

DEPARTEMENT DES MATHÉMATIQUES

**Master Mathématiques et Applications au Calcul Scientifique
(MACS)**

MEMOIRE DE FIN D'ETUDES

Pour l'obtention du Diplôme de Master Sciences et Techniques (MST)

**Théorie d'approximation des Réseaux de Neurones Artificiels,
Application à la résolution des équations différentielles
fractionnaires**

Réalisé par: Salim El azami El idrissi

Encadré par: Pr.ETTAOUIL Mohamed

Soutenu le :20/06/2019

Devant le jury composé de:

- | | |
|--------------------------------------|-----------------|
| - Pr. ETTAOUIL Mohamed | FST Fès |
| - Pr. EZZAKI Fatima | FST Fès |
| - Pr. YOUSSEFI Ahmed | ENSA Fès |
| - Pr. EL KHAOULANI EL IDRISSE Rachid | FST Fès |
| - Pr. HADDOUCH Khalid | ENSA AL-Hoceïma |
| -Dr. RAMCHOUN Hassan | Invité |

Année Universitaire : 2018 / 2019

Dédicaces

À mes parents :

Aucun mot ne pourra exprimer tous mes sentiments d'amour et de gratitude, je vous remercie beaucoup pour vos efforts le long de ces années, pour votre présence rassurante et pour tout l'amour que vous m'avez réservé, vous étiez et vous resterez pour moi les parents idéals qui m'ont toujours conseillé et orienté dans le bon sens, peut-être que j'ai encore des réserves sur des décisions que vous avez préféré de les prendre à ma place, mais soyez sûr que je n'ai jamais perdu ma croyance en votre sagesse que je la considère jusqu'à l'instant indiscutable.

chers parents, je vous remercie une autre fois pour vos sacrifices , que Dieu vous prête la longue vie et la bonne santé afin que je puisse vous rendre la pareille à mon tour, je vous aime très fort ♡.

À mes chers frères et sœurs :

Qui ont toujours été là pour moi, avec leur support leurs encouragements.
Qu'ils trouvent dans ce travail, toute ma reconnaissance et tout mon grand amour envers eux.
Je vous aime beaucoup

À tout les membres de ma famille :

En témoignage de profonds liens qui nous unissent, veuillez trouver à travers ce travail l'expression de mon grand amour, mon attachement et ma profonde reconnaissance.

Remerciements

Au terme de ce travail, je voudrais exprimer mes remerciements et ma profonde reconnaissance à tout ceux qui ont contribué de prêt ou de loin à sa réalisation.

Je voudrais, en tout premier lieu, exprimer ma profonde gratitude à mon encadrant, le professeur **ETTAOUIL MOHAMED**, d'avoir accepté de m'initier à la recherche et de diriger ce travail. Comme tous ceux qui ont eu la chance d'être sous sa direction, j'ai pu constater à quel point il est attachant, attentionné et profondément gentil, il restera pour moi un modèle, pour ses grandes compétences scientifiques, pour son exigence de clarté, sa franchise et ses qualités humaines. Je lui adresse mes remerciements les plus chaleureux, pour tout ce qu'il m'a appris durant la période pendant laquelle s'est déroulé ce travail.

j'aimerais bien adresser des chaleureux remerciements à tous les membres du jury :

- Madame **FATIMA EZZAKI** , Professeur et directrice du laboratoire LMCS à la Faculté des Sciences et Techniques de Fès.
- Monsieur **EL KHAOULANI EL IDRISSE RACHID** Professeur à la faculté des sciences et techniques de Fès.
- Monsieur **HADDOUCH RACHID** Professeur à l'Ecole Nationale des Sciences Appliquées d'Al-Hoceïma.
- et Monsieur **YOUSSEFI AHMED** Professeur à l'Ecole Nationale des Sciences Appliquées de Fès.

qui m'ont fait l'honneur de bien vouloir juger ce travail et de l'enrichir par leurs remarques et leurs critiques.

Je tiens à remercier également mon cher professeur le docteur **HASSAN RAMCHOUN**, et qui n'a épargné ni temps ni effort pour m'aider à concrétiser ce modeste travail.

Je tiens également à remercier le corps professoral de la Faculté des Sciences et Techniques de Fès pour leur contribution à ma formation et pour leurs enseignements précieux et multidisciplinaires.

Enfin, un merci très spécial à mes deux collègues **BRAHIM BEN DABIHI** et **HAKIM HABRI** pour leurs aides précieuses, leur gentillesse et leurs encouragements, tout au long de la préparation de ce mémoire.

Sommaire

Dédicaces.	i
Remerciements.	ii
Sommaire.	ii
Résumé.	vi
Liste des abréviations.	vii
Liste des figures.	ix
Liste des algorithmes.	x
Liste des tableaux.	xi
Introduction générale.	1
I Quelques outils d'analyse fonctionnelle pour la théorie d'approximation :	3
Introduction :	3
1 Passage de l'intégrabilité à la presque continuité :	4
1.1 Théorème d'Egorov :	4
1.2 Théorème de densité des fonctions continues dans les espaces L^p :	6
1.3 Premier théorème de Lusin :	9
2 Aperçu sur la densité dans l'espace des fonctions continues :	11
2.1 Théorème de Dini :	11
2.2 Théorème de Stone-Weierstrass :	11
3 Passage de la mesurabilité à la presque continuité :	15
3.1 Mesures régulières et espaces de Radon :	15
3.2 Deuxième théorème de Lusin :	17
4 C^0 -prolongeabilité des fonctions continues sur un fermé :	17
4.1 Théorème de prolongement de Dugundji :	17
4.2 Théorème de l'extension de Tietze :	18
Conclusion :	19
II Apprentissage Automatique :	
Réseaux de Neurones Artificiels.	20
Introduction :	20
1 De neurone biologique au neurone artificiel :	21
1.1 Neurone biologique :	21
1.2 Neurone artificiel :	22
2 Apprentissage d'un neurone artificiel :	30
2.1 Vue d'ensemble sur le concept d'apprentissage :	30
2.2 Optimisation et apprentissage supervisé d'un neurone formel :	34
2.3 Quelques algorithmes d'apprentissage supervisé pour un neurone formel :	36
3 Réseaux de neurones artificiels :	44

3.1	Nécessité d'introduire la structure de réseau :	44
3.2	Architectures des réseaux de neurones artificiels :	45
3.3	Quelques modèles de réseaux neuronaux artificiels :	46
3.4	Apprentissage supervisé du perceptron multicouches :	48
	Conclusion :	52
III Aptitude des Réseaux de Neurones Artificiels.		
	(Théorie d'approximation) :	53
	Introduction :	53
1	Théorèmes d'approximation au cas d'une fonction d'activation discriminatoire :	54
1.1	Vocabulaire et quelques notions de base :	54
1.2	Résultats principaux :	62
1.3	Application aux réseaux de neurones artificiels :	65
1.4	Résultats pour d'autres fonctions d'activation :	68
2	Théorèmes d'approximation au cas d'une fonction d'activation écrasante :	70
2.1	Définitions et notations :	70
2.2	Résultats fondamentaux :	74
3	Théorèmes d'approximation de Funahashi :	82
3.1	Mise en situation :	82
3.2	Résultats principaux :	82
4	Propriété de parcimonie chez les réseaux de neurones :	83
4.1	Mise en contexte :	83
4.2	Résultats principaux :	83
	Conclusion :	85
IV Application : Résolution numérique des équations différentielles fractionnaires par les Réseaux de Neurones Artificiels :		
	Introduction :	86
1	Initiation à l'analyse fractionnaire :	87
1.1	Quelques fonctions spéciales :	87
1.2	Intégrale fractionnaire de Riemann- Liouville :	91
1.3	Dérivée fractionnaire de Caputo :	93
2	Description et discrétisation du problème :	98
2.1	Description du problème :	98
2.2	Discrétisation du problème :	98
3	Principe de résolution par un processus d'apprentissage :	107
3.1	Passage à un problème d'optimisation :	107
3.2	Algorithme d'apprentissage proposé :	114
3.3	Exemple illustratif :	116
4	Quelques domaines d'applications des systèmes fractionnaires :	121
4.1	En Physique :	121
4.2	En automatique :	121
4.3	En Acoustique :	122
	Conclusion :	122
Conclusion générale et perspectives.		122

Annexe.	123
A Prérequis essentiel pour la lecture du troisième chapitre :	124
1 Théorèmes du Noyau fermé et de Hahn-Banach :	124
2 Espaces Hilbertiens et théorème de représentation de Riesz :	127
Index.	xiii
Bibliographie.	xviii

Résumé

Les équations différentielles fractionnaires (EDFs) apparaissent actuellement dans les différents domaines scientifiques, comme la physique, l'ingénierie, la médecine, l'électrochimie, la théorie du contrôle, ... etc, et l'efficacité remarquable qu'elles présentent, réside essentiellement dans leur grande capacité de modéliser les différents phénomènes du monde réel, ce fait va motiver et encourager, beaucoup de chercheurs à les étudier attentivement, et dans les deux côtés théorique et applicatif.

L'objectif de ce mémoire est de proposer une approche numérique convenable pour la résolution des équations différentielles fractionnaires ordinaires, et selon laquelle on essayera de franchir l'invalidité des schémas numériques classiques pour ce type d'équations, à travers un recours au domaine de l'intelligence artificielle.

Nous étudions tout d'abord et dans un cadre général, l'aspect d'apprentissage automatique, et au cours duquel on présentera les réseaux de neurones artificiels et leurs algorithmes d'apprentissage, puis on passera à la justification de la grande aptitude de ces réseaux, et ceci en faisant appel à la théorie d'approximation, qui va s'appuyer sur plusieurs théorèmes fondamentaux d'analyse fonctionnelle pour fournir des résultats sur leur capacité à approcher une certaine fonction donnée, enfin ces résultats vont nous servir pour approcher la solution de notre équation différentielle, et que l'on fera pratiquement à l'aide des algorithmes d'apprentissage, ce qui rend par conséquent ces réseaux, des générateurs de schémas itératifs pour la résolution numérique de ce type de problèmes.

Mots clés : Apprentissage automatique, Réseaux de Neurones Artificiels, Perceptron, Algorithme de Rétropropagation, Théorie d'approximation, Dérivation fractionnaire, Équations différentielles fractionnaires.

Liste des abréviations

-
- **"ADALINE"** : ADaptive LInear NEuron.
 - **"C.P-neurone"** : neurone de Mc.Culloch et Pitts.
 - **"EDF"** : Equation Différentielle Fractionnaire
 - **"LMS"** : Least Mean Squares.
 - **"MADALINE"** : Many ADALINE.
 - **"MV"** : Maximum de Vraisemblance.
 - **"PMC"** : Perceptron Multi-Couches.
 - **"RNA"** : Réseaux de Neurones Artificiels.
 - **"SVM"** : Support Vector Machine.
-

Liste des figures

I.1	Illustration de la distance de x à Y	6
I.2	Une suite de polynômes convergeant uniformément vers $x \mapsto \exp(x)$ sur $[-1, 1]$	14
I.3	Principe du théorème de Tietze pour $X = \mathbb{R}$	19
II.2	Schéma d'un neurone biologique.	22
II.3	Synapses entre deux neurones.	22
II.4	Comparaison entre la structure d'un neurone biologique et d'un autre formel.	23
II.5	La représentation simple d'un C.P-neurone.	25
II.6	Réalisation de la fonction "AND" par un C.P-neurone.	26
II.7	Réalisation de la fonction "OR" par un C.P-neurone.	26
II.8	Réalisation de la fonction "OR" à trois entrées par un C.P-neurone.	26
II.9	Modèle du perceptron de Rosenblatt.	27
II.10	Une solution possible de la fonction "OR" par un perceptron simple.	28
II.11	L'incapacité du perceptron simple face au problème "XOR".	28
II.12	Exemple d'un Adaline à deux entrées.	29
II.13	La courbe de la fonction $x \mapsto \frac{1}{2}(1 + \tanh(x))$	29
II.14	Quelques choix possibles pour la fonction d'activation	30
II.15	La différence entre l'apprentissage supervisé et non supervisé.	31
II.16	Apprentissage par renforcement.	32
II.17	Apprentissage par transfert face à l'apprentissage usuel.	32
II.18	La fonction (b) est régulière et elle peut bien jouer le rôle de la fonction (a).	38
II.19	Une illustration de la règle de Hebb.	39
II.20	La différence entre l'algorithme de Perceptron et celui de Widrow-Hoff.	43
II.22	La forme d'un réseau feedforward.	45
II.23	La forme d'un réseau récurrent.	46
II.24	Deux représentations différentes d'un réseau de Hopfield simple.	47
II.25	Principe général de la modélisation par une carte de Kohonen.	48
II.26	Schéma d'un perceptron multicouches.	48
II.27	Les étapes de l'algorithme de rétropropagation du gradient.	51
III.1	La courbe de la fonction φ	61
III.2	L'idée générale de la preuve.	67
III.3	La courbe du cosinus écraseur.	71
III.4	La courbe de la fonction Rampe.	71
III.5	L'architecture d'un réseaux de neurones $\Sigma\Pi^r(G)$	73
III.6	Indication de la démonstration	78
III.7	La courbe de la fonction $x \mapsto 2\Gamma\left(x - \frac{3\pi}{2}\right) - 1$	80

III.8 Exemple d'un réseau non parcimonieux.	84
IV.1 La courbe de la fonction Gamma.	89
IV.2 La courbe de la fonction bêta.	91
IV.3 La dérivée fractionnaire de $x \mapsto x$ pour certains ordres $\alpha \in [-1, 1]$	97
IV.4 Le réseau feed-forward réalisant $N(\cdot)$	99
IV.5 Le $\Sigma\Pi$ réseau réalisant $\tilde{V}(\cdot)$	99
IV.6 Le comportement d'une série entière de rayon de convergence R	100
IV.7 La courbe de la fonction $x \mapsto \varphi(x)$	103
IV.8 Maillage proposé du domaine d'étude $\Omega = [0, T]$	107
IV.9 La courbe de la fonction "Logsig".	116
IV.10 L'erreur moyenne en fonction du nombre d'itérations τ pour : $\tau \in \{0, \dots, 100\}$	118
IV.11 L'erreur moyenne en fonction du nombre d'itérations τ pour : $\tau = 100, 200$ et 500	119
IV.12 La courbe de l'erreur locale absolue pour des différents nombres de neurones cachés.	120
IV.13 L'effet des paramètres N et I sur la qualité d'apprentissage.	120
A.1 L'idée du théorème de Hahn-Banach pour $E = \mathbb{R}^2$, $F = \mathbb{R} \times \{0\}$ et $f(x) = 2x$	126

Liste des algorithmes

1	Algorithme d'apprentissage du perceptron simple par la descente de gradient. . . .	37
2	Algorithme de réajustement des poids de Hebb.	39
3	Algorithme d'apprentissage du perceptron.	40
4	Algorithme d'apprentissage de Widrow-Hoff.	42
5	Algorithme de rétropropagation du gradient.	52

Liste des tableaux

IV.1	Les solutions exactes et approchées dans le cas où : $\alpha = 0.25, I = 6$ et $\tau_{\max} = 500$	118
IV.2	Les valeurs de l'erreur locale absolue pour des différents nombres de neurones cachés.	119
IV.3	Les valeurs de l'erreur locale absolue pour des différents choix de α	121

Introduction générale

La classe des équations différentielles d'ordre fractionnaire s'est marquée par leur capacité de modéliser les différents problèmes et phénomènes du monde réelle, et aussi par le comportement de ses solutions qui sont très difficile à trouver, à décrire et à comprendre, pour cela on va proposer, et au cadre de ce mémoire, une approche qui se base sur des modèles d'apprentissage, pour estimer cette solution sur un domaine d'étude donné.

Dans ce contexte le modèle le plus classique qui peut être proposé, sera celui des réseaux de neurones artificiels, la particularité de ces derniers et qui occupent depuis environ soixante ans une place notable dans la recherche scientifique réside dans le fait que chaque structure régulière peut être considérée comme un approximateur universel, ce qui fait, que ce modèle est bien approprié au problème mentionné précédemment et qu'il peut estimer avec une précision approximative élevée la fonction inconnue sur le domaine dont on dispose.

Notre objectif dans la suite est de justifier d'abord le fait que ce modèle est un approximateur universel, et ensuite l'utiliser pour résoudre les problèmes d'équations différentielle fractionnaires

Ce mémoire s'articule principalement autour de quatre chapitres :

Le premier chapitre est consacré à une étude détaillée de quelques théorèmes d'analyse fonctionnelle, comme celui d'Egorov, de Lusin, de Dini et de Stone-Weierstrass, ... etc, et qui vont nous servir par la suite pour faire l'étude de la théorie d'approximation dans le cadre des réseaux de neurones artificiels.

Le deuxième est une présentation générale de l'apprentissage automatique, et qui va englober dans un premier temps une introduction aux réseaux de neurones artificiels, et au cours de laquelle on traitera leur histoire, leurs définitions, leurs types, leurs différentes structures, ainsi que leurs procédures d'apprentissage, et dans un second nous ferons notre premier contact avec le modèle neuronal le plus célèbre, et qui sera nommé "perceptron multicouches", nous allons définir sa nature, sa topologie, ses variantes, et enfin son algorithme d'apprentissage et qui sera appelé l'algorithme de rétro-propagation de gradient.

Le troisième concerne la théorie d'approximation, et qui va établir que les réseaux de neurones artificiels feed forward à une seule couche cachée, sont des approximateurs universels qui peuvent approcher et, au degré souhaité de précision, plusieurs types de fonctions, et on verra en outre, que la capacité de ces réseaux, dépend essentiellement et d'un point de vue théorique, du type de la fonction d'activation choisie et du nombre des nœuds adoptés.

Le dernier chapitre va se baser sur les résultats de la théorie d'approximation, pour appliquer les réseaux de neurones artificiels à la résolution des équations différentielles fractionnaires ordinaires.

Dans une première partie nous allons évoquer le cadre théorique de la dérivation fractionnaire au sens de Caputo , et dans une deuxième on présentera la forme de notre problème, la façon de sa discrétisation, et la procédure de sa résolution par un processus d'apprentissage , enfin, un exemple de test sera présenté pour mieux illustrer la méthodologie proposée.

Quelques outils d'analyse fonctionnelle pour la théorie d'approximation :

Résumé :

Dans ce premier chapitre on va représenter et étudier d'une façon bien détaillée, *quelques théorèmes généraux d'analyse fonctionnelle*, surtout ceux qui *vont nous servir pour la suite* dans le cadre des réseaux neuronaux artificiels (théorème de Lusin, de Stone-Weierstrass, de Dini... etc), et pour cela nous incitons le lecteur d'être un peu *courageux*, et *curieux* lors de sa lecture de cette partie, car elle est très importante, et elle rendra la théorie d'approximation, qui sera traitée ultérieurement dans le chapitre III, beaucoup *plus claire et compréhensible*.

Mots clés : Densité, approximation, la presque continuité, prolongeabilité.

Introduction :

Dans la pratique, on se pose souvent face à des problèmes où l'inconnue est une fonction appartenant à un certain espace fonctionnel déterminé, par exemple celui des fonction continues, dérivables, de classe C^∞ ...etc, et dans ce cas là l'analyse usuelle qui porte essentiellement sur l'étude des \mathbb{R} ou \mathbb{C} espaces vectoriels topologiques^b de dimension finie ou dénombrable devient complètement inutile.

Cette incapacité face à la manipulation des fonctions et qui sont généralement des objets mathématiques d'une nature très complexe et indescriptible, va nous obliger de faire le recours au domaine d'analyse fonctionnelle, et qui peut être vu tout simplement comme une extension naturelle à la dimension infinie de la géométrie euclidienne usuelle qui est en dimension finie, ce passage dimensionnel du fini à l'infini n'est pas du tout évident car ça peut causer parfois la perte de validité d'une grande partie de propriétés topologiques^c, comme il peut créer des contradiction avec l'intuition géométrique^d.

Dans les problèmes d'équations aux dérivées partielles par exemple, ou juste d'équations différentielles ordinaires, ou même d'ordre fractionnaires comme on le verra après au cours du dernier chapitre, le calcul explicite des solutions est souvent hors de portée, ce qui va nous pousser d'aller *chercher à décrire leur structure par des sous ensembles denses dans l'espace fonctionnel adap-*

b. Comme son nom l'indique c'est un espace vectoriel munit d'une topologie qui peut être définie abstraitement comme elle peut provenir d'une norme, d'une distance, d'un produit scalaire ou même d'un ordre.

c. La continuité de toutes les formes linéaires par exemple.

d. Par exemple dans le cas de dimension infinie, on peut trouver des sous espaces vectoriels non fermés.

tés au problème posé, et dont les éléments sont d'une forme descriptible.

afin d'arriver à ce but on présentera dans ce chapitre quelques résultats généraux qui vont nous aider après à justifier l'importance des fonctions réalisables par des réseaux de neurones artificiels, et qui formeront à priori par la suite cet ensemble dense dont on a parlé .

1 Passage de l'intégrabilité à la presque continuité :

1.1 Théorème d'Egorov :

Théorème 1.1.1 :(d'Egorov,[31])

Soit $(\Omega, \mathcal{F}, \mu)$ un espace mesuré *de mesure finie* , et $(f_n)_{n \in \mathbb{N}}$ une suite de fonctions $\mathcal{F} - \mathcal{B}_{\mathbb{R}}$ *mesurables* convergeant *μ -presque partout* vers une fonction f réelle et mesurable à son tour sur Ω .

Alors, $(\forall \varepsilon > 0)$ $(\exists A \in \mathcal{F})$ tel que : $\mu(A) < \varepsilon$ et f_n converge *uniformément* vers f sur $E \setminus A$.

Preuve :

Pour tout $(k, n) \in \mathbb{N}^*$ on considère la suite des ensembles $(B_k^n)_{(k,n) \in (\mathbb{N}^*)^2}$ définie par :

$$B_k^n = \bigcap_{i \geq n} \left\{ x \in E / |f_i(x) - f(x)| \leq \frac{1}{k} \right\}.$$

Pour tout $k \geq 1$ **fixé** , la suite $(B_k^n)_{n \in \mathbb{N}^*}$ est croissante pour l'inclusion ,donc et grâce à *la continuité croissante de la mesure* :

$$\mu \left(\bigcup_{n \geq 1} B_k^n \right) = \lim_{n \rightarrow \infty} \mu(B_k^n)$$

or,et comme la suite des fonctions $(f_n)_{n \in \mathbb{N}^*}$ converge simplement μ -p.p. vers f , on aura :

$$(\forall k \in \mathbb{N}^*) : \mu \left(\bigcup_{n \geq 1} B_k^n \right) = \lim_{n \rightarrow \infty} \mu(B_k^n) = \mu(E). \quad \star$$

donc,et car : $\mu(E) < +\infty$ on peut formuler \star ainsi :

$$(\forall k \in \mathbb{N}^*) (\forall \delta > 0) (\exists N_{(k,\delta)} \in \mathbb{N}^*) \text{ tel que : } (\forall n \geq N_{(k,\delta)}) : \underbrace{|\mu(E) - \mu(B_k^n)|}_{\text{positif}} = \mu(E) - \mu(B_k^n) < \delta. \quad \star \star$$

On **fixe** $\varepsilon > 0$, grâce à $\star \star$,et en posant $\delta := \left(\frac{\varepsilon}{2}\right)^k$ on peut trouver pour chaque $k \geq 1$ un entier $n_k := N_{\left(k, \left(\frac{\varepsilon}{2}\right)^k\right)}$ tel que : $\mu(B_k^{n_k}) \geq \mu(E) - 2^{-k} \varepsilon$.

Alors, l'ensemble : $A = \bigcup_{k \geq 1} (E \setminus B_k^{n_k})$ convient car :

$$\begin{aligned} \mu(A) &= \mu \left(\bigcup_{k \geq 1} (E \setminus B_k^{n_k}) \right) \\ &\leq \sum_{k \geq 1} \mu(E \setminus B_k^{n_k}) \quad (\text{grâce à la sous } \sigma\text{-additivité de la mesure}) \\ &\leq \sum_{k \geq 1} \mu(E) - \mu(B_k^{n_k}) \\ &\leq \sum_{k \geq 1} 2^{-k} \varepsilon = \varepsilon. \end{aligned}$$

et : $A = \bigcup_{k \geq 1} (E \setminus B_k^{n_k}) = E \setminus \left(\bigcap_{k \geq 1} B_k^{n_k} \right)$, donc : $E \setminus A = \left(\bigcap_{k \geq 1} B_k^{n_k} \right) = \bigcap_{k \geq 1} \bigcap_{i \geq n_k} \left\{ x \in E / |f_i(x) - f(x)| \leq \frac{1}{k} \right\}$.
d'où ce qu'il faut prouver. ■

Remarque 1.1.1 : (l'importance de finitude de la mesure) :

1. Considérons les fonctions mesurables f_n définies sur l'ensemble des réels ,muni de la tribu borélienne et la mesure de Lebesgue, par : $(\forall n \in \mathbb{N}) : f_n = \mathbb{1}_{[n, +\infty[}$.
(où $\mathbb{1}_{[n, +\infty[}$ désigne la fonction indicatrice de l'ensemble $[n, +\infty[$), alors, il est clair que la suite (f_n) converge simplement (donc μ -presque partout) vers la fonction nulle, **mais** il n'existe aucun borélien de mesure finie sur le complémentaire duquel cette convergence est uniforme.
donc l'hypothèse que μ **soit fini est essentiel**.
2. Ce résultat est d'une **grande importance**, et ça se voit à travers son apparence dans plusieurs démonstrations de théorèmes axiaux de la théorie de mesure et d'intégration, voyez par exemples les corollaires ci dessous ,mais dans la cadre d'apprentissage : il va nous servir essentiellement pour prouver le fameux théorème de Lusin.

Corollaire 1.1.1 :

On suppose μ finie. La convergence $f_n \xrightarrow[n \rightarrow \infty]{p.p.} f$ implique la convergence $f_n \xrightarrow[n \rightarrow \infty]{\mu} f$.

Preuve :

Fixons $\alpha > 0$, grâce au théorème d'Egorov :
pour tout $\varepsilon > 0$ il existe A mesurable tel que : $\mu(E \setminus A) < \varepsilon$ et $f_n \xrightarrow[n \rightarrow \infty]{} f$ **uniformément** sur A .
On aura donc et pour n assez grand de sorte que : $\sup_{x \in A} (|f_n - f| \leq \alpha)$:^a

$$\mu(|f_n - f| > \alpha) \leq \mu(\mathbb{C}_E^A) + \mu(\sup_{x \in A} |f_n - f| > \alpha) < \varepsilon.$$

ce qui achève la preuve. ■

Corollaire 1.1.2 : (théorème de convergence bornée)

Soit $(f_n)_{n \in \mathbb{N}^*}$ une suite de fonctions mesurables $(f_n : \mathbb{R}^d \rightarrow \mathbb{R})$ satisfaisant :

- il existe une constante $M > 0$ avec $|f_n| \leq M$ pour tout $n \geq 1$.
- il existe $E \subset \mathbb{R}^d$ mesurable avec $\mu(E) < \infty$ et $supp(f_n) \subset E$ pour tout $n \geq 1$.
- $f_n(x) \xrightarrow[n \rightarrow \infty]{} f(x)$ pour presque tout $x \in E$.

alors , la fonction limite f est **mesurable**, satisfait : $supp(f) \subset E$, et en plus :

$$\int_{\mathbb{R}^d} |f_n - f| \, d\mu \xrightarrow[n \rightarrow +\infty]{} 0$$

ce qui fait que : $\lim_{n \rightarrow +\infty} \int_{\mathbb{R}^d} f_n \, d\mu = \int_{\mathbb{R}^d} \lim_{n \rightarrow +\infty} f_n \, d\mu = \int_{\mathbb{R}^d} f \, d\mu$.

a. un tel n existe bien évidemment car sur $A : f_n \xrightarrow[n \rightarrow \infty]{} f$ **uniformément**.

