)}$, which serve as a ensemble d'entraînement $\mathcal{D}^{(\text{train})}$. Generalization refers to an apprentissage automatique method's ability to perform well outside the ensemble d'entraînement. Any mathematical theory of generalization needs some mathematical concept for the "outside the ensemble d'entraînement." For example, statistical learning theory uses a probabilistic model such as the i.i.d. assumption for données generation: the point de données in the ensemble d'entraînement are i.i.d. réalisations of some underlying loi de probabilité $p(\mathbf{z})$. A probabilistic model allows us to explore the outside of the ensemble d'entraînement by drawing

additional i.i.d. réalisations from $p(\mathbf{z})$. Moreover, using the i.i.d. assumption allows us to define the risque of a trained modèle $\hat{h} \in \mathcal{H}$ as the expected perte $\bar{L}(\hat{h})$. What is more, we can use concentration bounds or convergence results for sequences of i.i.d. VAs to bound the deviation between the risque empirique $\hat{L}(\hat{h}|\mathcal{D}^{(\text{train})})$ of a trained modèle and its risque [55]. It is possible to study generalization also without using probabilistic models. For example, we could use (deterministic) perturbations of the point de données in the ensemble d'entraînement to study its outside. In general, we would like the trained modèle to be robust, i.e., its prédictions should not change too much for small perturbations of a point de données. Consider a trained modèle for detecting an object in a smartphone snapshot. The detection result should not change if we mask a small number of randomly chosen pixels in the image [56].

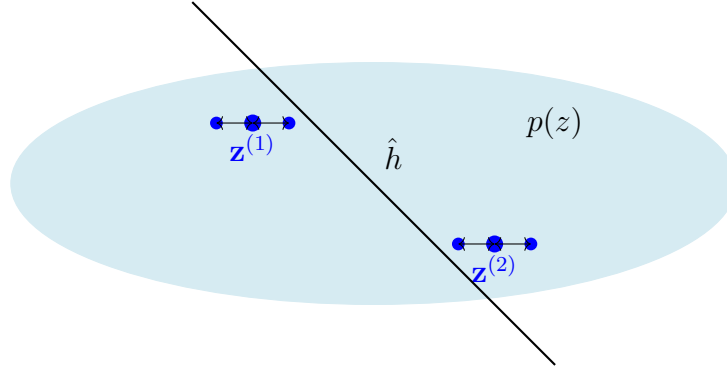


Figure 10: Two point de données $\mathbf{z}^{(1)}, \mathbf{z}^{(2)}$ that are used as a ensemble d'entraînement to learn a hypothèse \hat{h} via MRE. We can evaluate \hat{h} outside $\mathcal{D}^{(\text{train})}$ either by an i.i.d. assumption with some underlying loi de probabilité $p(\mathbf{z})$ or by perturbing the point de données.

generalized total variation (GTV) GTV is a measure of the variation of trained local models $h^{(i)}$ (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 [57].

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

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

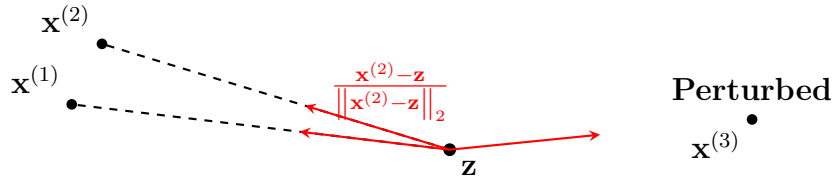


Figure 11: Consider a solution \mathbf{z} of (3) that does not coincide with any of the input vectors. The optimality condition for (3) requires that the unit vectors from \mathbf{z} to the input vectors sum to zero.

Figure 11 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 [58, 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].

gradient step Given a dérivable real-valued function $f(\cdot) : \mathbb{R}^d \rightarrow \mathbb{R}$ and a vector $\mathbf{w} \in \mathbb{R}^d$, the gradient step updates \mathbf{w} by adding the scaled negative gradient $\nabla f(\mathbf{w})$ to obtain the new vector (see Figure 12)

$$\hat{\mathbf{w}} := \mathbf{w} - \eta \nabla f(\mathbf{w}). \quad (4)$$

Mathematically, the gradient step is a (typically non-linear) operator $\mathcal{T}^{(f,\eta)}$ that is parametrized by the function f and the taille de pas η .

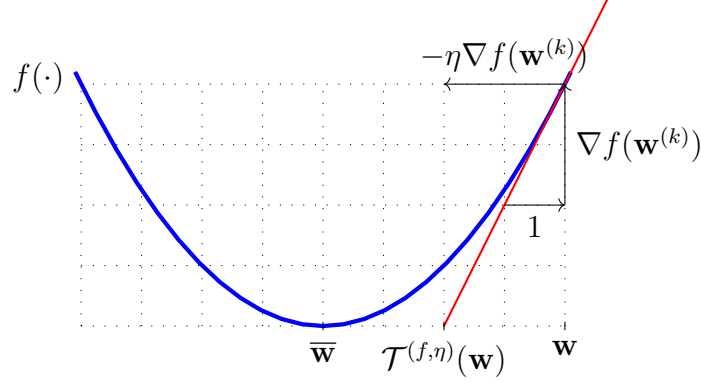


Figure 12: The basic gradient step (4) maps a given vector \mathbf{w} to the updated vector \mathbf{w}' . It defines an operator $\mathcal{T}^{(f,\eta)}(\cdot) : \mathbb{R}^d \rightarrow \mathbb{R}^d : \mathbf{w} \mapsto \hat{\mathbf{w}}$.

Note that the gradient step (4) optimizes locally - in a voisinage whose size is determined by the taille de pas η - a linear approximation to the function $f(\cdot)$. A natural generalization of (4) is to locally optimize the function itself - instead of its linear approximation - such that

$$\hat{\mathbf{w}} = \operatorname{argmin}_{\mathbf{w}' \in \mathbb{R}^d} f(\mathbf{w}') + (1/\eta) \|\mathbf{w} - \mathbf{w}'\|_2^2. \quad (5)$$

We intentionally use the same symbol η for the parameter in (5) as we used for the taille de pas in (4). The larger the η we choose in (5), the more progress the update will make towards reducing the function value $f(\hat{\mathbf{w}})$. Note that, much like the gradient step (4), also the update (5) defines a (typically non-linear) operator that is parametrized by the function $f(\cdot)$ and the parameter η . For a convexe function $f(\cdot)$, this operator is known as the proximal operator of $f(\cdot)$ [59].

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 deep nets 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 labeled datapoints 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.

graphe 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 edge weight [54, 61].

graphe A graph $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ is a pair that consists of a node set \mathcal{V} and an edge set \mathcal{E} . In its most general form, a graph is specified by a map that assigns each edge $e \in \mathcal{E}$ a pair of nodes [62]. One important family of graphs is simple undirected graphs. A simple undirected graph is obtained by identifying each edge $e \in \mathcal{E}$ with two different nodes $\{i, i'\}$. Weighted graphs also specify numeric poids A_e for each edge $e \in \mathcal{E}$.

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.

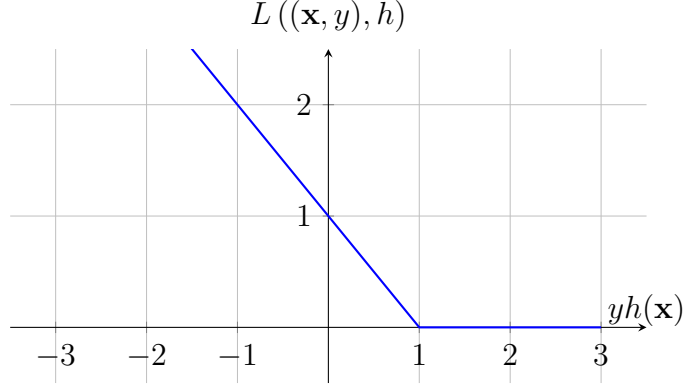
high-dimensional regime The high-dimensional regime of MRE is characterized by the dimension effective of the modèle being larger than the sample size, i.e., the number of (labeled) point de données in the ensemble d'entraînement. For example, linear regression methods operate in the high-dimensional 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 Euclidean space \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 hypoth ese map $h(\mathbf{x})$ is defined as

$$L((\mathbf{x}, y), h) := \max\{0, 1 - yh(\mathbf{x})\}. \quad (6)$$



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

histogram Consider a jeu de donn ees \mathcal{D} that consists of m point de donn eess $\mathbf{z}^{(1)}, \dots, \mathbf{z}^{(m)}$, 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 Δ . The histogram of \mathcal{D} assigns each elementary cell to the corresponding fraction of point de donn eess in \mathcal{D} that fall into this elementary cell. A visual example of such a histogram is provided in Figure 13.

horizontal federated learning (HFL) HFL uses jeu de donn ees locals constituted by different point de donn eess but uses the same caract eristiques to characterize them [67]. For example, weather forecasting uses a network of spatially distributed weather (observation) stations. Each

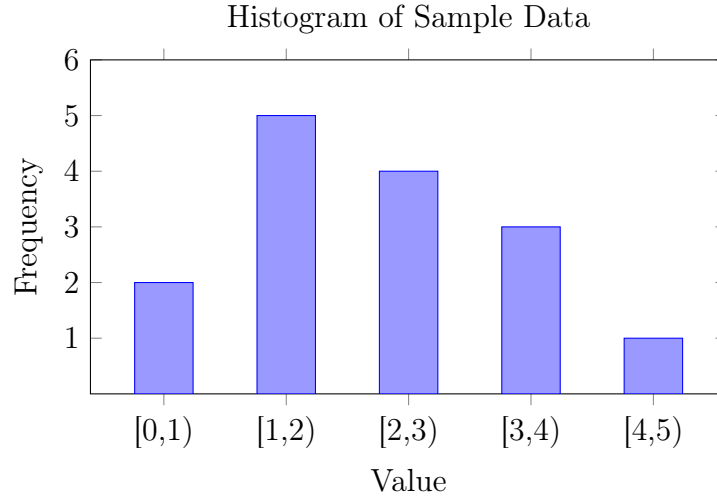


Figure 13: A histogram representing the frequency of data points falling within discrete value ranges (bins). Each bar height shows the count of samples in the corresponding interval.

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 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 squared error loss and the absolute error loss.

Huber regression Huber régression refers to MRE-based methods that use

the Huber loss as a measure of the prediction error. Two important special cases of Huber regression are least absolute deviation regression and linear regression. 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 lisse squared error loss.

