

ANNEE UNIVERSITAIRE 20__ / 20__

Niveau d'études (Ex : 2^{ème} année) :

Épreuve de :

Note de
l'épreuve

(1) Nom :
Prénoms :
Né(e) à :
le :

Il est interdit au candidat de signer sa copie ou d'y inscrire un signe quelconque pouvant en indiquer la provenance.

Ne pas mettre de colle sur la partie griseée

Ouvrir ici ▲

Nombre d'intercalaires ____ A ____ , le 16/01 20 24.

lolo@univ.fr

APPRENTISSAGE STATISTIQUE

Première

Apprentissage Statistique = Machine Learning (ML)
= Data Science

[Historique]

* 1930 - 1970 : Stat Inférentielle } stat
- Experiences plurifiles } paramétrique
- Test de modèles.

* 1970 - 1980 : 1ère période des stats non paramétriques
- 1^{ère} apparition des réseaux de neurones

* 1990 : - Jeux de données non plurifiles.
- Apprentissage des ML

* 2000 : Explosion du nombre de variables
(p = nbre de variable $\gg n$ = nombre de données)

- Sélection d'un modèle
- Sélection de dimension.

* 2010 : Explosion des ressources en données.

- Explosion des ressources en calcul
- Réseau des réseaux de neurones
- Théorie de l'optimisation

II/ Estimation VS Apprentissage

* Estimation : On cherche à approcher au mieux le "vrai" modèle.

• Si $y = f(x) + \epsilon \xrightarrow{\text{bruit}}$

On cherche \hat{f} très proche de f .

* Apprentissage : On cherche à prédire y .

$$\bullet \quad y = f(x) + \varepsilon.$$

On veut \hat{y} très proche de y .

II/ Apprentissage Supervisé VS non supervisé.

* Supervisé Soit le couple (X, Y) tel :

$$y = f(x) + \varepsilon$$

On observe un échantillon iid, $(x_i, y_i)_{i=1}^n$

On veut prédire y à partir de X

2 cas possibles.

Regression

$$Y \in \mathbb{R} (\text{rd})$$

Classification

$$Y \in \mathbb{N} (\text{z})$$

* Non Supervisé : On cherche à "comprendre" x à partir d'un échantillon $(x_i, y_i)_{i=1}^n$ iid.

- Estimation de densité
- Clustering
- Recherche de lien causal

APPRENTISSAGE

SUPERVISE

Soit (X, Y) un couple de n-a

* But: prédire Y à partir de X .

* Matériel: des observations $\{(X_i, Y_i)\}_{i=1}^n$

Choix possible de X et Y

X peut être	Y peut être
numérique	numérique (R): Régression
qualitatif	
mix des 2	qualitatif (ou discret): Classification

I-i) Le problème basique

À partir des observations $\{(X_i, Y_i)\}_{i=1}^n$,

On veut construire ~~un~~ un prédicteur (machine)

\hat{y}^n tel que $\hat{y}_{\text{new}} = \hat{y}^n(x_{\text{new}})$.

On apprend \hat{y}^n à partir des observations

Propriétés souhaitées

• Précision sur les données.

⇒ Les prédictions faites sur les données d'apprentissage doivent être proches des vraies valeurs (connues).

②

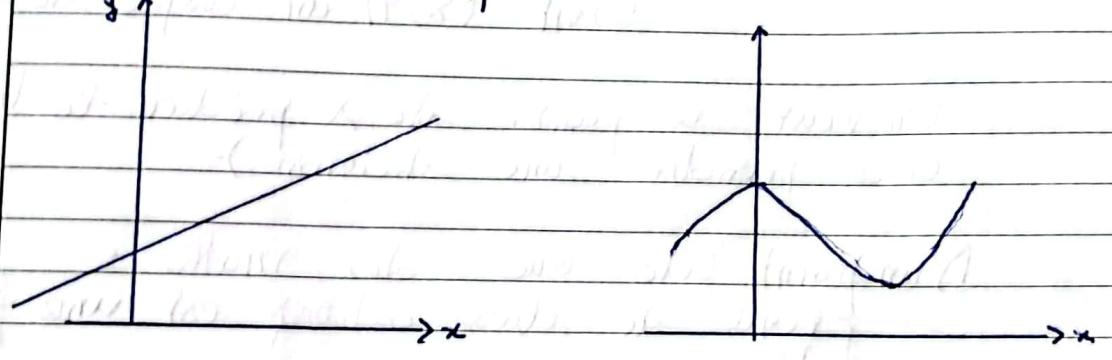
16.01.24

* Généralisabilité

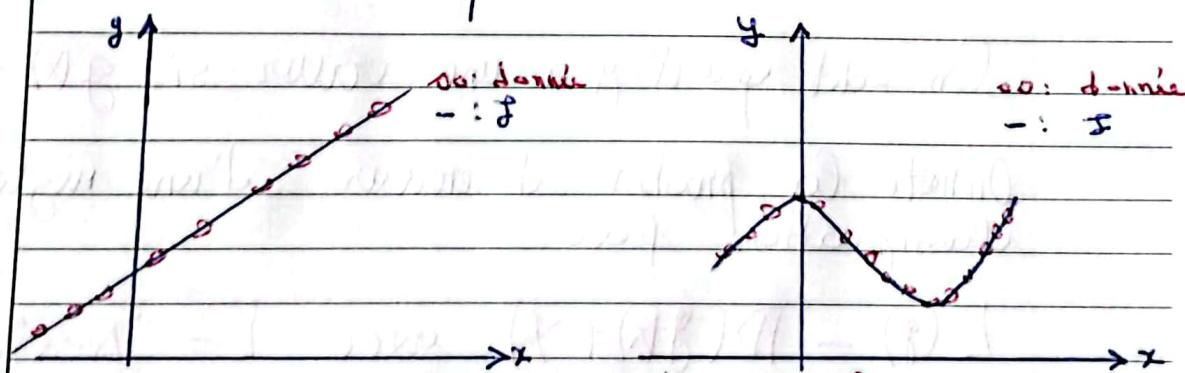
⇒ Les prédictions faites sur de nouvelles données doivent aussi être proches des vraies valeurs.

I. 2) Sous apprentissage et Sur apprentissage

- Fonction à apprendre

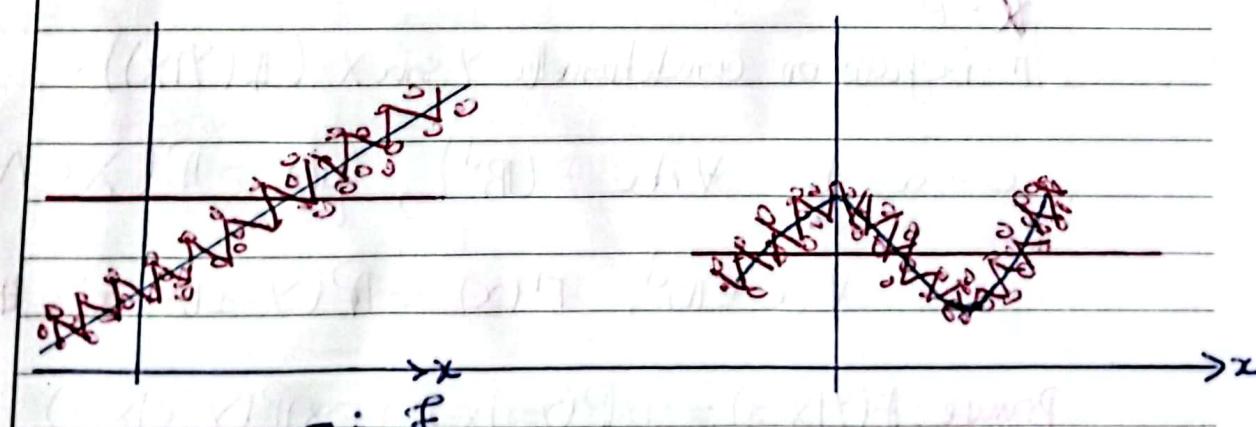


- Des données qu'on aimeraient avoir



NB: on a jamais cela en pratique.

- Des données qu'on a:



NN: Sur-apprentissage
—: sous-apprentissage

Pour avoir de bonnes performances

- * Avoir le meilleur algo (+)
- * Avoir de meilleures variables (+++)
- * Avoir de meilleures données (++++)

Classification Binaire

I) Généralités

Soit (x, y) un couple de $x \in \mathbb{R}^d$

On veut à partir de x prédire le label y (i.e prendre une décision).

D'un point de vue des maths:
→ la règle de classification est une fonction
une décision

$$g: \mathbb{R}^d \rightarrow \{0, 1\}.$$

On dit qu'il y a une erreur si $g(x) \neq y$.

On note la proba d'erreur d'une règle de classification par:

$$L(g) = P(g(x) \neq y) \text{ avec } L = \text{"Loss".}$$

On note (p, r) la loi du couple (x, y)

$\star \sim p$
région conditionnelle y sur x ($E(y|x)$).

$$\text{c.-à-d: } \forall A \in \mathcal{B}(\mathbb{R}^d), p(A) = P(x \in A)$$

$$\forall x \in \mathbb{R}^d, r(x) = P(y=1|x=x) = E(y|x=x)$$

$$\text{Prouve: } E(y|x=x) = 1 \cdot P(y=1|x=x) + 0 \cdot P(y=0|x=x)$$

③ 16.02.24

Question ? Existe-t-il une règle meilleure que toutes les autres ? oui

II Règle de Bayes (Bonne règle)

On appelle règle de Bayes, la fonction g^* définie par :

$$g^*(x) = \begin{cases} 1, & \text{si } P(Y=1|x=x) > P(Y=0|x=x) \\ 0, & \text{sinon} \end{cases} \quad (H(x) > 1/2)$$

Théorème :

$\forall g : \mathbb{R}^d \rightarrow \{0,1\}$ règle de décision

$$P(g^*(x) \neq y) \leq P(g(x) \neq y)$$

$$\text{avec } L^* = L(g^*) \text{ , } L^* = \inf_g L(g).$$

Preuve:

Soit g une règle de décision

$$P(g(x) \neq y | x) = 1 - P(g(x) = y | x)$$

$$= 1 - (P(Y=1, g(x)=1 | x) + P(Y=0, g(x)=0 | x))$$

$$= 1 - \left(\underset{g(x)=1}{P(Y=1|x)} + \underset{g(x)=0}{P(Y=0|x)} \right)$$

$$\text{or } H(x) = P(Y=1|x) \text{ et } P(Y=0|x) = 1 - H(x).$$

$$\text{D'où : } = 1 - \left(\underset{g(x)=1}{H(x)} + \underset{g(x)=0}{1-H(x)} \right) = 1 - (1 - H(x)).$$

$$\text{D'autre part, } P(g(x) \neq y | x) - P(g^*(x) \neq y | x)$$

$$= H(x) \left(\underset{g^*(x)=1}{1} - \underset{g(x)=1}{1} \right) + (1 - H(x)) \left(\underset{g^*(x)=0}{1} - \underset{g(x)=0}{1} \right)$$

$$\text{Donc : } = \underbrace{(2H(x)-1)}_A \underbrace{\left(\underset{g^*(x)=1}{1} - \underset{g(x)=1}{1} \right)}_B + \underbrace{(2(1-H(x))-1)}_{1-A} \underbrace{\left(\underset{g^*(x)=0}{1} - \underset{g(x)=0}{1} \right)}_B$$

2 cas possibles

$$* A > 0 \Leftrightarrow g^*(x) = 1$$

$$\Leftrightarrow \pi(x) > \frac{1}{2}$$

$$\Rightarrow B = 0 \text{ ou } 1 \Rightarrow AXB > 0$$

$$* A \leq 0 \Leftrightarrow g^*(x) = 0$$

$$\Rightarrow B = 0 \text{ ou } -1 \Rightarrow AXB > 0.$$

Dans tous les cas,

$$P(g(x) + y | x) - P(g^*(x) + y | x) > 0$$

$$\text{donc: } E[P(g(x) + y | x) - P(g^*(x) + y | x)] > 0$$

$$E(P(g(x) + y | x)) > E(P(g^*(x) + y | x))$$

$$P(g(x) + y) \geq P(g^*(x) + y)$$

D'où le résultat.

