

5. Generalization

VC-dimension des réseaux de neurones feed-forward Contrôle de l'erreur de généralisation

Intro

Sur les 3 termes : erreur d'approximation / de généralisation / d'optimisation, on va s'intéresser au 2^{ème}.

Plus précisément : on se restreint au problème de classification binaire :

- À partir d'un échantillon $(X_1, Y_1), \dots, (X_m, Y_m) \stackrel{i.i.d.}{\sim} P_{X,Y}$ où $X \in \{0,1\}$ proposer un classifieur $\hat{h}_m : X \rightarrow \{0,1\}$ tel que $P(\hat{h}_m(X) \neq Y)$ soit le plus petit possible.
- On va étudier des classifieurs \hat{h}_m de la forme $\hat{h}_m = \mathbb{1}_{\{f_m > 0\}} = \text{sgn}(f_m)$ où f_m est un réseau de neurones feed-forward.

On cherche à contrôler l'excès de risque

$$P(\hat{h}_m(X) \neq Y) - \inf_{h \in H} P(h(X) \neq Y)$$

$$\text{où } H = \text{sgn } F = \{ \text{sgn}(f) : f \in F \}$$

avec $F =$ toutes les fonctions qu'on peut obtenir à partir d'une architecture de réseau feed-forward fixée (on fait juste varier les biais et les poids sur les arêtes reliant les neurones).

1/ VC-dimension et borne de risque (quelques rappels)

cf cours de l'X de C. Giraud "Fondements mathématiques de l'apprentissage statistique".

Déf : Soit H un ensemble de classifieurs $h : X \rightarrow \{0,1\}$.

On appelle "coefficient d'éclatement" (ou : coefficient de pulvérisation, shattering coefficient, growth function) la quantité :

$$\left| \pi_H(m) = \max_{x_1, \dots, x_m \in \mathcal{X}} \text{card} \left\{ \underbrace{(h(x_1), \dots, h(x_m))}_{\in \{0,1\}^m} : h \in H \right\} \quad (m \in \mathbb{N}^*) \right.$$

Propriété: $\pi_H(m) \leq 2^m$ et $\pi_H(m) \leq \text{card } H$ si H fini.

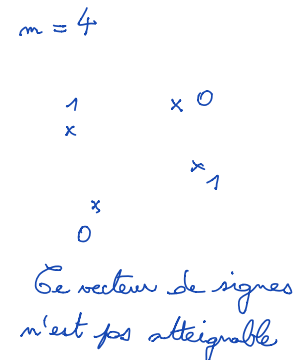
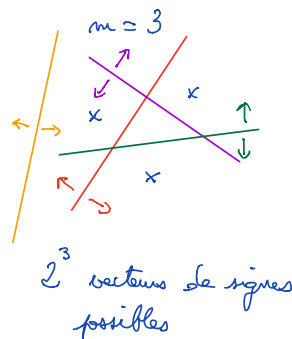
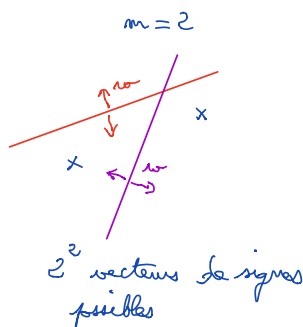
Def: On appelle dimension de Vapnik-Chervonenkis de H , notée $\text{VCdim}(H)$, la quantité (avec la convention $\pi_H(0) = 1$):
 $\text{VCdim}(H) := \sup \{ m \in \mathbb{N} : \pi_H(m) = 2^m \} \in \mathbb{N} \cup \{+\infty\}$.

$$\pi_H(m) = 2^m \iff \exists x_1, \dots, x_m \in \mathcal{X}, \forall \sigma \in \{0,1\}^m, \exists h \in H, \forall i \in \{1, \dots, m\}, h(x_i) = \sigma_i.$$

Interprétation: $\text{VCdim}(H)$ est la taille du plus grand échantillon (x_1, \dots, x_m) que H peut "écarter", i.e., tel qu'on puisse obtenir les 2^m vecteurs de signes possibles en appliquant des fonctions $h \in H$ à l'échantillon (x_1, \dots, x_m) .