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).

independent and identically distributed assumption (i.i.d. assumption)

The i.i.d. assumption interprets point de données of a jeu de données as the réalisations of i.i.d. VAs.

indépendantes et identiquement distribuées (i.i.d.) It can be useful to interpret point de données $\mathbf{z}^{(1)}, \dots, \mathbf{z}^{(m)}$ as réalisations of i.i.d. VAs with a common loi de probabilité. If these VAs are continuous-valued, their joint probability density function (pdf) is $p(\mathbf{z}^{(1)}, \dots, \mathbf{z}^{(m)}) = \prod_{r=1}^m p(\mathbf{z}^{(r)})$, with $p(\mathbf{z})$ being the common marginal pdf of the underlying VAs.

intelligence artificielle (IA) AI refers to systems that behave rationally in the sense of maximizing a long-term reward. The apprentissage

automatique-based approach to AI is to train a modèle for predicting optimal actions. These predictions are computed from observations about the state of the environment. The choice of fonction de perte sets AI applications apart from more basic apprentissage automatique applications. AI systems rarely have access to a labeled ensemble d'entraînement that allows the average perte to be measured for any possible choice of paramètres du modèle. Instead, AI systems use observed reward signals to obtain a (point-wise) estimate for the perte incurred by the current choice of paramètres du modèle.

interpretability An apprentissage automatique method is interpretable for a specific user if they can well anticipate the prédictions delivered by the method. The notion of interpretability can be made precise using quantitative measures of the uncertainty about the prédictions [45].

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 14 illustre un jeu de données utilisant des vaches comme points de données.



Figure 14: « Cows in the Swiss Alps » (*Des vaches dans les Alpes Suisses*) par User:Huhu Uet - sous licence [CC BY-SA 4.0](<https://creativecommons.org/licenses/by-sa/4.0/>)

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 14 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 [68], [69], [70], [71].

L'un des modèles de données les plus utilisés est le modèle relationnel, qui organise les données sous forme de tableau (ou relation) [35], [68].

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 14. 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 14.

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 trustworthy artificial intelligence (trustworthy AI).

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 [72].

jeu de données local The concept of a local jeu de données is in between the concept of a point de données and a jeu de données. A local jeu de données consists of several individual point de données, which are characterized by caractéristiques and étiquettes. In contrast to a single jeu de données used in basic apprentissage automatique methods, a local jeu de données is also related to other local jeu de données via different notions of similarity. These similarities might arise from probabilistic models or communication infrastructure and are encoded in the edges of an FL network.

kernel Consider point de données characterized by a vecteur de caractéristiques $\mathbf{x} \in \mathcal{X}$ with a generic espace des caractéristiques \mathcal{X} . A (real-valued) kernel $K : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ assigns each pair of vecteur de caractéristiques $\mathbf{x}, \mathbf{x}' \in \mathcal{X}$ a real number $K(\mathbf{x}, \mathbf{x}')$. The value $K(\mathbf{x}, \mathbf{x}')$ is often interpreted as a measure for the similarity between \mathbf{x} and \mathbf{x}' . Kernel methods use a kernel to transform the vecteur de caractéristiques \mathbf{x} to a new vecteur de caractéristiques $\mathbf{z} = K(\mathbf{x}, \cdot)$. This new vecteur de caractéristiques belongs to a linear espace des caractéristiques \mathcal{X}' which is (in general) different from the original espace des caractéristiques \mathcal{X} . The espace des caractéristiques \mathcal{X}' has a specific mathematical structure, i.e., it is a reproducing kernel Hilbert space [50, 66].

kernel method A kernel method is an apprentissage automatique method that uses a kernel K to map the original (raw) vecteur de caractéristiques

\mathbf{x} of a point de données to a new (transformed) vecteur de caractéristiques $\mathbf{z} = K(\mathbf{x}, \cdot)$ [50,66]. The motivation for transforming the vecteur de caractéristiquess is that, by using a suitable kernel, the point de données have a "more pleasant" geometry in the transformed espace des caractéristiques. For example, in a binary classification problem, using transformed vecteur de caractéristiquess \mathbf{z} might allow us to use modèle linéaires, even if the point de données are not linearly separable in the original espace des caractéristiques (see Figure 15).

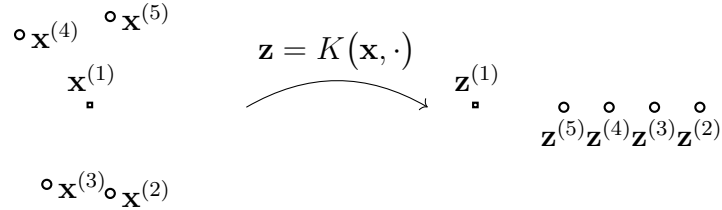


Figure 15: Five point de données characterized by vecteur de caractéristiques $\mathbf{x}^{(r)}$ and étiquettes $y^{(r)} \in \{\circ, \square\}$, for $r = 1, \dots, 5$. With these vecteur de caractéristiquess, there is no way to separate the two classes by a straight line (representing the frontière de décision of a linear classifier). In contrast, the transformed vecteur de caractéristiquess $\mathbf{z}^{(r)} = K(\mathbf{x}^{(r)}, \cdot)$ allow us to separate the point de données using a linear classifier.

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 [73].

labeled datapoint A point de données whose étiquette is known or has

been determined by some means which might require human labor.

Laplacian matrix The structure of a graphe \mathcal{G} , with nodes $i = 1, \dots, n$, can be analyzed using the properties of special matrices that are associated with \mathcal{G} . One such matrix is the graphe Laplacian matrix $\mathbf{L}^{(\mathcal{G})} \in \mathbb{R}^{n \times n}$, which is defined for an undirected and weighted graphe [61, 74]. It is defined element-wise as (see Figure 16)

$$L_{i,i'}^{(\mathcal{G})} := \begin{cases} -A_{i,i'} & \text{for } i \neq i', \{i, i'\} \in \mathcal{E}, \\ \sum_{i'' \neq i} A_{i,i''} & \text{for } i = i', \\ 0 & \text{else.} \end{cases} \quad (7)$$

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

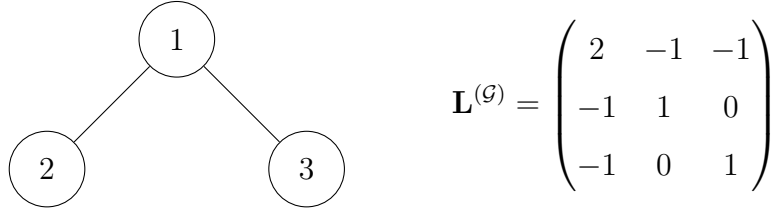


Figure 16: Left: Some undirected graphe \mathcal{G} with three nodes $i = 1, 2, 3$. Right: The Laplacian matrix $\mathbf{L}^{(\mathcal{G})} \in \mathbb{R}^{3 \times 3}$ of \mathcal{G} .

law of large numbers The law of large numbers refers to the convergence of the average of an increasing (large) number of i.i.d. VAs to the moyenne of their common loi de probabilité. Different instances of the law of large numbers are obtained by using different notions of convergence [75].

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 linear regression by adding the scaled ℓ_1 -norme $\alpha \|\mathbf{w}\|_1$ to the average squared error loss incurred on the ensemble d'entraînement.

linear classifier Consider point de données characterized by numeric caractéristiques $\mathbf{x} \in \mathbb{R}^d$ and a étiquette $y \in \mathcal{Y}$ from some finite espace des étiquettes \mathcal{Y} . A linear classifier is characterized by having région de décisions that are separated by hyperplanes in \mathbb{R}^d [8, Ch. 2].

linear regression Linear régression aims to learn a linear hypothèse map to predict a numeric étiquette based on the numeric caractéristiques of a point de données. The quality of a linear hypothèse map is measured using the average squared error loss incurred on a set of labeled datapoints, which we refer to as the ensemble d'entraînement.

lisse (ou régulière) A real-valued function $f : \mathbb{R}^d \rightarrow \mathbb{R}$ is smooth if it is dérivable and its gradient $\nabla f(\mathbf{w})$ is continuous at all $\mathbf{w} \in \mathbb{R}^d$ [76, 77]. A smooth function f is referred to as β -smooth if the gradient $\nabla f(\mathbf{w})$ is Lipschitz continuous with Lipschitz constant β , i.e.,

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

The constant β quantifies the amount of smoothness of the function f :

the smaller the β , the smoother f is. Optimization problems with a smooth objective function can be solved effectively by méthodes basées sur le gradient. Indeed, méthodes basées sur le gradient approximate the objective function locally around a current choice \mathbf{w} using its gradient. This approximation works well if the gradient does not change too rapidly. We can make this informal claim precise by studying the effect of a single gradient step with taille de pas $\eta = 1/\beta$ (see Figure 17).

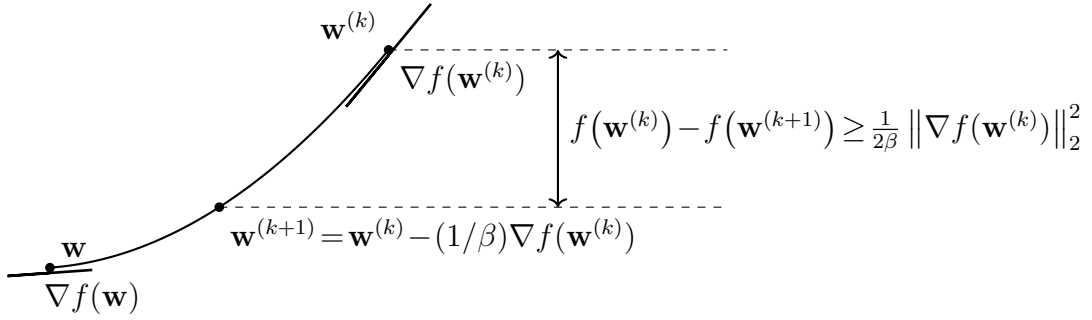


Figure 17: Consider an objective function $f(\mathbf{w})$ that is β -smooth. Taking a gradient step, with taille de pas $\eta = 1/\beta$, decreases the objective by at least $\frac{1}{2\beta} \|\nabla f(\mathbf{w}^{(k)})\|_2^2$ [76–78]. Note that the taille de pas $\eta = 1/\beta$ becomes larger for smaller β . Thus, for smoother objective functions (i.e., those with smaller β), we can take larger steps.

Local Interpretable Model-agnostic Explanations (LIME) Consider a trained modèle (or learnt hypothèse) $\hat{h} \in \mathcal{H}$, which maps the vecteur de caractéristiques of a point de données to the prédiction $\hat{y} = \hat{h}$. Local Interpretable Model-agnostic Explanations (LIME) is a technique for explaining the behaviour of \hat{h} , locally around a point de données with

vecteur de caractéristiques $\mathbf{x}^{(0)}$ [51]. The explanation is given in the form of a local approximation $g \in \mathcal{H}'$ of \hat{h} (see Fig.). This approximation can be obtained by an instance of MRE with carefully designed ensemble d'entraînement. In particular, the ensemble d'entraînement consists of point de données with vecteur de caractéristiques \mathbf{x} close to $\mathbf{x}^{(0)}$ and the (pseudo-)label $\hat{h}(\mathbf{x})$. Note that we can use a different modèle \mathcal{H}' for the approximation than the original modèle \mathcal{H} . For example, we can use a arbre de décision to approximate (locally) a deep net. Another widely-used choice for \mathcal{H}' is the modèle linéaire.