Rq: On note L^* le risque de Bayes

$$L^* = E(\min(\pi(0), 1 - \pi(0)))$$

Problème : g^* dépend de la loi de (X, Y) qui est inconnue.

Et la place on dispose d'un échantillon

$$S_n = \{(X_i, Y_i)\}_{i=1}^n \text{ iid de même loi que } (X, Y)$$

On va chercher une règle $g_n(x, S_n)$ pour estimer Y à partir de X et S_n .

On note $g_n(\omega) = g_n(x, S_n)$

③ 16.01.24

Le processus de création de g_n est appelé apprentissage supervisé.

La proba d'erreur :

$$L(g_n) = P(g_n(x) \neq g_n(x, s_n) | s_n),$$

On va trouver g_n qui minimise L .

II | Consistance faible/universelle/lente

On veut approcher L^* ,

On voudrait que d'une certaine façon, $L(g_n) \xrightarrow{n \rightarrow +\infty} L^*$.

On parle de règle de consistance

Def 1: Consistance faible/lente

Etant donnée une distribution de couple (X, Y) , une règle g_n est dite faiblement consistante/convergente

* Si $E(L(g_n)) = P(g_n(x) \neq Y) \xrightarrow{n \rightarrow +\infty} L^*$

elle est fortement consistante

* Si $L(g_n) \xrightarrow{n \rightarrow +\infty} L^*$ p.s

① Consistance forte \Rightarrow Consistance faible

② Consistance faible \Leftrightarrow Consistance en probab.

⚠ Ces membres sont relatifs à
distribution de (X, Y) fixée.

Def 2: Consistance universelle

g_n est universellement (fortement)

consistante si elle l'est pour toute
distribution de (X, Y)

Problème: Comment construire des
règles vérifiant ces
propriétés?

II. 1) Outils probabilistes

* Théorème de Borel-Cantelli (adapté)

Soit X une N.a., X_n une suite de N.a. Si $\forall \varepsilon > 0$,

$$\sum_{i=1}^n P(|X_i - X| > \varepsilon) < \infty \Rightarrow X_n \xrightarrow{\text{P.s}} X$$

Ainsi, pour avoir une consistante d'une règle g_n , on va chercher à avoir des inégalités du type

$$P(|\hat{L}_n - L^*| > \varepsilon) \leq k e^{-nC\varepsilon^2}$$

(k, C des cte > 0)

II.3) Méthode Plug-in

Soit (X, Y) couple de v.a, μ loi de X

Γ : régression conditionnelle de $Y|X$.

$$g^*(x) = \begin{cases} 1 & \text{si } \pi(x) > 1/2 \\ 0 & \text{sinon} \end{cases}$$

L'idée du plug-in, c'est de remplacer $\hat{M}(x)$ pour un estimateur $\hat{r}_n(x)$.

$$g_n(x) = \begin{cases} 1 & \text{si } \hat{r}_n(x) > 1/2 \\ 0 & \text{sinon} \end{cases}$$

L'espérance c'est que $\pi_i(r_i(x))$ est proche de $r_i(x)$, alors g_n sera proche de g et \hat{L}_n sera proche de L^*

*Théo:

Soit \hat{r}_n un estimateur de r , g_n la règle associée, alors

$$0 \leq g_n - L^* \leq 2 \int_{\mathbb{R}^d} |\hat{r}_n(x) - r(x)| \rho dx$$

Rq: La majoration par $2 \int_{\mathbb{R}^d} |\hat{r}_n(x) - r(x)| \rho dx$ n'est pas optimale.

Elle suffit pour montrer la consistante mais ne permet pas d'étudier les vitesses de convergences.

En fait, on a pas vraiment besoin que \hat{r}_n soit très proche de r sur tout \mathbb{R}^d , ce qui est critique c'est que :

$$\left\{ \begin{array}{l} \hat{r}_n(x) > \frac{1}{2} \text{ si } r(x) > \frac{1}{2} \\ \hat{r}_n(x) \leq \frac{1}{2} \text{ si } r(x) \leq \frac{1}{2}. \end{array} \right.$$

La zone critique est la zone où $r(x) \approx \frac{1}{2}$.

Une hypothèse utile est souvent

$$P(|r(x) - \frac{1}{2}| \leq t) \leq C t^\alpha, \forall 0 < t \leq t^*,$$

avec $d = \text{Cte} > 0$.

II.3) Théo de Stone (1977)

Le théo établit la consistante de certains estimateurs \hat{r}_n de r .

On considère des estimateurs de type "moyenne locale", i.e. qui aéviene

$$\hat{r}_n(x) = \sum_{i=1}^n w_{n,i}(x) Y_i, x \in \mathbb{R}^d \text{ avec}$$

$(W_{n,1}, \dots, W_{n,n})$ un vecteur de poids tel que chaque $W_{n,i}^{(x)}$ est une fonction mesurable de x et (X_1, \dots, X_n) , mais pas de Y_1, \dots, Y_n .

Chp.

intro

Les paires (X_i, Y_i) dont X_i est proche de x sont plus intéressantes que les autres.

Les poids devraient être plus importants pour les Y dont les X_i associés sont proches de x .

Ex: Estimateur à noyau naïf

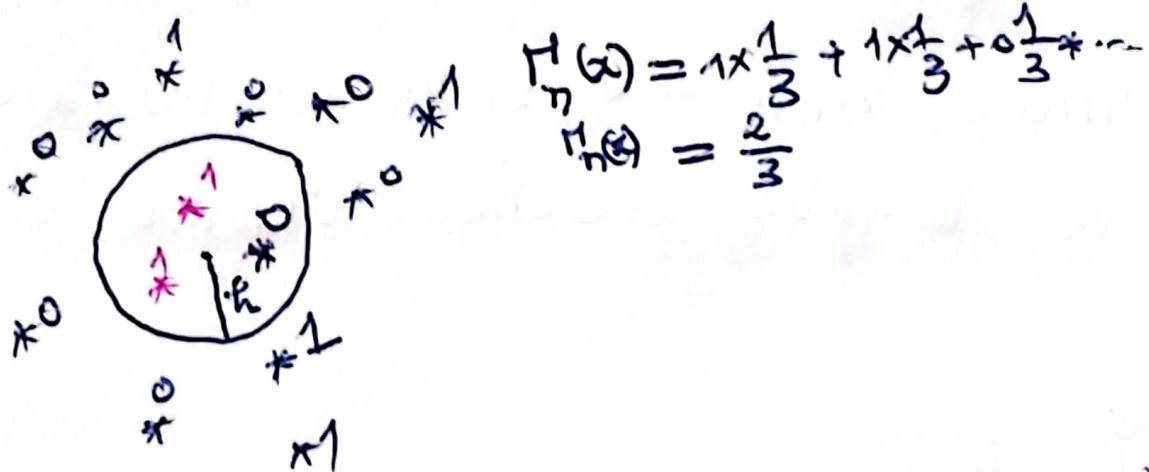
$$\hat{f}_n(x) = \frac{\sum_{i=1}^n W_{n,i}(x - x_i) y_i}{\sum_{i=1}^n W_{n,i}(x - x_i)} \quad ; \quad h > 0$$

$h = \text{"fenêtre"}$

$$w_{n,i} = \frac{\mathbb{I}_{\{||x-x_i|| \leq h\}}}{\sum_{i=1}^n \mathbb{I}_{\{||x-x_i|| \leq h\}}} = \begin{cases} \frac{1}{\# B(x, h)} & \text{si } x_i \in B(x, h) \\ 0 & \text{sinon} \end{cases}$$

En moyenne les y_i dont les x_i sont dans une boule de centre x et de rayon h .

Ex:



* Estimation à noyau (Nadaraya-Watson)

$$w_{n,i} = \frac{k \left(\frac{x - x_i}{h} \right)}{\sum_{i=1}^n k \left(\frac{x - x_i}{h} \right)}$$

$$H_n(x) = \frac{\sum_{i=1}^n y_i k \left(\frac{x - x_i}{h} \right)}{\sum_{i=1}^n k \left(\frac{x - x_i}{h} \right)}$$

où k = un noyau (une densité de proba)

Ex: Dans \mathbb{R} , k est un noyau Gaussien.

* Estimation des K plus proches voisins

- On fait la moyenne des \hat{y}_i dont les x_i sont les plus proches de x .
- On peut prendre des poids $\rightarrow \theta$ en fonction des distances $\|x - x_i\|$, mais $w_{n,i} = 0$ si on est pas dans les K -plus proches voisins.

23.01.24

- P) On considère les hypothèses suivantes
- i) $\exists c_1 > 0$ tq pour toute fonction $f: \mathbb{R}^d \rightarrow \mathbb{R}$ (borélienne) avec $E(\|f(x)\|) < \infty$.
 $E\left(\sum |w_{n,i}(x)| \cdot \|f(x)\|\right) \leq c_1 E(f(x)), \forall n \geq 1$.
 - ii) $\exists c_2 > 0$ tq $P\left(\sum |w_{n,i}(x)| \leq c_2\right) = 1, \forall n \geq 1$.
 - iii) $\forall \epsilon > 0, \sum |w_{n,i}(x)| \underset{\|x_i - x\| > \epsilon \sqrt{n \rightarrow \infty}}{\overset{P}{\rightarrow}} 0$.
 - iv) $\sum w_{n,i}(x) \underset{n \rightarrow \infty}{\overset{P}{\rightarrow}} 1$.
 - v) $\max_{1 \leq i \leq n} |w_{n,i}(x)| \underset{n \rightarrow \infty}{\overset{P}{\rightarrow}} 0$.

- Si i) - iv) sont vérifiés pour toute distribution de X , alors $\hat{\mu}_n$ est universellement LP consistant ($\forall p \geq 1$).
i.e. $\forall p \geq 1$, $E(|\hat{\mu}_n(X) - \mu(X)|^p) \rightarrow 0$

Rq: C'est presque une condition nécessaire et suffisante.