Ex: $H := \{ x \in \mathbb{R}^d \mapsto \text{sgn}(\langle w, x \rangle + b) ; w \in \mathbb{R}^d, b \in \mathbb{R} \}$ "perceptron"
 $\text{VCdim}(H) = d + 1$

Schéma pour $d=2$:



Lemme deauer : Soit H tel que $0 < V := V_{\dim}(H) < +\infty$, $\forall m \geq 1$,

$$\pi_H(m) \leq \sum_{i=0}^V \binom{m}{i} \begin{cases} = 2^m & \text{si } m \leq V \\ \leq \left(\frac{em}{V}\right)^V & \text{si } m \geq V \quad (\text{transition exp} \rightarrow \text{poly en } m) \end{cases}$$

Preuve par récurrence.

Proposition 1 (migration du risque de classification de l'ERM)

Soit $H \subset \{0,1\}^X$ tel que $V := V_{\dim}(H) \in \mathbb{N}^*$.

Alors :

$$(a) \mathbb{E} \left[\sup_{h \in H} \left| \mathbb{P}(h(X) \neq Y) - \frac{1}{m} \sum_{i=1}^m \mathbb{1}_{h(X_i) \neq Y_i} \right| \right] \leq 2 \sqrt{\frac{2 \lg(2\pi_H(m))}{m}}$$

\uparrow
sur $(X_i, Y_i)_{1 \leq i \leq m}$

Cela contrôle les déviations uniformes du risque empirique autour du vrai risque.
 \uparrow
sur H

(b) Le minimiseur \hat{h}_m du risque empirique (ERM) : $\hat{h}_m \in \arg \min_{h \in H} \frac{1}{m} \sum_{i=1}^m \mathbb{1}_{Y_i \neq h(X_i)}$

vérifie :

$$\mathbb{P}(\hat{h}_m(X) \neq Y) - \inf_{h \in H} \mathbb{P}(h(X) \neq Y) \leq 4 \sqrt{\frac{2 \lg(2\pi_H(m))}{m}}$$

Preuve : $\forall m \geq V$, $\lg \pi_H(m) \stackrel{\text{Lauer}}{\leq} V \lg\left(\frac{em}{V}\right)$ donc la borne de risque est en $\sqrt{\frac{V \lg(em/V)}{m}}$ (on pouvait enlever ici le facteur \lg avec la technique de "chargement").

• En combinant la preuve avec l'inégalité de Mc Diarmid, on obtient une borne en grande proba en $\sqrt{\frac{V \lg(em/V)}{m}} + \sqrt{\frac{\lg(1/\delta)}{m}}$ avec proba $\geq 1 - \delta$.

Proposition 2 (minoration de l'excès de risque dans le jeu des co, pour tout classifieur) [cf. annexe]

Soit $H \subset \{0,1\}^X$ de VC-dimension $V \in \mathbb{N}^*$.

Alors : $\forall m \geq c_1 V$, $\hat{g}_m^1(x, y)_{x \in \mathcal{X}, y \in \mathcal{Y}}$

$$\inf_{\hat{g}_m^1} \sup_{P \in \mathcal{P}_1^+(\mathcal{X} \times \{0,1\})} \left\{ P(\hat{g}_m^1(X) \neq Y) - \inf_{g \in H} P(g(X) \neq Y) \right\} \geq c_2 \sqrt{\frac{V}{m}}$$

où $c_1, c_2 > 0$ sont des constantes absolues et où l'infimum est pris sur tous les classifieurs $\hat{g}_m^1 : \underbrace{(\mathcal{X} \times \{0,1\})^m}_{\text{sur l'échantillon}} \times \mathcal{X} \rightarrow \{0,1\}$.