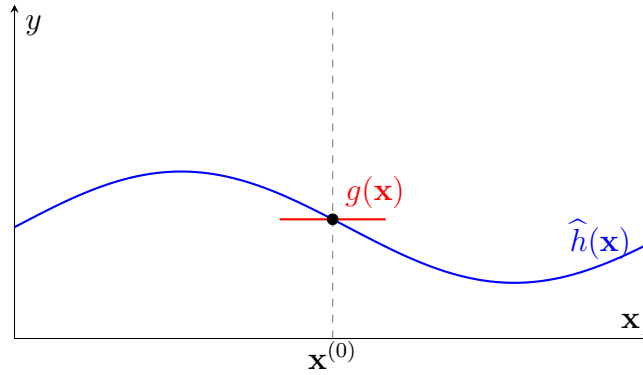


Figure 18: To explain a trained modèle $\hat{h} \in \mathcal{H}$, around a given vecteur de caractéristiques $\mathbf{x}^{(0)}$, we can use a local approximation $g \in \mathcal{H}'$.

local model Consider a collection of jeu de données locals that are assigned to the nodes of an FL network. A local modèle $\mathcal{H}^{(i)}$ is a espace des hypothèses assigned to a node $i \in \mathcal{V}$. Different nodes might be assigned different espace des hypothèsess, i.e., in general $\mathcal{H}^{(i)} \neq \mathcal{H}^{(i')}$ for different nodes $i, i' \in \mathcal{V}$.

logistic loss Consider a point de données characterized by the caractéristiques \mathbf{x} and a binary étiquette $y \in \{-1, 1\}$. We use a real-valued hypothèse h to predict the étiquette y from the caractéristiques \mathbf{x} . The logistic perte incurred by this prédiction is defined as

$$L((\mathbf{x}, y), h) := \log(1 + \exp(-yh(\mathbf{x}))). \quad (8)$$

Carefully note that the expression (8) for the logistic perte applies only for the espace des étiquettes $\mathcal{Y} = \{-1, 1\}$ and when using the thresholding rule (1).

logistic regression Logistic régression learns a linear hypothèse map (or classifier) $h(\mathbf{x}) = \mathbf{w}^T \mathbf{x}$ to predict a binary étiquette y based on the numeric vecteur de caractéristiques \mathbf{x} of a point de données. The quality of a linear hypothèse map is measured by the average logistic loss on some labeled datapoints (i.e., 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 de probabilité de cette VA. La loi de probabilité d'une VA binaire $y \in \{0, 1\}$ est entièrement déterminée par les probabilités $p(y = 0)$ et $p(y = 1) = 1 - p(y = 0)$. La loi de probabilité d'une VA à valeurs réelles $x \in \mathbb{R}$ peut être spécifiée par une pdf $p(x)$ telle que $p(x \in [a, b]) \approx p(a)|b - a|$. Dans le cas le plus général, une loi de probabilité est définie par une mesure de probabilité [6, 19].

loi normale multivariée La loi normale multivariée $\mathcal{N}(\mathbf{m}, \mathbf{C})$ est un probabilistic model 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], [79]. 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 covariance La matrice de covariance d'une VA $\mathbf{x} \in \mathbb{R}^d$ est définie comme $\mathbb{E} \left\{ (\mathbf{x} - \mathbb{E}\{\mathbf{x}\})(\mathbf{x} - \mathbb{E}\{\mathbf{x}\})^T \right\}$.

maximum The maximum of a set $\mathcal{A} \subseteq \mathbb{R}$ of real numbers is the greatest element in that set, if such an element exists. A set \mathcal{A} has a maximum if it is bounded above and attains its supremum (or least upper bound) [2, Sec. 1.4].

maximum likelihood Consider point de données $\mathcal{D} = \{\mathbf{z}^{(1)}, \dots, \mathbf{z}^{(m)}\}$ that are interpreted as the réalisations of i.i.d. VAs with a common loi de probabilité $p(\mathbf{z}; \mathbf{w})$ which depends on the paramètres du modèle $\mathbf{w} \in \mathcal{W} \subseteq \mathbb{R}^n$. Maximum likelihood methods learn paramètres du modèle \mathbf{w} by maximizing the probability (density) $p(\mathcal{D}; \mathbf{w}) = \prod_{r=1}^m p(\mathbf{z}^{(r)}; \mathbf{w})$ of

the observed jeu de données. Thus, the maximum likelihood estimator is a solution to the optimization problem $\max_{\mathbf{w} \in \mathcal{W}} p(\mathcal{D}; \mathbf{w})$.

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 $\Delta \mathbf{w} := \hat{\mathbf{w}} - \bar{\mathbf{w}}$. Here, $\bar{\mathbf{w}}$ denotes the true paramètres du modèle of the loi de probabilité of \mathbf{z} . The moyenne squared estimation error is defined as the espérance $\mathbb{E}\{\|\Delta \mathbf{w}\|^2\}$ of the squared Euclidean norm of the estimation error [17, 38].

minimisation du risque empirique (MRE) Risque empirique minimization is the optimization problem of finding a hypothèse (out of a modèle) with the minimum average perte (or risque empirique) on a given jeu de données \mathcal{D} (i.e., the ensemble d'entraînement). Many apprentissage automatique methods are obtained from risque empirique via specific design choices for the jeu de données, modèle, and perte [8, Ch. 3].

minimum Given a set of real numbers, the minimum is the smallest of those numbers.

missing data Consider a jeu de données constituted by point de données collected via some physical dispositif. 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 [80]. 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], [55]. Cependant, ce terme est également utilisé dans d'autres domaines avec des significations différentes. Par exemple, un probabilistic model désigne un ensemble paramétré de lois de probabilité.

modèle linéaire Consider point de données, each characterized by a numeric vecteur de caractéristiques $\mathbf{x} \in \mathbb{R}^d$. A linear modèle is a space des hypothèses which consists of all linear maps,

$$\mathcal{H}^{(d)} := \{h(\mathbf{x}) = \mathbf{w}^T \mathbf{x} : \mathbf{w} \in \mathbb{R}^d\}. \quad (9)$$

Note that (9) defines an entire family of espace des hypothèsess, which is parametrized by the number d of caractéristiques that are linearly combined to form the prédiction $h(\mathbf{x})$. The design choice of d is guided by computational aspects (e.g., reducing d means less computation),

statistical aspects (e.g., increasing d might reduce prédiction error), and interpretability. A linear modèle using few carefully chosen caractéristiques tends to be considered more interpretable [47, 51].

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

multi-armed bandit A multi-armed bandit (MAB) problem models a repeated decision-making scenario in which, at each time step k , a learner must choose one out of several possible actions, often referred to as arms, from a finite set \mathcal{A} . Each arm $a \in \mathcal{A}$ yields a stochastic reward $r^{(a)}$ drawn from an unknown loi de probabilité with moyenne $\mu^{(a)}$. The learner’s goal is to maximize the cumulative reward over time by strategically balancing exploration (gathering information about uncertain arms) and exploitation (selecting arms known to perform well). This balance is quantified by the notion of regret, which measures the performance gap between the learner’s strategy and the optimal strategy that

always selects the best arm. MAB problems form a foundational model in online learning, reinforcement learning, and sequential experimental design [81].

mutual information (MI) The MI $I(\mathbf{x}; y)$ between two VAs \mathbf{x}, y defined on the same espace probabilisé is given by [73]

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

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.

méthodes basées sur le gradient Gradient-based methods are iterative techniques for finding the minimum (or maximum) of a dérivable objective function of the paramètres du modèle. These methods construct a sequence of approximations to an optimal choice for paramètres du modèle that results in a minimum (or maximum) value of the objective function. As their name indicates, gradient-based methods use the gradients of the objective function evaluated during previous iterations to construct new, (hopefully) improved paramètres du modèle. One important example of a gradient-based method is 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 caractéristiques, 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 dispositifs.

networked exponential families (nExpFam) A collection of exponential families, each of them assigned to a node of an FL network. The paramètres du modèle are coupled via the network structure by requiring them to have a small GTV [82].

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 FL network $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ assigns a local model (i.e., a espace des hypothèses) to each node $i \in \mathcal{V}$ of the FL network \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-smooth We refer to a function as non-smooth if it is not lisse [76].

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 [83].

objective function An objective function is a map that assigns each value of an optimization variable, such as the paramètres du modèle \mathbf{w} of a hypothèse $h^{(\mathbf{w})}$, to an objective value $f(\mathbf{w})$. The objective value $f(\mathbf{w})$ could be the risque or the risque empirique of a hypothèse $h^{(\mathbf{w})}$.

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 [41], [42]. Unlike an offline algorithm, which has the entire input available from the start, an online algorithm must handle uncertainty 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., 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 out_k , and then subsequently receives the next input (point de données) in_{k+1} . A notable example of an online algorithm in apprentissage automatique is online gradient descent

(online GD), which incrementally updates paramètres du modèle as new point de données arrive.

See also: online learning, online GD.

online gradient descent (online GD) Consider an apprentissage automatique method that learns paramètres du modèle \mathbf{w} from some parameter space $\mathcal{W} \subseteq \mathbb{R}^d$. 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 objective function 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 descente de gradient updates $\mathbf{w}^{(t)}$ via a (projected) gradient step

$$\mathbf{w}^{(t+1)} := P_{\mathcal{W}}(\mathbf{w}^{(t)} - \eta_t \nabla_{\mathbf{w}} L(\mathbf{z}^{(t)}, \mathbf{w})). \quad (10)$$

Note that (10) is a gradient step for the current component $L(\mathbf{z}^{(t)}, \cdot)$ of the risque. The update (10) 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 η_t , online descente 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 descente 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 [42, 84].