Discussions des hypothèses

- i) hypothèse technique pour éviter de disposer que μ est continue.
- ii) et iii) disent que la somme des poids est bornee et tend $P_0 < 1$.
- iv) \Rightarrow que le poids global des valeurs hors d'un voisinage de $X \rightarrow 0$. Asymptotiquement, l'estimation ne dépend que des données proba du lui-même (moyenne local)

v) \Rightarrow que tous les poids deviennent petits quand n grandit
 \Rightarrow une donnée ne peut décider seule)

Cas classique

$(w_{n,1}(x), \dots, w_{n,n}(x))$ est un vecteur de \mathbb{P} ,
 (i.e. $w_{n,i} \geq 0, \sum w_{n,i} = 1$) .

ii) et iv) sont vérifiées.

Les autres sont une condition nécessaire et suffisante pour la convergence.

- On va utiliser ce théo pour notre problème de classification

Rappel: $\hat{g}(x) = \begin{cases} 1 & \text{si } f(x) > 1/2 \\ 0 & \text{sinon} \end{cases}$

Théo de Stone pour la classification

Si les hypo du théo de Stone sont vérifiées pour toute distribution de (x, y) , alors \hat{g} est universellement consistante.

i.e. $E(L(\hat{g})) \rightarrow L^* + \text{distrib}(x, y)$

RÈGLES NON-PARAMÉTRIQUES ET DATA-SPLITTING

I) Règle des plus proches voisins (k-NN)

Soit (x, y) couple de v.a $\rightarrow \mathbb{R}^d \times \{0, 1\}$

$$\Omega_n = \{(X_i, Y_i)\}_{i=1}^n \text{ iid.}$$

Pour x fixé ($x \in \mathbb{R}^d$), on réordonne l'échantillon $\Omega_n(x) = \{(X_{(1)}^{(x)}, Y_{(1)}^{(x)}) \dots (X_{(n)}^{(x)}, Y_{(n)}^{(x)})\}$ suivant les distances croissantes $\|X_{(1)} - x\| \dots \|X_{(n)} - x\|$.

$$\text{On a: } \|X_{(1)} - x\| \leq \|X_{(2)} - x\| \leq \dots \leq \|X_{(n)} - x\|$$

Déf 1: Soit $(w_{n,1}, \dots, w_{n,n})$ un vecteur à poids de somme 1. La règle des k-NN s'écrit $\forall x \in \mathbb{R}^d$.

$$\hat{g}_n^{(k)} = \begin{cases} 1 & \text{si } \sum w_{n,i} \cdot y_i = 1 \\ 0 & \text{sinon} \end{cases}$$

Pré: On prend $w_{n,i} \searrow 0$ en i.

\Rightarrow On vote fait en vote pondéré des plus proches voisins.

La choix classique est de prendre

$$\{w_{n,i}\}_{i=1}^n = \underbrace{\left\{ \frac{1}{k}; \frac{1}{k}; \dots; \frac{1}{k}; 0; \dots; 0; 0 \right\}}_{k \text{ fois}}$$

On obtient la règle classique des kppv.

$$\hat{g}_{n,k}(x) = \begin{cases} 1 & \text{si } \sum_{i=1}^k y_i = 1 \\ 0 & \text{sinon} \end{cases}$$

► On privilégi^e k: impaire

I. 2) Théo de Convergence

* Théo de Stone

Si $k_n \xrightarrow{n \rightarrow \infty} +\infty$ et $\frac{k_n}{n} \xrightarrow{n \rightarrow \infty} 0$, alors
pour toute loi (X, Y) , $E\{f(g_{n,k})\} \rightarrow L$
c'est la consistance faible universelle.

- $k_n \rightarrow +\infty \Rightarrow$ On prend suffisamment de voisins pour ne pas rester dans un voisinage temps de ∞ .
Pela évite le sur-apprentissage.
- $\frac{k_n}{n} \rightarrow 0 \Rightarrow$ On prend relativement peu de voisins pour ne pas considérer un voisinage temps grand.
Pela évite le sous-apprentissage.

23.01.24

⑥

⚠ Difficulté : Robibren t.

Theo I.31/lim $\lim_{n \rightarrow +\infty} L(g_n, \mu) = L_{\text{RNN}}$.

On connaît L_{RNN} en fait de Γ (espace)

Theo I.3.21

$$L^* \leq \dots \leq L_{(2k+1)m} \leq L_{(2k-1)m} \leq L_{3m} \leq L_{1m} \leq 2L^*$$

Cela est intéressant pour ce que ça dit que si $L^* \approx 0$, les règles R_m sont bonnes.

Theo I.33 : Pour toute distribution de (X, Y) et $R \geq 3$

$$\bullet L_{\text{RNN}} \leq L^* \left(1 + \frac{\text{cte}}{\sqrt{R}}\right) \text{ à peu près!}$$

30.04.2019

Solutions des paramètres

Datto Splitting V-fold

Toutes les méthodes précédentes dépendent des paramètres à calibrer (R, h, \dots).

I | Risque Empirique

On décompose l'erreur

\exists une famille de règles possible
(par exemple $\exists = \text{RPPV}$)

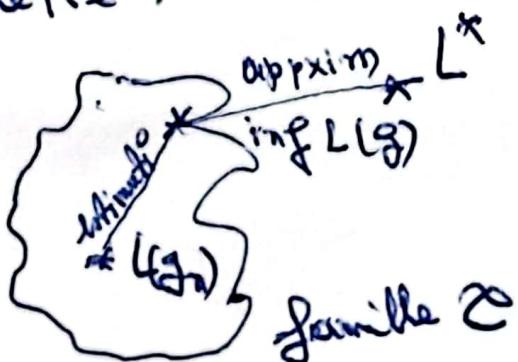
(b) = error of estimation

(2) = erreur d'approximation

↳ erreur structurelle que l'on

↳ exercice structuré que l'on connaît en se restreignant à une famille de règles.

Elle est indépendante du choix du paramètre.



On se focalise pour l'heure d'estimation.

On voudrait choisir \hat{g}_n^* la meilleure règle dans \mathcal{G} .
 $\hat{g}_n^* \in \operatorname{argmin}_{g \in \mathcal{G}} L(g)$.

Cela est impossible parce qu'on ne connaît pas la loi de (X, Y) et donc on ne peut pas calculer $L(g)$.

À la place, on minimise le risque empirique :

$$\hat{L}_n(g) = \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{\{g(x_i) \neq y_i\}}$$

approximati
Monte-Carlo

$$\hat{g}_n^* \in \operatorname{argmin}_{g \in \mathcal{G}} \hat{L}_n(g)$$

I-i) Cas des KPPV

$g_{n,k}$ règle des KPPV $\{(x_{(i)}, y_{(i)})\}_{i=1}^n$

échantillon trié selon les distances croissantes $\|x_i - x\|$.

$$g_{n,k}^{(i)} = \begin{cases} 1 & \text{si } \sum_{j=1}^k b_j y_j = 1 \\ 0 & \text{sinon.} \end{cases}$$

$$g_{n,k}^* \in \operatorname{argmin}_g L_n(g)$$

$$\Leftrightarrow k^* \in \operatorname{argmin}_{\substack{k \\ k \in \{1, \dots, n\}}} L(g_{n,k})$$

Ainsi, $g_{n,k}^* = g_{n,k^*}$

Pourquoi ça ne marche pas?

- On cherche à optimiser l'estimateur sur les données qui ont servi à la construction \rightarrow Construire.

L_n n'est pas consistant car L_n et $g_{n,k}$ dépendent des mêmes données.

$$\frac{1}{n} \sum_{i=1}^n \eta_{g_{n,k}}(x_i) + y_i \rightarrow L^* \text{ si } g_{n,k} \perp \!\!\! \perp (x_i, y_i).$$

- On sur-apprend

II. Validation croisée

II-1) Hold-out (Data-splitting)

II. 3) constance

* Data splitting

On peut calculer / contrôler

$$L(g_n^*) - \inf_{g_m \in \mathcal{G}_m} L(g_m)$$

\downarrow
taille de A_m

$$L(g_n^*) = P[g_n^*(x) \neq y | S_n]$$

$$L(g_m) = P[g_m(x) \neq y | A_m] \quad \forall A_m \subset S_n$$

Il est facile de montrer que

$$L(g_n^*) - \inf_{g_m \in \mathcal{G}_m} [L(g_m)] \leq \sup_{g_m \in \mathcal{G}_m} |\hat{L}_{m,\ell}(g_m) - L(g_m)|$$

on contrôle facilement les choses

si $\# \mathcal{G}_m$ est fini.

$$P \left\{ \sup_{g_m \in \mathcal{G}_m} |\hat{L}_{m,\ell}(g_m) - L(g_m)| > \varepsilon | S_n \right\} \leq 2 \times \# \mathcal{G}_m^{-\frac{-2\varepsilon^2}{\ell}}$$

$$\text{si } m = \ell = \frac{n}{2}$$

$$\mathbb{E}(L(g_m^*) - \inf_{g_m \in \mathcal{G}_m} L(g_m)) \leq \delta \sqrt{\frac{\log(2e/\delta)}{n}}$$

Si $\#\mathcal{G}_m$ est infini, on utilise la théorie de Vapnik-Chervonenkis

Dimension V.C

I | Error Estimation | Approximation

Soit C une famille de règles et
On veut choisir une règle dans C .

L_f :

$L_f(g) = P(g(x) \neq Y)$ soit proba de
 $\inf_{g \in C} \{P(g(x) \neq Y)\}$.

On minimise le risque empirique

$$g_n^* \in \arg \min g \underbrace{\frac{1}{n} \sum_{i=1}^n \mathbb{1}_{\{f(x_i) \neq y_i\}}}_{L_n(g)}$$

(en utilisant une procédure de validation croisée)

$$L(g_m^*) = P(g_m^*(x) \neq Y | S_n)$$

on espère que

$$L(g_m^*) \approx \inf_{g \in C} L(g) \quad [\text{vrai si } C \text{ est fini}]$$

Rappel: $L(g_m^*) - l^*$

$$= \underbrace{\{L(g_m^*) - \inf_{g \in C} L(g)\}}_{\text{Erreur d'estimation}} + \underbrace{\{\inf_{g \in C} L(g) - l^*\}}_{\text{Erreur d'approximation}}$$

On se concentre sur l'erreur d'estimation

Lemma (Important)

$$\text{1)} |L(g_n^*) - \inf_{g \in C} L(g)| \leq \sup_{g \in C} |\hat{L}_n(g) - L(g)|$$

$$\text{2)} |\hat{L}_n(g_n^*) - L(g_n^*)| \leq \sup_{g \in C} |\hat{L}_n(g) - L(g)|$$

Preuve

2) \Rightarrow trivial

$$\begin{aligned} \text{1)} & |L(g_n^*) - \inf_{g \in C} L(g)| \\ & \leq \underbrace{|L(g_n^*) - \hat{L}_n(g_n^*)|}_{(1)} + \underbrace{|\hat{L}_n(g_n^*) - \inf_{g \in C} L(g)|}_{(2)} \end{aligned}$$

$$\text{D'où: } (1) \leq \sup_{g \in C} |\hat{L}(g) - L(g)|$$

$$\textcircled{2} = \left| \inf_{g \in C} \hat{L}_n(g) - \inf_{g \in C} L(g) \right| \leq \sup_{(x)} \left| \hat{L}_n(x) - L(x) \right|$$

Exo

Prouver q: $\left| \inf_{x \in A} (f(x) - g(x)) \right| \leq \sup_{x \in A} |f(x) - g(x)|$

qq preuve (*) .