Cette proposition signifie que, quel que soit le classifieur \hat{g}_m^1 considéré, il existe une loi jointe P sur $\mathcal{X} \times \{0,1\}$ qui rend l'excès de risque de \hat{g}_m^1 au moins de l'ordre de $\sqrt{V/m}$. Cette borne inférieure s'appelle une "borne inférieure minimale". Elle permet d'identifier l'ordre de grandeur de l'excès de risque dans le jeu des co du meilleur classifieur.

2/ Contrôle de la VC-dimension d'un réseau de neurones feed-forward.

"Nearly-tight VC-dimension and pseudodimension bounds for piecewise linear neural networks", Bartlett, Hanneke, Liaw and Mohri, COLT 2017.

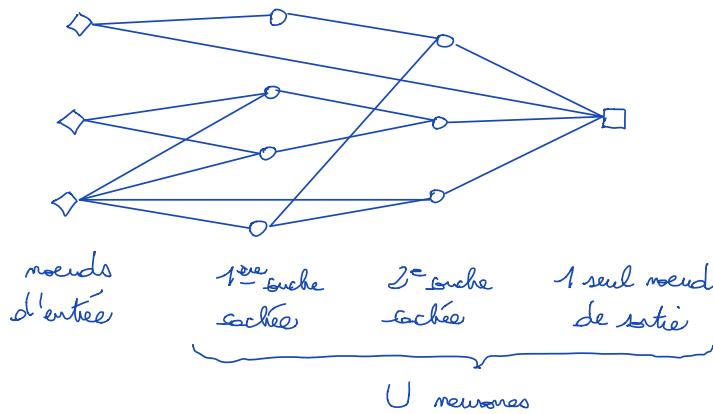
On va majorer la VC-dimension d'un réseau de neurones feed-forward (architecture fixée) en fonction de

- L : nb de couches (layers)
- U : nb de neurones (computation units)
- W : nb de poids (weights)

Définition formelle d'un réseau de neurones feed-forward :

- fonction d'activation $\sigma : \mathbb{R} \rightarrow \mathbb{R}$
- graphe orienté acyclique (DAG) $G = (S, A)$

- vecteur de poids : un réel par arête et un réel par neurone



Le graphe est orienté de gauche à droite.

couche 0 couche 1 couche 2 couche $L=3$

$L :=$ longueur maximale d'un chemin dans G .

$$W = |A| + U$$

\uparrow réels sur les arêtes \uparrow biais pour chaque neurone

Sur l'exemple ci-dessus :
 $L=3, U=8, W=24$

Les numéros de couche sont définis formellement ainsi :

- couche 0 = { nœuds de degré entrant nul }
- couche $l = \{$ nœuds qui admettent au moins un prédécesseur de la couche $l-1$, éventuellement d'autres prédécesseurs des couches $0, 1, \dots, l-2$; et aucun autre prédécesseur $\}$

On suppose qu'il existe un unique nœud de degré sortant nul (nœud de sortie, couche L).

Rem : on autorise des connexions entre couches non-consécutives.

Fonctionnement de chacun des $U-1$ neurones cachés : en entrée : $x \in \mathbb{R}$ ^{degré entrant}
 en sortie : $\sigma(\langle w, x \rangle + b)$ ^{biais}

Neurone de sortie : $x \mapsto \langle w, x \rangle + b$ ^{poids sur les arcs entrants}

On se restreint au cas "polynomial par morceaux": on suppose que $\sigma: \mathbb{R} \rightarrow \mathbb{R}$ est polynomiale sur chacun de $p+1$ morceaux ($p \geq 1$) I_1, I_2, \dots, I_{p+1} intervalles d'intérieurs non-vides, $\mathbb{Z} \geq 2$ disjointes, d'union \mathbb{R} .