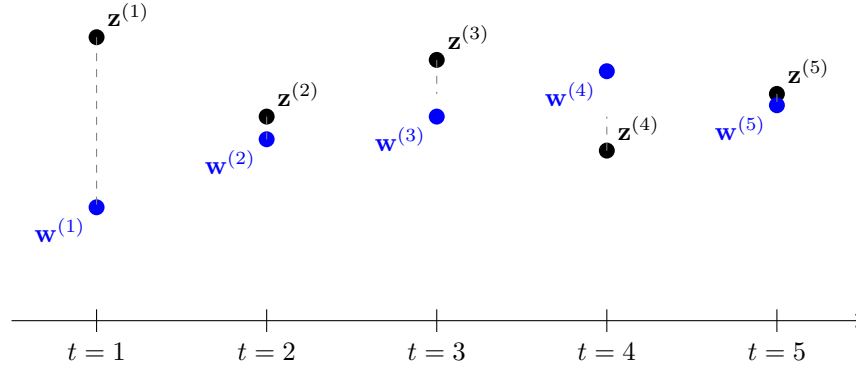


Figure 19: An instance of online descente de gradient that updates the paramètres du modèle $\mathbf{w}^{(t)}$ using the point de données $\mathbf{z}^{(t)} = x^{(t)}$ arriving at time t . This instance uses the squared error loss $L(\mathbf{z}^{(t)}, w) = (x^{(t)} - w)^2$.

online learning Some apprentissage automatique methods are designed to process données in a sequential manner, updating their paramètres du modèle as new point de données become available—one at a time. A typical example is time series data, such as daily minimum and maximum temperatures recorded by a FMI weather station. These values form a chronological sequence of observations. In online learning, the hypothèse (or its paramètres du modèle) is refined incrementally with each newly observed point de données, without revisiting past

données.

See also online GD, online algorithm.

optimism in the face of uncertainty apprentissage automatique methods learn paramètres du modèle \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 perte on a given jeu de données (i.e., the ensemble d'entraînement) as an estimate for the risque of a hypothèse. Using a probabilistic model, one can construct a confidence interval $[l^{(\mathbf{w})}, u^{(\mathbf{w})}]$ for each choice \mathbf{w} for the paramètres du modèle. One simple construction is $l^{(\mathbf{w})} := f(\mathbf{w}) - \sigma/2$, $u^{(\mathbf{w})} := f(\mathbf{w}) + \sigma/2$, with σ 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 paramètres du modèle according to the most favourable - yet still plausible - value $\tilde{f}(\mathbf{w}) := l^{(\mathbf{w})}$ of the performance criterion. Two examples of apprentissage automatique methods that use such an optimistic construction of an objective function are SRM [55, Ch. 11] and upper confidence bound (UCB) methods for sequential decision making [81, Sec. 2.2].

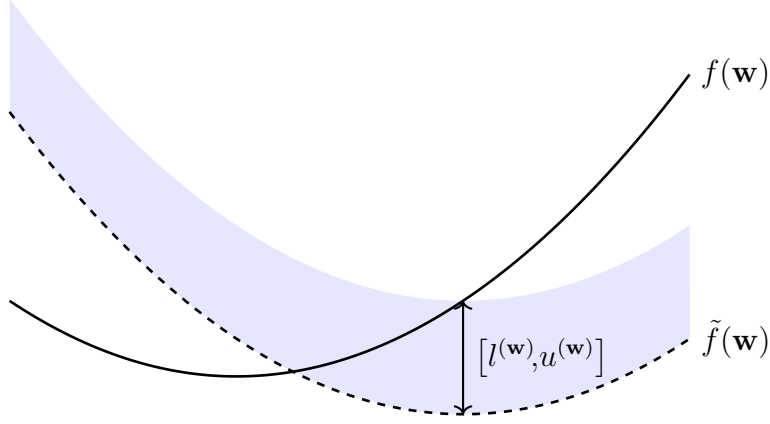


Figure 20: apprentissage automatique methods learn paramètres du modèle \mathbf{w} by using some estimate of $f(\mathbf{w})$ for the ultimate performance criterion $\bar{f}(\mathbf{w})$. Using a probabilistic model, one can use $f(\mathbf{w})$ to construct confidence intervals $[l^{(\mathbf{w})}, u^{(\mathbf{w})}]$ which contain $\bar{f}(\mathbf{w})$ with high probability. The best plausible performance measure for a specific choice \mathbf{w} of paramètres du modèle is $\tilde{f}(\mathbf{w}) := l^{(\mathbf{w})}$.

outlier Many apprentissage automatique methods are motivated by the i.i.d. assumption, which interprets point de données as réalisations of i.i.d. VAs with a common loi de probabilité. The i.i.d. assumption is useful for applications where the statistical properties of the données generation process are stationary (or time-invariant) [85]. However, in some applications the données consists of a majority of regular point de données that conform with an i.i.d. assumption 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 [86,87].

parameter space The parameter space \mathcal{W} of an apprentissage automatique modèle \mathcal{H} is the set of all feasible choices for the paramètres du modèle (see Figure 21). Many important apprentissage automatique methods use a modèle that is parametrized by vectors of the Euclidean space \mathbb{R}^d . Two widely used examples of parametrized modèles are modèle linéaires and deep nets. The parameter space is then often a subset $\mathcal{W} \subseteq \mathbb{R}^d$, e.g., all vectors $\mathbf{w} \in \mathbb{R}^d$ with a norme smaller than one.

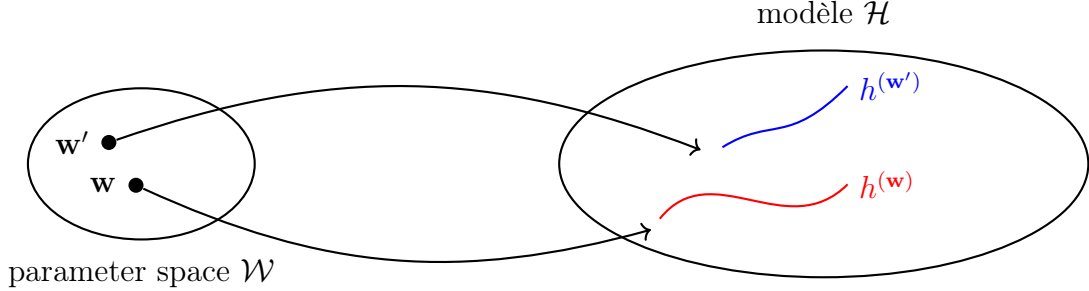


Figure 21: The parameter space \mathcal{W} of an apprentissage automatique modèle \mathcal{H} consists of all feasible choices for the paramètres du modèle. Each choice \mathbf{w} for the paramètres du modèle selects a hypothèse map $h^{(\mathbf{w})} \in \mathcal{H}$.

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 Modèle paramètres are quantities that are used to select a specific hypothèse map from a modèle. We can think of a list of modèle paramètres as a unique identifier for a hypothèse map, similar to how a social security number identifies a person in Finland.

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.

perte apprentissage automatique methods use a fonction de perte $L(\mathbf{z}, h)$ to measure the error incurred by applying a specific hypothèse to a specific point de données. With a slight abuse of notation, we use the term loss for both the fonction de perte L itself and the specific value $L(\mathbf{z}, h)$, for

a point de données \mathbf{z} and hypothèse h .

poids Consider a parametrized espace des hypothèses \mathcal{H} . We use the term weights for numeric paramètres du modèle that are used to scale caractéristiques or their transformations in order to compute $h^{(\mathbf{w})} \in \mathcal{H}$. A modèle linéaire uses weights $\mathbf{w} = (w_1, \dots, w_d)^T$ to compute the linear combination $h^{(\mathbf{w})}(\mathbf{x}) = \mathbf{w}^T \mathbf{x}$. Weights are also used in RNAs to form linear combinations of caractéristiques or the outputs of neurons in hidden layers.

point de données Un point de données correspond à tout objet qui transmet de l'information [73]. 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.

polynomial regression Polynomial régression aims at learning a polynomial hypothèse map to predict a numeric étiquette based on the numeric

caractéristiques of a point de données. For point de données characterized by a single numeric caractéristique, polynomial régression uses the espace des hypothèses $\mathcal{H}_d^{(\text{poly})} := \{h(x) = \sum_{j=0}^{d-1} x^j w_j\}$. The quality of a polynomial hypothèse map is measured using the average squared error loss incurred on a set of labeled datapoints (which we refer to as the ensemble d’entraînement).

predictor A predictor is a real-valued hypothèse 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.

privacy funnel The privacy funnel is a method for learning privacy-friendly caractéristiques of point de données [88].

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 probabilistic model 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 [89]).

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 model A probabilistic modèle interprets point de données as réalisations of VAs with a joint loi de probabilité. This joint loi de probabilité typically involves paramètres which have to be manually chosen or learned via statistical inference methods such as maximum likelihood estimation [17].

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

probability density function (pdf) The probabilité density function $p(x)$ of a real-valued VA $x \in \mathbb{R}$ is a particular representation of its loi de probabilité. If the probabilité density function exists, it can be used to

compute the probabilité that x takes on a value from a (measurable) set $\mathcal{B} \subseteq \mathbb{R}$ via $p(x \in \mathcal{B}) = \int_{\mathcal{B}} p(x') dx'$ [7, Ch. 3]. The probabilité density function of a vector-valued VA $\mathbf{x} \in \mathbb{R}^d$ (if it exists) allows us to compute the probabilité of \mathbf{x} belonging to a (measurable) region \mathcal{R} via $p(\mathbf{x} \in \mathcal{R}) = \int_{\mathcal{R}} p(\mathbf{x}') dx'_1 \dots dx'_d$ [7, Ch. 3].

probabilité We assign a probability value, typically chosen in the interval $[0, 1]$, to each event that might occur in a random experiment [6, 7, 37, 90].

projected gradient descent (projected GD) Consider an MRE-based method that uses a parametrized modèle with parameter space $\mathcal{W} \subseteq \mathbb{R}^d$. Even if the objective function of MRE is lisse, we cannot use basic descente de gradient, as it does not take into account constraints on the optimization variable (i.e., the paramètres du modèle). Projected descente de gradient extends basic descente de gradient to handle constraints on the optimization variable (i.e., the paramètres du modèle). A single iteration of projected descente de gradient consists of first taking a gradient step and then projecting the result back onto the parameter space.

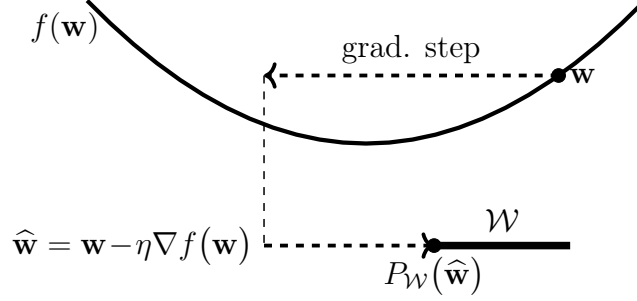


Figure 22: Projected descent de gradient augments a basic gradient step with a projection back onto the constraint set \mathcal{W} .

projection Consider a subset $\mathcal{W} \subseteq \mathbb{R}^d$ of the d -dimensional Euclidean space.