Indication:

$$|f(y) - g(y)| \leq |f(y) - f(z)| + |f(z) - g(z)| + |g(z) - g(y)|$$

$$\text{min}, |f(y) - f(z)| \leq \dots \leq \sup_x |f(x) - g(x)| \forall y$$

$$\Rightarrow |f(y) - g(y)| \leq \dots + \sup_x |f(x) - g(x)|$$

$$|f(y) - g(y)| \leq \dots + \sup_x |f(x) - g(x)|$$

Prendre \inf_x intelligemment

II. Classe C de taille finie

Tout repose sur l'inégalité de Hoeffding.

$$P(\hat{E} \neq E) \leq 2e^{-2n\delta^2}$$

13.02.2014

APP - STAT

②

Théo II.1

X_1, \dots, X_n r.v. a ~~IID~~ bornées $f_{q,f}(a_i, b_i)_{i=1}^n$

$P(X_i \in [a_i, b_i]) = 1$. On pose $S_n = \sum_{i=1}^n X_i$

$\forall \varepsilon > 0$

$$\textcircled{1} P\{S_n - E(S_n) \geq \varepsilon\} \leq e^{-\frac{-2\varepsilon^2}{\sum(b_i - a_i)^2}}$$

$$\textcircled{2} P\{S_n - E(S_n) \leq -\varepsilon\} \leq e^{-\frac{-2\varepsilon^2}{\sum(b_i - a_i)^2}}$$

$$\textcircled{3} P\{|S_n - E(S_n)| \geq \varepsilon\} \leq 2e^{-\frac{-2\varepsilon^2}{\sum(b_i - a_i)^2}}$$

Preuve

① Admis

② Découle de ① (Exo).

③ Découle de ① et ② (Exo).

Rq: X_1, \dots, X_n n'ont pas besoin d'avoir la même loi.

Exemple fondamental (Exo)

X_1, \dots, X_n iid $\sim B(p)$, alors $-2n\varepsilon^2$

$$P\left(\left|\frac{S_n}{n} - p\right| \geq \varepsilon\right) \leq 2e^{-2n\varepsilon^2}$$

on 1.

$$x \sim \mathcal{B}(p) \Rightarrow [a_i, b_i] = [0, 1],$$
$$p = \mathbb{E}(X) = \mathbb{E}\left(\frac{s_n}{n}\right).$$

with $\varepsilon > 0$.

$$\left| \frac{s_n}{n} - p \right| \geq \varepsilon \Leftrightarrow |s_n - np| \geq n\varepsilon$$

~~$$\Rightarrow P(|s_n - \mathbb{E}(s_n)| \geq n\varepsilon)$$~~

$$\Leftrightarrow P(|s_n - \mathbb{E}(s_n)| \geq n\varepsilon) \leq 2e^{-\frac{2n^2\varepsilon^2}{2(n-\varepsilon)^2}} = n$$

$$\Leftrightarrow P(|s_n - \mathbb{E}(s_n)| \geq n\varepsilon) \leq 2e^{-2n\varepsilon^2}.$$

Theo II.2

e c

Si $\#C \leq N$, alors $\forall \varepsilon > 0$

$$P\left(\sup_{g \in C} |\hat{L}_n(g) - L(g)| > \varepsilon\right) \leq 2N e^{-2n\varepsilon^2}.$$

Preuve

s i

Soit $Z_i = \mathbb{1}_{\{g(x_i) \neq y_i\}}$?

Dès lors que $\sum Z_i$ (exemple fonctionnelle)

$$P\left(\sup_{g \in C} |\hat{L}_n(g) - L(g)| > \varepsilon\right) \leq P\left(\left|\sum Z_i\right| > \varepsilon\right)$$
$$\leq \sum_{g \in C} P\left(\left|\hat{L}_n(g) - L(g)\right| > \varepsilon\right)$$

On a: $z_i \sim P(\underbrace{P(g(x_i) \neq y)}_{L(g)})$

$$S_n = \sum z_i$$

$$\frac{S_n}{n} - p = \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{\{g(x_i) \neq y_i\}} = \hat{L}_n(g)$$

$$p = L(g).$$

Donc

$$\frac{S_n}{n} - p = \hat{L}_n(g) - L(g) \quad \text{et donc,}$$

$$\forall g, P(|\hat{L}_n(g) - L(g)| > \varepsilon) \leq e^{-2n\varepsilon^2}$$

$$= P\left(|\frac{S_n}{n} - p| > \varepsilon\right) \leq e^{-2n\varepsilon^2}$$

avec S_n somme de Bernoulli. D'après

$$\text{Aussi, } P\left(\sup_{g \in C} |\hat{L}_n(g) - L(g)| > \varepsilon\right)$$

$$= P\left(\forall g \in C, |\hat{L}_n(g) - L(g)| > \varepsilon\right)$$

$$\leq \sum_{g \in C} P(|\hat{L}_n(g) - L(g)| > \varepsilon)$$

$$\leq \sum_{g \in C} e^{-2n\varepsilon^2} \leq N e^{-2n\varepsilon^2}$$

$$\leq \sum_{g \in C} 2e^{-2n\varepsilon^2} \leq 2N e^{-2n\varepsilon^2}$$

(parce que $\# C \leq N$)

NB: Cela est bien ce que nous pouvons intégrer $\ln P$.

Problème: gérer le cas $\# C = \infty$

→ théorie du V. C

III | La dimension est V. C

III.1) Poser $\sup_{y \in C} \mu \sup_{A \in \mathcal{A}} \lambda$

• r la loi du couple (X, Y) .

$\forall A \in \text{Borelliens } (\mathbb{R}^d) \times \{0, 1\}$

$$r(A) = P((X, Y) \in A)$$

• V_n la mesure empirique

$$V_n(A) = \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{\{(X_i, Y_i) \in A\}}$$

$$L(g) = P(g(X) \neq Y) = P((x, y) \in \{(x, y) : g(x) \neq y\})$$

$$= r(\underbrace{\{(x, y) : g(x) \neq y\}}_{A_g})$$

$$\hat{L}_n(g) = V_n(A_g)$$

13.02.24

APP. STAT

(3)

On peut réécrire

$$A_g = (\{x : g(x) = 1\} \times \{0\}) \cup (\{x : g(x) = 0\} \times \{1\})$$

On note: $\mathcal{W} = \{A_g, g \in C\}$ ~~et A_g~~

$$\sup_{g \in C} |\hat{L}_n(g) - L(g)| = \sup_{A \in \mathcal{W}} |V_n(A) - V(A)|.$$

On se ramène à l'étude de la convergence de la mesure empirique

Propriétés:

1) Si A mesurable

$$|V_n(A) - V(A)| \xrightarrow[n \rightarrow +\infty]{P.S.} 0 \quad (\text{LGN})$$

2) $P(|V_n(A) - V(A)| > \varepsilon) \leq 2e^{-2n\varepsilon^2}$ (Hoeffding)

3) Si A est de # finie (cardinal fini),

$P(\sup_{A \in \mathcal{W}} |V_n(A) - V(A)| > \varepsilon) \leq 2 \# \mathcal{W} e^{-2n\varepsilon^2}$ (Tho II.2)

4) Si \mathcal{W} est l'ensemble de tous les boreliens de $\mathbb{R}^d \times \{0,1\}$ (ENORME)

alors $\sup_{A \in \mathcal{A}} |\bar{V}_n(A) - V(A)|$ ne peut être
pas $\sqrt{\epsilon}$.

Conclusion: Il faut contrôler la taille de \mathcal{A}

- On veut \mathcal{A} grande pour l'erreur d'approximation
- On veut \mathcal{A} pas trop grande pour contrôler $\sup_{A \in \mathcal{A}} |\bar{V}_n(A) - V(A)|$

III.2) coefficient de généralisabilité et dimension V.C.

On se donne \mathcal{A} une famille d'ensemble de \mathbb{R}^d , z_1, \dots, z_n iid $\in \mathbb{R}^d$

$$\phi V(A) = P(z \in A)$$

$$\bar{V}_n(A) = \frac{1}{n} \sum H_{\{z_i \in A\}}$$

On veut des résultats sur la convergence de $\sup_{A \in \mathcal{A}} |\bar{V}_n(A) - V(A)|$

Def i)

Sont A une famille d'ensembles mesurables de \mathbb{R}^2 et $(z_1, \dots, z_n) \in (\mathbb{R}^2)^n$

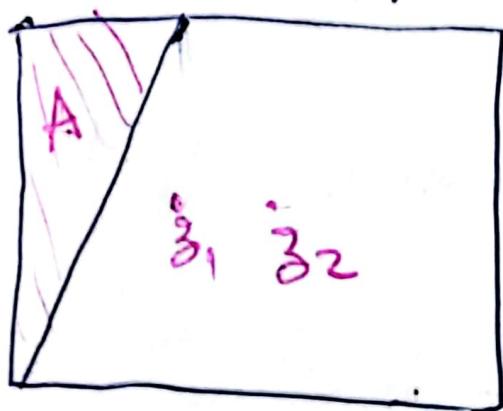
$N_A(z_1, z_n)$ cardinal de $\{(z_1, z_n) \cap A, A \in \mathcal{A}\}$

Lorsque $N_A(z_1, z_n) = 2^n$, on dit

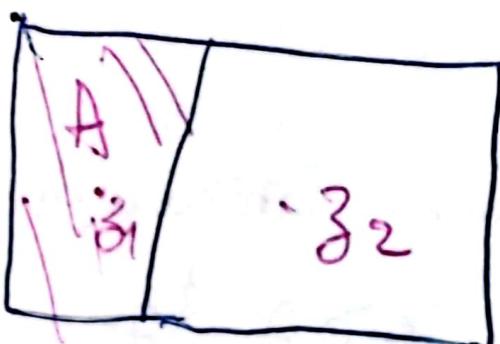
que A pulvérise $\{z_1, z_n\}$.