Ex: $\sigma = \text{ReLU}$ $I_1 = (-\infty, 0)$ $I_2 = [0, +\infty)$

$$\sigma(x) = \begin{cases} 0 & \text{sur } I_1 \\ x & \text{sur } I_2 \end{cases}$$

Théorème 1 (Bartlett, El Karoui, Liaw et Mehta, COLT 2017)

Soit $L \geq 1$, $U \geq 3$, $d \geq 0$, $p \geq 1$ et $W \geq U \geq L$.

Soit un réseau feed-forward avec W paramètres, U neurones, L couches tel que décrit ci-dessus. On note k_i le nb de neurones sur la i -ème couche ($i=1, \dots, L$). On suppose que les neurones cachés ont une fonction d'activation σ polynomiale sur $p+1$ morceaux ($p \geq 1$) de degré majoré par $d \in \mathbb{N}$. En sortie: fonction d'activation = identité

On pose: pour $i \in \{1, \dots, L\}$,

- Si $d=0$, $W_i =$ nb de paramètres utiles au calcul de tous les neurones de la couche i
 $= (\text{nb d'arcs entrants vers la couche } i) + k_i$ ↙ poids sur les arcs ↘ biais
- Si $d \geq 1$, $W_i =$ nb de paramètres (poids et biais) utiles au calcul de tous les neurones des couches 1 à i .

On pose $\bar{L} := \frac{1}{W} \sum_{i=1}^L W_i \in [1; L]$

→ égal à 1 si $d=0$

→ proche de L si $d \geq 1$ et si les neurones sont concentrés sur les premières couches (ou même uniformément répartis entre les couches)

$$\text{et } R := \sum_{i=1}^L k_i (1+(i-1)d)^{i-1} \leq U + U(L-1)d^{L-1} \left. \begin{array}{l} \} = U \text{ si } d=0 \\ \} \leq ULd^{L-1} \text{ si } d \geq 1 \end{array} \right\}$$

Alors : la classe \mathcal{F} de toutes les fonctions $f_{a \in \mathbb{R}^W} : \mathbb{R}^{\text{entrées}} \rightarrow \mathbb{R}$
 vérifie : $\forall m \geq W$,
 \uparrow paramètres du réseau

$$\mathcal{TC}_{\text{sgn}(\mathcal{F})}(m) \leq \prod_{i=1}^L 2 \left(\frac{2em k_i p (1+(i-1)d)^{i-1}}{W_i} \right)^{W_i} \quad (1)$$

$$\leq \left(4emp (1+(L-1)d^{L-1}) \right)^{\sum_{i=1}^L W_i} \quad (2)$$

Par ailleurs :

$$V\text{dim}(\text{sgn}(\mathcal{F})) \leq L + \bar{L}W \log_2 \left(4epR \log_2(2epR) \right) \quad (3)$$

$$= O(\bar{L}W \log(pU) + \bar{L}LW \log d) \text{ si } d \geq 1$$

\nwarrow négligeable

En particulier :

- Si $d=0$, $V\text{dim}(\text{sgn}(\mathcal{F})) \leq L + W \log_2(4epU \log_2(2epU))$
 $= O(W \log(pU)) = O(W \log(pW))$
- Si $d=1$, $V\text{dim}(\text{sgn}(\mathcal{F})) = O(\bar{L}W \log(pU))$

meilleure que la borne
 $O(\min\{W^2, W \log W + \underline{L^2 W}\})$
 connue jusqu'à alors. \uparrow inutile quand $d=1$.

NB : Cette borne est presque optimale dans le pire des cas : pour $d=1$,

$W \geq cL$ et $L \geq c$ il existe un réseau ReLU à $\leq L$ couches et $\leq W$ paramètres tel que $V\text{dim}(\text{sgn}(\mathcal{F})) \geq WL \log(\frac{W}{L}) / c$.
 ($c > 0$ constante absolue).

Preuve: repose sur le résultat de géométrie algébrique suivant:

Lemme 1 (prouvé par ex dans Anthony et Bartlett '95, theorem P.3)

Soit p_1, \dots, p_m des polynômes de degré au plus $d \geq 1$
en $n \leq m$ variables.