We define the projection $P_{\mathcal{W}}(\mathbf{w})$ of a vector $\mathbf{w} \in \mathbb{R}^d$ onto \mathcal{W} as

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

In other words, $P_{\mathcal{W}}(\mathbf{w})$ is the vector in \mathcal{W} which is closest to \mathbf{w} . The projection is only well-defined for subsets \mathcal{W} for which the above minimum exists [27].

proximable A convex function for which the proximal operator can be computed efficiently is sometimes referred to as proximable or simple [91].

proximal operator Given a convex function $f(\mathbf{w}')$, we define its proximal operator as [59, 92]

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

As illustrated in Figure 23, evaluating the proximal operator amounts to minimizing a penalized variant of $f(\mathbf{w}')$. The penalty term is the scaled

squared Euclidean distance to a given vector \mathbf{w} (which is the input to the proximal operator). The proximal operator can be interpreted as a generalization of the gradient step, which is defined for a lisse convexe function $f(\mathbf{w}')$. Indeed, taking a gradient step with taille de pas η at the current vector \mathbf{w} is the same as applying the proximal operator of the function $\tilde{f}(\mathbf{w}') = (\nabla f(\mathbf{w}))^T (\mathbf{w}' - \mathbf{w})$ and using $\rho = 1/\eta$.

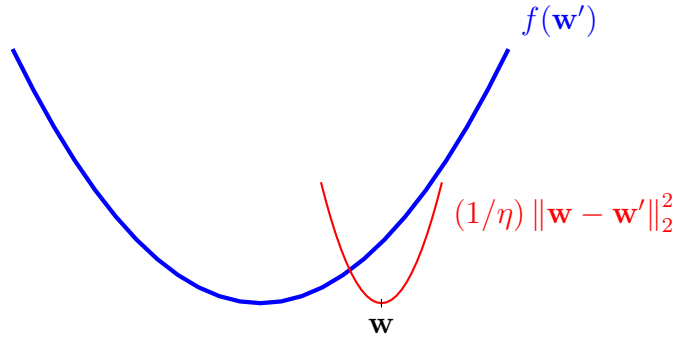


Figure 23: A generalized gradient step updates a vector \mathbf{w} by minimizing a penalized version of the function $f(\cdot)$. The penalty term is the scaled squared Euclidean distance between the optimization variable \mathbf{w}' and the given vector \mathbf{w} .

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 .

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 The regret of a hypothèse h relative to another hypothèse h' , which serves as a baseline, is the difference between the perte incurred by h and the perte incurred by h' [41]. The baseline hypothèse h' is also referred to as an 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 $\widehat{L}(h|\mathcal{D})$ 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),

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

The parameter $\alpha \geq 0$ controls the régularisation strength. For $\alpha = 0$, we recover standard MRE without régularisation. As α increases,

the learned hypoth ese is increasingly biased toward small values of $\mathcal{R}\{h\}$. The component $\alpha\mathcal{R}\{h\}$ in the objective function of (12) can be intuitively understood as a surrogate for the increased average perte that may occur when predicting  tiquettes for point de donn eess outside the ensemble d’entra nement. This intuition can be made precise in various ways. For example, consider a mod le lin aire trained using squared error loss 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 centr e r duites to the vecteur de caract ristiquess in the ensemble d’entra nement [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 [93, Sec. 7.2].

regularized loss minimization (RLM) See RERM.

regularizer A regularizer assigns each hypoth ese h from a espace des hypoth eses \mathcal{H} a quantitative measure $\mathcal{R}\{h\}$ for how much its pr diction error on a ensemble d’entra nement might differ from its pr diction errors on point de donn eess outside the ensemble d’entra nement. Ridge regression uses the regularizer $\mathcal{R}\{h\} := \|\mathbf{w}\|_2^2$ for linear hypoth ese 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 hypoth ese maps $h^{(\mathbf{w})}(\mathbf{x}) := \mathbf{w}^T \mathbf{x}$ [8, Ch. 3].

R nyi divergence The R nyi divergence measures the (dis)similarity between two loi de probabilit es [94].

reward A reward refers to some observed (or measured) quantity that al-

allows us to estimate the perte incurred by the prédiction (or decision) of a hypothèse $h(\mathbf{x})$. For example, in an apprentissage automatique application to self-driving vehicles, $h(\mathbf{x})$ could represent the current steering direction of a vehicle. We could construct a reward from the measurements of a collision sensor that indicate if the vehicle is moving towards an obstacle. We define a low reward for the steering direction $h(\mathbf{x})$ if the vehicle moves dangerously towards an obstacle.

ridge regression Ridge régression learns the poids \mathbf{w} of a linear hypothèse map $h^{(\mathbf{w})}(\mathbf{x}) = \mathbf{w}^T \mathbf{x}$. The quality of a particular choice for the paramètres du modèle \mathbf{w} is measured by the sum of two components. The first component is the average squared error loss incurred by $h^{(\mathbf{w})}$ on a set of labeled datapoints (i.e., the ensemble d'entraînement). The second component is the scaled squared Euclidean norme $\alpha \|\mathbf{w}\|_2^2$ with a régularisation parameter $\alpha > 0$. Adding $\alpha \|\mathbf{w}\|_2^2$ to the average squared error loss is equivalent to replacing each original point de données by the réalisation of (infinitely many) i.i.d. VAs centered around these point de données (see régularisation).

risque Consider a hypothèse h used to predict the étiquette y of a point de données based on its caractéristiques \mathbf{x} . We measure the quality of a particular prédiction using a fonction de perte $L((\mathbf{x}, y), h)$. If we interpret point de données as the réalisations of i.i.d. VAs, also the $L((\mathbf{x}, y), h)$ becomes the réalisation of an VA. The i.i.d. assumption allows us to define the risk of a hypothèse as the expected perte $\mathbb{E}\{L((\mathbf{x}, y), h)\}$. Note that the risk of h depends on both the specific choice for the

fonction de perte and the loi de probabilité of the point de données.

risque empirique The empirical risque $\widehat{L}(h|\mathcal{D})$ of a hypothèse on a jeu de données \mathcal{D} is the average perte incurred by h when applied to the point de données in \mathcal{D} .

robustness TBD

Règlement général sur la protection des données (RGPD) The GDPR was enacted by the European Union (EU), effective from May 25, 2018 [28]. It safeguards the privacy and données rights of individuals in the EU. The GDPR has significant implications for how données is collected, stored, and used in apprentissage automatique applications. Key provisions include the following:

- Data minimization principle: apprentissage automatique systems should only use the necessary amount of personal données for their purpose.
- Transparency and explainability: apprentissage automatique systems should enable their users to understand how the systems make decisions that impact the users.
- Données subject rights: Users should get an opportunity to access, rectify, and delete their personal données, as well as to object to automated decision-making and profiling.
- Accountability: Organizations must ensure robust données security and demonstrate compliance through documentation and regular audits.

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, 37]. Une réalisation de x est tout élément $a' \in \mathcal{N}$ pour lequel il existe un élément $\omega' \in \mathcal{P}$ tel que $x(\omega') = a'$.

réduction de dimension Dimensionality reduction methods map (typically many) raw caractéristiques to a (relatively small) set of new caractéristiques. These methods can be used to visualize point de données by learning two caractéristiques that can be used as the coordinates of a depiction in a scatterplot.

région de décision Consider a hypothèse map h that delivers values from a finite set \mathcal{Y} . For each étiquette value (category) $a \in \mathcal{Y}$, the hypothèse h determines a subset of caractéristique values $\mathbf{x} \in \mathcal{X}$ that result in the same output $h(\mathbf{x}) = a$. We refer to this subset as a decision region of the 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].

régularisation A key challenge of modern apprentissage automatique applications is that they often use large modèles, which have an dimension effective in the order of billions. Training a high-dimensional modèle using basic MRE-based methods is prone to surapprentissage: the learned hypothèse performs well on the ensemble d'entraînement but poorly

outside the ensemble d'entraînement. Regularization refers to modifications of a given instance of MRE in order to avoid surapprentissage, i.e., to ensure that the learned hypothèse performs not much worse outside the ensemble d'entraînement. There are three routes for implementing regularization:

- 1) Modèle pruning: We prune the original modèle \mathcal{H} to obtain a smaller modèle \mathcal{H}' . For a parametric modèle, the pruning can be implemented via constraints on the paramètres du modèle (such as $w_1 \in [0.4, 0.6]$ for the weight of caractéristique x_1 in linear regression).
- 2) Perte penalization: We modify the objective function of MRE by adding a penalty term to the erreur d'entraînement. The penalty term estimates how much larger the expected perte (or risque) is compared to the average perte on the ensemble d'entraînement.
- 3) Data augmentation: We can enlarge the ensemble d'entraînement \mathcal{D} by adding perturbed copies of the original point de données in \mathcal{D} . One example for such a perturbation is to add the réalisation of an VA to the vecteur de caractéristiques of a point de données.

Figure 24 illustrates the above three routes to regularization. These routes are closely related and sometimes fully equivalent: data augmentation using VA normale centrée réduites to perturb the vecteur de caractéristiquess in the ensemble d'entraînement of linear regression has the same effect as adding the penalty $\lambda \|\mathbf{w}\|_2^2$ to the erreur d'entraînement (which is nothing but ridge regression). The decision

on which route to use for regularization can be based on the available computational infrastructure. For example, it might be much easier to implement data augmentation than modèle pruning.

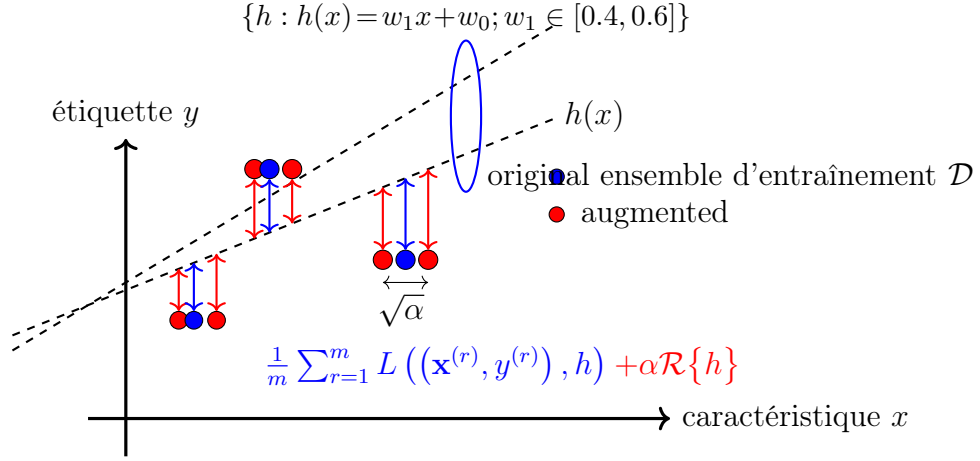


Figure 24: Three approaches to regularization: 1) data augmentation; 2) perte penalization; and 3) modèle pruning (via constraints on paramètres du modèle).

réseau de neurones artificiels (RNA) An ANN is a graphical (signal-flow) representation of a function that maps caractéristiques of a point de données at its input to a prédiction for the corresponding étiquette at its output. The fundamental unit of an ANN is the artificial neuron, which applies an fonction d'activation to its weighted inputs. The outputs of these neurons serve as inputs for other neurons, forming interconnected layers.