Preuve :

D'après les hypothèses, on voit immédiatement que la fonction limite f est bornée par la même constante M . On voit aussi que f s'annule hors de E , car : $supp(f) := \{x \in E / f(x) \neq 0\} \subset E$. soit $\varepsilon > 0$ le **théorème d'Egorov** nous permet de trouver un sous-ensemble mesurable $E_\varepsilon \subset E$ avec : $\mu(E \setminus E_\varepsilon) \leq \varepsilon$ en restriction auquel on a la convergence uniforme : $f_n/E_\varepsilon \xrightarrow{n \rightarrow \infty} f/E_\varepsilon$

Alors sur E_ε , nous pouvons trouver un entier N_ε assez grand pour lequel :

$$n \geq N_\varepsilon \implies (\forall x \in E_\varepsilon) : |f_n(x) - f(x)| \leq \varepsilon$$

donc, toujours pour $n \geq N_\varepsilon$:

$$\begin{aligned} \int_E |f_n - f| \, d\mu &= \int_{E_\varepsilon} |f_n - f| \, d\mu + \int_{E \setminus E_\varepsilon} |f_n - f| \, d\mu \\ &\leq \varepsilon \mu(E_\varepsilon) + 2M \mu(E \setminus E_\varepsilon) \\ &\leq \varepsilon \mu(E) + 2M \cdot \varepsilon \\ &\leq \varepsilon (\mu(E) + 2M) \end{aligned}$$

ce qui achève la preuve car : $\varepsilon > 0$ était arbitraire. ■

1.2 Théorème de densité des fonctions continues dans les espaces L^p :

Définition 1.2.1 : (la distance d'un point à une partie :)

Soit (X, d) un espace métrique et Y une de ses sous-partie, on définit la distance entre un élément $x \in X$ et Y par : $(\forall x \in X) : d(x, Y) := \begin{cases} \inf \{d(x, y) / y \in Y\} & \text{si } Y \neq \emptyset \\ +\infty & \text{sinon} \end{cases}$

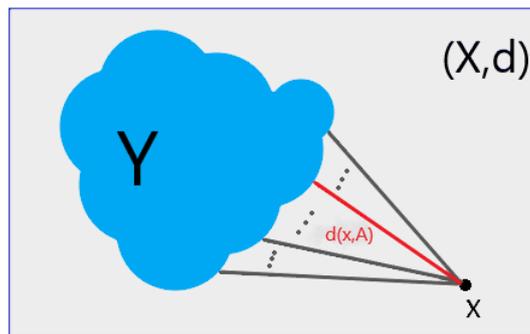


FIGURE I.1: Illustration de la distance de x à Y .

Lemme 1.2.1 :

Soit (X, d) un espace métrique, pour tout $Y \neq \emptyset$, un sous ensemble **fixé** de X , l'application : " $d(\cdot, Y)$ " qui à tout élément : x dans X associe l'image $d(x, Y) = \inf \{d(x, y) / y \in A\}$ est bien définie, **continue, et même 1-lipschitzienne**. (voir :[31])

Preuve :

Soit x et t dans X , montrons que : $d(\cdot, Y)$ est 1-lipschitzienne.
 prenons $\varepsilon > 0$, Par définition de $d(\cdot, Y)$ et qui est sous forme d'une borne inférieure, il existe $y_0 \in Y$ tel que : $d(x, y_0) \leq d(x, Y) + \varepsilon$, alors $d(t, Y) \leq d(t, x) + d(x, y_0) \leq d(t, x) + d(x, Y) + \varepsilon$, donc : $d(t, Y) \leq d(x, Y) + d(t, x) + \varepsilon$, et en faisant tendre ε vers 0 on obtient : $d(t, Y) \leq d(x, Y) + d(x, t)$.
 D'une façon analogue on aura $d(x, Y) \leq d(t, Y) + d(x, t)$, donc $|d(x, Y) - d(t, Y)| \leq d(x, t)$
 ce qui fait de $d(\cdot, Y)$ une application 1-lipschitzienne et par suite on déduit immédiatement qu'elle est continue.

Remarque 1.2.1 :

Cette proposition servira un peu partout, par exemple pour l'approximation des ouverts par des compacts dans un espace de dimension fini, ou pour démontrer le théorème de projection sur un **convexe complet**^a, et il sera encore très utile pour approcher des fonctions intégrables par des fonctions C^∞ ^b, et en plus, elle va apparaître pour nous sortir vraiment d'un grand pétrin dans la preuve du théorème III.1.3.1, portant sur l'approximation des fonctions de décision par des réseaux de neurones à une seule couche cachée.

Lemme 1.2.2 :

Soit (X, d) un espace métrique, \mathfrak{B} sa tribu **borélienne** et μ une **mesure finie** sur (X, \mathfrak{B}) ; pour tout borélien $B \in \mathfrak{B}$ et tout $\varepsilon > 0$, il existe un fermé F et un ouvert U de X , tels que : $F \subset B \subset U$, $\underbrace{\mu(U \setminus F)}_{\text{car la mesure est finie}} = \mu(U) - \mu(F) < \varepsilon$.

Preuve :

Notons \mathfrak{A} la classe de parties de X formée de tous les $A \in \mathfrak{B}$ tel que : pour tout $\varepsilon > 0$ **il existe** un fermé F et un ouvert U , de sorte que $F \subset A \subset U$ et $\mu(U \setminus F) < \varepsilon$; on va tout simplement essayer de montrer que la classe \mathfrak{A} **est une tribu qui contient tous les ouverts de X** : on en déduira que $\mathfrak{A} = \mathfrak{B}$, ce qui est le résultat voulu.

► **vérifions que \mathfrak{A} est une tribu :**

- il est évident que $X \in \mathfrak{A}$ car : $X \subset X \subset X$ et : $(\forall \varepsilon > 0) : \mu(X \setminus X) = \mu(\emptyset) = 0 < \varepsilon$.
- Il est clair d'abord que \mathfrak{A} **est stable par complémentaire** : or, pour $A \in \mathfrak{A}$ le complémentaire A^c vérifie l'encadrement $U^c \subset A^c \subset F^c$, et car : $F^c \setminus U^c = U \setminus F$, on déduit que : $\mu(U \setminus F) = \mu(F^c \setminus U^c) = \mu(U) - \mu(F) < \varepsilon$ ce qui fait : $A^c \in \mathfrak{A}$.
- **Pour la stabilité de la réunion dénombrable** : on prend une suite $(A_n)_{n \in \mathbb{N}} \subset \mathfrak{A}$, d'éléments de \mathfrak{A} on a : pour tout $\delta > 0$ et $n \in \mathbb{N}$ **il existe** un fermé F_n et un ouvert U_n , de sorte que $F_n \subset A_n \subset U_n$ et $\mu(U \setminus F) < \varepsilon$ on prend $\varepsilon > 0$ et on choisit $F_n \subset A_n \subset U_n$ pour tout $n \geq 0$ tels que $\mu(U_n \setminus F_n) < \varepsilon \left(\frac{1}{2}\right)^{n+2}$; on considère l'ouvert : $U = \bigcup_{n \in \mathbb{N}} U_n$, on pose $A = \bigcup_{n \in \mathbb{N}} A_n$ et $Y = \bigcup_{n \in \mathbb{N}} F_n$; bien sur Y **n'est pas fermé en général**, mais on a $Y \subset A \subset U$, et la relation :

$$U = \bigcup_{n \in \mathbb{N}} U_n \setminus \bigcup_{n \in \mathbb{N}} F_n \subset \bigcup_{n \in \mathbb{N}} (U_n \setminus F_n).$$

ce qui entraîne que : $\mu(U \setminus Y) \leq \sum_{n=0}^{+\infty} \mu(U_n \setminus F_n) \leq \sum_{n=0}^{+\infty} \varepsilon \left(\frac{1}{2}\right)^{n+2} = \frac{\varepsilon}{2}$.

Selon la continuité monotone de la mesure, on aura pour n assez grand :

$\mu(Y) - \mu(F_0, \dots, F_n) < \frac{\varepsilon}{2}$ et le fermé $F = F_0 \cup \dots \cup F_n$ fournit alors l'encadrement $F \subset A \subset U$ avec :

$$\mu(U) - \mu(F) = \mu(U) - \mu(Y) + \mu(Y) - \mu(F) = \underbrace{\mu(U \setminus Y)}_{\text{car : } Y \subset F} + \underbrace{\mu(Y \setminus F)}_{\text{car : } F \subset Y} < \frac{\varepsilon}{2} + \frac{\varepsilon}{2} = \varepsilon.$$

► **il reste à voir que \mathfrak{A} contient tous les ouverts de X :**

a. Attention : un convexe **complet** en général et pas juste fermé.
 b. Ce qui est bien notre cas.
 c. n'oubliez pas que **pour la topologie induite de sa distance**, X est **fermé** et **ouvert** au même temps.

soit U un ouvert de X , on pose pour tout $n \geq 0$ $F_n = \{x \in X / d(x, U^c) \geq 2^{-n}\}$, Il est clair que F_n est fermé ^a, et que $F_n \subset U$; de plus, comme U est ouvert, tout point $x \in U$ vérifie $d(x, U^c) > 0$ ^b donc appartient à F_n pour n assez grand. Il en résulte que $U = \bigcup_{n \in \mathbb{N}} F_n$. Si on donne $\varepsilon > 0$, on aura $\mu(U) - \mu(F_n) < \varepsilon$ pour n assez grand. On choisit alors l'encadrement $F_n \subset U \subset U$ montre que $U \in \mathcal{A}$.

donc et en guise de conclusion : \mathcal{A} est une tribu qui contient la topologie de X (c'est à dire l'ensemble de ses ouverts), alors il contiendra aussi la tribu engendrée par cette topologie et qui n'est autre que la tribu borélienne \mathcal{B} d'où $\mathcal{B} \subset \mathcal{A}$ et pour l'autre inclusion elle est triviale car \mathcal{A} ne contient que des ensembles mesurables, ce qui achève la preuve. ■

Théorème 1.2.1 :

De même que [60], on désignera par $C_c^0(\mathbb{R}^d)$, l'espace des fonctions continues à support compact définies sur \mathbb{R}^d . alors pour tout p tel que : $1 \leq p < +\infty$ l'espace $C_c^0(\mathbb{R}^d)$ est dense dans : $L^p(\mathbb{R}^d)$.

Preuve :

On veut approcher $f \in L^p(\mathbb{R}^d)$ par une fonction φ continue à support compact, au sens de norme L^p , et pour cela on va procéder à travers plusieurs réductions :

puisque $f = f^+ - f^-$, il suffit de le faire pour $f \geq 0$, or dans ce cas, f est une limite simple d'une certaine suite croissante $(f_n)_{n \in \mathbb{N}}$ de fonctions numériques étagées mesurables positives, $0 \leq f_n \leq f$.

La suite $(f - f_n)^p$ tend vers 0 en étant dominée. par la fonction intégrable f^p , donc $\|f - f_n\|_p \xrightarrow{n \rightarrow +\infty} 0$ par convergence dominée ^d

Il suffit donc de pouvoir approcher toute fonction étagée g . Puisque g est sous la forme d'une combinaison linéaire de fonctions de la forme $\mathbb{1}_B$ avec $B \in \mathcal{B}_{\mathbb{R}^d}$, il suffit d'approcher les fonctions indicatrices $\mathbb{1}_B$.

Enfin, si on pose $B_n = B \cap B(0, n[$, on montre comme avant que $\mathbb{1}_{B_n}$ tend vers $\mathbb{1}_B$ dans $L^p(\mathbb{R}^d)$, par convergence dominée. Finalement on veut approcher une fonction indicatrice $\mathbb{1}_B$ d'un borélien borné B .

on a $B \subset B(0, R[\subset K = \overline{B(0, R]} = B(0, R]$, introduisons aussi l'ouvert $V = B(0, R + 1[$; on a $B \subset K \subset V$, désignons par μ la mesure sur $(\mathbb{R}^d, \mathcal{B}_{\mathbb{R}^d})$, définie par : $(\forall A \in \mathcal{B}_{\mathbb{R}^d}) : \mu(A) = \lambda(A \cap V)$.

C'est une mesure finie à laquelle le lemme 1.1.2.2 s'applique.

On trouve ainsi, un fermé F et un ouvert U de \mathbb{R}^d tels que : $F \subset B \subset U$ et : $\mu(U) - \mu(F) < \varepsilon$, c'est à dire : $\lambda(U \cap V) - \lambda(F \cap V) < \varepsilon$ le fermé F est borné. puisque $F \subset K$, donc $K_1 = F$ est compact et

$$= \lambda(F) \text{ car : } F \subset B \subset K \subset V$$

l'ouvert $U_1 = U \cap V$ contient B . On a donc $K_1 \subset B \subset U_1$ et $\lambda(U_1 \setminus K_1) < \varepsilon$.

Maintenant, il reste à construire une fonction continue $\varphi \in C_c^1(\mathbb{R}^d)$, qui vérifie :

$$\|\mathbb{1}_B - \varphi\|_p := \left(\int_{\Omega} |\mathbb{1}_B - \varphi|^p d\mu \right)^{\frac{1}{p}} < \varepsilon$$

a. car et **selon le lemme 1.1.2.1** $d(x, U^c)$ est continue et l'image réciproque d'un fermé par une fonction continue reste toujours un fermé.

b. on rappelle que dans un espace métrique (X, d) on a : $(\forall A \in \mathcal{P}(X)) (\forall x \in X) : x \in \overline{A} \iff d(x, A) = 0$.

c. On appelle support de f , noté $\text{supp}(f)$, l'adhérence de l'ensemble des points en lesquels la fonction ne s'annule pas. $\text{supp}(f) = \{x \in X \mid f(x) \neq 0\}$.

d. je pense dans ce cas qu'on peut utiliser juste le théorème de la convergence monotone.

et pour cela, on va la choisir d'une façon, telle que : $\mathbb{1}_{K_1} \leq \varphi \leq \mathbb{1}_{U_1}$.
on pose d'abord :

$$\rho = d(K_1, U_1^c) := \inf \left\{ d(x, U_1^c) / x \in K_1 \right\}$$

. on a : $\rho \geq 0$ puisque cette quantité est la borne inférieure d'un ensemble d'éléments positifs, et en plus, on peut remarquer que : $\rho \neq 0$, car s'il n'est pas le cas, il existera forcément^a une suite $(x_n)_{n \in \mathbb{N}}$ d'éléments de K_1 qui vérifie : $\lim_{n \rightarrow +\infty} d(x_n, U_1^c) = 0$, or K_1 est compact, donc selon le théorème de **Bolzano-Weierstrass**, il existe une sous suite $(x_{\psi(n)})_{n \in \mathbb{N}}$ qui converge vers un élément : \bar{x} de K_1 . une autre fois, on fera appel au **lemme 1** qui va nous justifier l'égalité :

$$\lim_{n \rightarrow +\infty} d(x_{\psi(n)}, U_1^c) = d \left(\lim_{n \rightarrow +\infty} x_{\psi(n)}, U_1^c \right) = d(\bar{x}, U_1^c) = 0$$

et cela implique que : $\bar{x} \in \overline{U_1^c} = U_1^c$ ^b, ce qui est impossible car : $K_1 \subset U_1$ d'où $\rho > 0$; alors et par conséquent on peut aisément donner un sens à la fonction h définie par :

$$h(x) := 1 - \frac{d(x, K_1)}{\rho}.$$

de nouveau, on va se reposer sur le **lemme 1** pour mettre en évidence la continuité de cette fonction. et par suite on peut déduire que la fonction φ définie par :

$$(\forall x \in X) : \varphi(x) := \begin{cases} h(x) & \text{si : } h(x) \geq 0 \\ 0 & \text{sinon} \end{cases}$$

est continue à son tour et vérifie : $0 \leq \varphi(x) \leq 1$.

or, on remarque que : si $x \in K_1$, on a $d(x, K_1) = 0$ ce qui fait que : $\varphi(x) = 1$; et si $x \notin U_1$, on aura : $d(x, K_1) \geq \rho$ ce qui donne : $\varphi(x) = 0$; tout ceci montre que : $\mathbb{1}_{K_1} \leq \varphi \leq \mathbb{1}_{U_1}$ ^c, et comme on a aussi : $\mathbb{1}_{K_1} \leq \mathbb{1}_B \leq \mathbb{1}_{U_1}$ ^d, il en résulte que : $|\mathbb{1}_B - \varphi| \leq \mathbb{1}_{U_1} - \mathbb{1}_{K_1} = \mathbb{1}_{U_1 \setminus K_1}$, donc et enfin de compte :

$$\|\mathbb{1}_B - \varphi\|_p := \left(\int_{\Omega} |\mathbb{1}_B - \varphi|^p d\mu \right)^{\frac{1}{p}} = \mu(U_1 \setminus K_1)^{\frac{1}{p}} < \varepsilon$$

d'où ce qu'il faut prouver. ■

1.3 Premier théorème de Lusin :

Le théorème de Lusin a été cité pour la première fois en 1903 par Henri Lebesgue, qui a remarqué, dans un éclair de génie, que **toute fonction intégrable est presque continue**, mais cet énoncé restait considéré comme une conjecture, jusqu'au 1905 quand Giuseppe Vitali a réussi de l'établir, ce qui était redécouvert et complété en 1912 par Nicolas Lusin.

L'idée de base de ce théorème, est d'affirmer que toute fonction mesurable possède une restriction **continue** à une **grande** partie de son domaine de définition.

Lemme 1.3.1 : (La réciproque partielle de théorème de la convergence dominée :)

Soit $(\Omega, \mathcal{F}, \mu)$ un espace mesuré, on note : $\mathcal{L}^1(\Omega, \mathcal{F}, \mu)$ l'espace des fonctions **réelles** μ -intégrables.

Soit $(f_n)_{n \in \mathbb{N}}$ une suite de $\mathcal{L}^1(\Omega, \mathcal{F}, \mu)$ et $f \in \mathcal{L}^1(\Omega, \mathcal{F}, \mu)$ telles que : $\int_{\Omega} |f_n - f| d\mu \xrightarrow{n \rightarrow \infty} 0$, alors :

1. il existe une sous-suite (f_{n_k}) qui converge μ -p.p. vers la limite f .
2. il existe $h \in \mathcal{L}^1(\Omega, \mathcal{F}, \mu)$ telle que : $f_{n_k} \leq h$ μ -p.p.

a. selon la définition de la borne inférieure.

b. U est ouvert, ce qui fait que U^c est un fermé, d'où $\overline{U_1^c} = U_1^c$.

c. On peut montrer ça à travers une simple double inégalité.

d. Car : $K_1 \subset B \subset U_1$.

Preuve :

On observe que $(f_n)_{n \in \mathbb{N}^*}$ est une suite de Cauchy dans $\mathcal{L}^1(\Omega, \mathcal{F}, \mu)$, donc on peut définir une sous suite $(f_{n_k})_{k \in \mathbb{N}}$ par récurrence de telle sorte que :

$$(\forall k \geq 1) : \|f_{n_k} - f_{n_{k+1}}\|_1 \leq 2^{-k}.$$

on pose pour tout $x \in \Omega$ et $p \geq 1$:

$$g_p(x) := \sum_{k=1}^p |f_{n_{k+1}}(x) - f_{n_k}(x)| \text{ et } g(x) := \sum_{k=1}^{+\infty} |f_{n_{k+1}}(x) - f_{n_k}(x)| \in \overline{\mathbb{R}}^+.$$

il est clair que $\|g_p\|_1 \leq \sum_{k=1}^p 2^{-k} \leq 1$, donc d'après le **théorème de convergence dominée**, on en déduit que $g \in \mathcal{L}^1(\Omega, \mathcal{F}, \mu)$ et $g_n \xrightarrow[n \rightarrow \infty]{} g$ dans $\mathcal{L}^1(\Omega, \mathcal{F}, \mu)$, On observe aussi que pour $l > k$:

$$|f_{n_l} - f_{n_k}| \leq \sum_{i=k}^{l-1} |f_{n_{i+1}}(x) - f_{n_i}(x)| = g_l - g_k \leq g - g_k.$$

Il en résulte que (f_{n_k}) est $(\mu$ -p.p.) une suite de Cauchy donc $(\mu$ -p.p) une suite convergente ^a, en notant $f^* = \lim_{k \rightarrow \infty} f_{n_k}$ sa limite $(\mu$ -p.p) et en faisant tendre $l \rightarrow +\infty$ dans l'inégalité précédente, on obtient :

$$|f^* - f_{n_k}| \leq g - g_k \leq g.$$

En utilisant le théorème de convergence dominée pour la suite $f^* - f_{n_k}$, on déduit que $f_{n_k} \xrightarrow[k \rightarrow +\infty]{} f^*$, Par unicité de la limite, on a donc $f^* = f \mu$ p.p. et on conclut en posant $h := g + f^*$. ■

Théorème 1.3.1 : (de Lusin, version faible)

Soit $f : ([a, b], \mathcal{B}_{\mathbb{R}} \cap [a, b]) \rightarrow (\mathbb{R}, \mathcal{B}_{\mathbb{R}})$ une fonction **intégrable**, alors et selon [41], on a : pour tout $\varepsilon > 0$, il existe un **compact** $E_\varepsilon \subset [a, b]$ tel que : $\lambda([a, b] \setminus E_\varepsilon) \leq \varepsilon$ ^b et la restriction de f à E_ε **est continue au sens de la topologie induite sur E** , et pour résumer ça, on dit dans ce cas que **f est ε -presque continue sur $[a, b]$** .

Preuve :

Selon le théorème qu'on a traité précédemment portant sur la densité des fonctions continues à support compact dans les espaces L^p , on peut garantir dans ce cas, et en posant $p = 1$, qu'il existe une suite $(g_n)_{n \in \mathbb{N}}$ de fonctions continues, telle que : $g_n \xrightarrow[n \rightarrow +\infty]{} f$ dans L^1 , et via le **lemme précédent** qui a présenté une réciproque partielle de théorème de la convergence dominée, on peut de cette suite, extraire une sous-suite $(g_{n_k})_{k \in \mathbb{N}}$ telle que : $g_{n_k} \xrightarrow[k \rightarrow +\infty]{} f$ **presque partout** ensuite, et en utilisant le théorème d'Egorov, on trouvera que g_{n_k} converge **uniformément** sauf sur un ensemble de mesure inférieure à un certain ε choisi arbitrairement vers la fonction : f , et puisque la limite uniforme ^c d'une suite de fonctions continues reste encore continue, notre preuve sera bien finie. ■

Remarque 1.3.1 :

Ce résultat sera indispensable, pour montrer le théorème III.1.3.2, et qui va établir la capacité ou l'habileté des réseaux neuronaux, **pour l'apprentissage des fonction de décision**.

a. car \mathbb{R} est complet.
b. λ désigne la mesure de Lebesgue usuelle sur \mathbb{R} .
c. si elle existe bien sûr.

2 Aperçu sur la densité dans l'espace des fonctions continues :

2.1 Théorème de Dini :

Théorème 2.1.1 : (de Dini :)

Soit $(f_n)_{n \in \mathbb{N}}$ et f des fonctions à **valeurs réelles** définies et **continues** sur un espace topologique **compact** (X, Θ) , si $(f_n(x))_{n \in \mathbb{N}}$ est croissante pour tout $x \in X$, et si f_n converge vers f **simple-ment** alors : cette convergence sera uniforme.

Preuve :

On va reprendre pour établir ce résultat les idées de la preuve proposée, dans [47] page :272 en introduisant quelques modifications pour s'adapter au cas général des espaces compacts : on pose d'abord $(\forall x \in X) (\forall n \in \mathbb{N}) : \varphi_n(x) = f(x) - f_n(x) \geq 0$, on peut remarquer que la suite $(\varphi_n)_{n \in \mathbb{N}}$ est **décroissante et qu'elle converge simplement vers la fonction nulle**. fixons un nombre réel $\varepsilon > 0$ et considérons les ensembles $:\Lambda_n(\varepsilon) = \{x \in X / \varphi_n(x) < \varepsilon\}$, puisque les φ_n sont continues alors les $(\Lambda_n(\varepsilon))_{n \in \mathbb{N}, \varepsilon > 0}$ sont des ouverts.

on a :

$$\begin{aligned} (\forall x \in X) (\forall n \in \mathbb{N}) : f_n(x) \xrightarrow{n \rightarrow +\infty} 0 &\Leftrightarrow (\forall x \in X) (\forall \varepsilon > 0) (\exists \eta_{x,\varepsilon} > 0) \text{ tq : } (\forall n > \eta_{x,\varepsilon}) : |\varphi_n(x)| = \varphi_n(x) < \varepsilon \\ &\Rightarrow (\forall \varepsilon > 0) (\forall x \in X) (\exists \eta_{x,\varepsilon} > 0) \text{ tq : } (\forall n > \eta_{x,\varepsilon}) : x \in \Lambda_n(\varepsilon) \\ &\Rightarrow (\forall \varepsilon > 0) (\forall x \in X) (\exists M = [\eta_{x,\varepsilon}] + 1 > 0) \text{ tq : } (\forall n > \eta_{x,\varepsilon}) : x \in \Lambda_M(\varepsilon) \\ &\Rightarrow (\forall \varepsilon > 0) : X \subset \bigcup_{n \in \mathbb{N}} \Lambda_n(\varepsilon) \end{aligned}$$

comme X est compact, on peut extraire de $(\Lambda_n(\varepsilon))_{n \in \mathbb{N}}$ un sous-recouvrement fini, alors pour tout $\varepsilon > 0$ il existe un entier N_ε tel que $: X \subset \bigcup_{n \leq N_\varepsilon} \Lambda_n(\varepsilon)$.

maintenant soit $:\varepsilon > 0$, il est facile de voir que $:(\forall (n, m) \in \mathbb{N}^2) : n \leq m \implies \Lambda_n(\varepsilon) \subset \Lambda_m(\varepsilon)$.

or $:(\forall x \in X) (\forall n \leq m) : x \in \Lambda_n(\varepsilon) \implies \varphi_n(x) < \varepsilon \implies \varphi_m(x) \leq \varphi_n(x) < \varepsilon \implies x \in \Lambda_m(\varepsilon)$.

donc $: X \subset \bigcup_{n \leq N_\varepsilon} \Lambda_n(\varepsilon) = \Lambda_{N_\varepsilon}(\varepsilon) \iff X = \Lambda_{N_\varepsilon}(\varepsilon)$.

en se basant une autre fois sur l'hypothèse de la monotonie nous obtenons :

$\forall n \geq N_\varepsilon, \varphi_n(x) \leq \varphi_{N_\varepsilon}(x) < \varepsilon$, on déduit donc que la convergence de (f_n) vers 0 est uniforme sur X ce qui implique que $f_n \xrightarrow[n \rightarrow +\infty]{\text{unifor}} f$. ■

Remarque 2.1.1 :

le but général de théorème de Dini est d'énoncer des conditions sous lesquelles **la convergence simple d'une suite de fonctions implique la convergence uniforme**, ce théorème peut prendre d'autres versions plus généralisées, mais dans notre cas cet énoncé est suffisant et il va nous servir par la suite.

2.2 Théorème de Stone-Weierstrass :

Soit (X, d) un espace métrique **compact**, et $C(X, \mathbb{R})$ l'espace des fonctions continues^a de X dans \mathbb{R} , muni de la norme uniforme $:\|f\|_\infty = \sup_{x \in X} |f(x)|$.

Si $f, g \in C(X, \mathbb{R})$, on note $: f \times g \in C(X, \mathbb{R})$ la fonction $x \mapsto f(x) \cdot g(x)$, **muni de cette loi de produit, l'espace vectoriel $(C(X, \mathbb{R}), +, \cdot)$ devient une algèbre.**^b

a. On rappelle que $(C(X, \mathbb{R}), +, \cdot, \|\cdot\|)$ est un espace de Banach.

b. Une algèbre sur un corps commutatif K n'est autre qu'un K -espace vectoriel $(E, +, \cdot)$ muni d'une loi de composition interne bilinéaire " \times ".

Définition 2.2.1 :

1. On dit qu'une partie $\mathfrak{S} \subset C(X, \mathbb{R})$ est une sous-algèbre si et seulement si :
 $(\forall (f, g) \in C(X, \mathbb{R})^2) (\forall \lambda \in \mathbb{R})$ on a : $f + g \in \mathfrak{S}$, $f \times g \in \mathfrak{S}$, et $\lambda \cdot f \in \mathfrak{S}$.
2. On dit que \mathfrak{S} est une sous-algèbre unitaire si la fonction constante $\mathbb{1}$ appartient à \mathfrak{S} .
3. On dit que \mathfrak{S} sépare les points de X si : $(\forall x \neq y \in X) (\exists f \in \mathfrak{S})$ tel que : $f(x) \neq f(y)$.

Exemple 2.2.1 :

L'ensemble des fonctions polynomiales sur un segment $[a, b] \subset \mathbb{R}$ (où $a < b$) est une sous-algèbre unitaire qui sépare les points de $[a, b]$.