$K := \text{card} \left\{ (\text{sgn}(p_1(x)), \dots, \text{sgn}(p_m(x))) : x \in \mathbb{R}^n \right\}$

le nb de vecteurs de signes possibles.

Alors: $K \leq 2(2em d/m)^m$.

Obtenez $f(x, a)$ la sortie du réseau pour l'entrée $x \in \mathcal{X}$ et le vecteur de paramètres $a \in \mathbb{R}^w$.

$\mathcal{X} = \mathbb{R}^{\uparrow \text{dim des entrées}}$

Soit $x_1, \dots, x_m \in \mathcal{X}$. Afin de majorer $\pi_{\text{sgn}(f)}(m)$, majorons

$\text{card} \left\{ (\text{sgn}(f(x_1, a)), \dots, \text{sgn}(f(x_m, a))) : a \in \mathbb{R}^w \right\}$

$\leq \sum_{i=1}^N \text{card} \left\{ (\text{sgn}(f(x_1, a)), \dots, \text{sgn}(f(x_m, a))) : a \in P_i \right\}$

si P_1, \dots, P_N est une partition de \mathbb{R}^w qui sera choisie de sorte que les m fonctions $a \mapsto f(x_j, a)$ soient polynomiales sur chaque cellule P_i . Il suffit alors d'appliquer le lemme 1.

C'est l'exercice se réduit à la construction d'une bonne partition.

Partitions construites par récurrence: $\mathcal{I}_0, \mathcal{I}_1, \dots, \mathcal{I}_{L-1}$ partitions de \mathbb{R}^W
 telles que: partition finale: $\mathcal{I}_{L-1} = \{P_1, \dots, P_N\}$

(a) les partitions sont emboîtées: chaque $S \in \mathcal{I}_i$ est une réunion d'un ou plusieurs $S' \in \mathcal{I}_{i+1}$ ($0 \leq i \leq L-2$)

(b) $\text{card}(\mathcal{I}_0) = 1$ ($\mathcal{I}_0 = \{\mathbb{R}^W\}$) et: $\forall i \in \{1, \dots, L-1\}$,

$$\frac{|\mathcal{I}_i|}{|\mathcal{I}_{i-1}|} \leq c \left(\frac{2emk_p (1 + (i-1)d^{i-1})}{W_i} \right)^{W_i}$$

(c) Pour tout $i \in \{0, \dots, L-1\}$, tout $S \in \mathcal{I}_i$, tout $j \in \{1, \dots, m\}$, la sortie d'un neurone de la i -ième couche (relatif à l'entrée x_j) est une fonction polynomiale de W_i variables de $a \in S$, de degré $\leq i d^i$.

si $d=0$: les paramètres des neurones de la couche i
 si $d \geq 1$: les couches $1, \dots, i$

Récurrence:

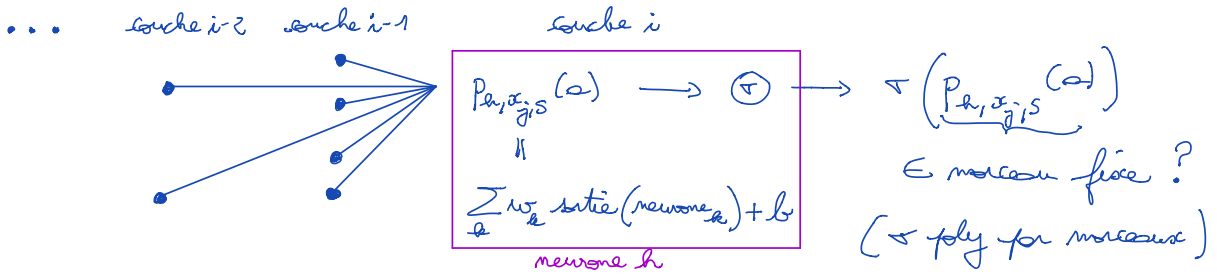
- $i=0$: $\mathcal{I}_0 = \{\mathbb{R}^W\}$ propriété (c) OK (en posant $W_0 = 0$)
 sortie d'un neurone d'entrée = fonction constante en $a \in \mathbb{R}^W$
- Pour $1 \leq i \leq L-1$. Supposons avoir construit des partitions emboîtées $\mathcal{I}_0, \mathcal{I}_1, \dots, \mathcal{I}_{i-1}$ vérifiant (b) et (c). Construisons \mathcal{I}_i .

