

I. Introduction

→ predictive maintenance, mobile health monitoring, drug discovery, recommendation Σ , object recognition

ChatGPT biased wtf qui s'y attendait ?!?

En gros ya du machine learning un peu partout

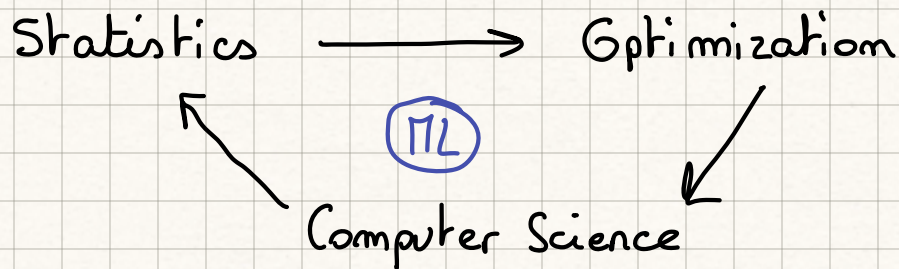
Def: type d'IA qui fait une tâche sans avoir été explicitement programmé pour.

Def by Tom Mitchell: learn from E with respect to class of task T and performance measure P

E = experience : data provided

T = task

P = accuracy on new data, ability to generalize



Supervised machine learning

- Predictive model: approximate a target function
- Conditional generative modeling approximate a target conditional distribution.

Unsupervised machine learning

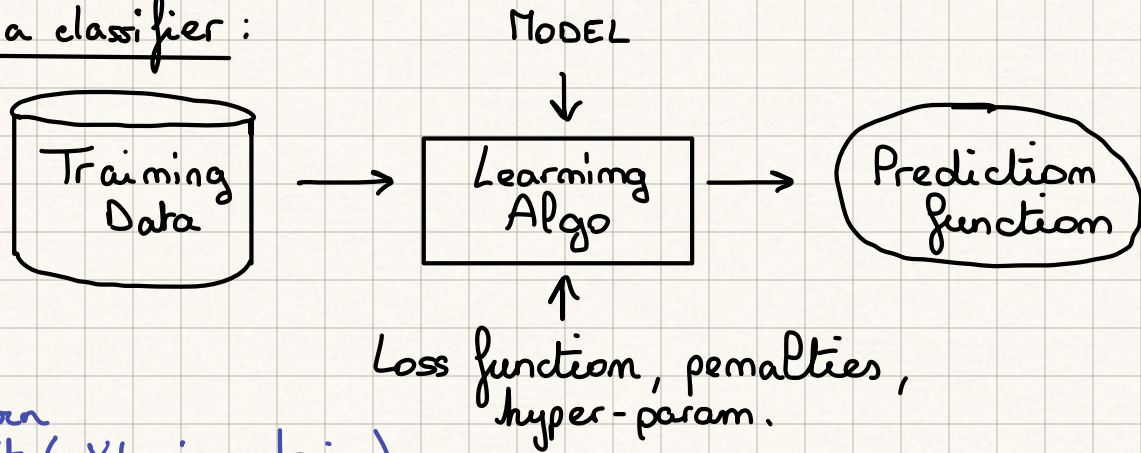
- Generative modeling (Gen AI): approximate a target dist.
- Clustering, Representation Learning, Dimension reduction

Learning paradigms

Customization learning VS Task-driven learning

II. Introduction to Supervised Learning with hands

Learning a classifier :



in sklearn
`clf.fit(Xtrain, ytrain)`

What do we need ?

- Data representation
- Output
- Hypothesis space
- Learning algorithm
- evaluation

III. Probabilistic and statistical setting of Supervised Learning

Risk of a predictive model :

- local loss function : $\ell : \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}$
- $R(g) = \mathbb{E}_{(x,y) \sim p} [\ell(y, g(x))]$

Min of true risk :

$$\operatorname{argmin}_{g: X \rightarrow Y} R(g)$$

Binary classification rule

$$P(g(x), y) = \mathbb{1}_{y \neq g(x)} \Rightarrow \arg \min_{g: X \rightarrow Y} P(Y \neq g(X))$$

Rappel: Bayes Rules

$$P(Y = k | x) = \frac{\overset{\text{likelihood}}{p(x | Y = k)} \overset{\text{prior probability}}{P(Y = k)}}{\underbrace{p(x | Y = -1)P(Y = -1) + p(x | Y = 1)P(Y = 1)}_{\text{proba density}}}$$

↑
posterior proba

Bayes classifier:

$$g_{\text{bayes}}(x) = \underset{\substack{\uparrow \\ \text{on dit que ça vaut } 1 \text{ ou } -1}}{\mathbb{1}}(P[Y = 1 | X = x] > 0,5)$$

bayes classifier achieves the minimal risk for the classification loss

!! bayes risk is characteristic of the "complexity" of the joint probability distribution P and the loss