Remarque 2.2.1 :

1. Si $\mathfrak{S} \subset C(X, \mathbb{R})$ est une sous-algèbre unitaire, alors pour tout $f \in \mathfrak{S}$ et tout polynôme $P \in \mathbb{R}[X]$ on a : $P(f) \in \mathfrak{S}$.
2. Si $\mathfrak{S} \subset C(X, \mathbb{R})$ est une sous-algèbre unitaire, alors l'adhérence $\overline{\mathfrak{S}} \subset C(X, \mathbb{R})$ est encore une sous-algèbre unitaire (démonstration facile).

Théorème 2.2.1 :(de Stone-Weierstrass algébrique, [58])

Soit (X, d) un **espace métrique compact** et $C(X, \mathbb{R})$ l'espace des fonctions continues de X dans \mathbb{R} , muni de la norme uniforme, si $\mathfrak{S} \subset C(X, \mathbb{R})$ est une sous-algèbre unitaire qui sépare les points de X , alors \mathfrak{S} est dense dans $C(X, \mathbb{R})$.

Preuve :

On va suivre la même démarche que : [14], où la démonstration comportera 5 étapes, dont 3 sont préliminaires. Le but des deux premières étapes est de montrer que, si $f_1, f_2, \dots, f_n \in \mathfrak{S}$, alors :
 $\max_{1 \leq i \leq n} (f_i) \in \overline{\mathfrak{S}}$, et $\min_{1 \leq i \leq n} (f_i) \in \overline{\mathfrak{S}}$.

► **Étape 1** : Il existe une suite de fonctions polynomiales (P_n) sur $[-1, 1]$ convergeant uniformément vers la fonction : $t \mapsto |t|$.

définissons cette suite par récurrence en posant :
$$\begin{cases} (\forall t \in [-1, 1]) : P_0(t) = 0 \\ (\forall t \in [-1, 1]) : P_{n+1}(t) = P_n(t) + \frac{1}{2} (t^2 - P_n(t)^2) \end{cases}$$

Il est facile de vérifier, par récurrence, que :

$$(\forall n \in \mathbb{N}) (\forall t \in [-1, 1]) : 0 \leq P_n(t) \leq |t| \quad (*)$$

en effet, c'est vrai pour $n = 0$ et l'hypothèse de récurrence nous fournit déjà que : $P_{n+1}(t) \geq 0$. D'autre part :

$$\begin{aligned} |t| - P_{n+1}(t) &= |t| - P_n(t) - \frac{1}{2} (|t| - P_n(t)) (|t| + P_n(t)) \\ (\text{et car : } |t| + P_n(t) &\leq 2|t| \leq 2 \text{ sur } [-1, 1]) = (|t| - P_n(t)) \left(1 - \frac{1}{2} (|t| + P_n(t)) \right) \geq 0 \end{aligned}$$

l'inégalité (*) montre aussi que $P_{n+1}(t) \geq P_n(t)$, donc pour tout $t \in [-1, 1]$ la suite de nombres réels $n \mapsto P_n(t)$ est croissante, et car elle est majorée alors, elle sera convergente. notons : $u(t) = \lim_{n \rightarrow +\infty} P_n(t)$. La formule de récurrence pour P_n donne à la limite :

$$(\forall t \in [-1, 1]) : u(t) = u(t) + \frac{1}{2} (t^2 - u(t)^2), \text{ donc : } u(t)^2 = t^2.$$

et comme $u(t) \geq 0$, on obtient que : $u(t) = |t|$, le théorème de Dini 1.2.1.1, nous permet de déduire la convergence uniforme de P_n vers la fonction $t \mapsto |t|$.

► **Étape 2** : Si $f, g \in \mathfrak{S}$, alors les fonctions $|f|, \max(f, g),$ et $\min(f, g)$ appartiennent à $\overline{\mathfrak{S}}$. montrons d'abord que, si $f \in \mathfrak{S}$, alors $|f| \in \overline{\mathfrak{S}}$. c'est évident si $f = 0$. Dans le cas contraire, on note $h = \frac{f}{\|f\|_\infty}$ et on observe que $h \in \mathfrak{S}$ et $h(x) \in [-1, 1]$ pour tout $x \in X$. Ainsi, d'après l'étape 1 :

$$\frac{|f|}{\|f\|_\infty} = |h| = \lim_{n \rightarrow \infty} P_n(h). \quad (\text{la convergence étant uniforme sur } X.)$$

Comme $P_n(h) \in \mathfrak{S}$ ^a pour tout $n \in \mathbb{N}$, on conclut que : $\frac{|f|}{\|f\|_\infty} \in \overline{\mathfrak{S}}$, donc $|f| \in \overline{\mathfrak{S}}$. soient maintenant $f, g \in \mathfrak{S}$, alors :

$$\max(f, g) = \frac{f + g + |f - g|}{2} \in \overline{\mathfrak{S}} \quad \text{et} \quad \min(f, g) = \frac{f + g - |f - g|}{2} \in \overline{\mathfrak{S}}$$

avec une simple itération de ce résultat on obtient que :

$$(\forall f_1, f_2, \dots, f_n) \in \mathfrak{S} : \max_{1 \leq i \leq n} (f_i) \in \overline{\mathfrak{S}}, \text{ et } \min_{1 \leq i \leq n} (f_i) \in \overline{\mathfrak{S}}.$$

► **Étape 3** : Pour tous les $x, y \in X$ tels que $x \neq y$ et tous les $\alpha, \beta \in \mathbb{R}$, il existe $g \in \mathfrak{S}$ tel que $g(x) = \alpha$ et $g(y) = \beta$.

En effet, comme \mathfrak{S} sépare les points de X , il existe $h \in \mathfrak{S}$ tel que $h(x) \neq h(y)$. on prend alors g la fonction définie par :

$$g : z \in X \mapsto \alpha + (\beta - \alpha) \frac{h(z) - h(x)}{h(y) - h(x)} = \lambda \cdot h(z) + \mu, \text{ comme } \mathfrak{S} \text{ contient les constantes, on a bien } g \in \mathfrak{S}.$$

► **Étape 4** : soit $f \in C(X, \mathbb{R})$ et $x_0 \in X$ on a :

$$(\forall \varepsilon > 0) (\exists g \in \mathfrak{S}) \text{ tel que : } g(x_0) = f(x_0), \text{ et } (\forall x \in X) : g(x) < f(x) + \varepsilon.$$

En effet, pour chaque $y \in X \setminus \{x_0\}$, l'étape 3 fournit l'existence d'une fonction $g_y \in C(X, \mathbb{R})$ telle que : $g_y(x_0) = f(x_0)$ et $g_y(y) = f(y)$. Comme g_y et f sont continues, l'ensemble :

$$\mathcal{U}_y = \{z \in X / g_y(z) < f(z) + \varepsilon\}$$

est un ouvert de X contenant x_0 et y . il s'ensuit que la famille $(\mathcal{U}_y)_{y \in X \setminus \{x_0\}}$ est un recouvrement ouvert de l'espace X **qui est compact**, donc :

il existe donc des points $y_1, \dots, y_n \in X \setminus \{x_0\}$ tels que $X = \bigcup_{i=1}^n \mathcal{U}_{y_i}$.

Soit à présent $g = \min(g_{y_1}, \dots, g_{y_n})$. alors : $g \in \overline{\mathfrak{S}}$ d'après la première étape, et $g(x_0) = f(x_0)$. en outre, pour tout $x \in X$, il existe $i \in 1, \dots, n$ tel que $x \in \mathcal{U}_{y_i}$, donc $g(x) \leq g_{y_i}(x) < f(x) + \varepsilon$.

► **Étape 5** : soit $f \in C(X, \mathbb{R})$ on a :

$$(\forall \varepsilon > 0) (\exists h \in \mathfrak{S}) \text{ tel que : } (\forall x \in X) : f(x) - \varepsilon < h(x) < f(x) + \varepsilon.$$

en effet, pour chaque $x \in X$, l'étape 4 fournit une fonction $h_x \in \overline{\mathfrak{S}}$ telle que $h_x(x) = f(x)$ et $h_x < f + \varepsilon$ sur tout X . l'ensemble : $V_x = \{z \in X / f(z) < h_x(z) + \varepsilon\}$ est un ouvert de X contenant x , et la famille

a. on a déjà signalé ça au premier point de la remarque :1.2.2.1

$(V_x)_{x \in X}$ constitue un recouvrement ouvert de l'espace X , **qui est compact**. donc il existe des points $x_1, \dots, x_m \in X$ tels que :

$$X = \bigcup_{i=1}^m V_{x_i}.$$

Soit maintenant $h = \max(h_{x_1}, \dots, h_{x_m})$. **d'après l'étape 1**, $h \in \overline{S}$, et il est clair que : $h(x) < f(x) + \varepsilon$ pour tout $x \in X$, puisque : $(\forall i \in \{1, \dots, m\}) : h_{x_i} < f + \varepsilon$. d'autre part, pour tout $z \in X$, il existe : $i \in \{1, \dots, m\}$ tel que : $z \in V_{x_i}$, donc : $f(z) - \varepsilon < h_{x_i}(z) \leq h(z)$.
d'où : $(\forall z \in X) : f(z) - \varepsilon < h_{x_i}(z) < f(z) + \varepsilon$. ■

Corollaire 2.2.1 : (Théorème d'approximation de Weierstrass :)

Si $f : [a, b] \rightarrow \mathbb{R}$ est une fonction continue, alors pour tout $\varepsilon > 0$ il existe une fonction polynomiale P telle que :

$$\|f - P\|_{\infty} = \sup_{x \in [a, b]} |f(x) - P(x)| \leq \varepsilon.$$

ce qui fait que :

$$\left(\forall f \in C^0([a, b]) \right) \left(\exists (P_n)_{n \in \mathbb{N}} \in (\mathbb{R}[X])^{\mathbb{N}} \right) \text{ tel que : } P_n \xrightarrow[n \rightarrow +\infty]{\text{uniforme}} f \text{ sur } [a, b].$$

Preuve :

il suffit de remarques que l'espace des fonctions polynomiales sur $[a, b]$ est une sous algèbre unitaire qui sépare les points de $[a, b]$ (car la fonction identité est en particulier polynomiale) et d'appliquer le théorème :I.2.2.1. ■

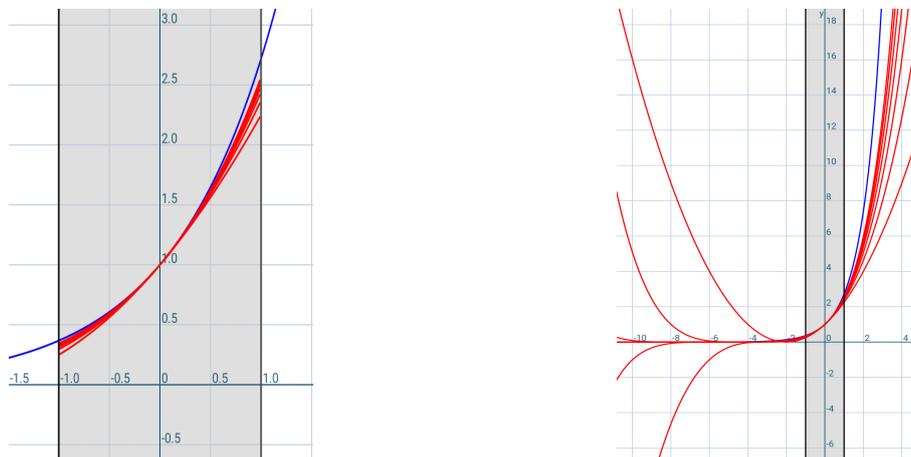


FIGURE I.2: Une suite de polynômes convergeant uniformément vers $x \mapsto \exp(x)$ sur $[-1, 1]$.

Remarque 2.2.2 :

Le théorème de Stone-Weierstrass, ne reste plus valable sans l'hypothèse de compacité, voir par exemple la fonction définie sur \mathbb{R}^+ par : $f(x) = e^{-x}$, et qui ne peut jamais être approchée uniformément par une suite des polynômes car :

$$\left(\forall P \in \mathbb{R}[X] \right) : \lim_{x \rightarrow +\infty} |P(x)| = +\infty \Rightarrow \lim_{x \rightarrow +\infty} |P(x) - f(x)| \geq \lim_{x \rightarrow +\infty} |P(x)| - \overbrace{\lim_{x \rightarrow +\infty} |f(x)|}^{=0} = +\infty.$$

3 Passage de la mesurabilité à la presque continuité :

3.1 Mesures régulières et espaces de Radon :

Définitions 3.1.1 :

Soit (X, Θ) un *espace topologique de Hausdorff* (appelé aussi espace séparé)^a, et ν une mesure sur la tribu Borélienne :

$$\mathfrak{B}_X := \tau(\Theta).$$

1. La mesure ν est appelé **régulière interne** ou serré si et seulement si :

$$(\forall B \in \mathfrak{B}_X) : \nu(B) = \sup \left\{ \nu(K) / K^c \in \Theta \text{ et } : B \supset K \right\}.$$

2. La mesure ν est appelée **régulière externe** si et seulement si :

$$(\forall B \in \mathfrak{B}_X) : \nu(B) = \inf \left\{ \nu(U) / U \in \Theta \text{ et } : B \subset U \right\}.$$

3. La mesure ν est dite **régulière** si et seulement si :elle vérifie les deux points précédents au même temps.

4. La mesure ν est appelée **localement finie** si et seulement si :

$$(\forall x \in X) (\exists U \in \Theta) : x \in U \text{ et } \nu(U) < +\infty.$$

5. La mesure ν est appelée **mesure du Radon** si elle est régulière et localement finie à la fois, et dans ce cas le triplet (X, \mathfrak{B}_X, ν) sera appelé un **espace mesuré de Radon**.

(voir pour ces définitions :[27]).

Proposition 3.1.1 :

Soit (X, d) un espace métrique, on note \mathfrak{B}_X sa tribu borélienne, et on considère une mesure μ définie sur cette tribu, si μ est **finie**, alors elle sera forcément régulière.

Preuve :

On défini l'ensemble $\mathcal{R} \subset \mathfrak{B}_X$, par :

$$(\forall A \in \mathcal{P}(X)) : A \in \mathcal{R} \iff \begin{cases} \mu(A) = \sup \left\{ \mu(K) / K \text{ est fermé et } A \supset K \right\} \\ \text{et :} \\ \mu(A) = \inf \left\{ \mu(U) / U \text{ est ouvert et } A \subset U \right\} \end{cases}$$

on veut établir que : $\mathcal{R} = \mathfrak{B}_X$, et pour cela ,on va montrer que : \mathcal{R} est une tribu et que : $\Theta \subset \mathcal{R}$.

► vérifions que \mathcal{R} est une tribu :

- il est évident que $X \in \mathcal{R}$ car :

$$\{U \text{ est ouvert et } X \subset U\} = \{X\} \implies \mu(X) = \inf \left\{ \mu(U) / U \text{ est ouvert et } A \subset U \right\}$$

$$\text{et } : X \in \{K \text{ est fermé et } A \supset K\} \implies \mu(X) = \sup \left\{ \mu(K) / K \text{ est fermé et } A \supset K \right\}^b$$

- Il est clair d'abord que \mathcal{R} **est stable par complémentaire** : or, soit $A \in \mathcal{R}$ et $\varepsilon > 0$, prenons un ouvert U et un fermé C de sorte que : $C \subset A \subset U$ et $\mu(C) - \varepsilon < \mu(A) < \mu(U) + \varepsilon$. alors, le complémentaire A^c vérifie l'encadrement $U^c \subset A^c \subset C^c$, et de plus :

$$\mu(U^c) - \varepsilon < \mu(A^c) < \mu(C^c) + \varepsilon$$

ce qui fait : $A^c \in \mathcal{R}$.

a. c'est un espace topologique dans lequel deux éléments différents admettent toujours des voisinages ouverts disjoints.

b. n'oubliez pas que **pour la topologie induite de sa distance**, X est **fermé** et **ouvert** à la fois.

- **Pour la stabilité de la réunion dénombrable :** on prend une suite $(A_n)_{n \in \mathbb{N}'} \subset \mathcal{R}$, d'éléments de \mathcal{R} et un certain $\varepsilon > 0$, prenons pour tout : $i \in \mathbb{N}$ un ouvert U_i un fermé : C_i tel que : $(\forall i \in \mathbb{N}) : C_i \subset A_i \subset U_i$ et $\mu(U_i) - 2^i \varepsilon < \mu(A_i) < \mu(C_i) + 2^{i+1} \varepsilon$.
donc : $\bigcup_{i \in \mathbb{N}} C_i \subset \bigcup_{i \in \mathbb{N}} A_i \subset \bigcup_{i \in \mathbb{N}} U_i$, et $\bigcup_{i \in \mathbb{N}} U_i$ est ouvert, (la réunion quelconque des ouverts reste un ouvert) avec :

$$\begin{aligned} \mu\left(\bigcup_{i \in \mathbb{N}} U_i\right) - \mu\left(\bigcup_{i \in \mathbb{N}} A_i\right) &\leq \mu\left(\bigcup_{i \in \mathbb{N}} U_i \setminus \bigcup_{i \in \mathbb{N}} A_i\right) \\ &\leq \mu\left(\bigcup_{i \in \mathbb{N}} (U_i \setminus A_i)\right) \leq \sum_{i \in \mathbb{N}} \mu(U_i \setminus A_i) \\ &= \sum_{i \in \mathbb{N}} (\mu(U_i) - \mu(A_i)) < \sum_{i \in \mathbb{N}} 2^i \varepsilon = \varepsilon. \end{aligned}$$

en outre : et par la continuité croissante de la mesure : $\mu\left(\bigcup_{i \in \mathbb{N}} C_i\right) = \lim_{k \rightarrow +\infty} \mu\left(\bigcup_{i=1}^k C_i\right)$

par conséquent :

il existe un certain k pour lequel : $\mu\left(\bigcup_{i \in \mathbb{N}} C_i\right) - \mu\left(\bigcup_{i=1}^k C_i\right) < \frac{\varepsilon}{2}$, donc : $C = \bigcup_{i=1}^k C_i \subset \bigcup_{i \in \mathbb{N}} A_i$ est fermé avec :

$$\begin{aligned} \mu\left(\bigcup_{i=1}^{+\infty} A_i\right) - \mu(C) &< \mu\left(\bigcup_{i=1}^{+\infty} A_i\right) - \mu\left(\bigcup_{i=1}^{+\infty} C_i\right) + \frac{\varepsilon}{2} \\ &\leq \mu\left(\bigcup_{i=1}^{+\infty} A_i\right) \setminus \bigcup_{i=1}^{+\infty} C_i + \frac{\varepsilon}{2} \\ &\leq \mu\left(\bigcup_{i=1}^{+\infty} (A_i \setminus C_i)\right) + \frac{\varepsilon}{2} \\ &\leq \sum_{i=1}^{+\infty} \mu(A_i \setminus C_i) + \frac{\varepsilon}{2} \\ &= \sum_{i=1}^{+\infty} (\mu(A_i) - \mu(C_i)) + \frac{\varepsilon}{2} \leq \frac{\varepsilon}{2} + \frac{\varepsilon}{2} \end{aligned}$$

donc on obtient : $\bigcup_{i=1}^{+\infty} A_i \in \mathcal{R}$, et ceci justifie la stabilité de \mathcal{R} par réunion dénombrable.

► **il reste à vérifier que \mathcal{R} contient tous les ouverts de X :**

Afin de montrer que : $\Theta \subset \mathcal{R}$ on va essayer d'établir que \mathcal{R} contient tous les fermés de X , et on se basera sur la stabilité par passage au complémentaire pour déduire ce résultat pour les ouverts.

soit $A \subset X$ un fermé, on note $U_n := \left\{x \in X \mid d(x, A) < \frac{1}{n}\right\} = \left\{x \in X \mid (\exists a \in A) : d(x, A) < \frac{1}{n}\right\}$

par le lemme I.1.2.1, la fonction $d(\cdot, A)$ est continue, ce qui fait que les U_n sont des ouverts, et de plus :

$$(\forall i \in \mathbb{N}^*) : U_{i+1} \subset U_i \text{ , et : } \bigcap_{i=1}^{+\infty} U_i = A \text{ (car } A \text{ est fermé.) }^a \text{ alors et}$$

car μ est finie on aura $\mu(U_1) < +\infty$ donc et selon **la monotonie décroissante de la mesure** : $\mu(A) = \lim_{n \rightarrow +\infty} \mu(U_n) = \inf_{n \in \mathbb{N}} \mu(U_n)$

donc : $\mu(A) \leq \inf \left\{ \mu(U) \mid U \text{ est ouvert et } A \subset U \right\} \leq \inf_{n \in \mathbb{N}} \mu(U_n) = \mu(A)$.

a. On rappelle que : $x \in \overline{A} \iff d(x, A) = 0$

donc et en guise de conclusion : \mathcal{R} est une tribu qui contient la topologie de X (c'est à dire l'ensemble de ses ouverts), alors il contiendra aussi la tribu engendrée par cette topologie et qui n'est autre que la tribu borélienne \mathcal{B} d'où $\mathcal{B} \subset \mathcal{R}$ et pour l'autre inclusion elle est triviale car \mathcal{R} ne contient que des ensembles mesurables, ce qui achève la preuve. ■

3.2 Deuxième théorème de Lusin :

Théorème 3.2.1 : (de Lusin, version généralisée :)

Soit (E, \mathcal{E}, ν) un espace mesuré de Radon, (H, Θ) un espace topologique, à base^a dénombrable muni de la tribu Borélienne et $f : E \rightarrow H$, une fonction **mesurable**, alors on a :
 $(\forall \varepsilon > 0) (\forall A \in \mathcal{E}) : \nu(A) < +\infty \Rightarrow (\exists F_{(\varepsilon, A)} \subset A)$ tq : $F_{(\varepsilon, A)}$ est fermé, $\nu(A \setminus F_{(\varepsilon, A)}) < \varepsilon$ et $f|_{F_{(\varepsilon, A)}}$ est continue.
 autrement dit, on peut restreindre toute les fonctions mesurables sur un espace de Radon à des autres qui sont continues sur une partie **fermée**, aussi grande que voulue, de leur domaine.

Preuve :

À son tour, cette version va se baser, sur les résultats présentés précédemment pour la démonstration du premier théorème de Lusin, (surtout celui d'Egorov), mais, le fait de passer au cas général des espaces de Radon et de remplacer l'hypothèse d'intégrabilité par celui de mesurabilité, nécessite l'intervention d'un arsenal théorique très solide, surtout en topologie, pour cela, on l'acceptera pour le moment, et n'ayez pas peur, cette forme généralisée du théorème de Lusin, n'apparaîtra pas beaucoup dans la suite. ■

Remarque 3.2.1 :

Quand même ce théorème va nous servir pour justifier un petit passage du lemme III.2.2.3 et qui sera d'une grande utilité pour la preuve du théorème III.2.2.2.

4 C^0 –prolongeabilité des fonctions continues sur un fermé :

4.1 Théorème de prolongement de Dungundji :

Définitions 4.1.1 : (Partition de l'unité :)

- On appelle partition de l'unité d'un espace topologique X , toute famille $(\phi_i)_{i \in I}$ de fonctions **continues**, définies de X à valeur dans $[0, 1]$, telles que :
 - pour tout point $x \in X$, il existe un voisinage de x tel que toutes les fonctions ϕ_i soient nulles sur ce voisinage à l'exception d'un nombre **fini** d'entre elles .
 - $(\forall x \in X) : \sum_{i \in I} \phi_i(x) = 1$.
- On dit qu'une partition de l'unité : $(\phi_i)_{i \in I}$ est **subordonnée à un recouvrement** $(U_i)_{i \in I}$ de X si et seulement si elle est indexée par le même ensemble I que le recouvrement, avec : pour tout $i \in I$, le support de ϕ_i est inclus dans un U_j .

a. On désigne ici la base au sens topologique, et qui est une classe d'ouverts tel que tout élément de Θ soit une réunion d'ensembles de cette classe.

Pensez par exemple à la classe des intervalles ouverts dans l'ensemble des réels \mathbb{R} .

Définition 4.1.2 : (Espace topologique paracompact :)

Un espace topologique **paracompact** est un espace topologique **séparé** dans lequel tout recouvrement ouvert \mathcal{U} admet un raffinement (un autre recouvrement dont chaque élément est inclus dans un autre de \mathcal{U}) ouvert **localement fini**^a.

Attention! : à ne pas confondre avec l'espace précompact dans lequel on peut extraire de toute suite une autre qui est de Cauchy.

Théorème 4.1.1 : (théorème de Dugundji :)

Soient (X, Θ) un espace topologique métrisable^b, $A \subset X$ un **fermé** de X et $(Y, +, \cdot, \|\cdot\|)$ un espace vectoriel normé alors et selon [18] :

toute application continue f de A dans Y admet un prolongement **continu** \tilde{f} de X dans Y , et on dit dans ce cas qu'elle est **C^0 -prolongeable sur X** .

Preuve :

Pour une distance \mathbf{d} fixée sur X , considérons, dans l'ouvert $X \setminus A$, le recouvrement constitué des boules ouvertes $\left(B \left(x, \frac{\mathbf{d}(x,A)}{2} \right) \right)_{x \in X \setminus A}$.

Selon le fameux **critère de métrisabilité de Smirnov**^c tout espace métrique est paracompact, ce qui fait qu'il existe un recouvrement ouvert localement fini $(U_i)_{i \in I}$ de $X \setminus A$ dont chaque ouvert est inclus dans l'une de ces boules : $U_i \subset B \left(x, \frac{\mathbf{d}(x,A)}{2} \right)$.

On choisit alors une partition de l'unité $(\phi_i)_{i \in I}$ subordonnée à ce recouvrement et pour tout i , un point a_i de A tel que : $\mathbf{d}(x_i, a_i) \leq 2 \cdot \mathbf{d}(x_i, A)$, et on prolonge f en posant :

$$\tilde{f}(x) = \begin{cases} f(x) & \text{si } x \in A \\ \sum_{i \in I} \phi_i(x) \cdot f(a_i) & \text{sinon.} \end{cases}$$

L'application \tilde{f} est clairement continue sur $X \setminus A$, il reste à montrer qu'elle l'est aussi sur A , c'est à dire :

$$(\forall a \in A) (\forall \varepsilon > 0) (\exists \delta_\varepsilon > 0), \text{ tel que : } f \left(B(a, \delta_\varepsilon) \cap A \right) \subset B(a, \varepsilon).$$

Pour affirmer que pour tout $x \in B \left(a, \frac{\delta_\varepsilon}{6} \right)$, $f(x) \in B(a, \varepsilon)$, il suffit d'utiliser que pour tout U_i contenant x , $a_i \in B(a, \delta_\varepsilon)$, d'après les inégalités suivantes :

$$\begin{aligned} \mathbf{d}(x, a_i) &\leq \mathbf{d}(x, x_i) + \mathbf{d}(x_i, a_i) \leq \mathbf{d}(x, x_i) + 2 \cdot \mathbf{d}(x_i, A) \leq 5 \cdot \mathbf{d}(x_i, A) - 5 \cdot \mathbf{d}(x, x_i) \leq 5 \cdot \mathbf{d}(x, A) \\ \mathbf{d}(a, a_i) &\leq \mathbf{d}(a, x) + \mathbf{d}(x, a_i) \leq 6 \cdot \mathbf{d}(a, x). \end{aligned}$$

■

4.2 Théorème de l'extension de Tietze :

Théorème 4.2.1 : (de Tietze :)

Toute fonction $f : X \rightarrow \mathbb{R}$ continue à valeurs réelles définie sur un fermé d'un espace topologique métrisable (X, θ) peut être prolongée continument sur tout l'espace X .

a. Une famille de parties d'un espace topologique est dite localement finie lorsque chaque point possède un voisinage qui ne rencontre qu'un nombre fini d'éléments de la famille.

b. C'est à dire que sa topologie peut provenir d'une norme.

c. Il affirme qu'un espace topologique est métrisable \iff il est **régulier** (dans lequel on peut séparer un point x et un fermé ne contenant pas x par deux ouverts disjoints) à base dénombrable localement finie.