On note $p_{h, x_j, S}(a)$ l'entrée (juste avant σ) du h -ième neurone ($h = 1, \dots, k_i$) de la couche i , pour l'entrée x_j , comme fonction de $a \in S$ avec $S \in \mathcal{I}_{i-1}$.

D'après l'hyp de récurrence (c), puisque $p_{h, x_j, S}(a)$ est de la forme $\sum_a w_a \text{sortie}(\text{neurone}_a) + b$, et puisque les partitions sont emboîtées, on a: $p_{h, x_j, S}$ est polynomiale sur S , de degré $\leq 1 + (i-1)d^{i-1}$, et dépend d'au plus W_i variables.

← NB: OK si $d=0$.

À cause de σ , la sortie du neurone h est polynomiale par morceaux sur S .
 On va découper S en sous-cellules pour que la sortie soit polynomiale sur chacune des sous-cellules, et ce \forall neurone h , \forall entrée x_j .



Soit $t_1 < t_2 < \dots < t_p$ les coupures des morceaux I_1, \dots, I_{p-1} de σ .
 Considérons les polynômes $\pm (P_{h, x_{j, S}}(a) - t_r)$, $h \in \{1, \dots, k\}$, $j \in \{1, \dots, m\}$, $r \in \{1, \dots, p\}$

$(\text{sgn}(u) = \begin{cases} 1 & u > 0 \\ -1 & u < 0 \end{cases})$

- + si I_{r+1} est ouvert en t_r
- si I_{r+1} est fermé en t_r

D'après le lemme 1, cet ensemble de polynômes sur \mathbb{R}^W atteint au plus

$$\pi := 2 \left(2e \underbrace{(k \cdot m \cdot p)}_{\text{nb de polynômes}} \underbrace{(1 + (p-1)d^{i-1})}_{\text{degré}} / W_i \right)^{W_i} \quad (\text{NB: } m \geq W \geq W_i)$$

$\uparrow \uparrow$
nb de variables effectives

vecteurs de signes différents $(\text{sgn}[\pm (P_{h, x_{j, S}}(a) - t_r)])_{h, j, r}$ quand $a \in \mathbb{R}^W$ et donc quand $a \in S$. On peut donc partitionner S en $\leq \pi$ sous-cellules de sorte que, sur chacune de ces sous-cellules, les $P_{h, x_{j, S}}(a)$ ne changent pas de morceau (de σ) lorsque a varie dans cette sous-cellule.

NB: Ces $\leq \pi$ sous-cellules de S sont les mêmes pour tous les neurones h et toutes les entrées x_j (utile pour la récurrence).

\Rightarrow On obtient une nouvelle partition \mathcal{I}_i de cardinal $\leq \pi \cdot \text{card}(\mathcal{I}_{i-1})$
 La propriété (b) est vérifiée.

Propriété (c) ?

$\forall S' \in \mathcal{J}_i$, la sortie du neurone $h \in \{1, \dots, k_i\}$:

$a \in S' \mapsto \nabla (P_{h, \alpha_j, S'}(a))$ est polynomiale de W_i variables
de degré $\leq d(1 + (i-1)d^{i-1}) \leq i d^i$

donc propriété (c) OK.

↑
multiplication du degré en entrée
par le degré de ∇ .

• Cela est la récurrence \rightarrow partitions emboîtées $\mathcal{J}_0, \mathcal{J}_1, \dots, \mathcal{J}_{L-1}$ vérifient (b) et (c).