Exemple :

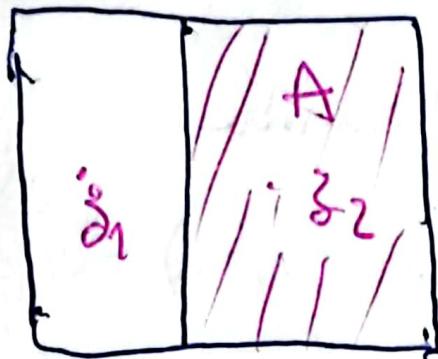
Soit \mathbb{R}^2 , A = ensemble des 1/2 espaces de \mathbb{R}^2



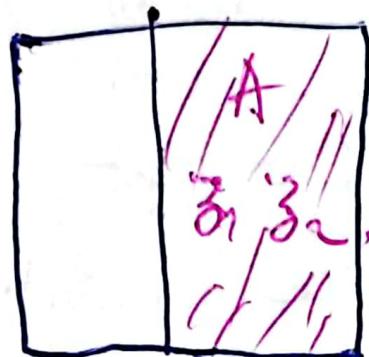
$$\rightarrow \{\{z_1, z_2\} \cap A\} = \emptyset$$



$$\rightarrow \{\{z_1, z_2\} \cap A\} = \{z_1\}$$



$$\text{d}\{z_1, z_2\}^n A = \{z_2\}$$



$$\{z_1, z_2\}^n A = \{z_1, z_2\}$$

Ainsi, $\mathcal{N}(z_1, z_2) = 2^2 = 4$

L'ensemble des 1/2 espaces joue le rôle de $\{z_1, z_2\}$.

Def 2

Le \mathcal{N}_{eff} de pulvérisation de (A) est

$$\mathcal{G}_A^{(n)} = \max_{(z_1, \dots, z_n) \in (\mathbb{R}^d)^n} \mathcal{N}(z_1, \dots, z_n)$$

Def 3

C'est une famille d'ensembles mesurables le plus grand entier $k \geq 1$ tel que $S_A^{(k)} = 2^k$ est appelé dimension H.C.

13.02.24

APP - STAT

43

il est noté V_A .

Si $S_A(k) = 2^k$ et $k > 1$, $V_A = \infty$.

Rq:

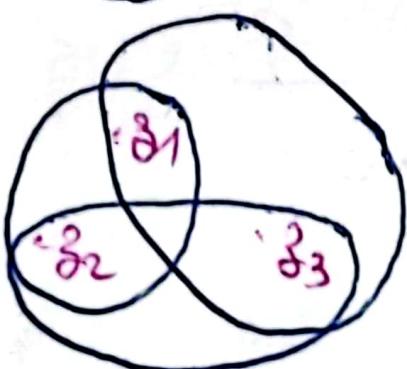
• Plus A est grande, et plus on va pouvoir pulvériser en grande nombre de points.

• Pulvériser n points, c'est être capable d'isoler nimporte quelle sous-partie de ces n points.

Exemple 1 $V_A = \text{disque de } R^2$.



$$S_A(2) = 2^2 = 4$$



$$S_A(3) = 2^3 = 8$$

$\cdot z_1$

$\cdot z_2$

$$\rightarrow S_{\text{ct}}(w) < 16.$$

$\cdot z_4$

$$\cdot z_3 \quad \text{Donc, } V_{\text{ct}} = 3$$

Exemple : $A = \text{ellipse de } \mathbb{R}^2$

$\cdot z_1$

$\cdot z_2$

$$S_{\text{ct}}(w) = 16$$



Rq : $S_{\text{ct}}(b) < 32$.

Donc $V_{\text{ct}} = 4$.

III-3 | Exemple important

1) Si A est une cardinal fini

$$S_{\text{ct}}(n) \leq \# A, \forall n \geq 4.$$

pour définition $S_A(V_A) = 2^{V_A}$,

(car $\forall k \in \mathbb{N} \text{ tel que } S_A(k) = 2^k$)

$\Rightarrow V_A \leq \# A$

$$\Rightarrow V_A \leq \frac{\ln(\# A)}{\ln(2)} = \log_2(\# A)$$

2) Dans \mathbb{R} , $A = \{x \in \mathbb{R}\}$

$$S_A(\emptyset) = \emptyset \rightarrow \{\emptyset\}$$

$$\cancel{\overbrace{B_1 \cup B_2 \cup \dots \cup B_n}} \rightarrow \{\emptyset, \{B_1\}, \{B_2\}, \dots, \{B_n\}\}$$

$$\Rightarrow S_A(\omega) < 4$$

$$\Rightarrow V_A = 1$$

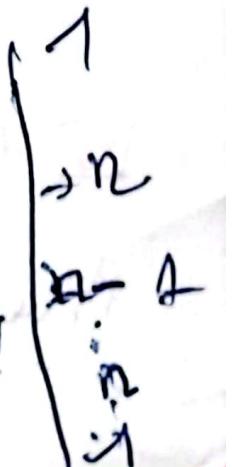
$$S_A(n) = n+1$$

3) Dans \mathbb{R} , $A = [a, b]$, $a, b \in \mathbb{R}$.

$$B_1, \dots, B_n \rightarrow \emptyset$$

$$\overbrace{B_1 \cup B_2 \cup \dots \cup B_n} \rightarrow \{\{B_1\}, \{B_2\}, \dots, \{B_n\}\}$$

$$\{B_1, B_2, B_3\}, \{B_2, B_3, B_4\}, \dots$$



$$S_n(1) = 2^1 = 2$$

$$S_n(2) = 2^2 = 4$$

$$S_n(3) = 7 < 2^3 = 8$$

$$V_{\text{eff}} = 2$$

4) Résultat dans \mathbb{R}^d :

- $\mathcal{A} = \{]-\infty; a_1] \times \dots \times]-\infty; a_d] \}$

$$V_{\text{eff}} = d$$

- $\mathcal{A} = \{ \text{rectangles de } \mathbb{R}^d \}, V_{\text{eff}} = d$

- $\mathcal{A} = \text{famille des } 1/2 \text{ espaces linéaires}$

$$\{ a_1x_1 + a_2x_2 + \dots + a_dx_d + b \geq 0 \}$$

$$V_{\text{eff}} \leq d+1.$$

20/01/24: IV: Convergence uniforme de la mesure empirique

IV-1: Le théorème de Vapnik-Chervonenkis

Théo: IV-1 (1971)

Théo

Soit une mesure de proba sur \mathbb{R}^d ,

$$N_n = \frac{1}{n} \sum_{i=0}^n \mathbf{1}_{\{X_i \in A\}},$$

et une famille d'ensembles mesurables.

~~\Rightarrow~~ $\forall n \geq 1, \forall \varepsilon > 0,$

$$P \left(\sup_{A \in \mathcal{A}} |\sqrt{N_n}(A) - \sqrt{\nu}(A)| > \varepsilon \right) \leq 8 \int_A (n) e^{-\frac{n\varepsilon^2}{32}}$$

①

Ref: ① Si $\nu(A)$ est de

$\int_A (n) \leq \# A$ et donc

la série $\int_A (n) e^{-\frac{n\varepsilon^2}{32}}$ est une série convergente ($\sum_{n=1}^{\infty} 2^\infty$).

et donc $\sup_{A \in \mathcal{A}} |\sqrt{N_n}(A) - \sqrt{\nu}(A)| \xrightarrow{ps} 0$

② si $\int_A (n)$ est un polynôme

en n degré fini, la série converge.

$$S_A(n) = \sum_{i=0}^n a_i A^{i+1}$$

③ Si $S_A(n) = 2^n$, la série ne converge pas.

IV.2 Application

* Si $A = [-\infty; 3]$, $z \in \mathbb{R}$, $S_A(n) = n+1$

$$\Pr \left\{ \sup_{A \in A} |V_n(A) - V(A)| > \varepsilon \right\} \leq 8(n+1) e^{-\frac{n\varepsilon^2}{32}}$$

$$\Pr \left\{ \sup_{z \in \mathbb{R}} |F_n(z) - F(z)| > \varepsilon \right\}.$$

Convergence uniforme de la fonction de répartition empirique
(Ghvenko-Cantelli)

$$\sup |F_n - F| \rightarrow 0 \text{ p.s.}$$

* Cas général

Lemma IV.4

X une v.a positive tq $\mathbb{E}X > 0$

$$P(X > \varepsilon) \leq C e^{-2n\varepsilon^2}$$

Alors $E[X^2] \leq \frac{\log(Ce)}{2n}$

$$E[X] \leq \sqrt{\frac{\log(Ce)}{2n}}$$

Preuve

$$\begin{aligned} E[X^2] &= \int_0^{+\infty} P(P(X > \varepsilon)) d\varepsilon \\ &= \int_0^u \underbrace{P(X > \varepsilon)}_{\leq 1} d\varepsilon + \int_u^{+\infty} P(X > \varepsilon) d\varepsilon \\ &\leq u + \int_u^{+\infty} P(X > \sqrt{\varepsilon}) d\varepsilon \\ &\leq u + \int_u^{+\infty} Ce^{-2n\varepsilon} d\varepsilon \end{aligned}$$

Il faut trouver u qui minimise cette borne.

On travaille sur $E[X] \leq \frac{\log(Ce)}{2n}$.

Rq: $E[f(z)] \geq f(E[z])$ si f est convexe

$E[f(z)] \leq f(E[z])$ si f est concave

c'est l'inégalité de JENSEN

$Z = x^2$, $f(x) = \sqrt{x}$ concave

$$E(X) = E(\sqrt{X^2}) \leq \sqrt{E(X^2)}$$

Corollaire IV:1 (IMPORTANT)

$$|E\left(\sup_{A \in \mathcal{A}} |V_n(A) - V(A)|\right)| \leq 8\sqrt{\frac{\log(8eS_{\lambda}(n))}{n}}$$

IV:3 Coefficient de pulvérisation et dimension Vopnick - C

$S_A(n) = \text{coef de pulvérisati}$

$$= \max_{(z_1, \dots, z_n) \in (\mathbb{R}^d)^n} K_d(z_1, \dots, z_n)$$

~~V_A~~ = dimension V.C de \mathcal{L}

= le plus grand R tq

$$S_A(k) = 2^k$$

$$= \infty \text{ si } S_A(k) = 2^k \forall k.$$

Propriété (Importante)

* Si $\sqrt{A} < \infty$, $S_A(n) \leq (n+1)^{\sqrt{A}}$
 soit $S_A(n) = 2^n$ si $n \geq 1$ ($\sqrt{A} \neq \infty$)
 $\Rightarrow \left\{ \begin{array}{l} \text{soit } S_A(n) \leq (n+1)^{\sqrt{A}}, \forall n \geq 1 (\sqrt{A} < \infty) \end{array} \right.$

\Leftrightarrow soit $S_A(n) = 2^n$, soit il est polynomial en n .

Résultat final

$$\text{si } \sqrt{A} < \infty, \exists c > 0 \text{ tq } F(\sup_{A \in A} |N_A| - \sqrt{A}) \leq c \frac{\sqrt{A \log(n)}}{n}$$

• si $\sqrt{A} \infty$, on a une convergence en $\sqrt{\frac{\log(n)}{n}}$

II/ Problème de Robustesse

III. 1 Existence d'adversaires

Principal contributeur: Ian Goodfellow

2013: "Intriguing properties of Neural Networks"

Possible de perturber très fortement les prédictions d'un réseau en perturbant très faiblement les données à classer.