Preuve :

Ce n'est autre qu'un cas particulier du théorème de Dugundji dans lequel l'espace Y est la droite réelle. ■

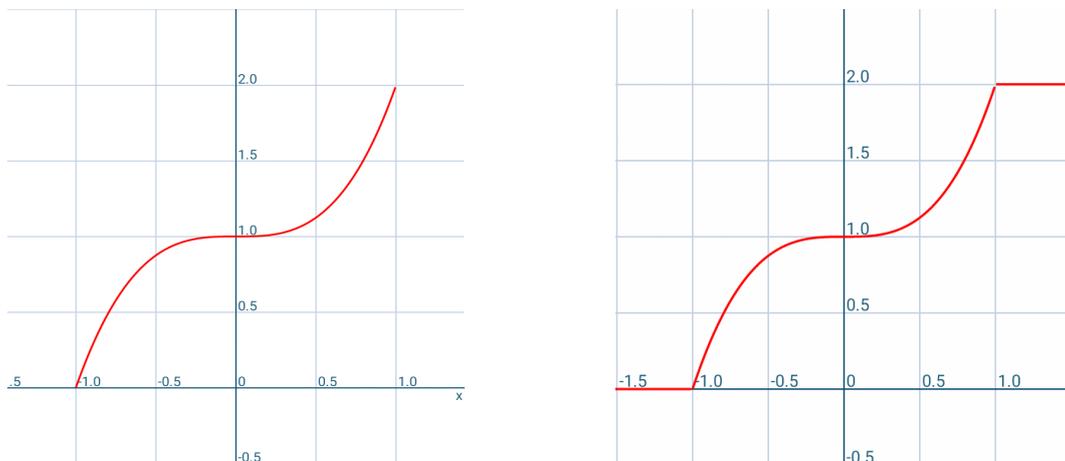


FIGURE I.3: Principe du théorème de Tietze pour $X = \mathbb{R}$.

Remarque 4.2.1 :

Ce résultat a été généralisée par le mathématicien russe Pavel Urysohn qui a remplacé l'espace métrisable X par n'importe quel espace normal ^a.

Conclusion :

Dans ce chapitre, et dans notre tentative de recouvrir une base théorique relativement solide et qui peut nous permettre de faire par la suite une lecture sans heurt des parties restantes, on a vu plusieurs résultats importants d'analyse fonctionnelle, comme les **deux théorèmes de Lusin**, **le théorème de Stone-Weierstrass** et **le théorème de Tietze**, ces derniers sont exceptionnellement distingués par leurs énoncés qui vont justifier ultérieurement la grande aptitude d'une machine d'apprentissage appelée **réseaux de neurones artificiels (RNA)**, et qui fera effectivement l'objet du chapitre suivant.

a. C'est un espace séparé dans lequel on peut trouver pour tous fermés disjoints F_1 et F_2 , deux ouverts disjoints O_1 et O_2 tels que : $F_1 \subset O_1$ et $F_2 \subset O_2$.

Apprentissage Automatique : Réseaux de Neurones Artificiels.

Résumé :

Dans ce chapitre, on va présenter une famille très intéressante de machines d'apprentissage et qui permet de faire le traitement et l'analyse de plusieurs classes de problèmes, que les outils ordinaires et les machines d'apprentissage classiques ne peuvent pas les réaliser ou juste ont des difficultés et du mal à résoudre.

le modèle de machines qu'on va proposer sera capable de reproduire et d'*imiter le plus fidèlement possible certains aspects de l'intelligence humaine*, et on verra que ceci revient à l'origine de l'inspiration de ce modèle du système nerveux, en outre l'introduction d'une telle machine qui portera le nom de "*réseau de neurones artificiels*" va causer l'émergence de plusieurs sous types remarquables, comme par exemple : le *perceptron*, les *réseaux de Hopfield*, *réseaux à compétition...*, etc, et même elle donnera naissance à des nouvelles gammes d'algorithmes appelés *algorithmes d'apprentissages*.

dans la suite, nous allons citer les caractéristiques de ces réseaux, voir de proche leur fonctionnement et essayer d'exposer au lecteur leurs principaux algorithmes d'apprentissage, tout en faisant de notre mieux pour lever toute équivoque possible.

Mots clés : Réseaux de neurones, apprentissage, algorithme, intelligence artificielle.

Introduction :

Afin d'analyser et de comprendre le *phénomène d'apprentissage*, l'Homme s'intéressa à l'étude du cerveau comme étant l'organe central de commande et de pensée ce qui va l'amener enfin de compte à faire ces études sur un système très compliqué *qui est le système nerveux*, cette complexité au niveau de l'étude va le pousser à se restreindre aux éléments de base qui le constituent et qui seront dans ce cas : *les neurones biologiques*, auxquels il proposera une représentation mathématique-informatique appelée : *Le neurone formel ou artificiel*.

L'introduction du premier exemple de neurones formels a eu lieu en 1943, et c'était dû à Warren Mc.Culloch et Walter Pitts, qui se sont basés sur des observations neurophysiologiques, *récentes à l'époque*, pour nous proposer un modèle *très simple* du neurones formels, et qui s'agissait, en ce temps là, d'une unité d'exécution binaire, retournant une sortie qui vaut 0 ou 1, le calcul de cette dernière, était par l'effectuation d'une somme pondérée des entrées, puis l'application d'une fonction de décision à seuil, qui donne comme résultat 1 lorsque cette somme pondérée dépasse une

certaines valeurs, et 0 au cas contraire.

Après avoir constaté l'insuffisance du modèle qu'ils ont proposé et son incapacité de résoudre certains problèmes (le fameux problème "XOR" à titre d'exemple.) Mc.Culloch et Pitts ont décidé de faire l'analogie avec le cerveau humain en rendant leur neurone comme une unité élémentaire construisant une structure plus compliquée et puissante, nommée dans ce cadre : "**réseaux de neurones artificiels**".

Cette idée va marquer une étape charnière dans l'histoire de l'intelligence artificielle, et il développera par la suite tout le domaine de la bio-informatique, et surtout à travers l'introduction d'une famille d'algorithmes capables d'améliorer leurs performances à résoudre des tâches sans être explicitement programmés pour chacun, et ceci fera des réseaux de neurones artificiels des systèmes habiles d'apprendre, et de mettre en œuvre le principe de l'induction, et spécialement avec l'émergence d'autres variantes de neurones qui ont été proposées, après les travaux de deux neurologues Mc.Culloch et Pitts.

1 De neurone biologique au neurone artificiel :

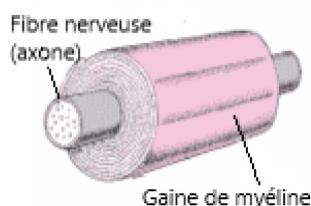
1.1 Neurone biologique :

1.1.1 Définitions générales et quelques statistiques :

Un neurone^a biologique se considère comme une cellule nerveuse élémentaire qui constitue l'unité de base du système nerveux, son rôle principal est de recevoir des stimulations extérieures, les transformer en impulsions électriques pour les transmettre^b enfin de compte à une autre cellule qui peut être nerveuse à son tour ou pas.

Le nombre total de ces neurones dans le corps humain est estimé à 100 milliards, dont 89% se concentre dans le cerveau, et les 11% restantes se répartissent sur d'autres membres comme la moelle spinale qui compte presque 4 milliards, l'intestin qui compte plus de 450 millions et le cœur qui contient un nombre dans les environs de 60 mille...etc.

Les neurones sont les composantes principales les plus importantes du tissu nerveux, et elles possèdent des caractéristiques très remarquables et particulières, d'abord **elles vivent très longtemps** et **elles nous accompagnent de la naissance à la mort**, et en plus elles **ne peuvent jamais être remplacées quand elles sont détruites**, et pour cette raison, on trouve dans le système nerveux, un type spécial de cellules appelées **les cellules gliales**, ou parfois **les neuroglies**, ou tout simplement **les glies**, qui s'occupent du soutien et de la nutrition des neurones ainsi que la production de **la myéline^c** qui protège le tissu nerveux.



(a) Une gaine de myéline.



(b) Gaine de myéline normale face à une autre lésée.

a. Ce terme était introduit dans le lexique médical pour la première fois en 1879 par le médecin et le biologiste allemand Heinrich Wilhelm Waldeyer (1836-1921).

b. Ce mécanisme de transmission reflète deux propriétés physiologiques très importantes chez les neurones et qui sont : l'**excitabilité** et la **conductivité**.

c. C'est une sorte de substance qui sert à **isoler** et à protéger **les fibres nerveuses**, exactement comme le fait du plastique autour des fils électriques, et on peut trouver certains troubles qui provoquent une démyélinisation qui touche directement le système nerveux central, ou les nerfs d'autres parties du corps, par exemple : **La polyneuropathie inflammatoire démyélinisante chronique « PIDC »**.

Le nombre des nervoglies est beaucoup plus important que celui des neurones ,et selon les estimations des scientifiques, leur nombre dépasse trois fois celui des neurones, et pour cela on trouve qu'elles représentent à peu près 50% du volume cérébral et plus de 65% des cellules du cerveau.

1.1.2 La structure générale d'un neurone biologique :

Un neurone biologique peut être décrit tout simplement comme une cellule ordinaire prolongée par des ramifications,et sa structure complexe se caractérise essentiellement par les composantes principales suivante :

- **les dendrites** : d'une moyenne de 10 mille par neurone , elles font la réception de l'information sous forme des impulsions électriques, d'autres neurones auxquels elles sont attachées via des *synapses*^a.
- **une membrane externe** : elle s'occupe de la propagation de cette information le long du neurone, et on peut dire qu'elle est la seule responsable de la *conductivité* chez les cellules nerveuses.
- **un corps cellulaire(ou Soma)** : il fait le traitement de toutes les informations provenant des dendrites pour produire une sortie qui se propage via l'axone enfin de compte.
- **un axone** : d'un diamètre compris entre 10 et 15 μm et d'une longueur qui varie d'un millimètre à plus d'un mètre, cette fibre nerveuse fait parvenir le potentiel d'action résultante à une autre cellule.

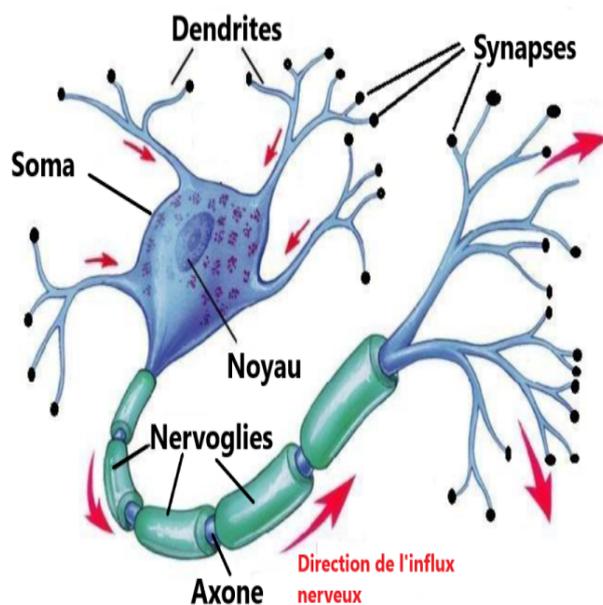


FIGURE II.2: Schéma d'un neurone biologique.

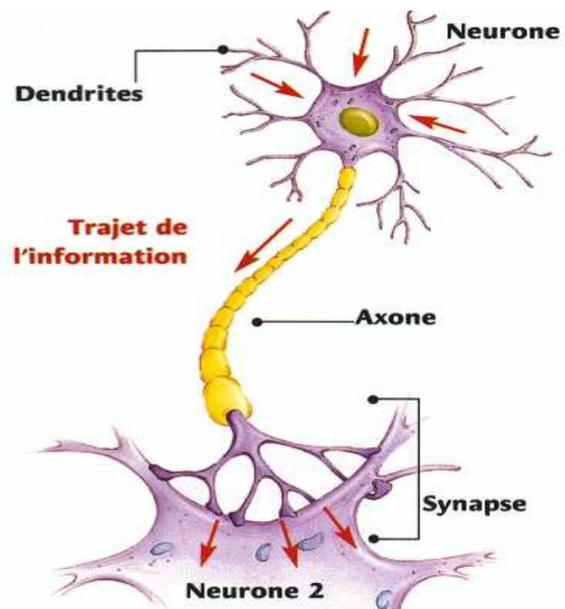


FIGURE II.3: Synapses entre deux neurones.

Donc,et après cette petite mise en situation ,on peut et sans problème parler d'une manière concrète *d'un neurone artificiel*, ou encore *neurone formel*.

1.2 Neurone artificiel :

1.2.1 Une description générale du neurone artificiel :

Un neurone artificiel ou formel n'est autre qu'une représentation mathématique informatique du neurone biologique et qui s'inspire exactement de ses composantes générales, et *il se caracté-*

a. C'est une zone de contact fonctionnelle qui s'établit entre deux neurones, ou entre un neurone et une autre cellule.

rise par :

- **des entrées** : qui correspondent aux dendrites.
- **des poids** : qui représentent les actions excitatrices des synapses, ce qui justifie le fait de les appeler parfois *les poids synaptiques*.
- **une fonction d'entrée totale** : qui définit le pré-traitement effectué sur les entrées du neurone, et à la plupart des cas, cette fonction est choisie comme celle de la somme pondérée des entrées par les poids.
- **une fonction d'activation** : elle agit sur l'entrée totale produite par la fonction précédente pour nous donner ce qu'on appelle l'état du neurone.
- **une fonction de sortie** : tout simplement elle calcule la sortie du neurone en agissant sur son état, et habituellement elle est omise par une identification avec la fonction identité.
- **la sortie** : c'est le résultat final du traitement effectué par le neurone et il représente l'information qui se propage au long de l'axone.
- **et enfin les types d'entrées et de sortie** : il faut quand-même mentionner que *ces types* restent parmi les facteurs les plus importants du neurone artificiel, et leur modification peut produire des changements radicaux sur sa capacité de résolution face à un problème précis ou à une classe de problèmes bien déterminée.

donc et d'une manière grossière, le neurone formel n'est autre que la composition des trois fonctions numériques définie précédemment, et de ce fait on pourra l'exprimer ainsi :

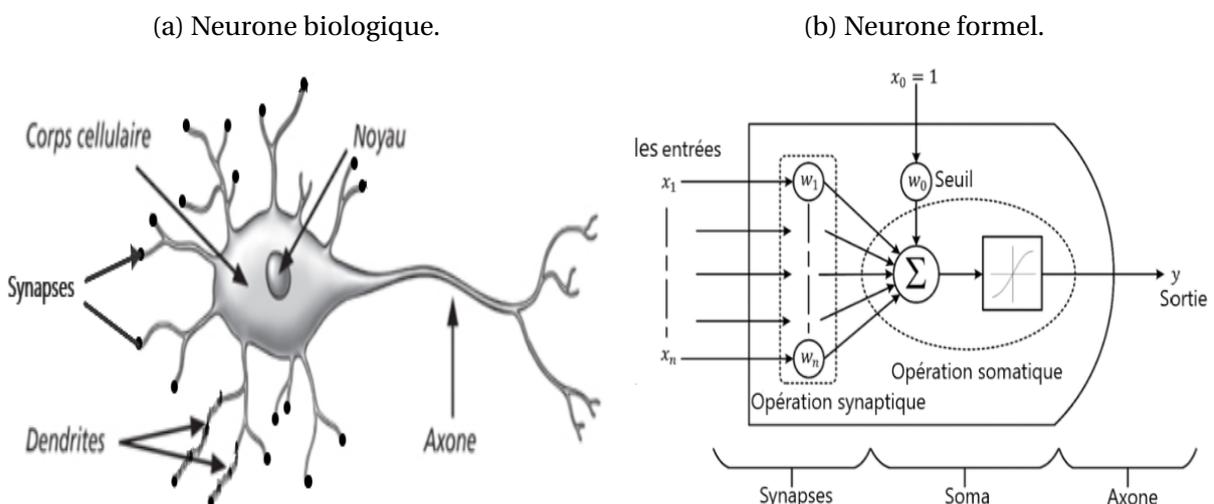
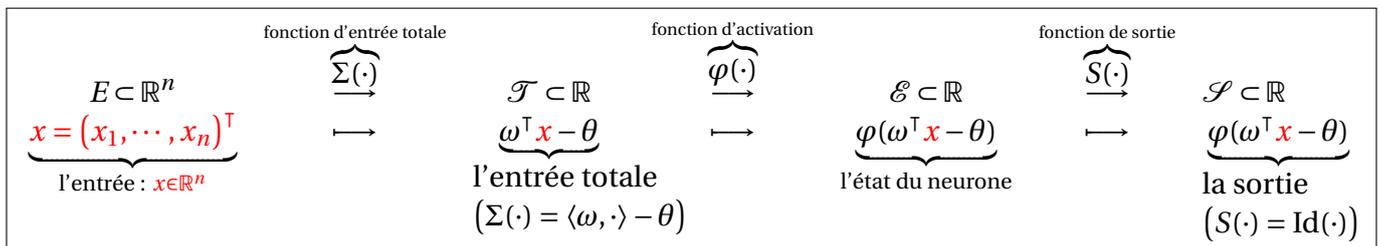


FIGURE II.4: Comparaison entre la structure d'un neurone biologique et d'un autre formel.

1.2.2 Le neurone artificiel d'un point de vue mathématique :

Définition : On appelle neurone formel ou artificiel, toute fonction numérique de n variables réelles x_1, \dots, x_n , et qui prend la forme suivante :

$$\begin{array}{l} \Psi : E \subset \mathbb{R}^n \longrightarrow \mathcal{S} \subset \mathbb{R} \\ x \longmapsto \Psi \left(\omega_0 + \sum_{j=1}^n \omega_j x_j \right) \end{array}$$

avec : $n \in \mathbb{N}^*$ est **un paramètre** entier non nul indiquant **la taille des entrées**, $\omega = (\omega_1, \dots, \omega_n)^\top \in \mathbb{R}^n$, est un vecteur fixé à l'avance nommé : **le vecteur des poids** et $\omega_0 \in \mathbb{R}$ est une certaine constante appelée **le biais**, et dans ce cas là, l'image : $y = \Psi \left(\omega_0 + \sum_{j=1}^n \omega_j x_j \right) \in \mathbb{R}$, portera le nom de **la sortie du neurone associée au vecteur** : x et de plus les ensembles : $E \subset \mathbb{R}^n$ et $\mathcal{S} \subset \mathbb{R}$, désignent respectivement **l'espace d'entrée et de sortie** du neurone

Remarque : Dans la définition précédente on a adopté les choix habituels pour les fonctions d'entrée^a et de sortie, or la première était confondue avec celle des somme pondérées des poids avec les entrées, et la deuxième avec l'identité, mais **faites attention!**, ceci peut être modifié sans aucun problème dans certains cas, à savoir les $\Sigma\Pi$ neurones par exemple.

1.2.3 Comment on est arrivé à cette description du neurone formel? :

On a bien vu dans ce qui précède que le neurone formel n'était qu'une simple présentation du neurone biologique, et même on a réussi de réduire et de simplifier sa description pour la rendre sous la forme d'une fonction mathématique possédant une expression un peu particulière et spécifique, mais **cette simple formulation n'était pas obtenue d'un seul coup, ou d'une façon miraculeuse**, cependant **elle était le fruit d'une longue histoire de recherche et de travail**, dans la suite, on va présenter les grands repères historiques derrière l'évolution de cette description en exposant les différentes représentations adoptées avant son émergence.

★ **Le neurone de Culloch-pitts (1943) :** Comme on a déjà dit, le premier neurone formel a été proposé en 1943 par les deux chercheurs américains **Mc.Culloch** et **Walter Pitts** (logicien) dans leur article [42], où il était présenté comme suit :

► **Description :**

- **les entrées :** pour ce cas elles sont de type binaire, et elles peuvent être vues comme un vecteur : $x = (x_1, \dots, x_n)^\top \in \{0, 1\}^n$, en plus on pourra dans ce cadre faire la distinction entre deux types d'entrées :

- **des entrées inhibitrices :** sont celles qui **ont un effet maximum sur la prise de décision** quelles que soient les autres entrées.
- **des entrées excitatrices :** elles ne sont pas du tout celles qui vont déclencher le feu de neurone mais **elles pourraient le déclencher quand elles sont combinées**.

× **les poids :** en ce temps là cette notion n'était pas encore abordée, et pour cela on va les négliger en leur attribuant tous la valeur 1.

a. Certaines sources appellent cette fonction **la fonction d'agrégation**.

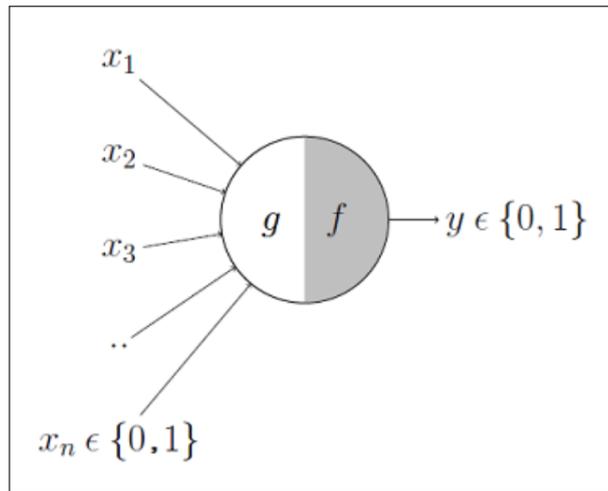


FIGURE II.5: La représentation simple d'un C.P-neurone.

- **la fonction d'entrée totale:** à l'époque elle était juste la somme des entrées, c'est à dire :

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

$$x = \begin{pmatrix} x_1 \\ \vdots \\ x_n \end{pmatrix} \mapsto \sum_{j=1}^n x_j$$

× **la fonction d'activation:** selon la toute première définition proposée par Culloch et pitts dans leur article [42], c'était la fonction de décision définie par : $f(x) = \mathbb{1}_{[\theta, +\infty[}(x) := \begin{cases} 1 & \text{si } : x \geq \theta \\ 0 & \text{sinon.} \end{cases}$ avec θ est un réel appelé **paramètre de seuillage**.

× **la fonction de sortie:** exactement comme les poids, cette notion n'avait pas de sens à l'époque, ce qui fait qu'on peut la confondre avec la fonction identité pour le moment.

- **la sortie:** clairement elle est de type binaire^a avec :

$$y = \begin{cases} 1 & \text{si } : \sum_{j=1}^n x_j \geq \theta \\ 0 & \text{si } : \sum_{j=1}^n x_j < \theta \end{cases}$$

Remarque : cette expression de la sortie montre bien et justifie le fait d'appeler θ par paramètre de seuillage, et on signale que pour améliorer la capacité de discrimination chez notre neurone et l'approcher d'un comportement désiré on doit régler et changer **à la main** la valeur prise par ce paramètre.

► **Mise en œuvre :** voyons maintenant comment ce modèle peut être utilisé pour représenter quelques fonctions booléennes.

1. **La fonction logique "AND" :**

a. c'est à dire qu'elle prend ses valeurs dans : $\{0, 1\}$.

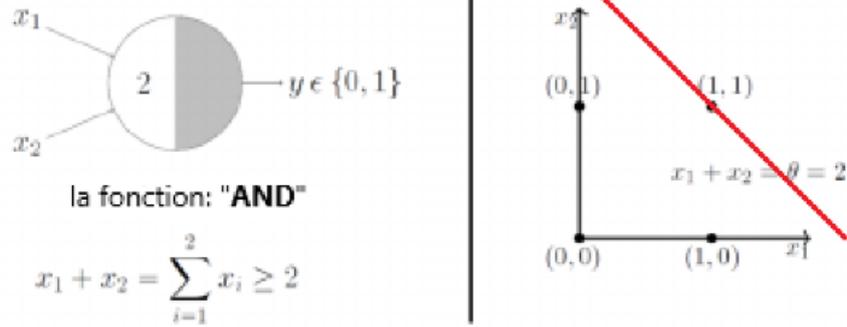


FIGURE II.6: Réalisation de la fonction "AND" par un C.P-neurone.

2. La fonction logique "OR" :

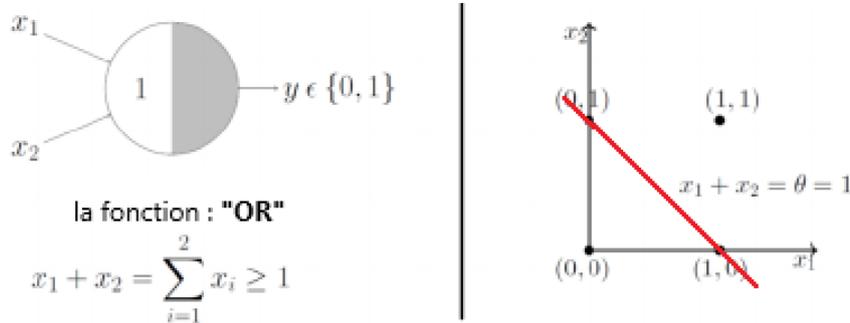


FIGURE II.7: Réalisation de la fonction "OR" par un C.P-neurone.

3. La fonction logique "OR" au cas de trois entrées :

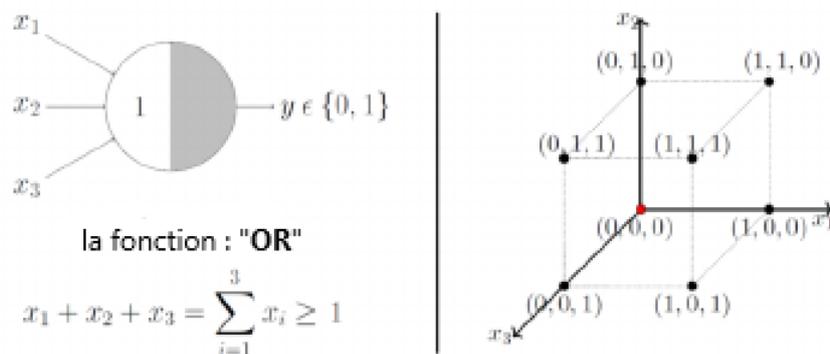


FIGURE II.8: Réalisation de la fonction "OR" à trois entrées par un C.P-neurone.

► **Les limites du modèle de Culloch-pitts :**

1. Que peut on faire au cas des entrées non booléennes (par exemple, réelles)?
2. Que faire si nous voulons attribuer plus d'importance à certaines entrées?
3. Avons-nous toujours besoin de coder *manuellement* le seuil?
4. Comment résoudre des fonctions qui ne sont pas linéairement séparables?(la fonction "XOR" à titre d'exemple).

★ **Le Perceptron simple de Rosenblatt (1958)** : Le psychologue américain **Frank Rosenblatt** va se baser sur les travaux de Culloch et Pitts pour introduire en 1958 une machine de **classification linéaire** très puissante appelée "**Perceptron**" et dans son article [50] , cette machine a été décrite ainsi :

► **Description :**

- **les entrées** : dans ce cas sont des réels qui peuvent être vus comme un vecteur : $x = (1, x_1, \dots, x_n)^T \in \mathbb{R}^{n+1}$ où la première composante $x_0 = 1$ représente une **entrée fictive** qui revient à l'introduction du biais ω_0 .
- **les poids** : sont aussi des réels qui sont vus également comme un vecteur des pondérations : $\omega = (\omega_0, \omega_1, \dots, \omega_n)^T \in \mathbb{R}^{n+1}$, possédant une première composantes $\omega_0 = -\theta$ qui représente **le biais** et qui est tout simplement **l'opposé du seuil**.
- **la fonction d'entrée totale** : c'est tout simplement celle des sommes pondérées définie par.

$$\Sigma : \mathbb{R}^n \longrightarrow \mathbb{R}$$

$$x = \begin{pmatrix} x_1 \\ \vdots \\ x_n \end{pmatrix} \longmapsto \left\langle \omega, \begin{pmatrix} 1 \\ x \end{pmatrix} \right\rangle = \omega_0 + \sum_{j=1}^n \omega_j x_j = \underbrace{\sum_{j=0}^n \omega_j x_j}_{\text{ici on a : } x_0=1.}$$

- **la fonction d'activation** : selon la toute première définition proposée par Culloch et pitts dans leur article [42],c'était la fonction Heavside définie par : $\varphi(x) = \mathbb{1}_{\mathbb{R}^+}(x) := \begin{cases} 1 & \text{si } :x \geq 0 \\ 0 & \text{sinon.} \end{cases}$

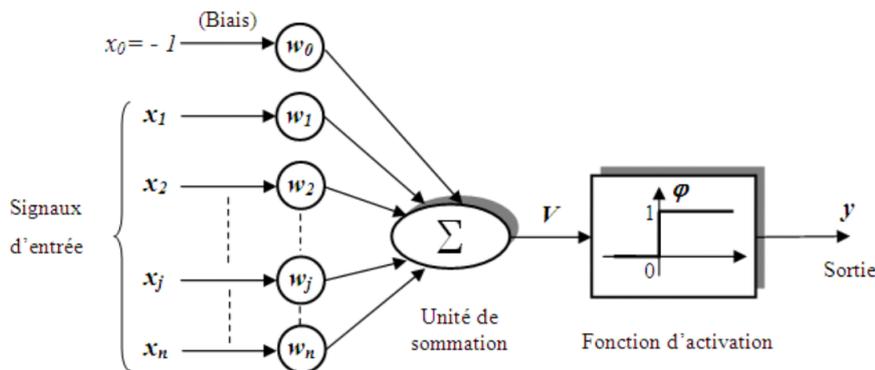


FIGURE II.9: Modèle du perceptron de Rosenblatt.