Sum up: Target function in

- supervised classification: Bayes classifier for the 0-1 loss
- regression: $h(x) = E[Y | x]$ for the square loss

en général, target funcⁿ depend bcp de la loss

Statistical supervised learning problem

find a classifier (regressor) in \mathcal{G} that minimizes

$$R(g) = E_{(x,y) \sim P} [P(g(x), y)]$$

only on a finite training sample $S_m = \{(x_i, y_i)_{i=1}^m\}$

We look for $g^* \in \arg \min_{g: X \rightarrow Y} R(g)$ by providing an estimate $\hat{h} \in \mathcal{G}$ of g^* from S_m

IV - Minimization of the empirical risk

Empirical risk : $R_m(g) = \frac{1}{m} \sum_{i=1}^m \ell(g(x_i), y_i)$

Statistical Learning by empirical risk minimization :

$$\hat{g} \in \arg \min_{g \in \mathcal{G}} R_m(g)$$

excess risk : $R(\hat{g}) - R^* = \underbrace{R(\hat{g}) - \inf_{g \in \mathcal{G}} R(g)}_{\text{estimation error}} + \underbrace{\inf_{g \in \mathcal{G}} R(g) - R^*}_{\text{approximation error}}$
 $R^* = R(g^*)$

Risk convexification :

Pb : empirical risk hard to minimize
 \Rightarrow on prend $u \mapsto \ell(u, v)$ convexe, différentiable
une autre loss function

Misclassification : $\ell(g(x), y) = \mathbb{1}(y g(x) < 0)$

Exponential : $\ell(g(x), y) = e^{-y g(x)}$

Bimomial : $\ell(g(x), y) = -\log(1 + e^{-2y g(x)})$

Squared error : $\ell(g(x), y) = (1 - y g(x))^2$

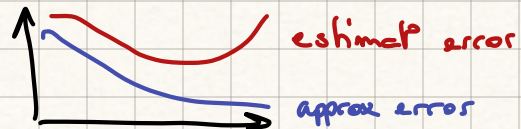
Summary: For (X, Y) et ℓ on veut un classifieur proche de $g^* = \operatorname{argmin} E(\ell(g(X), Y))$

Strategie: training sample of $(X, Y) \rightarrow$ on minimise le risque empirique

Method: Numerical optimization (ex: descente de grad)

II - Relevance of Empirical risk minimization

Compromise biais / variance:



- if model too simple \rightarrow large biais, no universality \Rightarrow Underfitting
- if model too complex \rightarrow large variance, no consistency \Rightarrow Overfitting

Vapnik and Chervomenkis' s result

$$\forall P, S_n \text{ drawn from } P, \forall h \in \mathcal{H}, R(h) \leq R_n(h) + \underbrace{B(d, n)}_{\substack{\text{biais} \\ \text{measure of complexity of } \mathcal{H}}}$$

Theorem: Soit \mathcal{H} une famille de fonctions prenant des valeurs dans $\{-1, 1\}$ de VC-dim d_{VC}

Alors $\forall \delta > 0, \forall h \in \mathcal{H}$ avec proba $1 - \delta$

$$R(h) \leq R_n(h) + \sqrt{\frac{8 d_{VC} (\ln \frac{2n}{d_{VC}} + 1) + 8 \log(\frac{4}{\delta})}{n}}$$

Idée: on veut contrôler la complexité de \mathcal{H} et réduire l'erreur empirique

\Rightarrow on remplace la minimisation du risk empirique par celle du structural risk

Shattering: \mathcal{H} is said to shatter a set of data points if, \forall 2^n possible assignments of binary labels to those points, $\exists h \in \mathcal{H}$ tq h me fait 0 erreur de prédiction

VC-dimension: size of the largest set that can be fully shattered by \mathcal{H}

$$d_{VC}(\mathcal{H}) = \max \{ m : \exists (x_1, \dots, x_m) \in \mathcal{X}^m \text{ shattered by } \mathcal{H} \}$$

mb: si $d_{VC} = d$ ça veut dire que ya un set de taille d pas que tous les sets de taille d ou \leq sont shatterables.

VC-dimension of hyperplanes

$$d_{VC}(\mathcal{H}_d) = d + 1$$

\uparrow
hyperplanes in \mathbb{R}^d

Regularisation: $+ \lambda \Omega(h)$

\uparrow role = control of the model complexity

\rightarrow imposition of some prior knowledge

$$\arg \min_{h \in \mathcal{H}} \frac{1}{n} \sum_i \ell(y_i, h(x_i)) + \lambda \Omega(h)$$

\uparrow reg. parameter (hyperparameter) \uparrow complexity penalty