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 sample size.

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

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

sample size The number of individual point de données contained in a jeu de données.

scatterplot A visualization technique that depicts point de données by markers in a two-dimensional plane. Fig. 25 depicts an example of a scatterplot.

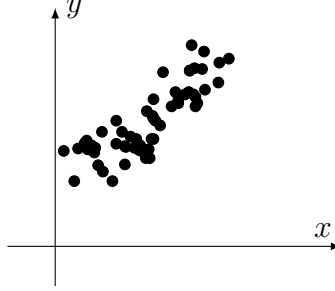


Figure 25: A scatterplot with circle markers, where the point de données represent daily weather conditions in Finland. Each point de données is characterized by its minimum daytime temperature x as the caractéristique and its maximum daytime temperature y as the étiquette. The temperatures have been measured at the FMI weather station Helsinki Kaisaniemi during 1.9.2024 - 28.10.2024.

A scatterplot can enable the visual inspection of points de données that are naturally represented by vecteurs de caractéristiques in high-dimensional spaces.

See also: réduction de dimension.

semi-définie positive A (real-valued) symmetric matrix $\mathbf{Q} = \mathbf{Q}^T \in \mathbb{R}^{d \times d}$ is referred to as psd if $\mathbf{x}^T \mathbf{Q} \mathbf{x} \geq 0$ for every vector $\mathbf{x} \in \mathbb{R}^d$. The property of being psd can be extended from matrices to (real-valued) symmetric kernel maps $K : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ (with $K(\mathbf{x}, \mathbf{x}') = K(\mathbf{x}', \mathbf{x})$) as follows: For any finite set of vecteur de caractéristiquess $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}$, the resulting

matrix $\mathbf{Q} \in \mathbb{R}^{m \times m}$ with entries $Q_{r,r'} = K(\mathbf{x}^{(r)}, \mathbf{x}^{(r')})$ is psd [50].

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 Consider an apprentissage automatique method that uses MRE to learn a hypothèse with the minimum risque empirique on a given ensemble d'entraînement. Such a method is underfitting the ensemble d'entraînement if it is not able to learn a hypothèse with a sufficiently small risque empirique on the ensemble d'entraînement. If a method is underfitting, it will typically also not be able to learn a hypothèse with a small risque.

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 Laplacian matrix $\mathbf{L}^{(\mathcal{G})}$ 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 Euclidean space-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 Euclidean space \mathbb{R}^d (see Figure 26).

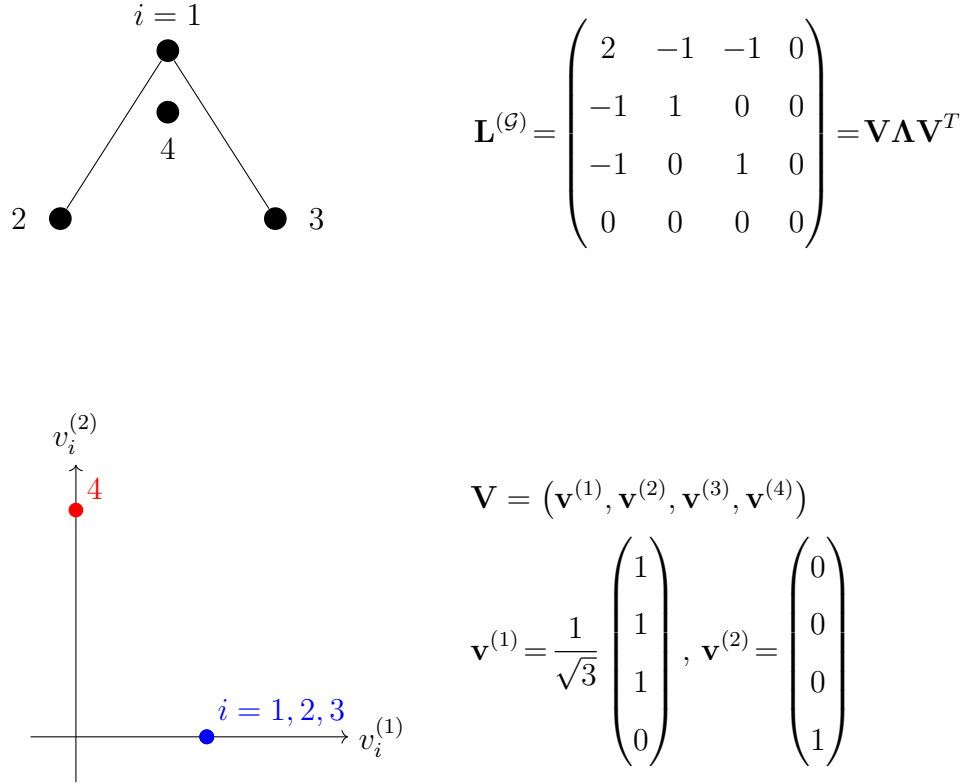


Figure 26: **Top.** Left: An undirected graphe \mathcal{G} with four nodes $i = 1, 2, 3, 4$, each representing a point de données. Right: The Laplacian matrix $\mathbf{L}^{(\mathcal{G})} \in \mathbb{R}^{4 \times 4}$ and its décomposition en éléments propres. **Bottom.** Left: A scatterplot of point de données using the vecteur de caractéristiquess $\mathbf{x}^{(i)} = (v_i^{(1)}, v_i^{(2)})^T$. Right: Two vecteur propres $\mathbf{v}^{(1)}, \mathbf{v}^{(2)} \in \mathbb{R}^d$ corresponding to the valeur propre $\lambda = 0$ of the Laplacian matrix $\mathbf{L}^{(\mathcal{G})}$.

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) [95]. Figure 27 depicts a time signal along with its spectrogram.

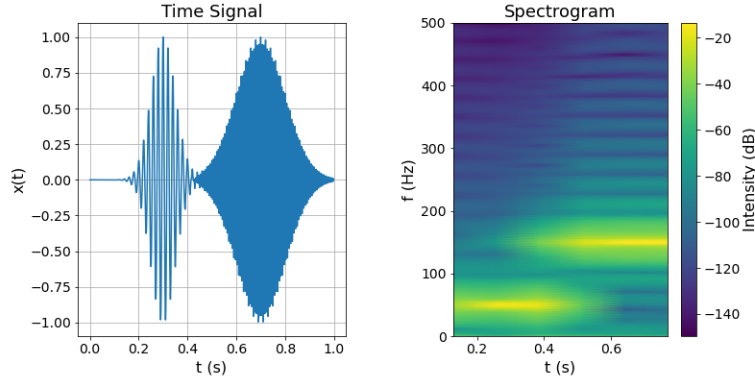


Figure 27: Left: A time signal consisting of two modulated Gaussian pulses. Right: An intensity plot of the spectrogram.

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 deep nets originally developed for image classification and object detection [96]. It is worth noting that, beyond the spectrogram, several alternative representations exist for the time-frequency distribution of signal energy [97, 98].

squared error loss The squared error perte measures the prédiction error of a hypothèse h when predicting a numeric étiquette $y \in \mathbb{R}$ from the

caractéristiques \mathbf{x} of a point de données. It is defined as

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

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 generalization error or risque of the method; see [55, Ch. 13]. To build intuition, consider the three datasets depicted in Fig. 28, 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.

statistical aspects By statistical aspects of an apprentissage automatique method, we refer to (properties of) the loi de probabilité of its output under a probabilistic model for the données fed into the method.

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} [99]. In its most basic variant, the stochastic

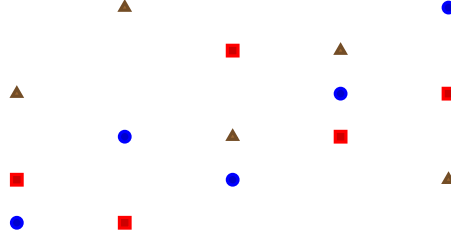


Figure 28: Three jeu de données $\mathcal{D}^{(*)}$, $\mathcal{D}^{(\square)}$, and $\mathcal{D}^{(\triangle)}$, each sampled independently from the same données-generating loi de probabilité. A stable apprentissage automatique method should return similar outputs when trained on any of these jeu de données.

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.

stochastic gradient descent (SGD) Stochastic descent de gradient is obtained from descente de gradient by replacing the gradient of the objective function with a stochastic approximation. A main application of stochastic descent de gradient is to train a parametrized modèle via MRE on a ensemble d'entraînement \mathcal{D} that is either very large or not readily available (e.g., when point de données are stored in a database distributed all over the planet). To evaluate the gradient of the risque

empirique (as a function of the paramètres du modèle \mathbf{w}), we need to compute a sum $\sum_{r=1}^m \nabla_{\mathbf{w}} L(\mathbf{z}^{(r)}, \mathbf{w})$ over all point de données in the ensemble d'entraînement. We obtain a stochastic approximation to the gradient by replacing the sum $\sum_{r=1}^m \nabla_{\mathbf{w}} L(\mathbf{z}^{(r)}, \mathbf{w})$ with a sum $\sum_{r \in \mathcal{B}} \nabla_{\mathbf{w}} L(\mathbf{z}^{(r)}, \mathbf{w})$ over a randomly chosen subset $\mathcal{B} \subseteq \{1, \dots, m\}$ (see Figure 29). We often refer to these randomly chosen point de données as a lot. The lot size $|\mathcal{B}|$ is an important parameter of stochastic descent de gradient. Stochastic descent de gradient with $|\mathcal{B}| > 1$ is referred to as mini-lot stochastic descent de gradient [100].

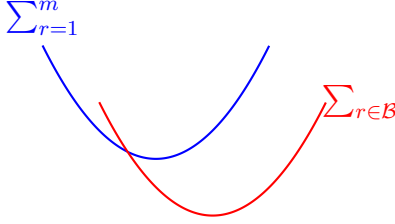


Figure 29: Stochastic descent de gradient for MRE approximates the gradient $\sum_{r=1}^m \nabla_{\mathbf{w}} L(\mathbf{z}^{(r)}, \mathbf{w})$ by replacing the sum over all point de données in the ensemble d'entraînement (indexed by $r = 1, \dots, m$) with a sum over a randomly chosen subset $\mathcal{B} \subseteq \{1, \dots, m\}$.

strongly convex A continuously dérivable real-valued function $f(\mathbf{x})$ is strongly convexe with coefficient σ if $f(\mathbf{y}) \geq f(\mathbf{x}) + \nabla f(\mathbf{x})^T (\mathbf{y} - \mathbf{x}) + (\sigma/2) \|\mathbf{y} - \mathbf{x}\|_2^2$ [76], [78, Sec. B.1.1].

structural risk minimization (SRM) Structural risk minimization (SRM) is an instance of RERM, which 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 [55, 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}' [101, 102].

subgradient descent Subgradient descent is a generalization of descente de gradient that does not require differentiability of the function to be minimized. This generalization 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 objective function. The objective function might be the risque empirique $\hat{L}(h^{(\mathbf{w})}|\mathcal{D})$ 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 linear regression and logistic regression, 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 30, 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 (6) [39, 66, 103].