- **la fonction de sortie** : comme d'habitude c'est la fonction identité.
- **la sortie** : clairement elle est de type binaire^a avec :

$$S = \begin{cases} 1 & \text{si : } \sum_{j=1}^n \omega_j x_j \geq \theta \\ 0 & \text{si : } \sum_{j=1}^n \omega_j x_j < \theta \end{cases}$$

► **Mise en œuvre** : Une autre fois on va voir comment ce modèle peut être utilisé pour représenter quelques fonctions booléennes , mais cette fois ci on se contentera par l'étude de la simple fonction "**OR**" en laissant au lecteur le soin de traiter les autres exemples.

1. **La fonction logique "OR"** :

a. c'est à dire qu'elle prend ses valeurs dans : $\{0,1\}$.

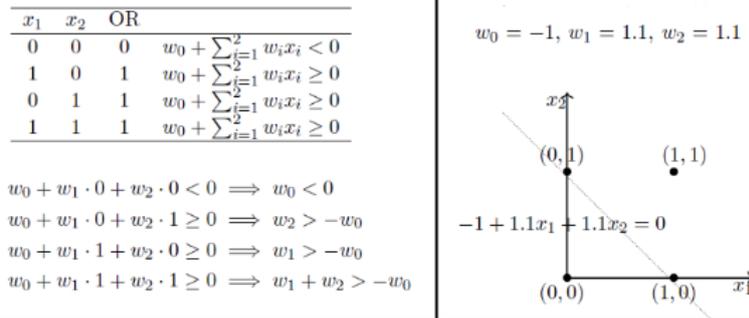


FIGURE II.10: Une solution possible de la fonction "OR" par un perceptron simple.

► **Les limites du Perceptron simple :** On ne peut pas nier que le perceptron est un modèle très puissant, qui a pu dépasser quelques limites du C.P-neurone en fournissant des solutions et des réponses à plusieurs questions, mais bien sûr *pas toutes les questions*, et pour cette raison, on trouve qu'il y en a encore des interrogations qui se considèrent comme des bornes de capacité chez ce type de machines, et on propose à titre d'exemples les questions suivantes :

1. Avons-nous toujours besoin de coder *manuellement* le seuil?
2. Comment résoudre des fonctions qui ne sont pas linéairement séparables? (la fonction "XOR" à titre d'exemple).

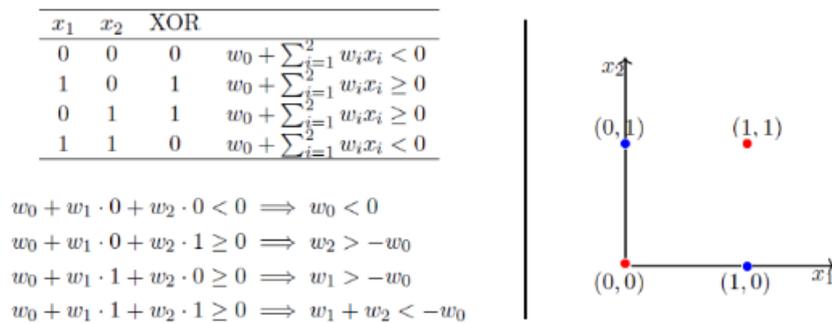


FIGURE II.11: L'incapacité du perceptron simple face au problème "XOR".

Remarque : en générale, un *seul* perceptron (simple) ne peut être utilisé que pour mettre en œuvre des fonctions linéairement séparables^a, et il reste totalement impuissant devant celles qui ne le sont pas, ce grand point faible sera à l'origine de l'apparence de l'idée de contribution entre les neurones, et de construction des *réseaux neuronaux*^b par la suite.

★ **L'Adaline de Widrow-Hoff (1960) :** le professeur américain **Bernard Widrow** et son étudiant, **Ted Hoff**, ont proposé en 1960 dans une partie de leur rapport technique [62], une machine d'apprentissage appelée "*neurone linéaire adaptatif*" en anglais : *Adaptive Linear Neuron* ou « **Adaline** » en abrégé.

en somme ce modèle n'était autre que le perceptron de Rosenblatt munit de la fonction d'activation linéaire pure, c'est à dire "l'identité", et la sortie dans ce cas sera sans surprise :

$$S = \sum_{i=0}^n \omega_i x_i \in \mathbb{R}.$$

a. C'est à dire qu'il est possible de tracer une ligne pour séparer les entrées positives des entrées négatives.
b. Nous reviendrons un peu plus tard à cette notion.

mais ce qui va marquer la différence complète de cette machine , c'est sa capacité **de régler automatiquement ses paramètres en se basant sur des expériences déterminées et des données correctes** , sans que ceci soit réalisé manuellement , ce fait va la distinguer de tous systèmes logique usuels à l'époque , et on verra par la suite que ce réglage se fait via un algorithme bien précis et qui était globalement inspiré d'une règle appelée la règle de Hebb.

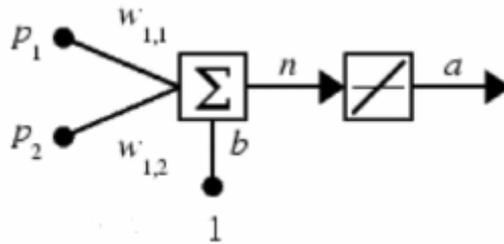


FIGURE II.12: Exemple d'un Adaline à deux entrées.

★ **Le perceptron de Minsky et Papert (1969)** :le modèle proposé par Rosenblatt va être analysé et amélioré par les deux neuropsychologues américains **Marvin Lee Minsky** et **Seymour Aubrey Papert** en 1969, dans leur article [43], au sein duquel ils ont modifié légèrement quelques détails au niveau de la structure de cette machine en donnant plus d'opportunités à sa fonction d'activation (**la fonction signe à titre d'exemple**), et ils se sont également inspirés des travaux de leurs collègues Widrow et Hoff, et qui ont été récentes à l'époque, pour attribuer à leur modèle un algorithme permettant de réaliser "le réglage" automatique des poids synaptiques, on va revenir par la suite à l'étude de cet algorithme d'une façon plus détaillée.

La plupart des neurones formels proposés après 1958 et jusqu'à maintenant ne sont que des variantes du perceptron de Rosenblatt dont la fonction d'activation est modifiée, par exemple :

★ **Les perceptrons sigmoïdes** : les neurones artificiels usuels qu'on a proposé jusqu'à cet instant disposent tous d'**une fonction d'activation à seuil**, qui **produit un changement soudain** au niveau de la décision de 0 à 1 lorsque la valeur de l'entrée totale : $S = \sum_{i=0}^n \omega_i x_i \in \mathbb{R}$ dépassera un certain seuil $\theta = -\omega_0$, mais pour la plupart des applications réelles , **cette logique de seuillage sévère est vraiment déconseillée** , et **surtout à cause des difficultés algorithmiques qui peuvent se produire** à cause de la non différentiabilité de ces fonctions (l'invalidité des algorithmes basés sur la descente de gradient par exemple) et pour cela nous nous attendons à une fonction de décision plus fluide, **qui passe progressivement** de 0 à 1 , ce qui va nous pousser à introduire les neurones formels ou les perceptrons sigmoïdes où la fonction d'activation est beaucoup plus lisse que la fonction du pas ce qui semble être une chose logique et évidente à faire. rappelez-vous qu'une fonction sigmoïde est une fonction mathématique d'une courbe en forme de "S" et qui passe de 0 à 1, Il existe de nombreuses fonctions qui ont cette caractérisation, par exemple : $\Psi : x \mapsto \frac{1}{2} (1 + \tanh(x))$.

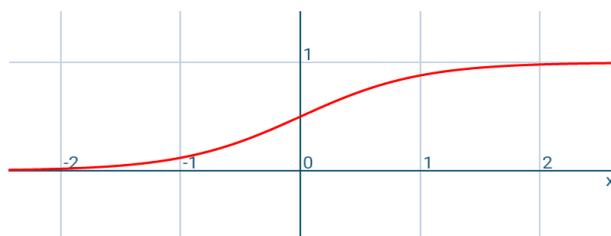


FIGURE II.13: La courbe de la fonction : $x \mapsto \frac{1}{2} (1 + \tanh(x))$.

(a) Les fonctions d'activation les plus courantes.

Nom de la fonction	Relation d'entrée/sortie	Icône
seuil	$a = 0$ si $n < 0$ $a = 1$ si $n \geq 0$	
seuil symétrique	$a = -1$ si $n < 0$ $a = 1$ si $n \geq 0$	
linéaire	$a = n$	
linéaire saturée	$a = 0$ si $n < 0$ $a = n$ si $0 \leq n \leq 1$ $a = 1$ si $n > 1$	
linéaire saturée symétrique	$a = -1$ si $n < -1$ $a = n$ si $-1 \leq n \leq 1$ $a = 1$ si $n > 1$	
linéaire positive	$a = 0$ si $n < 0$ $a = n$ si $n \geq 0$	
sigmoïde	$a = \frac{1}{1+\exp^{-n}}$	
tangente hyperbolique	$a = \frac{e^n - e^{-n}}{e^n + e^{-n}}$	
compétitive	$a = 1$ si n maximum $a = 0$ autrement	

(b) Les fonctions d'activation les plus utilisées.

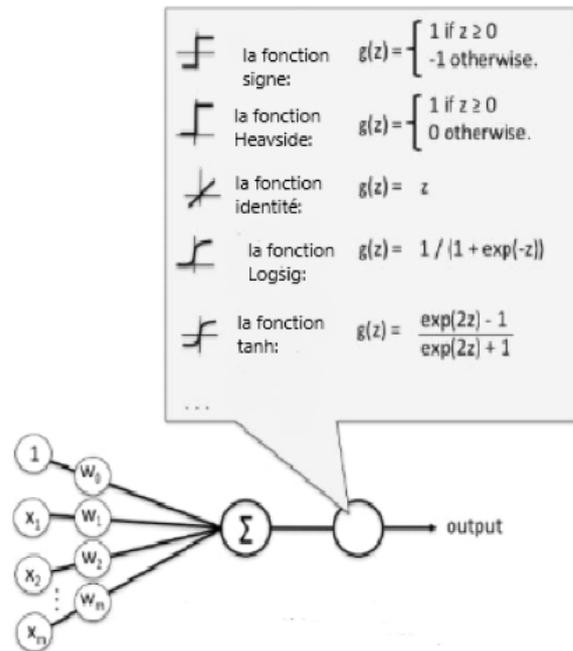


FIGURE II.14: Quelques choix possibles pour la fonction d'activation .

Remarque : on peut trouver d'autres applications numériques qui peuvent être proposées pour jouer le rôle de la fonction d'activation , **mais faites attention**, les seules applications qui peuvent bénéficier de ce privilège sont celles qui garantissent la réalisation de la propriété **d'approximation universelle** chez les réseaux de neurones artificiels, ce qui rend la tâche du choix d'une fonction convenable un peu plus délicate.

pour le moment contentez vous des simples exemples qu'on a présenté, et ne vous inquiétez pas on reviendra dans le **chapitre III** pour voir en profondeur quelques modèles de fonctions qui peuvent bien nous servir dans ce sens.

2 Apprentissage d'un neurone artificiel :

2.1 Vue d'ensemble sur le concept d'apprentissage :

2.1.1 La notion d'apprentissage :

Dans un cadre général **l'apprentissage** désigne le fait d'**acquérir des informations** ou des propriétés relativement stables ,et on peut le décrire dans le cadre humain comme étant un **mélange** complexe entre **les contraintes innées** et le **hasard acquis**, mais dans le cadre de l'intelligence artificiels, l'apprentissage désigne une phase durant laquelle une machine ,comme le perceptron simple par exemple, **se change** et se développe,**en modifiant** son comportement et **ses paramètres pour arriver à une attitude désirée**.

le terme « **apprentissage** » a été introduit dans le cadre des machines pour la première fois en 1959 par l'informaticien américain **Arthur Samuel** dans son article [53] , et ceci à la suite de sa création d'un programme qui est capable de jouer au fameux Jeu de Dames et de s'améliorer en jouant.

2.1.2 Les buts de l'apprentissage :

L'apprentissage possède plusieurs objectifs mais en somme on les résume comme suit :

1. l' **extraction des informations** pertinentes à partir des données.
2. la **classification** des données **et la prédiction** si une certaine donnée appartient à une telle ou telle classe.
3. le **regroupement** des données en classes homogènes.
4. la **mémorisation des informations** afin d' **identifier** toutes données bruitée ou incomplète à une autre déjà connue .
5. la **généralisation** qui nous permet d' **étendre nos résultats** et de faire le passage **à des nouvelles données**.

2.1.3 Les types principaux d'apprentissage :

Il existe beaucoup de types d'apprentissage, et dans la littérature on trouve souvent qu'ils sont au nombre de six (voir :[20]) mais principalement , on peut distinguer juste entre quatre grands types :

-**l'apprentissage supervisé** : il se fait sur une base de données particulière appelée **base d'apprentissage** et qui est constituée d'une collection finie de couples ayant la forme suivante :

$$(entrée, sortie désirée)$$

ces sorties désirées forment un **étiquetage** qui permet de subdiviser les entrées en classes, ce qui fait que nos exemples sont connus, nos classes sont prédéterminées, et le système n'a que d' **apprendre à les classer selon un modèle de classification bien précis** ainsi ,la machine peut aussi apprendre à **généraliser** et à prédire la classification de futures données en s'appuyant sur les différentes techniques de **la régression** statistique .

-**l'apprentissage non supervisé** : comme son nom l'indique ,cet apprentissage se fait sur un autre type de base,appelé **base de test** et dans laquelle **on ne dispose plus de sorties désirées**,juste des entrées ,et ce type d'apprentissage est voué à des utilisations particulières comme le **clustering** ou la répartition des données, or,on souhaite en l'occurrence de mettre des exemples hétérogènes sous forme de collections dont les éléments sont liés par des caractéristiques communes , et c'est à la machine de détecter les traits en communs sans aucune intervention extérieure , on mentionne que la similarité entre les données est souvent déterminée via une fonction de distance entre paires d'exemples.

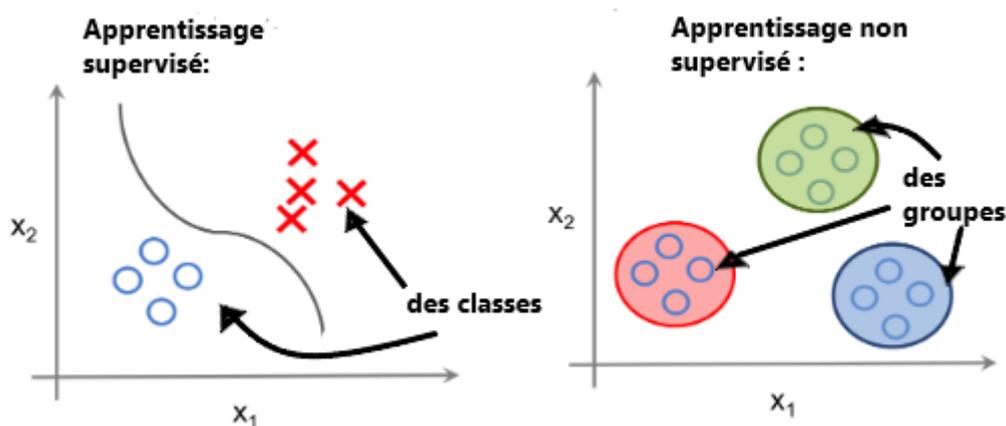


FIGURE II.15: La différence entre l'apprentissage supervisé et non supervisé.

-**l'apprentissage par renforcement** : ce type est basé sur trois notions *l'agent*, *l'environnement* et *le système de récompense*.

l'agent va réagir en fonction d'un état de *l'environnement*, pour renvoyer une action en fonction de celui-ci, le *système de récompense* va évaluer positivement ou négativement l'agent, en fonction de l'action qu'il a pris.

le but de l'agent est de **collecter le maximum de points possible**, ce qui va l'obliger de comprendre la différence entre une bonne et une mauvaise action, et donc au fur et à mesure de favoriser les plus convenables.

cette façon d'apprentissage vise en fait de reproduire le mécanisme naturel d'acquisition des connaissances, et ce qui la rend extrêmement puissante, c'est le fait qu'elle n'a pas besoin de bases de données, contrairement aux deux types d'apprentissages précédents.

-**l'apprentissage par transfert** : comme son nom l'indique ce type d'apprentissage est basé sur la capacité d'un système à **faire l'extension des connaissances** et des compétences, acquises à partir des tâches antérieures **sur des nouvelles tâches** ou domaines **partageant des propriétés communes**, (par exemple faire le passage de la résolution d'un cube de Rubik de dimensions $2 \times 2 \times 2$ à un autre de taille $3 \times 3 \times 3$) mais ce type d'apprentissage est rarement adopté en pratique, surtout à cause des difficultés qu'il présente au niveau de la détection des similitudes et le transfert des connaissances entre les tâches.

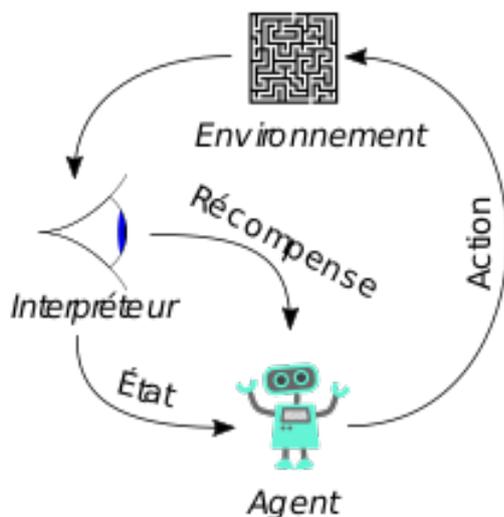


FIGURE II.16: Apprentissage par renforcement.

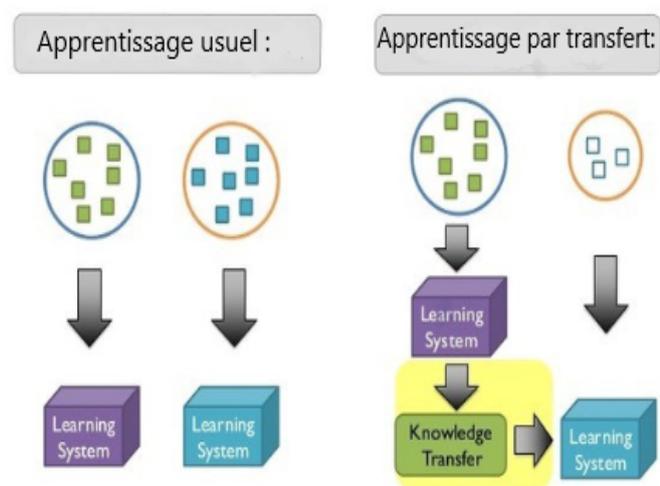


FIGURE II.17: Apprentissage par transfert face à l'apprentissage usuel.

2.1.4 Les facteurs qui influencent la qualité d'apprentissage :

La qualité de l'apprentissage dépend en générale de plusieurs facteurs, et qui sont essentiellement liés **soit aux performances de la machine** d'apprentissage utilisée, **soit à la nature de la base de donnée** fournie à cette machine (bien sûr lorsque son apprentissage nécessite la présence d'une telle base) :

► **les facteurs liés à la base des données** : et qui sont généralement les suivants :

- **la taille de la base de données** : lorsque le nombre d'exemplaires est **acceptablement élevé**, la machine sera capable de se doter d'un nombre suffisant d'informations et de renseignements, ce qui va favoriser **en principe** ses chances de réaliser un apprentissage de bonne qualité.
- **la proportion des données incomplètes** : dans une base de taille importante la présence d'un tel type de données est fort probable, et ceci peut perturber notre machine

et par conséquent réduire sa qualité d'apprentissage , pour garder cette réduction dans le cadre de l'acceptable ,**on doit** tout simplement **se poser dans des cas où le pourcentage de cette catégorie de données est faible** par rapport à la taille totale de la base .

- **la proportion des données bruité** : la présence des valeurs douteuses , par exemple des éléments erronés , mal classés ,ou non-conformes au mode de distribution générale des exemples,ou même ayant des valeurs bizarres et aberrantes, peut obscurcir la vision de notre machine , et par la suite affaiblir sa capacité d'apprendre , à l'instar des données incomplètes on va franchir ce problème en se mettant **dans un cadre où la proportion de cette famille de données est négligeable** par rapport à la taille totale de la base .
- **le nombre et le type des attributs** : les attributs qui décrivent nos exemples déterminent entièrement le nombre de leurs classes, ce qui montre qu'ils jouent quand même un rôle au niveau de l'amélioration de la qualité d'apprentissage, surtout au cas du supervisé, or, si ce nombre est très élevé la tâche d'apprendre sera plus compliquée, et spécialement pour une machine qui adopte un modèle de classification simple comme le modèle linéaire par exemple, en plus le type de ces attributs et même des exemples peut apparaître comme un facteur dominant dans le processus d'apprentissage surtout pour le semi et le non-supervisé car dans lesquels on aura besoin de définir une **distance** entre deux exemples , une chose qui est simple dans le cas où ce type est quantitative (poids,âge,taille, ...) mais en revanche il ne l'est pas du tout pour le cas est qualitative (couleur,humeur, ...).

► **les facteurs liés à la machine d'apprentissage** : la machine est un élément clé dans le processus d'apprentissage et sa capacité reste **le paramètre le plus important** qui peut parfaire et perfectionner la qualité de ce processus, et pour les facteurs principaux qui sont liés à cette machine, on trouve :

- **le modèle** : on distingue souvent entre trois grands modèles de machines d'apprentissage et dont chacune possède des propriétés et des particularités propres à elle :
 - ★ **les machines à vecteurs de support** : également appelées réseaux de vecteurs de support, ou les séparateurs à vaste marge «**SVM**», ces machines constituent un ensemble de méthodes d'apprentissage supervisé **consacré aux problèmes de classification et de régression**, or, et à partir d'une collection des **exemples d'apprentissage**, dont chacun est identifié à une classe déterminée ,un algorithme d'apprentissage **SVM** fournit un modèle qui **prédit** si un nouvel exemple tombe dans une classe ou dans une autre.
cette machine s'est développée grâce aux travaux de **Vladimir Vapnik** sur la théorie d'apprentissage statistique (voir :[61]), et elle reste jusqu'à présent **un exemple flagrant d'une machine qui a été créée à base d'une étude théorique pure**.
 - ★ **les réseaux bayésiens** : ce sont des **machines calculatoires des probabilités conditionnelles**, et plus précisément des modèles graphique probabiliste qui représentent un ensemble de variables aléatoires et leur indépendance conditionnelle à l' aide d'**un graphe acyclique dirigé de causalité** (voir : [44]).
Un réseau bayésien peut nous servir par exemple, pour représenter les relations probabilistes entre les maladies et les symptômes, et compte tenu de ces derniers, le réseau peut calculer les probabilités de présence de diverses maladies.
pour l'apprentissage de ces réseaux, il est basé sur la recherche d'une **estimation des distributions de probabilités ou des paramètres des lois correspondantes**, à partir de données disponibles.
L'estimation de ces distributions est un sujet très vaste et complexe et dans ce sens

Il existe des méthodes et des algorithmes efficaces qui l'effectuent ^a, et si vous êtes intéressés par ce sujet nous vous invitons d'aller voir : [39].

- ★ **les réseaux de neurones artificiels** : cette machine fera l'objet de la section suivante, mais pour le moment, et sans entrer dans les détails, on peut la définir comme étant un groupe de neurones formels interconnectés entre eux, pour former un système semblable au réseau neuronal du cerveau.

Cette machine est très puissante et peut être utilisée pour résoudre plusieurs classes de problèmes par exemple les problèmes de **classification**, de **régression**, de **prédiction**, de **détection d'anomalies**, de **clustering**, ...etc, et elle était notamment adoptée pour réaliser pas mal de tâches, comme, la reconnaissance de la parole, de l'image, la traduction automatique, le diagnostic médical...etc.

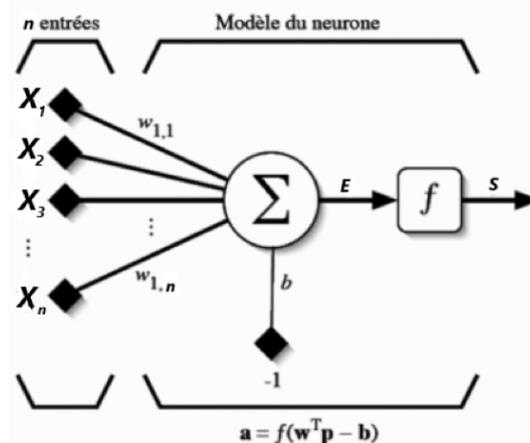
- **le type** : si on a spécifié le modèle de la machine qu'on va utiliser **pour résoudre un problème précis**, alors, on doit ensuite déterminer le type de ce modèle, par exemple : si on a pris comme modèle **les réseaux de neurones artificiels** alors on doit signaler est ce qu'il est un perceptron multicouche, un madaline, un réseau de Hopfield, une carte auto adaptative de Kohonen ...etc?.

(on aura après l'opportunité de traiter quelques-uns de ces types.)

- **l'algorithme d'apprentissage** : à ce stade le modèle et le type de la machine sont déjà déterminés, mais la façon par laquelle elle apprend ne l'est pas encore, pour cela on doit toujours indiquer l'algorithme d'apprentissage adopté par notre machine car il peut vraiment marquer la différence au niveau de sa qualité.

2.2 Optimisation et apprentissage supervisé d'un neurone formel :

Une mise en situation : on considère un perceptron simple (de Rosenblatt) qui se présente comme suit :



avec :

— " $\Sigma(\cdot)$ " est **la fonction d'entrée totale** qui fait la somme pondéré des x_i , et elle est définie par :

$$\Sigma(x_1, \dots, x_n) = \sum_{i=1}^n \omega_i x_i + \overbrace{\omega_0}^{:= -b} = \sum_{i=0}^n \omega_i x_i = \left\langle \begin{pmatrix} 1 \\ x \end{pmatrix}, \omega \right\rangle.$$

a. Par exemple dans le cas où toutes les variables sont observées, la méthode la plus simple et la plus utilisée est l'**estimation statistique** qui consiste à estimer la probabilité d'un événement par la fréquence son apparition dans la base de données, cette approche est appelée **le maximum de vraisemblance (MV)** et nous écrivons :

$$\hat{P}(X_i = x_k / \text{parents}(X_i) = x_j) = \hat{\theta}_{i,j,k} = \frac{N_{i,j,k}}{\sum_k N_{i,j,k}}$$

où $N_{i,j,k}$ est le nombre d'événements dans la base de données pour lesquels la variable X_i est dans l'état x_k et ses parents sont dans la configuration x_j .

- f est **la fonction d'état ou d'activation** qu'on peut la prendre comme : $\tanh(\cdot)$, $\text{sign}(\cdot)$, $\text{logsig}(\cdot)$, ...
- $\omega = (\omega_0, \omega_1, \dots, \omega_n)^\top \in \mathbb{R}^{n+1}$ est **le vecteur des pondération** associé à notre neurone.
- et enfin, $u = \sum_{i=1}^n \omega_i x_i + \omega_0$ et $S = f(u)$ sont respectivement l'entrée totale et la sortie fournie par le perceptron pour un vecteur d'entrée : $X = (x_1, \dots, x_n) \in \mathbb{R}^n$.

D'abord, et au cours de cette section il faut savoir qu'**on se pose dans le cadre d'apprentissage supervisé**, dans lequel on dispose d'une base d'apprentissage : $\mathcal{B} = (X^h, d^h)_{1 \leq h \leq \ell}$, avec :

- ★ ℓ : est la taille de la base \mathcal{B} .
- ★ et $(\forall h \in \{1, \dots, \ell\}) : X^h = (x_1^h, \dots, x_n^h)^\top \in \mathbb{R}^n$, et $d^h \in \mathbb{R}$ sont respectivement les composantes du couple (X^h, d^h) qui est un élément de notre base \mathcal{B} qu'on l'appel "**exemplaire**", le vecteur X^h est un **vecteur exemple** et d^h est tout simplement **la sortie désirée associée à cet exemple**.