En particulier, \mathcal{J}_{L-1} est une partition de \mathbb{R}^W telle que la sortie de chaque neurone des couches $0, \dots, L-1$ est polynomiale de degré $\leq (L-1)d^{L-1}$ sur chaque $S \in \mathcal{J}_{L-1}$, pour tout $j \in \{1, \dots, m\}$.
car les \mathcal{J}_i sont emboîtées

Ainsi, pour chaque cellule $S \in \mathcal{J}_{L-1}$ et chaque entrée $x_j \in X$, la fonction $a \in S \mapsto f(x_j, a)$ en sortie de réseau est polynomiale de degré $\leq \textcircled{1} + (L-1)d^{L-1}$ et ne dépend que de W_L variables.

↑
à cause de la combinaison
linéaire calculée par le
neurone de sortie

D'autre, d'après le lemme 1 :

$$\text{card} \left\{ \left(\text{sgn}(f(x_1, a)), \dots, \text{sgn}(f(x_m, a)) \right) : a \in S \right\}$$

$$\leq 2 \left(\frac{2em(1 + (L-1)d^{L-1})}{W_L} \right)^{W_L} \quad (\text{NB: } m \geq W \geq W_L)$$

Des lors :

$$\text{card} \left\{ \left(\text{sgn}(f(x_1, a)), \dots, \text{sgn}(f(x_m, a)) \right) : a \in \mathbb{R}^W \right\}$$

$$\leq \sum_{S \in \mathcal{J}_{L-1}} \text{card} \left\{ \left(\text{sgn}(f(x_1, a)), \dots, \text{sgn}(f(x_m, a)) \right) : a \in S \right\}$$

$$\leq \text{ord}(\mathcal{J}_{L-1}) \times e^{\left(\frac{2em(1+(L-1)d^{L-1})}{W_L} \right)^{W_L}} \quad (4)$$

Or, d'après la propriété (b) :

$$\text{ord}(\mathcal{J}_{L-1}) \leq \prod_{i=1}^{L-1} e^{\left(\frac{2em k_i p (1+(i-1)d^{i-1})}{W_i} \right)^{W_i}}$$

et donc (car (4) est valable pour tous $x_1, \dots, x_m \in \mathcal{X}$) :

$$\pi_{\text{sgn}(F)}(m) \leq \prod_{i=1}^L e^{\left(\frac{2em k_i p (1+(i-1)d^{i-1})}{W_i} \right)^{W_i}} \rightarrow (1) \text{ OK}$$

Aparté (moyenne géométrique \leq moyenne arithmétique)

• $\forall y_1, \dots, y_k > 0, \forall \alpha_1, \dots, \alpha_k \geq 0$ tq $\sum_{i=1}^k \alpha_i = 1,$

$$\prod_{i=1}^k y_i^{\alpha_i} \leq \sum_{i=1}^k \alpha_i y_i \quad (\text{preuve par Jensen})$$

• D'où: $\forall y_1, \dots, y_k > 0, \forall a_1, \dots, a_k \geq 0$ tq $\sum_{i=1}^k a_i > 0,$

$$\prod_{i=1}^k y_i^{a_i} \leq \left(\frac{\sum_{i=1}^k a_i y_i}{\sum_{i=1}^k a_i} \right)^{\sum_{i=1}^k a_i}$$

$$\pi_{\text{sgn}(F)}(m) \leq e^L \left(\frac{2emp \sum_{i=1}^L k_i (1+(i-1)d^{i-1})}{\sum_{i=1}^L W_i} \right)^{\sum_{i=1}^L W_i}$$

$$= e^L \left(\frac{2empR}{\sum_{i=1}^L W_i} \right)^{\sum_{i=1}^L W_i} \text{ par définition de } R \quad (5)$$

$$\leq \left(\frac{4emp(1+(L-1)d^{L-1}) \sum_{i=1}^L k_i}{\sum_{i=1}^L W_i} \right)^{\sum_{i=1}^L W_i} \text{ car } L \leq \sum_{i=1}^L W_i$$

$$\leq \left(4emp(1+(L-1)d^{L-1}) \right)^{\sum_{i=1}^L W_i} \text{ car } \sum_{i=1}^L k_i \leq \sum_{i=1}^L W_i \rightarrow (2) \text{ OK}$$