La perturbation dépend des données à classer.

La perturbation dépend de la donnée d'entrée.

2017: "Universal Adversarial perturbation"
→ perturbation d'une donnée d'entrée.

La principale méthode : Goodfellow 2015
 "Adversarial training"

Principe général :

On considère une fonction de coût J

(par $J = \sum_i J(w, x_i, y_i)$, $J(w, x, y)$
 poids entre les labels.

on cherche à minimiser J en

w , on va regarder une
 version perturbée.

$$\tilde{J}(w, x, y) = \alpha J(w, x, y) + \\ (1-\alpha) J(w, \text{perturb}(x), y)$$

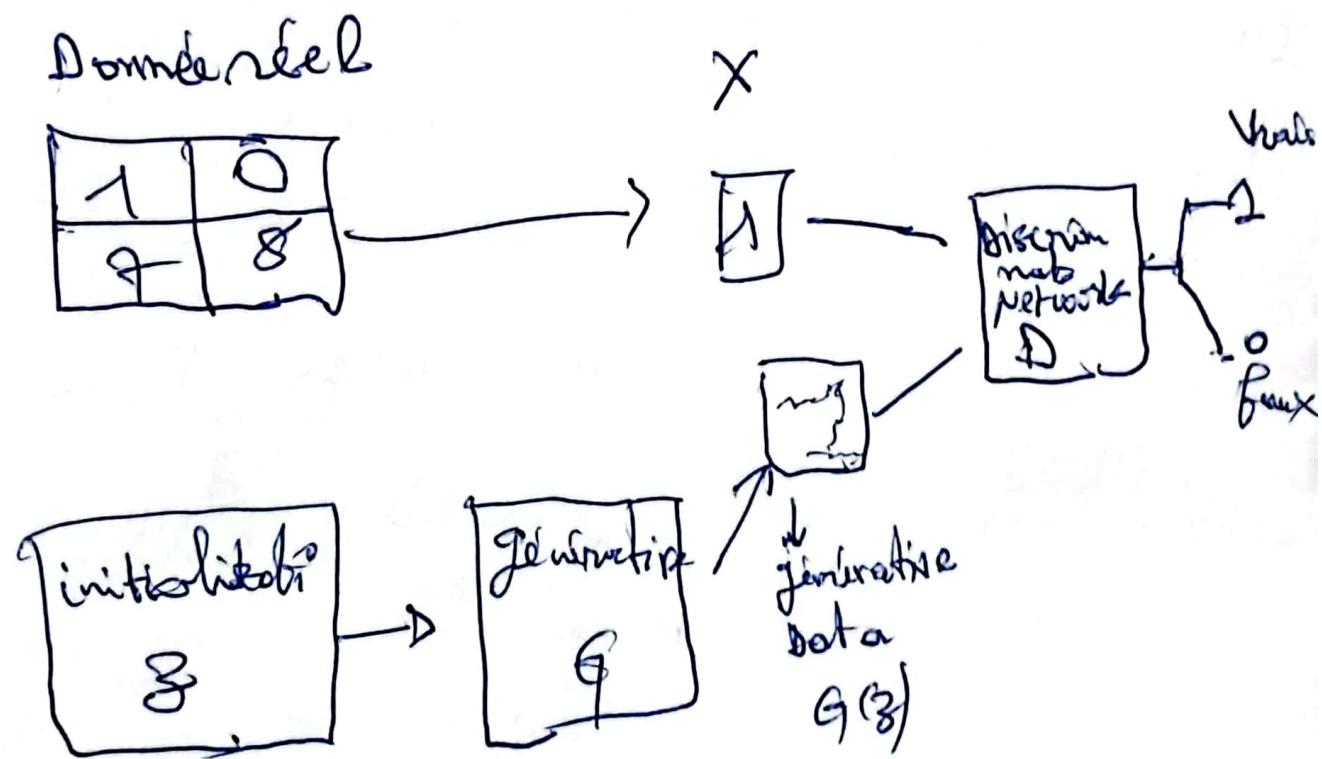
$x + \delta$

$\delta \in [0, 1]$.

en essayant de prendre des
 perturbation qui "entraînent" le
 plus le réseau.

N.1. Generative adversarial Network

Le but est de générer des fausses données "les plus "vraies possible".



La fctⁱ cible $\nabla(G, D)$

$$\nabla(G, D) = \mathbb{E}_x [\log(D(x))] + \mathbb{E}_z [\log(1 - D(G(z)))]$$

On optimise le F.A.N en cherchant min max $\nabla(G, D)$

En pratique

- ① on part d'un jeu de données réelles et on entraîne D (on vise $D(x) = 1$) .
- ② on génère avec G et on entraîne D (on vise $D(G(z)) = 0$)
- 3) on fixe D et on entraîne G (on vise $D(G(z)) = 1$)
- On itérera (1) + (2), et 3.

On appelle ce sur-apprentissage

On part d'un échantillon

$\{(x_i, y_i)\}_{i=1}^n$ et on veut

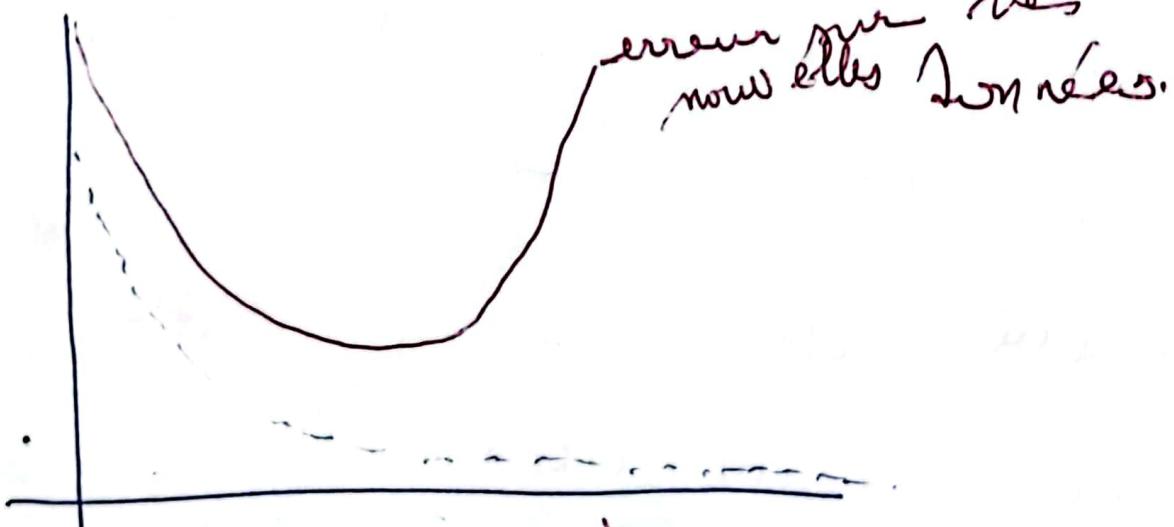
construire une règle g_n

$R^d \rightarrow \mathbb{R}^2$ minimisant $P(g_n(x) \neq y)$

On a vu que selon la complexité
de la feuille de règles
(et le choix des paramètres)

On peut aller du sous-
apprentissage (trop simple et on
prédit la même chose)
au sur-apprentissage

(temps spécifique aux données
d'apprentissage)



Pour le projet
cooler moyen naif \rightarrow moyenne
mobile.

$$\frac{1}{B(x, h)} \sum_{z \in B(x, h)} \mathbb{1}_{\{y_i = g\}}$$

Apprentissage flou

Mardi 21/04/24
Pew
Cours

Pref.: * Michel Grabisch, 2020

* Louis Leogane, "Éléments de logique

* Fuzzy models for pattern recognition

Intro

La théorie des ensembles flous développée en 1965 par Lotfi Zadeh. C'est une théorie de l'algèbre abstraite qui a pour but de représenter l'imprécision relative à l'appartenance à certaines classes d'objets. Cette théorie sert de fondement à la logique floue.

Tres rapidement en apprentissage on n'a pas appris que la notion de classes

(notamment pour la reconnaissance de fleurs) pourraient bénéficier de cette théorie.

En effet, une classe est souvent représentée comme un groupe d'individus partageant certaines similitudes. Ces similitudes peuvent être fortes entre les individus d'une même classe, et un individu peut présenter des similitudes avec des individus d'autres classes.

Ainsi, l'appartenance à un individus une classe peut être réduite à une seule classe mais "distribuée" sur plusieurs classes. Cela rejoint le formalisme des ensembles.

L'ensemble, où les ensembles flora peuvent sembler constituer un cadre naturel pour la notion de classes, force est de constater que les algo classiques considèrent implicitement les classes comme l'nette.

Les méthodes Bayésiennes donnent bien une proba d'appartenance à une classe, mais cette notion de proba (associé à la notion de fréquence d'appartenance à un événement basé sur répétition de l'expérience) est ~~fortement~~ fondamentalement différente de la notion d'appartenance à une classe (sur laquelle repose la théorie "flore").

- Exemple: L l'ensemble des liquides
- L le sous-ensemble des liquides potables
- L le sous-ensemble des jus.

On est dans le désert, sans eau depuis longtemps et on trouve 2 bouteilles A et B. Il faut choisir laquelle boire.

On a 2 infos:

- * Le degré d'appartenance de B à L est 0,9 ($\text{Echelle} = [0, 1]$)
- * On proba que A \in L est 0,9.
Si on choisit B, on voit que le liquide est "probablement" potable (ne contiendra pas les substances les plus toxiques),
si on choisit A, on a 9 chance sur 100 d'avoir un liquide parfaitement

postable, mais 1 chance de mourir

II - Ensemble flows - Relation flows

II.1) Ensemble flows (1965)

Dans le cadre de la théorie des ensembles conventionnelles, les éléments appartiennent ou non (de manière binaires) à des ensembles.

Dans le cadre de la théorie flows, l'appartenance à un ensemble n'est pas si tranchée.

Pour représenter ça mathématiquement, en introduisant une fonction "d'appartenance".
Avec élément à un ensemble, les ensembles flows sont définis de la

méthode vivante.

Def II-1 : Ensemble flou

Un ensemble flou A est un ensemble, caractérisé par une fonction d'appartenance $f_A: E \rightarrow [0,1]$, où E est un espace quelconque.

Si $f_A(x)$ est fraction de 1, x a un fort degré d'appartenance à A .

Si $f_A(x)$ est proche de 0, x a un faible degré d'appartenance à A .

Si $f_A: E \rightarrow \{0,1\}$, on retrouve sur la théorie conventionnelle.

Def II-2 : Ensemble nishe

A (ensemble flou) est nishe si $f_A = 0$

Def II-3 égalité

Les ensembles A et B sont égaux si
 $f_A = f_B$ pour tout $x, f_A(x) = f_B(x)$.

Def III-4 complémentaire

A ensemble flou, \bar{A} son complémentaire
 défini par $f_{\bar{A}}(x) = 1 - f_A(x)$.