L'apprentissage de ce perceptron consiste en fait à trouver une façon ou une méthode pour le rendre "optimal" plus clairement on veut qu'il approche le mieux notre base d'apprentissage. une manière traditionnelle pour attribuer un sens à ce terme de "**mieux**", est d'adopter le critère des moindres carrés « **L·M·S** », (Least Mean Squares), et qui est inspiré essentiellement de la statistique.

ce critère est **basé sur la minimisation d'une certaine erreur** E comme étant **une fonction des pondérations et de biais** $(\omega_0, \omega_1, \dots, \omega_n) \in \mathbb{R}^{n+1}$, et plus explicitement, pour trouver **le perceptron optimal**, il faut et il suffit de déterminer les poids :

$$(\omega_0^*, \dots, \omega_n^*)^\top \in \mathbb{R}^{n+1}$$

qui minimisent l'erreur totale moyenne :

$$C(\omega) = \frac{1}{\ell} \sum_{h=1}^{\ell} \frac{1}{2} (d^h - g(X^h, \omega))^2.$$

Autrement dit, **le problème d'apprentissage va se ramener enfin de compte au problème d'optimisation** suivant :

$$(\mathcal{P}) \left\{ \begin{array}{l} \min C(\omega) = \frac{1}{2\ell} \sum_{h=1}^{\ell} (d^h - g(X^h, \omega))^2 \\ \omega \in \mathbb{R}^{n+1} \end{array} \right.$$

Remarques 2.2.1 :

1. Certain ouvrages ([16] par exemple) font la distinction entre trois types d'erreurs :

-**l'erreur locale** : ou instantané c'est tout simplement l'erreur pour un seul exemple, et il est défini par :

$$e_h(\omega) = \frac{1}{2} (d^h - g(X^h, \omega))^2.$$

-**l'erreur totale** : c'est l'erreur provenant de toute la base d'apprentissage et ce n'est autre que la somme de tout les erreurs locales de ses exemples, et on écrit donc :

$$E(\omega) = \sum_{h=1}^{\ell} e_h(\omega) = \sum_{h=1}^{\ell} \frac{1}{2} (d^h - g(X^h, \omega))^2.$$

-**l'erreur moyenne** : comme son nom l'indique c'est l'erreur en moyenne pour tout les exemples d'apprentissage, et on note : $C(\omega) = \frac{1}{\ell} E(\omega) = \frac{1}{\ell} \sum_{h=1}^{\ell} \frac{1}{2} (d^h - g(X^h, \omega))^2$.

2. on peut voir facilement que : $E(\omega) = 0 \iff (\forall h \in \{1, \dots, \ell\}) : e_h(\omega) = 0$.

3. pour le problème d'optimisation (\mathcal{P}) , l'adjonction du coefficient $\frac{1}{\ell}$ dans l'expression de la fonction "objectif" est essentiellement due à l'étude du cas aléatoire, où sa présence était indispensable pour légitimer le passage de la minimisation de l'erreur théorique à la minimisation de l'erreur empirique en faisant appel au théorème central limite, mais dans notre cas ici on se pose dans un cadre déterministe et pour cela on pourra et sans aucun souci l'éliminer^a de la syntaxe du problème (\mathcal{P}) .

2.3 Quelques algorithmes d'apprentissage supervisé pour un neurone formel :

2.3.1 Algorithme de descente de gradient :

À ce stade, on a bien vu que le principe général d'apprentissage supervisé pour un neurone formel consiste en fait de trouver un élément neuronal optimal qui approche le mieux notre base de données, ce qui est revenu comme on a déjà dit précédemment, à la résolution du problème d'optimisation :

$$(\mathcal{P}) \left\{ \begin{array}{l} \min E(\omega) = \frac{1}{2} \sum_{h=1}^{\ell} \left(d^h - g(X^h, \omega) \right)^2 \\ \omega \in \mathbb{R}^{n+1} \end{array} \right.$$

et pour le résoudre on va adopter des algorithmes ou des méthodes adaptatives, qui vont se baser sur la génération d'une suite $(\omega^k)_{k \in \mathbb{N}}$ des poids qui convergent vers le minimum de l'erreur locale $e_h(\omega)$, c'est à dire que ces algorithmes vont minimiser l'erreur locale et non pas l'erreur total ou moyen.

la famille des algorithmes qui peut se présenter dans ce cas comme un candidat idéal est tout simplement la famille de **descente de gradient**.

maintenant **on suppose que f est différentiable** ou plus simplement dérivable (car c'est une fonction d'une seule variable réelle).

$$\begin{aligned} \text{on a : } (\forall i \in \{0, \dots, n\}) : \frac{\partial e_h(\omega)}{\partial \omega_i} &= \frac{\partial}{\partial \omega_i} \left(\frac{1}{2} (S^h - d^h)^2 \right) = \frac{\partial}{\partial S^h} \left(\frac{1}{2} (S^h - d^h)^2 \right) \times \frac{\partial S^h}{\partial \omega_i} \\ &= (S^h - d^h) \times \frac{\partial f(u)}{\partial u} \times \frac{\partial u}{\partial \omega_i} = (S^h - d^h) \times f'(u) \times \frac{\partial \sum_{k=0}^n \omega_k x_k^h}{\partial \omega_i} \\ &= (S^h - d^h) \times f' \left(\sum_{k=0}^n \omega_k x_k^h \right) \times x_i^h \\ &= \underbrace{-(d^h - S^h) \times f' \left(\sum_{k=0}^n \omega_k x_k^h \right)}_{:= \delta} \times x_i^h \\ &= -\delta \cdot x_i^h \end{aligned}$$

alors le vecteur gradient de la fonction e_h s'écrira ainsi :

$$\nabla_{\omega} e_h(\omega) := \begin{pmatrix} \frac{\partial e_h(\omega)}{\partial \omega_0} \\ \vdots \\ \frac{\partial e_h(\omega)}{\partial \omega_n} \end{pmatrix} = \begin{pmatrix} -\delta \cdot x_0^h \\ \vdots \\ -\delta \cdot x_n^h \end{pmatrix} = -\delta X^h.$$

avec X^h est un exemple (un vecteur d'entrée) de notre classe.

a. bien sûr car le nombre ℓ est positif, en fait et pour la même raison on peut éliminer à fois le "2" de l'expression de la fonction à minimiser mais, il est préférable de la garder pour alléger l'expression du gradient par la suite.

donc dans ce cas l'algorithme d'apprentissage qu'on aura sera de la forme :

Algorithme 1 Algorithme d'apprentissage du perceptron simple par la descente de gradient.

Entrée: une base de données : $(X^h, d^h)_{1 \leq h \leq \ell}$, une fonction d'activation dérivable f , des poids de départ : $\omega = (\omega_0, \dots, \omega_n) \in \mathbb{R}^{n+1}$ et un nombre maximal d'itérations N_{\max} .

Sortie: les poids optimaux : $\omega = (\omega_0, \dots, \omega_n)$

```

1:  $k \leftarrow 1$ ;
2: tant que  $k \leq N_{\max}$  faire:
3:   pour  $h$  allant de 1 à  $\ell$  faire:
4:     pour  $i$  allant de 0 à  $n$  faire:
5:        $u^h \leftarrow \sum_{j=0}^n \omega_j^{(k)} x_j^h$ ;
6:        $S^h \leftarrow f(u^h)$ ;
7:        $e^h \leftarrow d^h - S^h$ ;
8:        $\delta \leftarrow e^h \times f'(u^h)$ ;
9:        $\omega_i^{(k+1)} \leftarrow \omega_i^{(k)} + \eta_k \times \delta \times x_i^h$ ;
10:    fin pour
11:  fin pour
12: fin tant que
13: retourner:  $\omega = (\omega_0, \dots, \omega_n)$ 

```

Remarques 2.3.1 :

1. Comme son nom l'indique cet algorithme est inspiré de celui de descente de gradient usuel (Cauchy 1847, Hadamard 1908), appliqué à la fonction d'erreur locale $e_h(\omega)$.
2. Le critère d'arrêt peut être choisi autrement, par exemple : $\|\nabla E(\omega)\| \leq \varepsilon, \|\omega^{(k+1)} - \omega^{(k)}\| \leq \varepsilon, \dots$
3. η_k est tout simplement le pas de descente qu'on l'appelle dans ce cas **le taux d'apprentissage**, et dans la littérature il y en a plusieurs façons pour le choisir, par exemple le fixer dès le début (algorithme de gradient à pas fixé) ou avant chaque itération (à pas prédéterminé) ou le prendre par une procédure d'optimisation unidimensionnelle ^a (à pas optimal)...
4. En général rien ne prouve que la diminution d'erreur pour un point va produire une diminution pour les autres aussi, mais cette façon de procéder est commune entre tous les algorithmes adaptatifs.

Théorème 2.3.1 :

1. Si on adopte pour l'algorithme précédent un choix prédéterminé du pas d'apprentissage et qui vérifie les deux conditions : $\lim_{k \rightarrow +\infty} \eta_k = +\infty$ et $\sum_{k=0}^{+\infty} \eta_k = +\infty$, alors il convergera globalement, c'est à dire pour n'importe quel point de départ, vers les pondérations optimales.
2. si la fonction d'erreur $\omega \mapsto E(\omega)$ est de classe $C^\infty(\mathbb{R})$, avec : $\lim_{\|\omega\| \rightarrow +\infty} E(\omega) = +\infty$, et si notre actualisation du pas η est faite par une procédure d'optimisation unidimensionnelle, alors pour toute initialisation des poids, l'algorithme d'apprentissage par descente de gradient convergera mais vers des pondérations critiques et pas forcément optimales.

a. c'est à dire prendre : $\eta = \underset{\lambda > 0}{\text{Argmin}} E(\omega^{old} - \lambda \cdot \nabla E(\omega^{old}))$.

Preuve :

La preuve des deux critères précédents se découle nécessairement des théorèmes usuels de convergence pour l'algorithme de descente de gradient, voir pages :15, 17 dans [19] , et en plus on verra ultérieurement que cet algorithme n'est autre qu'un cas particulier de celui de la rétropropagation de gradient ce qui fait que ces critères vont être généralisés dans ce cas voir pour cela [10] page 215. ■

malgré la simplicité et l'efficacité de la procédure d'apprentissage par descente de gradient, **la particularité du choix de la fonction d'activation entrave sa capacité**, or et dans plusieurs cas pratiques , par exemple les problèmes de décision, on peut se trouver face à des fonctions d'activation non dérivables, et dans ce contexte cet algorithme devient invalide, et pour cette raison on sera obligé d'introduire des modifications soit au niveau du problème, soit au niveau de l'algorithme.

▷ **modifier le problème** : le principe de cette approche est facile mais elle peut nous sortir de plusieurs pétrins, qui sont souvent rencontrés dans les problèmes de décision, et en somme il consiste à remplacer ou substituer la fonction d'activation f et qui est par hypothèse non dérivable par une autre \tilde{f} qui est dérivable et très proche de f .

cette idée va apparaître dans la preuve du théorème : III.1.3.2 concernant l'approximation des fonctions de décision par des réseaux neuronaux à une seule couche cachée.

Exemple 2.3.1 :

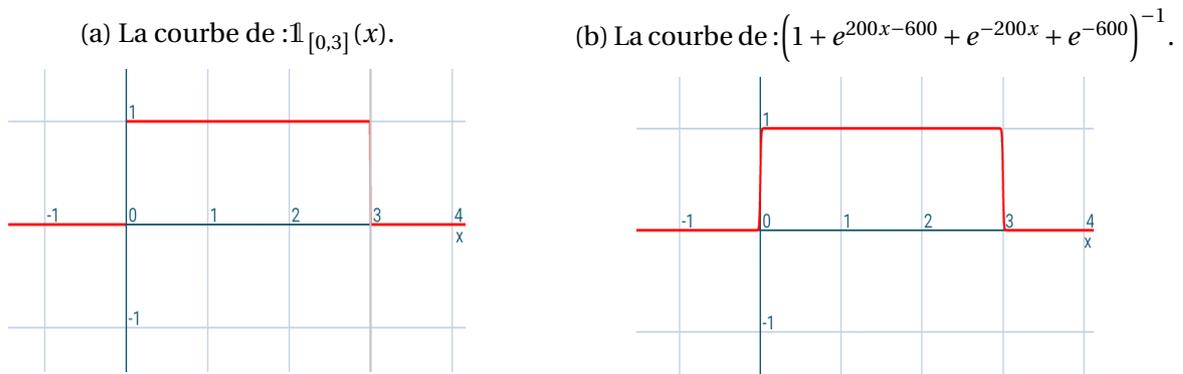


FIGURE II.18: La fonction (b) est régulière et elle peut bien jouer le rôle de la fonction (a).

Remarque 2.3.2 :

Pour l'approche de modification du problème, toutes les techniques permettant de faire l'approximation régulière ^a, comme le développement limité, l'interpolation lagrangienne, l'interpolation Hermitienne ^b ... etc, peuvent être évoquées pour nous servir dans ce cas, mais malgré cela la réalisation d'une telle approximation est parfois très difficile et moins évidente à trouver (voir l'exemple précédent), ce qui ne nous laissera aucun choix que de modifier l'algorithme.

a. C'est à dire qui approchent notre fonction f par une autre régulière. (par exemple l'interpolation linéaire n'est pas acceptée dans ce cas).

b. C'est une extension de l'interpolation de Lagrange, qui consiste, pour une fonction dérivable en un nombre fini de points donnés à construire un polyôme interpolateur dont les valeurs de la dérivée à ces points coïncident avec celles de la dérivée de la fonction, l'avantage de cette méthode c'est qu'elle permet de donner une meilleure approximation lorsque le nombre des points augmentent, c'est à dire qu'elle peut franchir le phénomène de Runge (voir : [15]).

▷ **modifier l'algorithme** : contrairement au principe précédent, cette approche consiste à garder le problème tel qu'il est mais de modifier l'algorithme pour qu'il s'adapte à la non différentiabilité de la fonction d'activation, cette approche va donner naissance à deux fameux algorithmes d'apprentissage et qui sont *l'algorithme d'apprentissage du perceptron* et *l'algorithme d'adaline de Widrow-Hoff*.

2.3.2 L'algorithme d'apprentissage du perceptron :

La règle de Hebb : le neuropsychologue canadien **Donald Hebb** (1904 – 1985) s'est basé sur des conjectures biologiques pour affirmer que les neurones se lient entre eux s'ils s'excitent ensemble, c'est à dire que lorsque deux neurones sont excités conjointement, il se crée ou renforce un lien qui les unissent, cette idée sera la base de la notion d'apprentissage automatique, et elle fournira le tout premier algorithme de réajustement des poids, et qui était apparu en 1949 dans l'article de Hebb[29] sous la formulation moderne suivante :

Algorithme 2 Algorithme de réajustement des poids de Hebb.

Entrée: une collection d'exemples $(X^h)_{1 \leq h \leq \ell}$, des poids de départ $\omega = (\omega_0, \dots, \omega_n) \in \mathbb{R}^{n+1}$, un facteur de correction η , et un nombre maximal d'itérations N_{\max} .

Sortie: des poids plus convenables $\omega = (\omega_0, \dots, \omega_n)$

- 1: $k \leftarrow 1$;
- 2: **tant que** $k \leq N_{\max}$ **faire**:
- 3: **pour** h allant de 1 à ℓ **faire**:
- 4: **pour** i allant de 0 à n **faire**:
- 5: $u^h \leftarrow \sum_{j=0}^n \omega_j^{(k)} x_j^h$;
- 6: $S^h \leftarrow \mathbb{1}_{\mathbb{R}^+}(u^h)$; (la fonction d'activation dans cas est la fonction Heavside.)
- 7: $\omega_i^{(k+1)} \leftarrow \omega_i^{(k)} + \eta_k \times (S^h \cdot x_i^h)$;
- 8: **fin pour**
- 9: **fin pour**
- 10: **fin tant que**
- 11: **retourner**: $\omega = (\omega_0, \dots, \omega_n)$

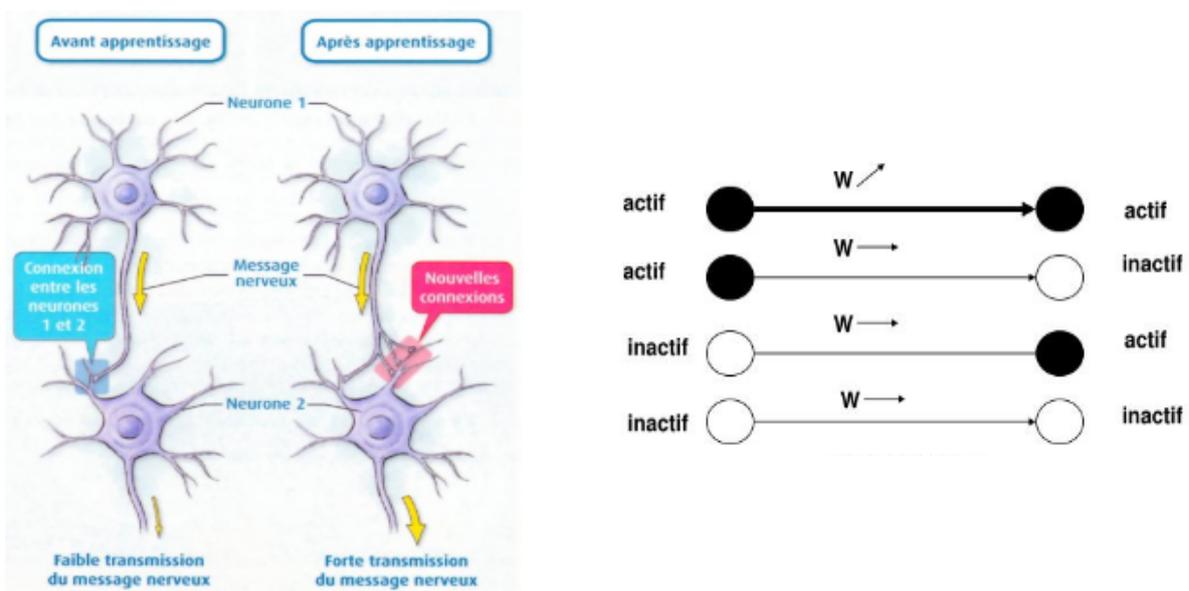


FIGURE II.19: Une illustration de la règle de Hebb.

Remarques 2.3.3 :

1. l'idée de cette règle était juste une conjecture et une vision théorique jusqu'à 2000, quand le médecin et le chercheur américain **Eric Kandel** a reçu le prix Nobel en fournissant une preuve expérimentale qui assure l'implication des mécanismes d'apprentissages de Hebb dans les synapses réels.
2. cette règle ne dispose pas des théorèmes de convergence, et reste inapplicable dans certains cas malgré l'existence de la solution, ce qui fait qu'elle n'est pas encore à la hauteur d'être considérée comme un algorithme d'apprentissage.

l'algorithme du perceptron cet algorithme est exactement analogue à la règle précédente, juste il a ajouté une idée un peu maline dans la correction des poids en s'inspirant de la descente de gradient, or, à la place de corriger comme suit : $\omega_i^{(k+1)} \leftarrow \omega_i^{(k)} + \eta_k \times (x_i^h \cdot S^h) \times x_i^h$ on va le faire ainsi : $\omega_i^{(k+1)} \leftarrow \omega_i^{(k)} + \eta_k \times (d^h - S^h) \times x_i^h$, et de ce fait l'algorithme se présente ainsi :

Algorithme 3 Algorithme d'apprentissage du perceptron.

Entrée: une base de données : $(X^h, d^h)_{1 \leq h \leq \ell}$, une fonction d'activation dérivable f , des poids de départ : $\omega = (\omega_0, \dots, \omega_n) \in \mathbb{R}^{n+1}$ et un nombre maximal d'itérations N_{\max} .

Sortie: les poids optimaux : $\omega = (\omega_0, \dots, \omega_n)$

- 1: $k \leftarrow 1$;
 - 2: **tant que** $k \leq N_{\max}$ **faire:**
 - 3: **pour** h allant de 1 à ℓ **faire:**
 - 4: **pour** i allant de 0 à n **faire:**
 - 5: $u^h \leftarrow \sum_{j=0}^n \omega_j^{(k)} x_j^h$;
 - 6: $S^h \leftarrow \mathbb{1}_{\mathbb{R}^+}(u^h)$; (la fonction d'activation dans cas est la fonction Heavside.)
 - 7: $e^h \leftarrow d^h - S^h$;
 - 8: $\omega_i^{(k+1)} \leftarrow \omega_i^{(k)} + \eta_k \times e^h \times x_i^h$;
 - 9: **fin pour**
 - 10: **fin pour**
 - 11: **fin tant que**
 - 12: **retourner:** $\omega = (\omega_0, \dots, \omega_n)$
-

Théorème 2.3.2 : (Novikov 1962 [1])

Si la base $\mathcal{B} = (X^h, d^h)_{1 \leq h \leq \ell}$ est linéairement séparable par un hyperplan passant par l'origine, la procédure d'apprentissage précédente *convergera globalement*, c'est à dire pour n'importe quelles pondérations initiales, *et en un nombre fini d'itérations* vers un perceptron sans seuil résolvant le problème de décision déterminé par \mathcal{B} .

Preuve :

Soit $n + 1$ le nombre d'entrées du perceptron, on suppose que chaque exemple possède une $n + 1$ -ème coordonnée égale à 1 (qui correspond au biais).

on note \mathcal{E}_0 et \mathcal{E}_1 respectivement l'ensemble *des exemples* de \mathcal{B} pour lesquels le perceptron doit retourner la valeur 1 et la valeur 0.

par hypothèse, il existe un hyperplan de l'espace \mathbb{R}^{n+1} qui sépare \mathcal{E}_0 de \mathcal{E}_1 , ce qui fait qu'il existe

un vecteur $V = (v_0, \dots, v_n)^T$ tel que :

$$\begin{cases} \langle V, X \rangle < 0 & \text{si } X \in \mathcal{E}_0 \\ \langle V, X \rangle \geq 0 & \text{si } X \in \mathcal{E}_1 \end{cases}$$

et puisque \mathcal{B} est finie, on trouvera forcément un réel $\varepsilon > 0$ strictement positif tel que :

$$(\forall X \in \mathcal{E}_0) : \langle V, X \rangle \leq -\varepsilon \quad \text{et} : \quad (\forall X \in \mathcal{E}_1) : \langle V, X \rangle \geq \varepsilon$$

alors pour tout X de \mathcal{E} , on a : $|\langle V, X \rangle| \geq \varepsilon$.

soit M un majorant ^a de $\{X^2 = \langle X, X \rangle / X \in \mathcal{E}\} \subset \mathbb{R}^+$.

Soit $\omega^{(0)} = (\omega_0^{(0)}, \omega_1^{(0)}, \dots, \omega_n^{(0)}) \in \mathbb{R}^{n+1}$ les valeurs d'initialisation des poids synaptiques. on procédera par l'absurde, en supposant que l'algorithme d'apprentissage ne s'arrête pas et en montrant par la suite que cette hypothèse va nous ramener à une contradiction.

on note $(\omega^{(k)})_{k \in \mathbb{N}} \in (\mathbb{R}^{n+1})^{\mathbb{N}}$ les états successifs différents ^b des coefficients synaptiques (c'est-à-dire que pour tout $k \in \mathbb{N}$, on a $\omega^{(k)} \neq \omega^{(k+1)}$), comme on a supposé que l'algorithme ne s'arrête pas, cela implique en particulier qu'il y a un nombre infini des vecteurs ω_i .

on va déterminer des bornes pour ; ω_i^2 et $\omega.V$.

1. **Premier cas** : on présente un exemple X de \mathcal{E}_1 :

- si $\langle \omega^{(k)}, X \rangle \geq 0$, alors la réponse fournie par le neurone sera juste, ce qui fait que l'algorithme ne modifiera pas les poids.
- si $\langle \omega^{(k)}, X \rangle < 0$, alors la réponse est fautive ce qui fait que les poids seront modifiés, donc : $\langle V, X \rangle > \varepsilon$ et : $\omega^{(k+1)} = \omega^{(k)} + \eta(d - S).X$, d'où :

$$\begin{aligned} (\omega^{(k+1)})^2 &= (\omega^{(k)} + \eta(d - S).X)^2 = (\omega^{(k)} + \eta(1 - 0).X)^2 \\ &= (\omega^{(k)})^2 + 2\eta \cdot \langle \omega^{(k)}, X \rangle + (\eta.X)^2 \\ &\leq (\omega^{(k)})^2 + \eta^2.M \end{aligned}$$

$$\text{d'une autre part} : \langle \omega^{(k+1)}, V \rangle = \langle \omega^{(k)} + \eta.X, V \rangle = \langle \omega^{(k)}, V \rangle + \eta \cdot \langle X, V \rangle \geq \langle \omega^{(k)}, V \rangle + \eta \cdot \varepsilon$$

2. **deuxième cas** : on présente maintenant un exemple X de \mathcal{E}_0 :

- si $\langle \omega^{(k)}, X \rangle < 0$, alors la réponse fournie par le neurone sera juste, ce qui fait que l'algorithme ne fera aucune modification sur les poids.
- si $\langle \omega^{(k)}, X \rangle \geq 0$, alors la réponse est fautive ce qui fait que les poids seront modifiés, donc : $\langle V, X \rangle < -\varepsilon$ et : $\omega^{(k+1)} = \omega^{(k)} + \eta(d - S).X$, d'où :

$$\begin{aligned} (\omega^{(k+1)})^2 &= (\omega^{(k)} + \eta(d - S).X)^2 = (\omega^{(k)} + \eta(0 - 1).X)^2 \\ &= (\omega^{(k)})^2 - 2\eta \cdot \langle \omega^{(k)}, X \rangle + (\eta.X)^2 \\ &\leq (\omega^{(k)})^2 + \eta^2.M \end{aligned}$$

$$\text{d'une autre part} : \langle \omega^{(k+1)}, V \rangle = \langle \omega^{(k)} - \eta.X, V \rangle = \langle \omega^{(k)}, V \rangle - \eta \cdot \langle X, V \rangle \geq \langle \omega^{(k)}, V \rangle + \eta \cdot \varepsilon$$

a. bien sûr cet ensemble est majoré car il est fini.

b. en fait pour cet algorithme si deux poids successifs sont identiques, alors il va aboutir son état de stabilité et il ne modifiera aucun poids par la suite.

finalement, on déduit que dans tout les cas, et pour n'importe quel $k \in \mathbb{N}$, on a :

$$\langle \omega^{(k+1)}, V \rangle \geq \langle \omega^{(k)}, V \rangle + \eta \cdot \varepsilon \text{ et : } \left(\omega^{(k+1)} \right)^2 \leq \left(\omega^{(k)} \right)^2 + \eta^2 \cdot M.$$

on en déduit que pour tout : $k \geq 1$:

$$\langle \omega^{(k)}, V \rangle \geq \langle \omega^{(0)}, V \rangle + k \times \eta \cdot \varepsilon \text{ et : } \left(\omega^{(k)} \right)^2 \leq \left(\omega^{(0)} \right)^2 + k \times \eta^2 \cdot M$$

comme : $\varepsilon > 0$, on trouvera forcément un $k_0 \in \mathbb{N}$, prenez par exemple : $k_0 = \left\lceil \frac{-\langle \omega^{(0)}, V \rangle}{\varepsilon} \right\rceil + 1$, tel

que : $\langle \omega^{(0)}, V \rangle + k \times \eta \cdot \varepsilon > 0$.

donc pour : $k \geq k_0$ on a :

$$\left(\langle \omega^{(0)}, V \rangle + k \times \eta \cdot \varepsilon \right)^2 \leq \left(\langle \omega^{(k)}, V \rangle \right)^2 \leq V^2 \times \left(\omega^{(0)} \right)^2 + k \times \eta^2 \cdot M.$$

ce qui est impossible, car le terme de gauche est un polynôme de deuxième degré en k alors que le terme de droite est linéaire en k .

d'où ce qu'il faut prouver. ■

Remarque 2.3.4 :

On peut voir via la démonstration précédente que dans le cas d'un choix fixé du taux d'apprentissage, ce dernier n'a pas d'effet sur la convergence de l'algorithme, car il restera convergent en tout cas, mais malgré cela, un bon choix de ce paramètre peut réduire le nombre de corrections exécutées, en outre cet algorithme nous ramène (dans le cas d'un choix fixé) vers des pondérations convenables mais qui ne sont pas forcément optimales.