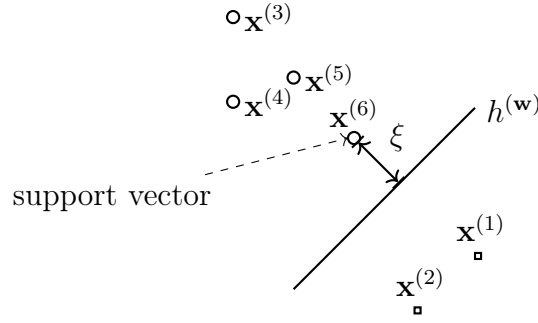


Figure 30: The SVM learns a hypoth se (or classifier) $h^{(w)}$ with minimal average soft-margin hinge loss. Minimizing this perte is equivalent to maximizing the margin ξ between the fronti re de d cision de $h^{(w)}$ and each class of the ensemble d'entra nement.

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

surapprentissage Consider an apprentissage automatique method that uses MRE to learn a hypoth se with the minimum risque empirique on a given ensemble d'entra nement. Such a method is overfitting the ensemble d'entra nement if it learns a hypoth se with a small risque empirique on the ensemble d'entra nement but a significantly larger perte outside the ensemble d'entra nement.

taille de pas See taux d'apprentissage.

taux d'apprentissage Consider an iterative apprentissage automatique method for finding or learning a useful hypothèse $h \in \mathcal{H}$. Such an iterative method repeats similar computational (update) steps that adjust or modify the current hypothèse to obtain an improved hypothèse. One well-known example of such an iterative learning method is descente de gradient and its variants, SGD and projected gradient descent (projected GD). A key parameter of an iterative method is the learning rate. The learning rate controls the extent to which the current hypothèse can be modified during a single iteration. A well-known example of such a parameter is the taille de pas used in descente de gradient [8, Ch. 5].

total variation See GTV.

transparency Transparency is a fundamental requirement for trustworthy AI [104]. In the context of apprentissage automatique methods, transparency is often used interchangeably with explainability [45, 105]. However, in the broader scope of IA systems, transparency extends beyond explainability and includes providing information about the system's limitations, reliability, and intended use. In medical diagnosis systems, transparency requires disclosing the confidence level for the prédictions delivered by a trained modèle. In credit scoring, IA-based lending decisions should be accompanied by explanations of contributing factors, such as income level or credit history. These explanations allow humans (e.g., a loan applicant) to understand and contest automated decisions. Some apprentissage automatique methods inherently offer

transparency. For example, logistic regression provides a quantitative measure of classification reliability through the value $|h(\mathbf{x})|$. Arbres de décisions are another example, as they allow human-readable decision rules [47]. Transparency also requires a clear indication when a user is engaging with an IA system. For example, IA-powered chatbots should notify users that they are interacting with an automated system rather than a human. Furthermore, transparency encompasses comprehensive documentation detailing the purpose and design choices underlying the IA system. For instance, modèle datasheets [72] and IA system cards [106] help practitioners understand the intended use cases and limitations of an IA system [107].

trustworthy artificial intelligence (trustworthy AI) Besides the computational aspects and statistical aspects, a third main design aspect of apprentissage automatique methods is their trustworthiness [108]. The EU has put forward seven key requirements (KRs) for trustworthy IA (that typically build on apprentissage automatique methods) [109]:

- 1) KR1 - Human agency and oversight;
- 2) KR2 - Technical robustness and safety;
- 3) KR3 - Privacy and data governance;
- 4) KR4 - Transparency;
- 5) KR5 - Diversity, non-discrimination and fairness;
- 6) KR6 - Societal and environmental well-being;
- 7) KR7 - Accountability.

tâche d'apprentissage Consider a jeu de données \mathcal{D} constituted by several point de données, each of them characterized by caractéristiques \mathbf{x} . For example, the jeu de données \mathcal{D} might be constituted by the images of a particular database. Sometimes it might be useful to represent a jeu de données \mathcal{D} , along with the choice of caractéristiques, by a loi de probabilité $p(\mathbf{x})$. A learning task associated with \mathcal{D} consists of a specific choice for the étiquette of a point de données and the corresponding espace des étiquettes. Given a choice for the fonction de perte and modèle, a learning task gives rise to an instance of MRE. Thus, we could define a learning task also via an instance of MRE, i.e., via an objective function. Note that, for the same jeu de données, we obtain different learning tasks by using different choices for the caractéristiques and étiquette of a point de données. These learning tasks are related, as they are based on the same jeu de données, and solving them jointly (via apprentissage multitâche methods) is typically preferable over solving them separately [110], [111], [112].

uncertainty Uncertainty refers to the degree of confidence—or lack thereof—associated with a quantity such as a model prediction, parameter estimate, or observed data point. In apprentissage automatique, uncertainty arises from various sources, including noisy data, limited training samples, or ambiguity in model assumptions. Probability theory offers a principled framework for representing and quantifying such uncertainty.

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 reward signal $r^{(a_k)}$. A widely used probabilistic model for this type of sequential decision-making problem is the stochastic multi-armed bandit setting [81]. 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 reward $\mu^{(a)}$, 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 [81].

valeur propre We refer to a number $\lambda \in \mathbb{R}$ as an eigenvalue of a square matrix $\mathbf{A} \in \mathbb{R}^{d \times d}$ if there is a non-zero vector $\mathbf{x} \in \mathbb{R}^d \setminus \{\mathbf{0}\}$ such that $\mathbf{A}\mathbf{x} = \lambda\mathbf{x}$.

validation Consider a hypothèse \hat{h} that has been learned via some apprentissage automatique method, e.g., by solving MRE on a ensemble d’entraînement \mathcal{D} . Validation refers to the practice of evaluating the perte incurred by the hypothèse \hat{h} on a set of point de données that are not contained in the 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 [55]) 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'Euclidean space \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 pdf est donnée par [7], [19], [75]

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

Étant donnée une VA normale centrée réduite x , on peut construire une VA normale x' ayant pour moyenne μ et variance σ^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 [113].

Les VA normales sont largement utilisées comme probabilistic models 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 [36].

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 Euclidean space 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 kernel method).

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 $\mathbf{Ax} = \lambda\mathbf{x}$ with some valeur propre λ .

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 [114]. 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 The neighborhood of a node $i \in \mathcal{V}$ is the subset of nodes constituted by the voisins of i .

voisins The neighbors of a node $i \in \mathcal{V}$ within an FL network are those nodes $i' \in \mathcal{V} \setminus \{i\}$ that are connected (via an edge) to node i .

zero-gradient condition Consider the unconstrained optimization problem $\min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w})$ with a lisse and convexe objective function $f(\mathbf{w})$. 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 classifier $h(\mathbf{x})$ 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 ou au-dessus de sa courbe :

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

Une fonction est convexe si et seulement si son épigraphe est un ensemble convexe [27, 101].

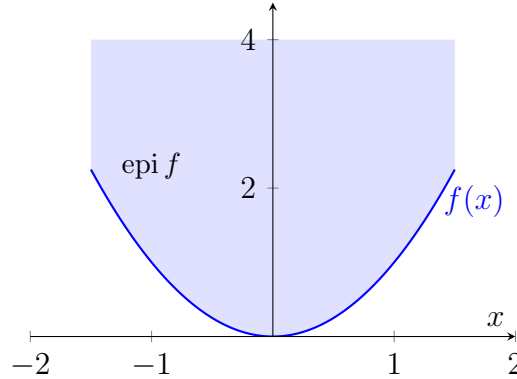


Figure 31: Épigraphe de la fonction $f(x) = x^2$ (zone grisée).

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, 117
- k -fold cross-validation (k -fold CV), 16
- k -moyennes, 16
- absolute error loss, 16
- accuracy, 16
- algebraic connectivity, 17
- algorithme, 17
- analyse en composantes principales (ACP), 18
- application programming interface (API), 18
- apprentissage automatique, 19
- apprentissage fédéré, 19
- apprentissage multitâche, 19
- apprentissage semi-supervisé, 20
- arbre de décision, 20
- autoencoder, 21
- backdoor, 21
- bagging, 22
- baseline, 22
- Bayes estimator, 24
- Bayes risk, 24
- biais, 25
- boosting, 25
- bootstrap, 26
- borne supérieure, 27
- caractéristique, 27
- classification, 27
- classification multi-classe, 27
- classifier, 28
- cluster, 28
- clustered federated learning (CFL), 28
- clustering assumption, 30
- computational aspects, 30
- condition number, 31
- confusion matrix, 31
- connected graph, 31
- convex clustering, 31
- convexe, 32
- Courant–Fischer–Weyl min-max characterization, 32
- critère d’arrêt, 33
- data augmentation, 33
- data minimization principle, 34

data normalization, 34
 data poisoning, 35
 deep net, 35
 degree of belonging, 35
 denial-of-service attack, 35
 density-based spatial clustering of
 applications with noise
 (DBSCAN), 35
 descente de gradient, 36
 differential privacy (DP), 37
 dimension effective, 37
 discrepancy, 38
 dispositif, 38
 distributed algorithm, 38
 données, 39
 décomposition en valeurs
 singulières, 39
 décomposition en éléments propres,
 40
 dérivable, 40
 edge weight, 40
 ensemble d'entraînement (ou
 d'apprentissage), 40
 ensemble de test (ou jeu de test),
 41
 ensemble de validation (ou jeu de
 validation), 41
 erreur d'entraînement, 41
 erreur de validation, 41
 espace des caractéristiques, 42
 espace des étiquettes, 42
 espace probabilisé, 43
 espérance, 43
 estimation error, 44
 Euclidean space, 44
 expectation-maximization (EM),
 44
 expert, 45
 explainability, 45
 explainable empirical risk
 minimization (EERM), 46
 explainable machine learning
 (explainable ML), 46
 explanation, 46
 feature learning, 47
 feature map, 47
 feature matrix, 48
 FedAvg, 48
 federated learning network (FL
 network), 48