Def II-5 : inclusion

A, B flous : A est inclus dans B ($A \subset B$)
 si $f_A(x) \leq f_B(x)$ pour tous x

Théo classeur

$\forall x \in A \quad \exists \alpha \in B$

$$x \in A \Rightarrow f_A(x) = 1 \Rightarrow f_B(x) \geq f_A(x) \Rightarrow \underbrace{f_B(x) = 1}_{\Rightarrow x \in B}.$$

Def II. 6 Union

A, B flous

$$C = A \cup B \text{ est un ensemble}$$

$$\text{flou. Donc } f_C(x) = \max(f_A(x), f_B(x))$$

$\forall x \in E$.

Théorème classique

- Si $x \in A, x \notin B$, $f_A(x) = 1, f_B(x) = 0, f_C(x) = 1$
- Si $x \in B, x \notin A$, $f_A(x) = 0, f_B(x) = 1, f_C(x) = 1$
- Si $x \in A, x \in B$, $f_A(x) = 1, f_B(x) = 1, f_C(x) = 1$
- Si $x \notin A, x \notin B$, $f_A(x) = 0, f_B(x) = 0, f_C(x) = 0$

Proposition II. 1

$$\cancel{AB} \quad A \cup (B \cup C) = (A \cup B) \cup C$$

$$A \subset A \cup B$$

$$B \subset A \cup B$$

12-03. 24

Covers App. STAT

(3)

$$\begin{aligned} f_{A \cup (B \cup C)}^{(x)} &= \max \{f_A^{(x)}, f_{B \cup C}^{(x)}\} \\ &= \max \left\{ f_A^{(x)}, \max \left\{ f_B^{(x)}, f_C^{(x)} \right\} \right\} \\ &= \max \left\{ f_A^{(x)}, f_B^{(x)}, f_C^{(x)} \right\} \\ &= \max \left\{ \max \left\{ f_A^{(x)}, f_B^{(x)} \right\}, f_C^{(x)} \right\} \\ &= f_{(A \cup B) \cup C}^{(x)} \end{aligned}$$

$$f_{A \cup B}^{(x)} = \max \{f_A^{(x)}, f_B^{(x)}\} \geq f_A^{(x)}$$

$\Rightarrow A \subset A \cup B.$

Def: If 7 intersection
A, B flows, $c = A \cap B$ define for
 $f_c^{(x)} = \min \{f_A^{(x)}, f_B^{(x)}\}$

P.II.2

$$\cdot A \cap (B \cap C) = (A \cap B) \cap C \quad \left. \begin{array}{l} \\ \\ \end{array} \right\} \text{Exo.}$$

$\cdot A \cap B \subset A$
 $\cdot A \cap B \subset B$

P.II.3

$$\cdot \overline{(A \cup B)} = \overline{A} \cap \overline{B}$$

$$\cdot \overline{A \cap B} = \overline{A} \cup \overline{B}$$

~~$\cdot \overline{C \cap (A \cup B)} = (\overline{C} \cap A) \cup (\overline{C} \cap B)$~~

~~$\cdot \overline{C \cap (A \cup B)} = (\overline{C} \cap A) \cup (\overline{C} \cap B)$~~

$$\cdot C \cap (A \cup B) = (C \cap A) \cup (C \cap B) \quad \left. \begin{array}{l} \\ \end{array} \right\} \text{Exo}$$

$$\cdot C \cup (A \cap B) = (C \cup A) \cap (C \cup B) \quad \left. \begin{array}{l} \\ \end{array} \right\} \text{Exo}$$

II.2) Relation floues

$D \subset \mathbb{R}^d$, $k > 1$ (cashe concurrent pour commencer)

Une partition de D en \mathbb{R}

clusters / classes sont donnée par des fonctions indicaterices

μ_1, \dots, μ_k

$$\mu_i(x) = \begin{cases} 1 & \text{si } x \in \text{classe } i, i \in \{1, \dots, k\} \\ 0 & \text{sinon} \end{cases}$$

Une relation $\begin{pmatrix} \text{"sure"} \\ \text{"nette"} \end{pmatrix}$ sur D

est définie comme une fonction

$$r: D \times D \rightarrow \{0, 1\}$$

On dit que x et y ($\in D$) partagent une relation si

$$r(x, y) = 1.$$

formalisation matricielle

$$D = \{x_1, \dots, x_n\} \subset \mathbb{R}^d$$

$$N_{ij}(x_j) = N_{ij} \begin{cases} i \in \{1, \dots, k\} \\ j \in \{1, \dots, n\} \end{cases}$$

$$R(x_i, x_j) = R_{ij} \begin{cases} i \in \{1, \dots, n\} \\ j \in \{1, \dots, n\} \end{cases}$$

\mathcal{D}_k = ensemble de toutes les matrices de partition U .

$$\dim(U) = k \times 2$$

$$U[i,j] = p_{i,j} + \lambda \cdot 1^T$$

$$\sum_{i=1}^k p_{i,j} = 1 \quad (\text{chq. donnée à une unique classe})$$

$$\sum_{j=1}^k p_{i,j} \geq 0 \quad \forall i \quad (\text{pas de classe vide})$$

Pour $U \in \mathcal{D}_k$, on associe une matrice de relation $R = \{r_{ij}\}_{i,j}$ avec $\dim(R) = n \times n$ définie par:

$$r_{ii,l} = \begin{cases} 1 & \text{si } p_{i,l} = p_{i,i} = \frac{1}{k}, \text{ pour un certain } i \\ 0 & \text{sinon.} \end{cases}$$

~~Exemple~~

Exemple 3 données ($n = 3$)
 2 classes ($k = 2$)

$$U = \begin{pmatrix} x_1 & x_2 & x_3 \\ 0 & 0 & 1 \\ 1 & 1 & 0 \\ & \frac{1}{p_{22}} & \frac{0}{p_{23}} \end{pmatrix}$$

$\xrightarrow{\text{2ème classe}}$ $\xrightarrow{\text{2ème donnée}}$

$x_1 \in \text{classe } 2$
 $x_2 \in \text{--- } 2$
 $x_3 \in \text{--- } 1$

→ Matrice de relation R produite

$$R = \begin{pmatrix} x_1 & x_2 & x_3 \\ 1 & 1 & 0 \\ 1 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix}$$

NB: se poser la qti si on est au
point ensemble.

* Cache flou

Def II-8 Matrice flou

$U \in M_{f, k}$ (ensemble des partitions
flous à la classe)

Si $\dim(U) = k \times n$ ($n = \text{nbr de données}$
 $\bar{a} \text{ classes}$)

$U[i, j] = p_{ij} \in [0, 1]$ (proto = 4 si
la donnée j
appartient à la
classe i)

$$\sum_{i=1}^k p_{ij} = 1$$

$$\sum_{j=1}^n p_{ij} > 0$$

On peut définir une matrice de
probabilités flow associée à U .

Def II.3 - Matrice de relati^o flow

$U \in M_{k,n}$, la matrice R_f
associée à U pour

$$\begin{cases} 1 \leq j \leq n \\ 1 \leq l \leq m \end{cases} \text{ par}$$

$$R_f[j, l] = R_{jl} = \max_{1 \leq i \leq k} (\min(p_{ij}, p_{il}))$$

Exemple 3 données ($n = 3$)
2 colonnes ($k = 2$)

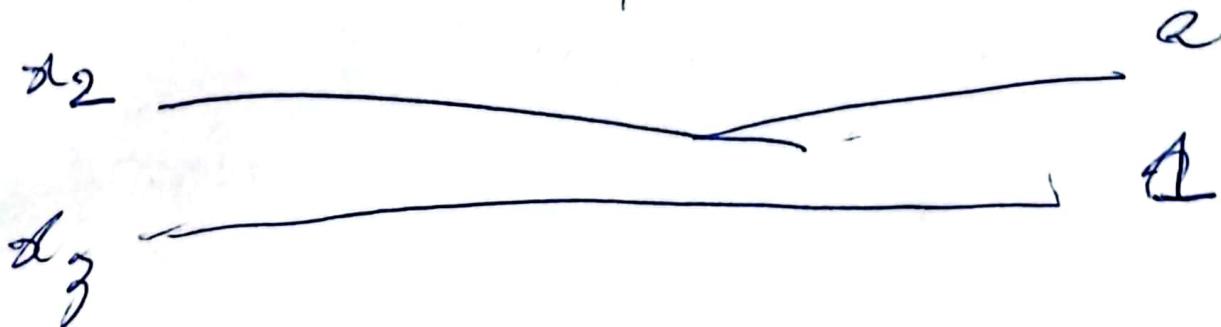
$$U = \begin{pmatrix} x_1 & x_2 & x_3 \\ 0,3 & 0,2 & 0,9 \\ 0,7 & 0,8 & 0,1 \end{pmatrix}$$

~~x_1 grande
 x_2 petite~~

\hookrightarrow degré d'appartenance de x_1 à

la classe 2.

x_1 ressemble plutôt à la classe 2



$$x = R_{1,2} = \max \left\{ \min \left\{ \mu_{1,1}, \mu_{1,2} \right\}, \min \left\{ \mu_{2,1}, \mu_{2,2} \right\} \right\}$$

$$R_{1,2} = \begin{pmatrix} x_1 & x_2 & x_3 \\ 0,7 & 0,7 & 0,3 \\ 0,2 & 0,8 & 0,2 \\ 0,3 & 0,2 & 0,9 \end{pmatrix}$$

x_1 et x_2 plutôt ensemble

x_1 et x_3 pas ensemble

x_2 et x_3 plutôt pas ensemble.