2.3.3 L'algorithme d'apprentissage de Widrow-Hoff :

Cet algorithme a été proposé en 1960 par les deux chercheurs américains **Bernard Widrow** et **Ted Hoff** dans leur rapport technique [62], son idée générale est de franchir la non dérivabilité de f par une annulation totale de cette fonction.

autrement dit on va prendre exactement l'algorithme de descente de gradient mais en supposant que la fonction f n'existe plus c'est à dire qu'on va l'appliquer juste à la partie Adaline de notre perceptron ce qui justifie le fait de l'appeler parfois l'algorithme d'apprentissage d'adaline et on aura dans ce cas :

Algorithme 4 Algorithme d'apprentissage de Widrow-Hoff.

Entrée: une base de données : $(X^h, d^h)_{1 \leq h \leq \ell}$, une fonction d'activation non dérivable f , des poids de départ : $\omega = (\omega_0, \dots, \omega_n) \in \mathbb{R}^{n+1}$ et un nombre maximal d'itérations N_{\max} .

Sortie: les poids optimaux : $\omega = (\omega_0, \dots, \omega_n)$

- 1: $k \leftarrow 1$;
 - 2: **tant que** $k \leq N_{\max}$ **faire:**
 - 3: **pour** h allant de 1 à ℓ **faire:**
 - 4: **pour** i allant de 0 à n **faire:**
 - 5: $u^h \leftarrow \sum_{j=0}^n \omega_j^{(k)} x_j^h$;
 - 6: $\omega_i^{(k+1)} \leftarrow \omega_i^{(k)} + \eta_k \times (d^h - u^h) \times x_i^h$;
 - 7: **fin pour**
 - 8: **fin pour**
 - 9: **fin tant que**
 - 10: **retourner:** $\omega = (\omega_0, \dots, \omega_n)$
-

On remarque que cet algorithme garde la même procédure de correction que l'algorithme d'apprentissage du perceptron, mais il modifie les pondérations à chaque itération et dans tous les cas, et pas seulement au cas de la présence d'une erreur^a.

l'avantage de cet algorithme par rapport au précédent est le fait qu'il ne s'intéresse pas à la sortie calculée par le neurone, ce qui fait que même si l'information fournie est erronée, (c'est-à-dire que certains exemples de l'échantillon sont mal classés par le perceptron), l'algorithme restera encore en cours d'exécution et de plus si on choisit convenablement le paramètre η , il va converger vers une solution optimale (voir [16] page : 244).

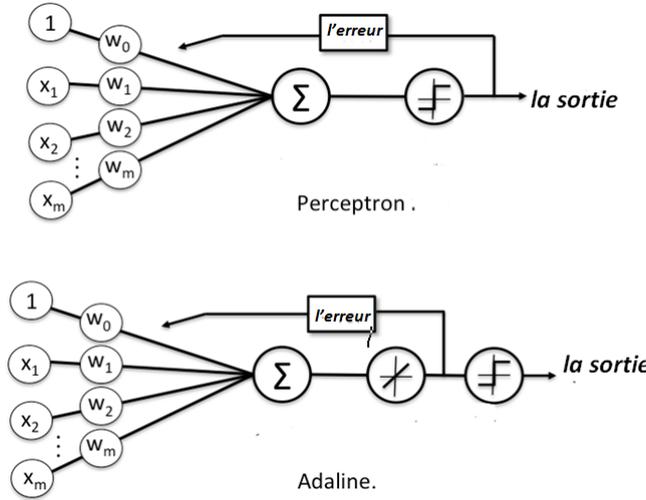


FIGURE II.20: La différence entre l'algorithme de Perceptron et celui de Widrow-Hoff.

Théorème 2.3.3 :

Si le problème déterminé par la base $\mathcal{B} = (X^h, d^h)_{1 \leq h \leq \ell}$ admet une solution optimale ω^* , et si on affecte au pas d'apprentissage une valeur positive très petite alors les pondérations fournies par l'algorithme de Widrow-Hoff s'approchent^b, de cette solution à chaque itération.

Preuve :

Pour montrer la convergence de cet algorithme, nous allons calculer la distance entre les termes $\omega^k \in \mathbb{R}^{n+1}$ qu'il engendre et les pondérations optimales $\omega^* \in \mathbb{R}^{n+1}$, et dans ce cadre on a :

$$\begin{aligned}
 \left\| \omega^{(k+1)} - \omega^* \right\|^2 - \left\| \omega^{(k)} - \omega^* \right\|^2 &= \sum_{i=0}^n \left(\left(\omega_i^{(k+1)} - \omega_i^* \right)^2 - \left(\omega_i^{(k)} - \omega_i^* \right)^2 \right) \\
 &= \sum_{i=0}^n \left(\left(\omega_i^{(k)} + \eta \times (d - u) x_i - \omega_i^* \right)^2 - \left(\omega_i^{(k)} - \omega_i^* \right)^2 \right) \\
 &= \sum_{i=0}^n \left(2 \times \eta \times (d - u) \left(\omega_i^{(k)} - \omega_i^* \right) x_i + \eta^2 \times (d - u)^2 x_i^2 \right) \\
 &= (d - u) \sum_{i=0}^n \left(2\eta \times \left(\omega_i^{(k)} - \omega_i^* \right) x_i + \eta^2 (d - u) x_i^2 \right) \\
 \left(\text{car : } u &= \sum_{i=0}^n \omega_i^{(k)} x_i \text{ et : } u^* = \sum_{i=0}^n \omega_i^* x_i \right) &= 2\eta (d - u) \sum_{i=0}^n \left(\left(\omega_i^{(k)} - \omega_i^* \right) x_i \right) + \eta^2 \times (d - u)^2 \sum_{i=0}^n x_i^2 \\
 &= \eta^2 (d - u)^2 \|x\|^2 - 2\eta (d - u) (u^* - u)
 \end{aligned}$$

a. Et pour ça on dit qu'il n'est pas un algorithme d'apprentissage par correction d'erreur

b. **Attention!**: s'approchent de cette solution et pas convergent vers elle.

si on prend maintenant un nombre M tel que : $\|X\| \leq M$ (par exemple : $M = \max_{h \in [1, \ell]} \|x^h\|$), alors on aura :

$$\begin{aligned}
\eta^2 (d-u)^2 \|x\|^2 - 2\eta (d-u) (u^* - u) &\leq \eta^2 M^2 (d-u)^2 - 2\eta (d-u) (u^* - u) \\
&\leq \eta^2 M^2 (d-u)^2 - 2\eta (d-u) (u^* - d + d - u) \\
&\leq \eta^2 M^2 (d-u)^2 - 2\eta (d-u)^2 - 2\eta (d-u) (u^* - d) \\
&\leq (\eta M^2 - 2) \times \eta (d-u)^2 - 2\eta (d-u^* + u^* - u) (u^* - d) \\
&\leq (\eta M^2 - 2) \times \eta (d-u)^2 + 2\eta (d-u^*)^2 + 2\eta (d-u^*) (u^* - u) \\
\left(\text{car : } a \times b \leq \frac{1}{4} (a+b)^2 \right) &\leq (\eta M^2 - 2) \times \eta (d-u)^2 + 2\eta (d-u^*)^2 + \frac{\eta}{2} (d-u)^2 \\
&\leq \left(\eta M^2 - \frac{3}{2} \right) \times \eta (d-u)^2 + 2\eta (d-u^*)^2
\end{aligned}$$

donc on peut affirmer que :

$$\|\omega^{(k+1)} - \omega^*\|^2 - \|\omega^{(k)} - \omega^*\|^2 \leq \left(\eta M^2 - \frac{3}{2} \right) \times \eta (d-u)^2 + 2\eta (d-u^*)^2$$

de plus on a :

$$\begin{aligned}
\left(\eta M^2 - \frac{3}{2} \right) \times \eta (d-u)^2 + 2\eta (d-u^*)^2 \leq 0 &\iff 2\eta (d-u^*)^2 \leq \left(\frac{3}{2} - \eta M^2 \right) \times \eta (d-u)^2 \\
&\iff 4(d-u^*)^2 \leq (3 - 2 \times \eta M^2) \times (d-u)^2 \\
&\iff \underbrace{(d-u^*)^2}_{\text{erreur commise par } \omega^*} \leq \frac{(3 - 2 \times \eta M^2)}{4} \times \underbrace{(d-u)^2}_{\text{erreur commise par } \omega}
\end{aligned}$$

donc si on prend η si petit, par exemple : $\eta \leq \frac{3}{2M^2}$, la suite $(\|\omega^{(k)} - \omega^*\|)_{k \in \mathbb{N}}$ sera décroissante, ce qui fait que les $\omega^{(k)}$ s'approchent à chaque itération de la solution optimale, d'où ce qu'il faut établir. ■

Remarque 2.3.5 :

Jusqu'à l'instant, il n'y a pas une règle qui peut nous garantir la convergence de l'algorithme d'apprentissage d'adaline vers la solution optimale, ou même vers un perceptron capable de classer correctement tout les vecteurs d'entraînement (voir [40] pages 10 et 11).

3 Réseaux de neurones artificiels :

3.1 Nécessité d'introduire la structure de réseau :

Parfois, on se pose dans des cas où un seul perceptron n'est plus capable de réaliser la tâche désirée, et que la réalisation de cette dernière demande la contribution de plusieurs neurones, à ce stade là on fera l'appel à ce qu'on appelle **les réseaux neuronaux**.

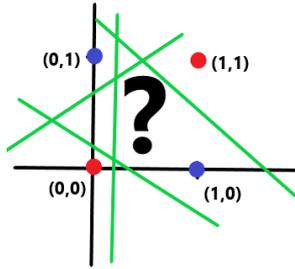
Définition 3.1.1 :

Un réseau neuronal est une collection de plusieurs neurones reliés et inter-connectés entre eux selon une architecture déterminée.

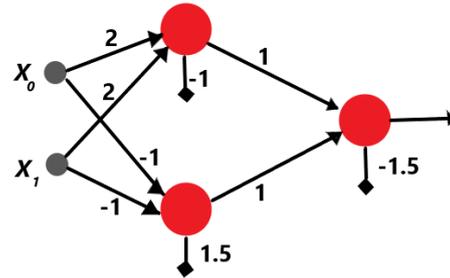
une architecture désigne en fait la topologie, la façon ou la méthode selon laquelle les neurones

de notre réseau sont inter-connectés, on verra par la suite quelques exemples célèbres d'architectures.

Exemple 3.1.1 :



(a) l'incapacité d'un seul perceptron face au problème "Xor".



(b) Un réseau résolvant le problème : "Xor".

3.2 Architectures des réseaux de neurones artificiels :

On distingue entre deux grands types d'architectures de réseaux de neurones : Les réseaux feedforward non récurrents et les réseaux dynamiques récurrents.

3.2.1 Les réseaux feedforward :

Dans un réseaux feedforward, appelé aussi statique ou non récurrent, la sortie courante d'un nœud ne peut jamais apporter un effet sur ses sorties futures, c'est à-dire qu'elle ne peut pas être injectée à son entrée ni directement ni indirectement via d'autres neurones, ce qui fait que, l'information circule dans un seule sens, de l'entrée vers la sortie, et ceci justifie bien l'utilisation du terme "feedforward".

généralement ce type de réseaux réalisent des transformations de la forme : $O = \Psi(x)$ où $x \in \mathbb{R}^n$ et $O \in \mathbb{R}^m$, sont respectivement les vecteurs d'entrée et de sortie de notre réseau, en plus dans son architecture la plus générale, l'entrée de chaque neurone est connectée à toutes les sorties des neurones précédents, mais, la plupart des réseaux neuronaux statiques utilisés, sont organisés en couches, et pour ça ils sont appelés des réseaux multicouches ou perceptrons multicouches. un réseau multicouche comporte généralement : une couche d'entrée, une ou plusieurs couches cachées et une couche de sortie, et dans un tel réseaux, la sortie de chaque neurone d'une couche est connectée seulement aux entrées de chaque neurone de la couche suivante.

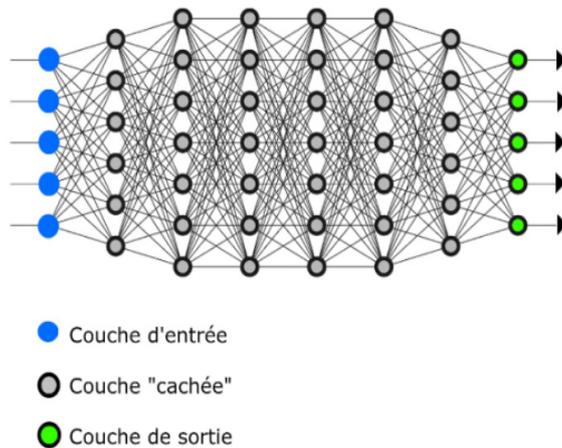


FIGURE II.22: La forme d'un réseau feedforward.

3.2.2 Les réseaux dynamiques :

Ces réseaux, qui sont appelés aussi réseaux récurrents, ou bouclés sont organisés de telle sorte que chaque neurone reçoit sur ses entrées en plus des informations externes, une partie ou la totalité de l'état du réseau c'est à dire des sorties d'autres neurones, ce qui s'interprète par le fait qu'il existe au moins un feedback d'une couche sur la précédente, autrement dit il existe une boucle dans l'ensemble de ses connexions.

dans ce cas l'état global du réseau dépend bien évidemment de ses états précédents, ce qui fait que l'équation du neurone, sera décrite par des suites récurrentes en temps discret, ou des équations différentielles en temps continu, et cela justifie bien l'emploi du terme "récurrent" .

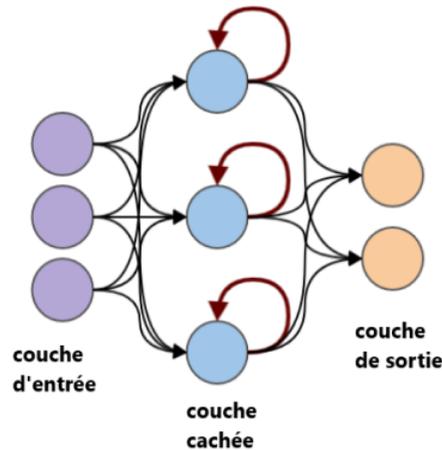


FIGURE II.23: La forme d'un réseau récurrent.

3.3 Quelques modèles de réseaux neuronaux artificiels :

Dans cette section on va essayer de présenter quelques modèles connexionnistes incontournables et qui décrivent les différentes topologies de la grande majorité des réseaux neuronaux classiques.

3.3.1 Le modèle de Hopfield :

Le physicien **John Hopfield** (1933, ...) a introduit en 1982 dans son article [32] un modèle connexionniste constitué de plusieurs neurones formels à sorties binaires, totalement connectés entre eux, et qui sont à la fois des unités d'entrée et de sortie, ce réseau récurrent, peut être décrit formellement comme un graphe complet non orienté valué $G = (N, [N]^2, f)$, doté d'une fonction de valuation symétrique f qui relie chaque paires (i, j) d'unités à la valeur de la pondération ω_{ij} correspondante.

les connexions dans un réseau Hopfield ont généralement les propriétés suivantes :

- i) $(\forall (i, j) \in \{1, \dots, n\}^2) : \omega_{ij} = \omega_{j,i} .$
- ii) $(\forall (i, j) \in \{1, \dots, n\}) : \omega_{ii} = 0 .$

La mise à jour de la sortie d'une unité (un nœud du graphe qui simule un neurone artificiel) du réseau de Hopfield est effectuée à l'aide de la règle d'actualisation suivante :

$$(\forall i \in \{1, \dots, n\}) : s_i^{(new)} = \begin{cases} +1 & \text{si : } \sum_{j=1}^n \omega_{ij} s_j^{(old)} \geq \theta_i, \\ 0 & \text{sinon.} \end{cases} .$$

cette procédure de mise à jour dans un réseau de Hopfield mis en valeur la notion du temps et montre qu'elle peut être effectuées de deux manières différentes :

Synchrone : toutes les unités sont mises à jour en même temps, ce qui nécessite une horloge centrale du système afin de maintenir la synchronisation, cette méthode est souvent considérée comme moins réaliste.

Asynchrone : chaque unité est mise à jour à la fois, et elle peut être choisie au hasard, ou à un ordre prédéfini, comme elle peut être imposé dès le début.

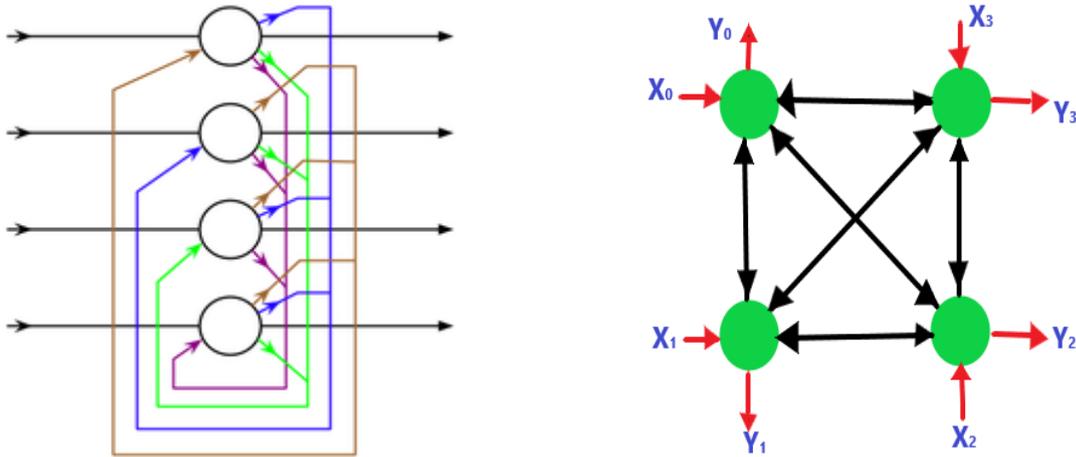


FIGURE II.24: Deux représentations différentes d'un réseau de Hopfield simple.

Remarque 3.3.1 :

Les modèles de Hopfield représentent une importance plus historique que pratique, et qui revient au fait de leur émergence pendant un tournant de l'histoire du connexionnisme, où ils ont été à la base de son redémarrage, en revanche ils ne sont plus utilisés maintenant dans leur version de base en raison de leur coût en terme de temps de calculs.

3.3.2 Le modèle de Kohonen :(la carte auto adaptative)

Le chercheur finlandais **Teuvo Kohonen** (1934, ...) a proposé en 1982 dans son article [38] de projeter l'espace des données \mathcal{D} , sur un autre espace de faible dimension, en général 1, 2 ou 3D au maximum, cet espace de projection sera appelé une , qu'on la note \mathcal{C} pour la suite .

cette carte est constitué d'un ensemble de neurones interconnectés selon une structure de graphe non orienté et qui peut prendre plusieurs formes, plan, cylindre, tore ...etc , mais en général on la réalise à partir d'un réseau neuronal à deux couches, une en entrée et une en sortie.

les neurones de la couche d'entrée sont entièrement connectés à ceux de la couche de sortie et qui sont placés dans un espace d'une ou de deux dimensions en général, chacun de ces neurones possède des voisins dans cet espace, et donc des connexions latérales récurrentes dans sa couche, ce qui va lui permettre de procéder en inhibant les neurones éloignés et laissant agir ceux qui sont voisins ^a.

en somme le réseau de Kohonen est un réseau de neurones dont la particularité est d'agir en tant que compresseur de données, qui se base sur la quantification vectorielle pour conserver uniquement les informations caractérisantes.

a. la notion de voisinage nécessite l'introduction d'une distance, et souvent on prend celle qui est définie comme étant la longueur du plus court chemin entre deux sommets du graphe \mathcal{C} .

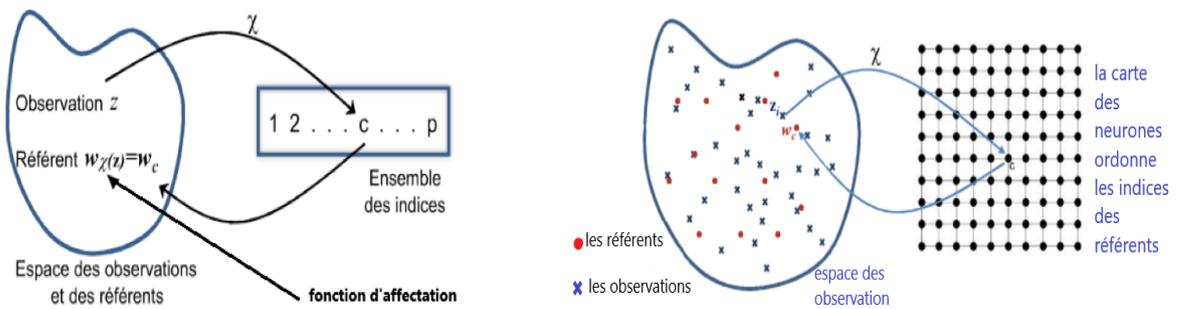


FIGURE II.25: Principe général de la modélisation par une carte de Kohonen.

3.3.3 Le Perceptron multicouche :

Le perceptron multicouche, et comme son nom l'indique est un réseau neuronal multicouche à propagation avant constitué de trois couches dont chacune a un rôle particulier.

La première couche s'appelle **la rétine** et son rôle est la réception de l'information de l'extérieur.

La deuxième s'appelle **une couche cachée** de pré-traitement : cette couche se situe entre la rétine et la couche de sortie ,elle est composée d'un nombre de neurones appelés neurones d'association et son rôle est de faire le pré traitement des informations provenantes de la rétine .

Et **la dernière couche** s'appelle une **couche de sortie** : tout simplement elle fait le traitement de ce qui provient de la couche cachée pour produire la sortie de notre réseau et il faut signaler qu'on plus de cet architecture particulière un perceptron est entièrement déterminé par ses pondérations .

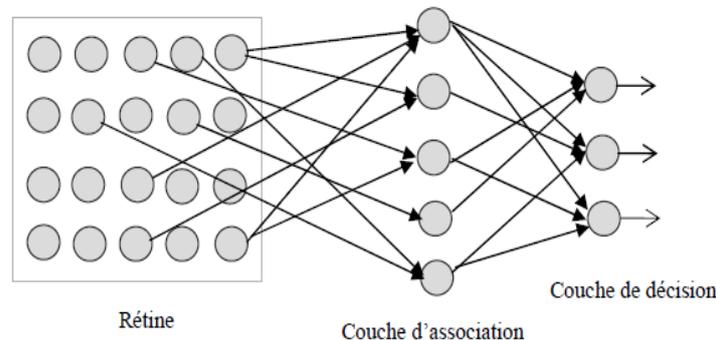


FIGURE II.26: Schéma d'un perceptron multicouches.

3.4 Apprentissage supervisé du perceptron multicouches :

Pour la suite, et afin d'éviter toute ambogüité, on signale qu'on se pose toujours dans le même cadre usuel d'apprentissage supervisé, et qui était déjà présenté pour le perceptron simple dans la page 35.

3.4.1 L'algorithme de Rétro-propagation du gradient :

Le problème de l'apprentissage pour les perceptrons multi-couches consiste à reconnaître l'influence de chaque poids sur l'erreur globale du réseau, et l'algorithme de rétro-propagation du gradient permet de faire ceci en trois étapes :

1. la propagation de l'information de l'entrée jusqu'à la sortie.
2. le calcul de l'erreur en sortie.

3. la rétro-propagation de l'erreur de la sortie jusqu'aux entrées.

nous décrirons dans la suite les différentes équations qui permettent de faire la mise à jour des pondérations de notre réseau.

★ **Pour la couche de sortie :**

prenons un neurone j de la couche de sortie, lorsque le $n^{\text{ème}}$ **exemple est lui présenté**, il produit une sortie $y_j(n)$ tandis que c'est la valeur $d_j(n)$ qui est attendue, ce qui fait qu'il commet une erreur :

$$e_j(n) = d_j(n) - y_j(n)$$

et de ce fait, l'erreur globale du réseau pour l'exemple n est :

$$\varepsilon(n) = \frac{1}{2} \sum_{j=1}^m e_j^2(n) = \|d(n) - y(n)\|_2 \quad (*)$$

où m est le nombre de neurones dans la couche de sortie.

comme d'habitude, le but de l'apprentissage est de minimiser l'erreur moyenne correspondante aux N exemples d'apprentissage, c'est à dire minimiser :

$$\varepsilon_{\text{moy}} = \frac{1}{N} \sum_{n=1}^N \varepsilon(n) \quad (**)$$

on note :

$$v_j(n) = \left(\sum_i \omega_{ij}(n) \times y_i(n) \right) - b_j(n) \quad (\blacklozenge)$$

l'entrée totale d'un neurone dont la fonction d'activation est f_j , sa sortie s'écrit alors sous la forme de

$$y_j(n) = f_j(v_j(n))$$

Le mécanisme de rétro-propagation est basé sur l'effectuation d'une correction $\Delta\omega_{ji}(n)$ et $\Delta b_j(n)$ aux poids et aux biais.

en ce qui concerne les , on utilisera la règle dite du delta (delta rule), c'est à dire que la modification sera proportionnelle au gradient suivant :

$$\frac{\partial \varepsilon(n)}{\partial \omega_{ji}(n)} = \frac{\partial \varepsilon(n)}{\partial e_j(n)} \times \frac{\partial e_j(n)}{\partial y_j(n)} \times \frac{\partial y_j(n)}{\partial v_j(n)} \times \frac{\partial v_j(n)}{\partial \omega_{ji}(n)}$$

qui détermine la direction dans laquelle on recherche les valeurs de $\omega_{ij}(n)$, nous appellerons taux d'apprentissage le facteur de proportionnalité et nous le noterons η par la suite.

d'après l'équation (**), on peut écrire :

$$\frac{\partial \varepsilon(n)}{\partial e_j(n)} = e_j(n)$$

de plus avec * :

$$\frac{\partial e_j(n)}{\partial y_j(n)} = f_j'(v_j(n))$$

finalement et à partir de (\blacklozenge) :

$$\frac{\partial v_j(n)}{\partial \omega_{ji}(n)} = y_i(n)$$

soit $\delta_j(n)$ le gradient local défini par :

$$\delta_j(n) = -\frac{\partial \varepsilon(n)}{\partial e_j(n)} \times \frac{\partial e_j(n)}{\partial y_j(n)} \times \frac{\partial y_j(n)}{\partial v_j(n)}$$

d'après ce qui précède ceci est égal à :

$$\delta_j(n) = e_j(n) \times f'_j(v_j(n)) \quad (\blacklozenge\blacklozenge)$$

La correction appliquée au poids $\omega_{ji}(n)$ sera donc la suivante (règle du delta) :

$$\Delta \omega_{ji}(n) = -\eta \times \frac{\partial \varepsilon(n)}{\partial \omega_{ji}(n)}$$

soit encore :

$$\Delta \omega_{ji}(n) = \eta \times \delta_j(n) y_i(n)$$

du même, la correction apportée au biais s'écrira sous la forme

$$\Delta b_j(n) = \eta \times \frac{\partial \varepsilon(n)}{\partial b_j(n)}$$

en suivant un raisonnement analogue à ce qui précède on obtient

$$\begin{aligned} \frac{\partial \varepsilon(n)}{\partial b_j(n)} &= \frac{\partial \varepsilon(n)}{\partial e_j(n)} \times \frac{\partial e_j(n)}{\partial y_j(n)} \times \frac{\partial y_j(n)}{\partial v_j(n)} \times \frac{\partial v_j(n)}{\partial b_j(n)} \\ &= -\delta_j(n) \times \frac{\partial v_j(n)}{\partial b_j(n)} \\ &= \delta_j(n) \end{aligned}$$

d'où :

$$\Delta b_j(n) = -\eta \cdot \delta_j(n)$$

ces deux équations de mise à jour ne sont valables que pour la couche de sortie, car on a utilisé l'erreur de sortie $e_j(n) = d_j(n) - y_j(n)$.

★ **Pour les couches cachées :**

En ce qui concerne les autres couches, on ne peut plus utiliser cette formule pour l'erreur, alors nous serons forcément obligés de déterminer les nouvelles équations de mise à jour.