Pour pousser la borne (3) sur $VCdim(\text{sgn}(F))$, on va combiner (5) et le lemme suivant :

Lemme 2: Soit $r \geq 16$ et $w \geq t > 0$
 Alors, pour tout $m > t + w \log_2(2r \log_2 r) =: x_0$, on a
 on a $2^m > 2^t \left(\frac{mr}{w}\right)^w$

Suite de la preuve du théorème 1: d'après (5) et le lemme 2 avec $t=L$, $w = \sum_{i=1}^L W_i$ et $r = 2epR \geq 2eU \stackrel{U \geq 3}{\geq} 16$, on a :

$$\forall m > V := L + \left(\sum_{i=1}^L W_i\right) \log_2(4epR \log_2(2epR)) \geq W, \quad \pi_H(m) < 2^m$$

et donc $VCdim(\text{sgn}(F)) \leq V$ par définition de la VC-dimension, ce qui prouve (3). ■

Preuve du lemme 2: Soit $x_0 := t + w \log_2(2r \log_2 r)$.

Montrons que, pour tout $x > x_0$,

$$2^x > 2^t \left(\frac{xr}{w}\right)^w \Leftrightarrow \underbrace{x - t - w \log_2\left(\frac{xr}{w}\right)}_{f(x)} > 0$$

Il suffit de montrer que $f(x_0) \geq 0$ et $f'(x) > 0 \forall x \geq x_0$.

$$\begin{aligned} \bullet f(x_0) \geq 0 &\Leftrightarrow x_0 - t - w \log_2\left(\frac{x_0 r}{w}\right) \geq 0 \\ &\Leftrightarrow w \log_2(2r \log_2 r) \geq w \log_2\left(\frac{x_0 r}{w}\right) \quad (\text{déf de } x_0) \end{aligned}$$

$$\Leftrightarrow 2 \lg_2(x) \geq \frac{x_0}{w} = \frac{t}{w} + \lg_2(2x \lg_2 x)$$

$$\Leftrightarrow \lg_2\left(\frac{x^2}{2x \lg_2 x}\right) \geq \frac{t}{w}$$

$$\Leftrightarrow \frac{x}{2 \lg_2 x} \geq 2^{t/w}$$

ce qui est vrai car $\frac{x}{2 \lg_2 x} \geq 2$ car $x \geq 16$ et $\frac{t}{w} \leq 1$.

• Pour tout $x > x_0$, $f'(x) = 1 - \frac{w}{x \lg_2 x}$

d'où $f'(x) > 0 \Leftrightarrow x > \frac{w}{\lg_2 x}$

ce qui est vrai car $x_0 \geq w \lg_2(2x \lg_2 x) \stackrel{x \geq 16}{>} w / \lg_2 x$
 On en déduit que $f(x) > 0$ pour tout $x > x_0$. \square

Remarque : cette borne de complexité, quasi-optimale dans le pire des cas, n'explique vraisemblablement pas les bonnes performances pratiques des réseaux de neurones, car elles sont permissives.

Ex de papiers proposant des bornes dépendant plus finement du réseau considéré :

- "Spectrally-normalized margin bounds for neural networks", Bartlett, Foster, Gelbovich, NIPS 2017.
- "Type-independent sample complexity of neural networks", Gelbovich, Rakhlin, Lohamir, COLT 2018.

La recherche est encore très active dans ce domaine ; cf. par ex :

- "Uniform convergence may be unable to explain generalization in deep learning", Nagarajan and Kolter, NeurIPS 2019.