III - Classification

III-1. Projet

$S = \{x_1, \dots, x_n\}$ dont on connaît les labels (pas flou)

$$\mu_{ij} \in \{0, 1\}$$

$$\Rightarrow \mu_{ij} \in \{0, 1\}$$

$S' = \{x'_1, \dots, x'_{n'}\}$ dont les labels sont inconnus.
 $1 \leq i \leq n'$ classe = c

On calcul $\mu'_{ij}, 1 \leq j \leq m$

$$\mu'_{ij} = \frac{\sum_{s=1}^k \mu_{is} \times \frac{1}{\|x'_i - x_s\|^2 / (\lambda - 1)}}{\sum_{j=1}^m \frac{1}{\|x'_i - x_j\|^2 / (\lambda - 1)}}$$

$$1 \leq i \leq n$$

qui est un paramètre ($1 < \lambda < \infty$) qui détermine le degré de flou de la classification

12.03.24

(Centres App - start)

⑤

Exemple : 3 données, 2 classe 2ppv, $\lambda = 2$

$$U = \begin{pmatrix} 0 & 0 & 1 \\ 1 & 1 & 0 \end{pmatrix} \quad \begin{aligned} x_1 &= (0, 0) \\ x_2 &= (1, 1) \\ x_3 &= (4, 4) \end{aligned}$$

2 nouvelles données $x'_1 = (0, 1)$
 $x'_2 = (4, 1)$

Les 2ppr sole x'_1 : x_1 et x_2 $\begin{cases} \|x'_1 - x_1\| = 1 \\ \|x'_1 - x_2\| = 1 \end{cases}$

— x'_2 : x_2 et x_3 $\begin{cases} \|x'_2 - x_2\| = 1 \\ \|x'_2 - x_3\| = 3 \\ \|x'_2 - x_1\| = 3 \end{cases}$

$$P_{1,1}^1 = \frac{\frac{0}{\lambda} \times \frac{1}{\|x'_1 - x_1\|^2} + \frac{1}{\lambda} \times \frac{1}{\|x'_1 - x_2\|^2}}{\frac{1}{\|x'_1 - x_1\|^2} + \frac{1}{\|x'_1 - x_2\|^2}}$$

$$P_{1,1}^1 = 0$$

$$P_{1,2}^1 = \frac{\frac{0}{\lambda} \times \frac{1}{\|x'_2 - x_2\|^2} + \frac{1}{\lambda} \times \frac{1}{\|x'_2 - x_3\|^2}}{\frac{1}{\|x'_2 - x_2\|^2} + \frac{1}{\|x'_2 - x_3\|^2}}$$

$$\cancel{N_{1,2}^1 = \frac{1}{9} X}$$

$$N_{1,2}^1 = \frac{1/9}{1/9 + 1/9} = \frac{1}{2} = 0.5$$

$$U^1 = \begin{pmatrix} 0 & 0.5 \\ 1 & 0.5 \end{pmatrix}$$

Create U

label = Data[, Nbcol] # recuperer le label.

unique(c(0, 0, 1, 1, 1, 2, 2))
 $\hookrightarrow (0, 1, 2)$

recuperer le nbr de label

Nblabel = length(unique.label)

Def une matrice contenant que des zeros

U = matrix(0, ncol = nbdata, nrow = Nblabel)

Remplir la matrice

for (i in 1:Nblabel) {

U[i,] = as.numeric(label == unique[\cancel{i}])

}

over fuzzyknn:

fuzzyknn = function(x, Data, k, lambda, U = 0) {

Nbcol = ncol(Data)

nData = nrow(Data)

Data = Data[order(Data[, Nbcol]),]

donnes = Data[, -Nbcol]

label = Data[, Nbcol]

uniquelabel = unique(label)

Nblabel = length(uniquelabel)

if (U == 0) { U = createU(Data)}

D = ... # calcul des distances ||X - xi||

ODk # indice des kppv

opp # le vecteur dans lequel on va mettre les degrés d'appartenance de x dans chaque groupe

for(i in 1:Nblabel) { # applique la formule p

return (opp)

si $D = (1, 0.8, 2, 0.6)$

$ODk = (4, 2, 1)$

ordre \Rightarrow donne les indice

du plus grand (nbr) ou
du plus petit (nbr)

fuzzypred = ($d = 4$)

hardpred = ($d = 1,05$)

Afficher les

comparer avec plusieurs Valeurs
de d .

26.03.24

App. STAT

①

Apprentissage séquentiel

Introduction

En stat, l'analyse séquentielle est une analyse où la taille de l'échantillon n'est pas connue à l'avance.
À la place, les données sont évaluées au fur et à mesure qu'elles sont recelées.

— Jusqu'à présent on avait un échantillon $s_n = \{(x_i, y_i)\}_{i=1}^n$ et on construisait une règle de décision \hat{f}_n à partir de s_n .

Si on récupère de nouvelles données, il faut tout refaire.

L'approche séquentielle est différente :
On veut pouvoir actualiser \hat{f}_{n+1}

à partir de \hat{y}_n et $(x_{n+1}; y_{n+1})$

Intérêts :

* (1) gain Computational : On peut avoir un faible coût d'actualisation. Il peut être ~~très~~ intéressant de traiter des jeux de données de manière séquentielle même si on les a depuis le début.

* (2) gain d'espace de stockage : une fois \hat{y}_n calculé, on ne va pas y revenir pour passer à S_n .

I - kppv

I-1) De quoi a-t-on besoin ?

* Retour au cas "Classique"

$\mathcal{D}_n = \{(x_i, y_i)\}_{i=1}^n, x \in \mathbb{R}^d, k$ fixé,
label

On réalise l'échantillon
 $\mathcal{Q}_n = \{(x_i^{(n)}, y_i^{(n)})\}_{i=1}^n$ et on
affète à x , la classe la plus
représentée dans ses kppv.

$$g_n(x) \in \operatorname{argmax}_{j \in \text{Label}} \left\{ \sum_{i=1}^n \mathbb{1}_{y_i^{(n)} = j} \right\}$$

* Si une nouvelle donnée $(x_{n+1}; y_{n+1})$
arrive, de quoi on t-on
besoin pour actualiser \hat{g}_n ?

Mais, on a besoin de savoir
si x_{n+1} est dans les kppv de x .

Pour cela, il faut pouvoir comparer
 $\|x_{n+1} - x\|$ avec $\|x_1 - x\|, \|x_2 - x\|, \dots, \|x_k - x\|$
de l'étape précédente.

\Rightarrow Il faut stocker d_1, \dots, d_K

Si $\|X_{n+1} - x\| > d_K$, on change rien.

Si $\|X_{n+1} - x\| < d_K$, il faut intégrer X_{n+1} dans les k_{ppv} , actualiser $\{d_1, \dots, d_K\}$, et il faut recalculer

$$\sum_{i=1}^{q_k} \left\{ \gamma_{(i)}^{(n)} = j \right\} = n_j .$$

\Rightarrow On va garder $\{n_1, \dots, n_m\}$ pt

$$Y_1 \mapsto Y_K$$

I.2) kpp v Segmentatio

Flopes 1: On a

$$S_n = \{(x_i, y_i)\}_{i=1}^n, x, k \text{ fixes}$$

m classes ($\{1, \dots, m\}$)

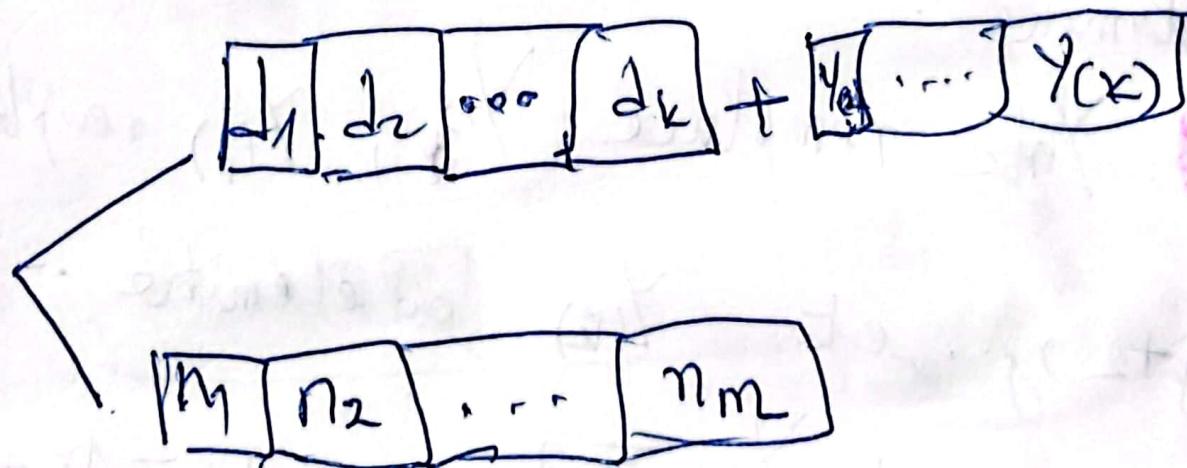
26.03.24

App - START

(2)

On calcule les distances
 $d_1 \leq d_2 \leq \dots \leq d_n$ et on garde
 $\{d_1, \dots, d_k\}, \{Y(1), \dots, Y(k)\}$

on obtient



On fait calculer
 $\hat{g}_m(x) \in \arg \max_{j \in \{1, \dots, m\}} n_j$

Etape 2: Nouvelle donnée
 (x_{n+1}, Y_{n+1})

Si $\|x_{n+1} - x\| \geq d_k \Rightarrow$ on fait rien.

si $\|x_{n+1} - x\| \in [d_i, d_j]$

($0 < i < j \leq R$, $d_0 = \infty$)

On remplace d_j par $\|x_{n+1} - x\|$,

d_{j+1} par d_j ; d_{j+2} par d_j ,

d_{j+2} par d_{j+1} ... et d_k est

éliminé.

② y_{n+2} remplace $y_{(i)}$, $y_{(k)}$ remplace

y_{j+2}, \dots et $y_{(k)}$ est éliminé.

③ si $\begin{cases} y_{n+2} = p \\ y_{(k)} = f \end{cases} \Rightarrow \begin{cases} n_s = n_s + 1 \\ n_b = n_{n+1} \end{cases}$

$j_{n+1} \in \arg \max \{y_j\}$

* Ajustage: trouer Chose si
Calculer, pour se
chose si pt bon.

* Inconvénient : R est fixé
comme VR

- ① Appr "classe" sur un 1^{er} nor de donnée (Etap 1)
- ② apprendre une fct^o qui actualise
 - ↳ f[backtrace, kLabels,
m, x, NomData,]

EXAMEN 24 VOT~~A~~ 24

Règles d'association

* Apprentissage non supervisé.

*→ Influence coussale

• Si quelqu'un achète du pain et du beurre est-ce qu'il va aussi acheter du lait.

Sac Portentante

1	Bearre, pain, lait
2	pain, viande
:	:
n	Poisson, pain

2 critères: "Support", "confidence"

* Support: % d'observat° qui contiennent les éléments de la règle.

100 sac, 30 avec au moins lait, pain, beurre.

Support de lait, pain, Beurre = 0.3 (30%)

* confidence: % d'observat° contenant y sachant qu'il y avait x.

26.03.24

App. STAT

3)

100 sacs, ~~30~~ 33 contiennent
Pain et beurre et dans les 33,
30 contiennent du lait.

$$\text{Confidence: } \frac{30}{33} \approx 0.91$$

\Rightarrow Règles intéressantes: gain support
gain confidence

Un algo: "A priori" (1997)

Etapes: On cherche des sous-ensembles de support assez élevé.

Etape 2: Dans ces sous-ensembles, on cherche les règles de grande confiance.

Exemple: Urrière { support > 0.4
confidence > 0.7 }

1	A, B, C
2	A, B, C
3	B, D, E

Etape 1 :

table 4 → utile pour éliminer les combinaisons de taille supérieure.

A	0.67
B	1
C	0.67
D	0.33
E	0.33

proportion
d'opposition.
support

table 2

AB	0.67
BC	0.67
AC	0.67

Taillé 2

$d(A, B, C)$	0.67
--------------	------

taillé 2

taillé 2

Confidence 70.7

$A \rightarrow B$	$\frac{2}{2} = 1$
$B \rightarrow A$	$2/3 = 0.67$
$A \rightarrow C$	1
$B \rightarrow A$	1
$B \rightarrow C$	0.67
$C \rightarrow B$	1

2 qui
contient
A et pas
C et 2, 2
contient B

2 groupes