- FedGD, 49
- FedProx, 49
- FedRelax, 49
- FedSGD, 49
- Finnish Meteorological Institute
(FMI), 49
- flow-based clustering, 49
- fonction d'activation, 50
- fonction de perte (ou de coût), 50
- frontière de décision, 51
- Gaussian mixture model (GMM),
51
- Gaussian Process (GP), 52
- generalization, 53
- generalized total variation (GTV),
55
- generalized total variation
minimization (GTVMin),
55
- geometric median (GM), 55
- gradient, 56
- gradient step, 56
- grand modèle de langage (GML),
57
- graph clustering, 58
- graphe, 58
- hard clustering, 58
- high-dimensional regime, 59
- Hilbert space, 59
- hinge loss, 59
- histogram, 60
- horizontal federated learning
(HFL), 60
- Huber loss, 61
- Huber regression, 61
- hypothèse, 62
- independent and identically
distributed assumption
(i.i.d. assumption), 62
- indépendantes et identiquement
distribuées (i.i.d.), 62
- intelligence artificielle (IA), 62
- interpretability, 63
- jeu de données, 63
- jeu de données local, 66
- kernel, 66
- kernel method, 66
- Kullback-Leibler divergence (KL
divergence), 67

labeled datapoint, 67	(MRE), 74
Laplacian matrix, 68	minimum, 74
law of large numbers, 68	missing data, 74
least absolute deviation regression, 69	model selection, 75
least absolute shrinkage and selection operator (Lasso), 69	modèle, 75
linear classifier, 69	modèle linéaire, 75
linear regression, 69	moyenne, 76
lisse (ou régulière), 69	multi-armed bandit (MAB), 76
Local Interpretable Model-agnostic Explanations (LIME), 70	mutual information (MI), 77
local model, 71	méthodes basées sur le gradient, 77
logistic loss, 72	nearest neighbor (NN), 77
logistic regression, 72	networked data, 78
loi de probabilité, 72	networked exponential families (nExpFam), 78
loi normale multivariée, 72	networked federated learning (NFL), 78
lot, 73	networked model, 78
matrice de covariance, 73	node degree, 78
maximum, 73	non-smooth, 78
maximum likelihood, 73	norme, 79
mean squared estimation error (MSEE), 74	objective function, 79
minimisation du risque empirique	online algorithm, 79
	online gradient descent (online GD), 80
	online learning, 81

- optimism in the face of uncertainty, 82
- outlier, 83
- parameter space, 84
- paramètres, 84
- paramètres du modèle, 85
- partitionnement de données, 85
- perte, 85
- poids, 86
- point de données, 86
- polynomial regression, 86
- predictor, 87
- privacy funnel, 87
- privacy leakage, 87
- privacy protection, 87
- probabilistic model, 88
- probabilistic principal component analysis (PPCA), 88
- probability density function (pdf), 88
- probabilité, 89
- projected gradient descent (projected GD), 89
- projection, 90
- proximable, 90
- proximal operator, 90
- prédiction, 91
- quadratic function, 92
- Rényi divergence, 93
- random forest, 92
- rectified linear unit (ReLU), 92
- regret, 92
- regularized empirical risk minimization (RERM), 92
- regularized loss minimization (RLM), 93
- regularizer, 93
- reward, 93
- ridge regression, 94
- risque, 94
- risque empirique, 95
- règlement général sur la protection des données (RGPD), 95
- réalisation, 96
- réduction de dimension, 96
- région de décision, 96
- régression, 96
- régularisation, 96
- réseau de neurones artificiels (RNA), 99

sample, 99	surapprentissage, 109
sample covariance matrix, 99	taille de pas, 110
sample mean, 99	taux d'apprentissage, 110
sample size, 99	total variation, 110
scatterplot, 100	transparency, 110
semi-définie positive, 100	trustworthy AI, 111
sensitive attribute, 101	tâche d'apprentissage, 112
similarity graph, 101	
soft clustering, 101	uncertainty, 112
sous-apprentissage, 102	upper confidence bound (UCB),
spectral clustering, 102	112
spectrogram, 104	
squared error loss, 104	valeur propre, 113
stability, 105	validation, 113
statistical aspects, 105	Vapnik–Chervonenkis dimension
stochastic block model (SBM), 105	(VC dimension), 113
stochastic gradient descent (SGD),	variable aléatoire (VA), 114
106	variable aléatoire normale centrée
	réduite, 114
strongly convex, 107	variance, 115
structural risk minimization	vecteur de caractéristiques, 115
(SRM), 107	vecteur propre, 116
subgradient, 108	vertical federated learning (vertical
subgradient descent, 108	FL), 116
support vector machine (SVM),	voisinage, 116
108	voisins, 116

weights, 86

zero-gradient condition, 116

épigraphe, 117

étiquette, 117

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*. New York, NY, USA: Springer Science+Business Media, 2001.
- [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] I. Goodfellow, Y. Bengio, and A. Courville, *Deep Learning*. Cambridge, MA, USA: MIT Press, 2016.
- [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] S. Ross, *A First Course in Probability*, 9th ed. Boston, MA, USA: Pearson Education, 2014.
- [37] P. R. Halmos, *Measure Theory*. New York, NY, USA: Springer-Verlag, 1974.
- [38] S. M. Kay, *Fundamentals of Statistical Signal Processing: Estimation Theory*. Upper Saddle River, NJ, USA: Prentice Hall PTR, 1993.
- [39] C. M. Bishop, *Pattern Recognition and Machine Learning*. New York, NY, USA: Springer Science+Business Media, 2006.
- [40] 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.
- [41] N. Cesa-Bianchi and G. Lugosi, *Prediction, Learning, and Games*. New York, NY, USA: Cambridge Univ. Press, 2006.

- [42] E. Hazan, “Introduction to online convex optimization,” *Found. Trends Optim.*, vol. 2, no. 3–4, pp. 157–325, Aug. 2016, doi: 10.1561/24000000013.
- [43] 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
- [44] 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.
- [45] 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.
- [46] 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>
- [47] 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.

- [48] 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/>
- [49] 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.
- [50] B. Schölkopf and A. J. Smola, *Learning with Kernels: Support Vector Machines, Regularization, Optimization, and Beyond*. Cambridge, MA, USA: MIT Press, 2002.
- [51] 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.
- [52] 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>
- [53] 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
- [54] 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.
 - [55] S. Shalev-Shwartz and S. Ben-David, *Understanding Machine Learning: From Theory to Algorithms*. New York, NY, USA: Cambridge Univ. Press, 2014.
 - [56] 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.
 - [57] 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.
 - [58] 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.
 - [59] N. Parikh and S. Boyd, “Proximal algorithms,” *Found. Trends Optim.*, vol. 1, no. 3, pp. 127–239, Jan. 2014, doi: 10.1561/24000000003.

- [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
- [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] 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/>
- [69] S. Abiteboul, R. Hull, and V. Vianu, *Foundations of Databases*. Reading, MA, USA: Addison-Wesley Publishing Company, 1995.
- [70] S. Hoberman, *Data Modeling Made Simple: A Practical Guide for Business and IT Professionals*, 2nd ed. Basking Ridge, NJ, USA: Technics Publications, 2009.
- [71] R. Ramakrishnan and J. Gehrke, *Database Management Systems*, 3rd ed. New York, NY, USA: McGraw-Hill, 2002.
- [72] T. Gebru et al., “Datasheets for datasets,” *Commun. ACM*, vol. 64, no. 12, pp. 86–92, Nov. 2021, doi: 10.1145/3458723.
- [73] T. M. Cover and J. A. Thomas, *Elements of Information Theory*, 2nd ed. Hoboken, NJ, USA: Wiley, 2006.
- [74] 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
- [75] A. Papoulis and S. U. Pillai, *Probability, Random Variables, and Stochastic Processes*, 4th ed. New York, NY, USA: McGraw-Hill Higher Education, 2002.

- [76] Y. Nesterov, *Introductory Lectures on Convex Optimization: A Basic Course*. Boston, MA, USA: Kluwer Academic Publishers, 2004.
- [77] S. Bubeck, “Convex optimization: Algorithms and complexity,” *Found. Trends Mach. Learn.*, vol. 8, no. 3–4, pp. 231–357, Nov. 2015, 10.1561/22000000050.
- [78] D. P. Bertsekas, *Convex Optimization Algorithms*. Belmont, MA, USA: Athena Scientific, 2015.
- [79] A. Lapidoth, *A Foundation in Digital Communication*. Cambridge, U.K.: Cambridge Univ. Press, 2009.
- [80] 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.
- [81] 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.
- [82] 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.
- [83] R. A. Horn and C. R. Johnson, *Matrix Analysis*, 2nd ed. New York, NY, USA: Cambridge Univ. Press, 2013.
- [84] 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>
- [85] P. J. Brockwell and R. A. Davis, *Time Series: Theory and Methods*, 2nd ed. New York, NY, USA: Springer-Verlag, 1991.
 - [86] 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.
 - [87] 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.
 - [88] 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.
 - [89] 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.
 - [90] O. Kallenberg, *Foundations of Modern Probability*. New York, NY, USA: Springer-Verlag, 1997.
 - [91] 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.

- [92] 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.
- [93] S. Shalev-Shwartz and A. Tewari, “Stochastic methods for ℓ_1 regularized loss minimization,” in *Proc. 26th Annu. Int. Conf. Mach. Learn.*, L. Bottou and M. Littman, Eds. Jun. 2009, pp. 929–936.
- [94] 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.
- [95] L. Cohen, *Time-Frequency Analysis*. Upper Saddle River, NJ, USA: Prentice Hall PTR, 1995.
- [96] 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.
- [97] B. Boashash, Ed. *Time Frequency Signal Analysis and Processing: A Comprehensive Reference*. Oxford, U.K.: Elsevier, 2003.
- [98] S. Mallat, *A Wavelet Tour of Signal Processing: The Sparse Way*, 3rd ed. Burlington, MA, USA: Academic, 2009.
- [99] 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>

- [100] 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.
- [101] D. P. Bertsekas, A. Nedic, and A. E. Ozdaglar, *Convex Analysis and Optimization*. Belmont, MA, USA: Athena Scientific, 2003.
- [102] D. P. Bertsekas, *Nonlinear Programming*, 2nd ed. Belmont, MA, USA: Athena Scientific, 1999.
- [103] 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.
- [104] 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>
- [105] C. Gallese, “The AI act proposal: A new right to technical interpretability?,” *SSRN Electron. J.*, Feb. 2023. [Online]. Available: <https://ssrn.com/abstract=4398206>
- [106] M. Mitchell et al., “Model cards for model reporting,” in *Proc. Conf. Fairness, Accountability, Transparency*, 2019, pp. 220–229, doi: 10.1145/3287560.3287596.
- [107] 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.
- [108] 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>
- [109] 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>
- [110] R. Caruana, “Multitask learning,” *Mach. Learn.*, vol. 28, pp. 41–75, Jul. 1997, doi: 10.1023/A:1007379606734.
- [111] 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.
- [112] 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.
- [113] C. E. Rasmussen and C. K. I. Williams, *Gaussian Processes for Machine Learning*. Cambridge, MA, USA: MIT Press, 2006.
- [114] 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.