▷ considérons un neurone j sur la dernière couche cachée, (celle qui se place juste avant celle de sortie.), nous pouvons et en utilisant une formule analogue à ($\blacklozenge\blacklozenge$) définir le gradient local par :

$$\begin{aligned} \delta_j(n) &= \frac{\partial \varepsilon(n)}{\partial y_j(n)} \times \frac{\partial y_j(n)}{\partial v_j(n)} \\ &= \frac{\partial \varepsilon(n)}{\partial y_j(n)} \times f'_j(v_j(n)) \end{aligned}$$

et d'après (**) on peut écrire :

$$\begin{aligned} \frac{\partial \varepsilon(n)}{\partial y_j(n)} &= \sum_k e_k(n) \times \frac{\partial e_k(n)}{\partial y_j(n)} \\ &= \sum_k e_k(n) \times \frac{\partial e_k(n)}{\partial v_k(n)} \times \frac{\partial v_k(n)}{\partial y_j(n)} \end{aligned}$$

or, le neurone k est sur la couche de sortie, d'où

$$\begin{aligned}\frac{\partial e_k(n)}{\partial v_k(n)} &= \frac{\partial(d_k(n) - y_k(n))}{\partial v_k(n)} \\ &= \frac{\partial(d_k(n) - f_k(v_k(n)))}{\partial v_k(n)} \\ &= f'_k(v_k(n))\end{aligned}$$

en outre, et d'après (♦) :

$$\frac{\partial e_k(n)}{\partial y_j(n)} = \omega_{kj}(n)$$

donc nous obtenons :

$$\begin{aligned}\frac{\partial \mathcal{E}(n)}{\partial y_j(n)} &= - \sum_k e_k(n) \times f'_k(v_k(n)) \times \omega_{kj}(n) \\ &= - \sum_k \delta_k(n) \times \omega_{kj}(n)\end{aligned}$$

le gradient local pour un neurone de la dernière couche cachée est alors donné par :

$$\delta_j(n) = \left(\sum_k \delta_k(n) \times \omega_{kj}(n) \right) \times f'_j(v_j(n))$$

ceci peut se généraliser et sans aucun souci aux autres couches internes, nous obtenons ainsi les règles de mise à jour des poids et des biais avec la méthode de rétro-propagation du gradient.

pour cette raison, les poids et les biais doivent être actualisés en utilisant :

$$\begin{aligned}\omega_{ji} &\leftarrow \omega_{ji} + \eta \times \delta_j(n) \times y_i(n) \\ b_j &\leftarrow b_j - \eta \times \delta_j(n)\end{aligned}$$

avec :

$$\delta_j(n) = e_j(n) \times f'_j(v_j(n))$$

et :

$$e_j(n) = \begin{cases} d_j(n) - y_j(n) & \text{si } j \text{ est un neurone de sortie.} \\ \sum_k \delta_k(n) \times \omega_{kj}(n) & \text{si } j \text{ est un neurone interne.} \end{cases}$$

Ainsi, l'erreur de sortie se trouve propagée vers l'entrée (c'est à dire dans le sens inverse de celui qui est utilisé pour propager un signal à travers le réseau) afin de corriger de couche en couche les valeurs des poids et des biais.

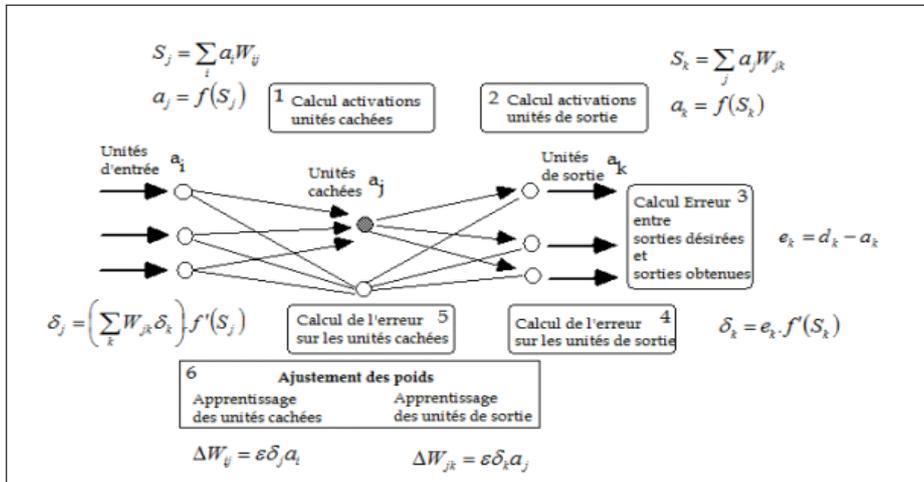


FIGURE II.27: Les étapes de l'algorithme de rétropropagation du gradient.

Algorithme 5 Algorithme de rétropropagation du gradient.

Entrée: une base de données : $(X^h, d^h)_{1 \leq h \leq \ell}$, une fonction d'activation non dérivable f .

Sortie: les poids optimaux : ω_{ij} .

1. **Initialisation :**

initialiser tous les poids et les biais du réseau de neurones.

2. **Présentation d'un exemple en entrée :**

présenter un exemple en entrée $X^n = (x_1(n), x_2(n), \dots, x_p(n))$ ainsi que la sortie désirée correspondante $d^n = (d_1(n), d_2(n), \dots, d_m(n))$.

3. **Calculer les sorties :**

calcules successivement les sorties des différentes couches pour l'entrée x .

4. La mise à jour des poids et des biais :

modifier récursivement les poids et les biais du réseau suivant les règles détaillées avant :

$$\omega_{ji} \leftarrow \omega_{ji} + \eta \times \delta_j(n) \times y_i(n)$$

$$b_j \leftarrow b_j - \eta \times \delta_j(n)$$

avec

$$\delta_j(n) = e_j(n) \times f'_j(v_j(n))$$

et

$$e_j(n) = \begin{cases} d_j(n) - y_j(n) & \text{si } j \text{ est un neurone de sortie} \\ \sum_k \delta_k(n) \times \omega_{kj}(n) & \text{si } j \text{ est un neurone interne} \end{cases}$$

5. **Critère d'arrêt :**

répéter les étapes 2, 3 et 4 jusqu'au nombre maximum d'itérations ou jusqu'à ce que l'erreur quadratique moyenne soit inférieure à un certain seuil.

retourner: ω_{ij}

Conclusion :

À partir du cerveau humain et du comportement des éléments neuronaux qui le constituent, on a pu construire des modèles artificiels très complexes appelés : "**réseaux de neurones artificiels**", dans ce chapitre on a bien traité ces réseaux, et on spécifié leurs définitions, leurs types, leurs concepts de base et aussi leurs mécanisme d'apprentissage, surtout pour le perceptron multicouches.

ces réseaux ont formé un outil très puissant, capable de traiter plusieurs familles de problèmes, que ça soit pour des problèmes de classification, de régression, de clustering, de prédiction ...etc, et cette capacité revient essentiellement à leurs propriétés d'approximations universelles, et qui feront effectivement l'objet du chapitre suivant.

Aptitude des Réseaux de Neurones Artificiels. (Théorie d'approximation) :

Résumé :

Dans ce chapitre, nous démontrons que les combinaisons linéaires *finies* de compositions d'une fonction *fixe* uni-variée (ou d'une seule variable réelle) et d'un ensemble de fonctionnelles affines peuvent approximer, *uniformément*, toute fonction continue de n variables réelles *avec un support inclus dans l'hypercube unité*; notons que, seules des conditions faibles sont imposées à notre fonction uni-variée.

les résultats de ce chapitre vont répondre à une question très importante ^b et qui porte sur la représentabilité de la classe des réseaux de neurones *à une seule couche cachée*.

En particulier, nous montrons que des régions de décision arbitraires peuvent être *bien approchées* ^c par des réseaux de neurones à propagation en avant (feedforward neural network) et avec une seule couche interne dont les nœuds ont tous une fonction d'activation *sigmoïdale continue* et *commune*.

on va discuter aussi des propriétés d'approximation d'autres types possibles de fonctions, et qui peuvent être mises en œuvre par des réseaux de neurones artificiels.

Mots clés : Réseaux de neurones, approximation, complétude, densité.

Introduction :

Un certain nombre de domaines d'application, sont basés sur l'utilisation des fonctions de vecteurs réelles à n dimensions, ($x \in \mathbb{R}^n$), et qui s'expriment par une combinaison linéaire *finie* de la forme :

$$f(x) = \sum_{j=1}^N \alpha_j \sigma(\omega_j^\top x + \theta_j) \tag{III.1}$$

où $\omega_j \in \mathbb{R}^n$, $\theta_j \in \mathbb{R}$, $\alpha_j \in \mathbb{R}$ et $\sigma \in \mathcal{F}(\mathbb{R}, \mathbb{R})$ sont *fixés* ^d. (ω^\top désigne ici la transposée de vecteur ω donc, et sans surprise $\omega^\top x$ ne va être que le produit scalaire usuel de x et ω).

La fonction $\sigma : \mathbb{R} \rightarrow \mathbb{R}$ dépend fortement du contexte de l'application, mais, notre principale

b. en fait elle était ouverte jusqu'au 1988, avant les travaux de **A.R.Gallant et H.White**.
 c. on verra après dans quel sens.
 d. c'est à dire qu'ils sont considérés comme des paramètres.

préoccupation concerne souvent ce qu'on appelle **les fonctions sigmoïdales**^a, car elles apparaissent naturellement dans la théorie des réseaux neuronaux en tant que fonctions d'activation d'un nœud neuronal (ou d'une unité, qui devient plutôt le terme le plus courant).

Le but principal de ce chapitre est de démontrer le fait que les sommes de la forme III.1, avec σ est sigmoïdale continue, (une condition qui peut être modifiée ou plus affaiblie par la suite.) sont **denses** dans l'espace des fonctions continues sur l'hypercube unitaire de dimension n .

Les utilisations possibles d'un réseau de neurones artificiels et surtout dans le traitement et le contrôle du signal, a généré une attention considérable depuis les années quatre-vingt-dix, et on peut constater en adoptant un point de vu un peu "mathématique", que sa description reste très facile, malgré sa puissance et ses différentes applications, or, un tel réseau est formé à partir de compositions et de superpositions d'une simple fonction d'activation^b, en conséquence, la sortie qu'il fournit n'est autre que la valeur de la fonction qui résulte de cette composition.

En particulier, la plus simple classe non triviale de réseaux, qu'on peut imaginer, sont ceux d'une seule couche interne et ils **implémentent** exactement **la classe de fonctions donnée par (1)**.

Donc, vous pouvez dire que ce chapitre vise de mettre en évidence le fait qu'un réseau de neurones mono-couche, **et d'un point de vu théorique**, possède une grande aptitude, et il est presque **capable de tout faire**.

1 Théorèmes d'approximation au cas d'une fonction d'activation discriminatoire :

*Pour une bonne lecture de cette section, il est indispensable de se référer à l'annexe A page : 124, qui va **actualiser et enrichir** le pré-acquis de notre lecteur, par quelques théorèmes d'analyse fonctionnelle et d'intégration, surtout ceux qui **vont apparaitre** pour légitimer et justifier quelques passages dans **le cadre de l'approximation des fonctions par des réseaux de neurones artificiels**. pour cette raison, nous vous invitons d'**aller voir cette annexe avant de traiter ce qui suit**.*

1.1 Vocabulaire et quelques notions de base :

Définition 1.1.1 : (mesure extérieure :)

On appelle une **mesure extérieure** ou **mesure de Crathéodory** sur un ensemble Ω toute application :

$$v : \mathcal{P}(\Omega) \longrightarrow \overline{\mathbb{R}}^+ = [0, +\infty]$$

$$A \longmapsto v(A)$$

qui vérifie les trois propriétés suivantes :

1. $v(\emptyset) = 0$.
2. $(\forall (A, B) \in \mathcal{P}(\Omega)) : A \subset B \implies v(A) \leq v(B)$. (monotonie.)
3. $(\forall (A_i)_{i \in \mathbb{N}} \in \mathcal{P}(\Omega)^{\mathbb{N}}) : v\left(\bigcup_{i \in \mathbb{N}} A_i\right) \leq \sum_{i \in \mathbb{N}} v(A_i)$. (sous additivité dénombrable.)

Exemple 1.1.1 :

Il est clair que toute mesure positive sur $(\Omega, \mathcal{P}(\Omega))$ est en particulier une mesure extérieure.

- a. tout simplement, ce sont des fonctions $\sigma : \mathbb{R} \longrightarrow \mathbb{R}$ qui vérifient les propriétés suivantes :

$$\lim_{x \rightarrow -\infty} \sigma(x) = 0 \text{ et } \lim_{x \rightarrow +\infty} \sigma(x) = 1.$$

- b. **Attention!**, on a supposé ici que les unités constituant notre réseaux sont **identiques**.

Définition 1.1.2 : (mesurabilité au sens de Carathéodory :)

Soit Ω un ensemble quelconque (même vide), et ν une mesure extérieure sur cet ensemble, une sous partie $E \subset \Omega$ est dite ν -**mesurable** ou mesurable au sens de Carathéodory par ν si et seulement si :

$$(\forall A \in \mathcal{P}(\Omega)) : \nu(A) = \nu(E^c \cap A) + \nu(E \cap A).$$

autrement dit, un sous ensemble est ν mesurable si et seulement si elle divise toute partie A de Ω d'une manière additive par rapport à la mesure extérieure ν .

Remarque 1.1.1 :

Notons qu'en vertu de la sous σ -additivité, l'inégalité $\nu(A) \leq \nu(E^c \cap A) + \nu(E \cap A)$ aura toujours lieu; donc pour vérifier qu'une partie E est ν -mesurable, il suffit de vérifier l'inégalité inverse.

Définition 1.1.3 : (mesure de Borel régulière :)

Une mesure extérieure ν sur \mathbb{R}^n est dite : mesure régulière de Borel si et seulement si :
 $(\forall B \in \mathcal{B}_{\mathbb{R}^n}) : B$ est ν -mesurable au sens de Carathéodorie, et $(\forall A \subset \mathbb{R}^n) (\exists B \in \mathcal{B}_{\mathbb{R}^n}) : A \subset B$ et $\nu(A) = \nu(B)$.

Définition 1.1.4 : (mesure signée :)

Dans un espace mesurable, (Ω, \mathcal{F}) une mesure signée est tout simplement une fonction $\mu : \mathcal{F} \rightarrow \overline{\mathbb{R}}$ (à valeurs dans $\overline{\mathbb{R}}$), qui vérifie :

1. $(\exists A \in \mathcal{F}) : \mu(A) = +\infty$ (resp $-\infty$) $\Rightarrow (\forall A \in \mathcal{F}) : \mu(A) \neq -\infty$ (resp $+\infty$).
2. $\mu(\emptyset) = 0$.
3. pour toute suite d'éléments deux à deux disjointes $(A_i)_{i \in \mathbb{N}} \in \mathcal{F}^{\mathbb{N}}$ on a : $\mu\left(\bigcup_{n=1}^{\infty} A_n\right) = \sum_{n=1}^{\infty} \mu(A_n)$.

c'est à dire qu'elle garde exactement la même définition qu'une mesure ordinaire sauf qu'on change l'espace d'arrivée, et qu'on ajoute la condition "1" pour qu'elle soit bien définie.

Exemples 1.1.2 :

- Considérons **une mesure positive** μ sur un espace mesuré (Ω, \mathcal{F}) , et prenons une fonction $(\Omega, \mathcal{B}_{\mathbb{R}})$ mesurable $f : \Omega \rightarrow \mathbb{R}$ telle que $f \in L^1(\Omega, \mathcal{F})$, c'est à dire que : $\int_{\Omega} |f(x)| d\mu(x) < +\infty$ alors :
l'application $\nu : \mathcal{F} \rightarrow \mathbb{R}$ donnée par : $\nu(A) = \int_A f(x) d\mu(x)$ pour tout $A \in \mathcal{F}$, est bien définie, et en plus, **c'est une mesure signée** à valeurs dans \mathbb{R} , et pour lui permettre de prendre $+\infty$ (ou $-\infty$) comme valeur, il suffit d'affaiblir l'hypothèse selon laquelle f est absolument intégrable et ceci en la remplaçant par :

$$\int_{\Omega} f^-(x) d\mu(x) < +\infty, \left(\text{ou} \int_{\Omega} f^+(x) d\mu(x) < +\infty \right)$$

avec : $f^+ = \max(f, 0)$ et $f^- = \max(-f, 0)$.

- Considérons cette fois ci **deux mesures positives** μ et ν sur un espace mesuré (Ω, \mathcal{F}) , et **supposons que l'une d'elles est finie**, alors, et sans surprise : l'application $\nu = \mu - \nu$ sera bien **une mesure signée** à valeurs dans $\overline{\mathbb{R}}$, mais en revanche, et ce qui est un peu surprenant est le fait que **n'importe quelle mesure signée est forcement de cette forme**, c'est un résultat très important, mais un peu compliqué à démontrer surtout qu'il se base essentiellement sur **le**

théorème de décomposition de Hahn^a, donc ne vous laissez pas tromper par la sympathie de l'énoncé.

Remarques 1.1.2 :

1. Malgré cette définition qui est simple d'apparence, **l'étude des mesures signées est très compliquée**, et même elle présente un champ de recherche riche dont pas mal de questions sont encore ouvertes^b, surtout que plusieurs propositions, comme le lemme des classes monotones ou le théorème de Dynkin-Sierpinski^c ou même celui de convergence monotones de Beppo-Levi ne restent plus valables, ce qui va poser, et même dans notre cadre de travail, beaucoup de problèmes^d.
2. Pour n'importe qu'elle notion analogue, comme mesure extérieure, mesure régulière, mesure de Borel...etc, l'attribution du terme "**signée**" signifie la considération de notre mesure comme une différence de deux mesures positives de même type, avec, bien sûr, une réalisation d'un changement de l'espace d'arrivée de $\overline{\mathbb{R}}^+$ à $\overline{\mathbb{R}}$ et l'ajout de la condition "1" pour éviter la forme indéterminée " $(+\infty) + (-\infty)$ ".

Définition 1.1.5 : (la transformée de Fourier d'une mesure finie :)

Soit μ une mesure signée **finie** sur $(\mathbb{R}^d, \mathcal{B}_{\mathbb{R}^d})$, on appelle **transformée de Fourier** de μ la fonction $\widehat{\mu} : \mathbb{R}^d \rightarrow \mathbb{C}$, définie par :

$$\left(\forall t \in \mathbb{R}^d \right) : \widehat{\mu}(t) := \int_{\mathbb{R}^d} \exp(i\langle t, x \rangle) d\mu(x).$$

Remarques 1.1.3 :

- On notera que c'est l'hypothèse de la finitude de la mesure ($\mu(\mathbb{R}^d) < +\infty$) qui nous a permis de donner un sens à la transformée de Fourier, car elle rend automatiquement la fonction $\Psi_t : x \mapsto \exp(i\langle t, x \rangle)$ intégrable, sur \mathbb{R}^d pour tout $t \in \mathbb{R}^d$, or :

$$\int_{\mathbb{R}^d} \exp(i\langle t, x \rangle) d\mu(x) \leq \int_{\mathbb{R}^d} |\exp(i\langle t, x \rangle)| d\mu(x) = \int_{\mathbb{R}^d} 1 d\mu(x) = \mu(\mathbb{R}^d) < +\infty.$$

mais faites attention, s'il n'est pas le cas, on ne parle plus d'une telle notion, par exemple : Ψ n'est pas λ_d -intégrable sur \mathbb{R}^d ce qui fait qu'on ne peut pas définir $\widehat{\lambda}_d$.

- Lorsque μ est la loi P_X d'un vecteur aléatoire X , $\widehat{\mu}$ n'est autre que **la fonction caractéristique** de X c'est à dire :

$$\mu = P_X \implies \widehat{\mu}(t) = \varphi_X(t) = \mathbb{E} \left(e^{i\langle t, X \rangle} \right).$$

a. ce théorème affirme que : pour tout espace mesurable (Ω, \mathcal{F}) et **toute mesure signée** μ définie sur \mathcal{F} , il existe deux ensembles mesurables $P, N \in \mathcal{F}$ tel que :

• $P \cup N = \Omega, P \cap N = \emptyset$ • $(\forall E \in \mathcal{F}) : E \subseteq P \implies \mu(E) \geq 0$ • $(\forall E \in \mathcal{F}) : E \subseteq N \implies \mu(E) \leq 0$.

de plus, **cette décomposition est unique presque sûrement**.

b. Par exemple quelles conditions on doit ajouter dans le cas des mesures signées, pour garder la validité du théorème de Levy, qui relie la convergence en loi d'une suite de variables aléatoires réelles, avec la convergence ponctuelle de leurs fonctions caractéristiques?.

c. Il énonce que deux mesures σ -finies confondues sur un π -système, seront aussi confondues sur la tribu qu'il engendre.

d. Par exemple au niveau de la démonstration que toute fonction sigmoïdal mesurable bornée est en particulier discriminatoire (voir la preuve du lemme 2.2.1 page 63).

- si la mesure μ est à densité par rapport à celle de Lebesgue, c'est à dire :

$$\left(\exists f \in L^1(\mathbb{R}^d) \right) \text{ tel que : } (\forall A \in \mathfrak{B}_{\mathbb{R}}) : \mu(A) = \int_A f d\lambda$$

alors, la transformée de Fourier de la mesure μ sera exactement la transformée de Fourier usuelle de la fonction f .

Proposition 1.1.1 :

Soit μ une mesure signée finie sur $(\mathbb{R}^d, \mathfrak{B}_{\mathbb{R}^d})$, on note $\hat{\mu}$ sa transformée de Fourier associée, alors on a :

1. $\hat{\mu}$ est toujours continue et bornée sur \mathbb{R}^d .
2. la **transformation de Fourier**^a \mathcal{F} définie sur l'ensemble $M(\mathbb{R}^d)$ des mesures signées finies sur \mathbb{R}^d par :

$$\begin{aligned} \mathcal{F} : M(\mathbb{R}^d) &\longrightarrow C^0(\mathbb{R}, \mathbb{C}) \\ \mu &\longmapsto \mathcal{F}(\mu) = \hat{\mu} \end{aligned}$$

est une application **bien définie, linéaire**, et en plus **injective**, et c'est ce qui est très important (voir :[37]).

Preuve :

1. **★ Pour la continuité** : on a d'une part : $(\forall t \in \mathbb{R}^d) : \hat{\mu}(t) := \int_{\mathbb{R}^d} \exp(i\langle t, x \rangle) d\mu(x)$ et d'une autre :
 - pour tout $(x \in \mathbb{R}^d), t \mapsto e^{i\langle t, x \rangle}$ est continue sur \mathbb{R}^d .
 - pour tout $(t \in \mathbb{R}^d), x \mapsto e^{i\langle t, x \rangle}$ est continue sur \mathbb{R}^d .
 - pour tout $: t, x \in \mathbb{R}^d \mid e^{i\langle t, x \rangle} \mid \leq 1$, et car μ **est fini**, on obtiendra immédiatement l'intégrabilité de la fonction constante $\mathbb{1}$, or : $\mu(\mathbb{R}^d) = \int_{\mathbb{R}^d} \mathbb{1} d\mu(x) < +\infty$.

donc , et selon **le théorème de continuité des intégrales paramétriques**^b on peut déduire que : $\hat{\mu}(t) := \int_{\mathbb{R}^d} e^{i\langle t, x \rangle} d\mu(x)$ est continue.

★ pour la bornitude ; il suffit de remarquer que :

$$(\forall t \in \mathbb{R}^d) : \int_{\mathbb{R}^d} e^{i\langle t, x \rangle} d\mu(x) \leq \int_{\mathbb{R}^d} \mid e^{i\langle t, x \rangle} \mid d\mu(x) = \int_{\mathbb{R}^d} \mathbb{1} d\mu(x) = \mu(\mathbb{R}^d) \implies \hat{\mu} \leq \mu(\mathbb{R}^d).$$

2. **★ montrons que l'application \mathcal{F} est bien définie** : d'abord, l'existence de l'image est un fait qui provient directement du premier point (la transformée de Fourier d'une mesure est toujours continue), et l'unicité est triviale, et on peut dire même qu'elle est vérifiée par définition.

★ maintenant pour la linéarité :

soit μ et ν deux mesures signées **bornées** sur \mathbb{R}^d , on a :

$$(\forall A \in \mathfrak{B}_{\mathbb{R}^d}) : \int_{\mathbb{R}^d} \mathbb{1}_A d(\mu + \nu) := [\mu + \nu](A) := \mu(A) + \nu(A) = \int_{\mathbb{R}^d} \mathbb{1}_A d\mu + \int_{\mathbb{R}^d} \mathbb{1}_A d\nu.$$

a. Attention, c'est la transformation de Fourier et pas la transformée, soyez prudents et faites la différence entre ces deux notions.

b. **un théorème classique de l'analyse** qui s'énonce ainsi : soit $(\Omega, \mathcal{F}, \mu)$ un espace mesuré et $f : (t, x) \mapsto f(t, x)$ une fonction de $I \times \Omega$ dans \mathbb{C} (où I est un intervalle de \mathbb{R}) on suppose que :

♣ $(\forall t \in I) : x \mapsto f(t, x)$ est mesurable. ♣ $(\forall x \in \Omega) : t \mapsto f(t, x)$ est continue sur I .

♣ $(\exists \phi \in L^1((\Omega, \mathcal{F}, \mu); \mathbb{R}^+))$ tel que : $(\forall t \in I) \mid f(t, x) \mid \leq \phi(x)$ μ p.p sur E .

alors la fonction $F : t \mapsto F(t) = \int_E f(t, x) d\mu$ est bien définie et continue.

donc pour toute fonction numérique étagée mesurable positive (appelée également fonction simple) φ de la forme $\varphi = \sum_{i=1}^n a_i \mathbb{1}_{A_i}$ où $(A_i)_{i \in \{1 \dots n\}} \subset \mathcal{B}_{\mathbb{R}^d}$ sont deux-à-deux disjoints et

$(a_i)_{i \in \mathbb{N}} \subset \mathbb{R}^+$, on a :

$$\begin{aligned} \int_A \varphi d(\mu + \nu) &:= \int_{\mathbb{R}^d} \sum_{i=1}^n a_i \mathbb{1}_{A_i} d(\mu + \nu) = \sum_{i=1}^n a_i \int_{\mathbb{R}^d} \mathbb{1}_{A_i} d(\mu + \nu) = \sum_{i=1}^n a_i \int_{\mathbb{R}^d} \mathbb{1} d\mu + \sum_{i=1}^n a_i \int_{\mathbb{R}^d} \mathbb{1} d\nu \\ &= \int_{\mathbb{R}^d} \varphi d\mu + \int_{\mathbb{R}^d} \varphi d\nu. * \end{aligned}$$

Soit maintenant f une fonction mesurable positive **bornée**, d'une part et selon le théorème d'approximation, f est la limite simple d'une suite **croissante** de fonctions étagées mesurables positives $(\varphi_n)_{n \in \mathbb{N}}$, et d'une autre, et comme on a déjà signalé au deuxième point de la remarque III.1.1.2, que n'importe quelle mesure signée peut s'écrire comme une différence de deux mesures positives, donc et en appliquant ce fait aux deux mesures μ et ν : on pourra bien écrire $\mu = \mu^+ - \mu^-$ et $\nu = \nu^+ - \nu^-$ avec μ^+, μ^- et ν^+, ν^- sont tous des mesures positives ^a.

en outre, on a selon le résultat * :

$$(\forall n \in \mathbb{N}) : \int_{\mathbb{R}^d} \varphi_n d(\mu + \nu) = \int_{\mathbb{R}^d} \varphi_n d\mu + \int_{\mathbb{R}^d} \varphi_n d\nu \quad \blacklozenge$$

et encore par le même résultat il est évident que :

$$(\forall n \in \mathbb{N}) : \int_{\mathbb{R}^d} \varphi_n d\mu = \int_{\mathbb{R}^d} \varphi_n d(\mu^+ + \mu^-) = \int_{\mathbb{R}^d} \varphi_n d\mu^+ + \int_{\mathbb{R}^d} \varphi_n d\mu^- = \int_{\mathbb{R}^d} \varphi_n d\mu^+ - \int_{\mathbb{R}^d} \varphi_n d(-\mu^-) \quad \text{par définition.}$$

or, et selon le théorème de la convergence monotone on a :

$$\begin{aligned} \lim_{n \rightarrow +\infty} \int_{\mathbb{R}^d} \varphi_n d\mu^+ &= \int_{\mathbb{R}^d} \lim_{n \rightarrow +\infty} \varphi_n d\mu^+ = \int_{\mathbb{R}^d} f d\mu^+ < +\infty \\ \lim_{n \rightarrow +\infty} \int_{\mathbb{R}^d} \varphi_n d(-\mu^-) &= \int_{\mathbb{R}^d} \lim_{n \rightarrow +\infty} \varphi_n d(-\mu^-) = \int_{\mathbb{R}^d} f d(-\mu^-) < +\infty \end{aligned}$$

donc et en faisant tendre $n \rightarrow +\infty$ on peut déduire que :

$$\begin{aligned} \lim_{n \rightarrow +\infty} \int_{\mathbb{R}^d} \varphi_n d\mu &= \lim_{n \rightarrow +\infty} \int_{\mathbb{R}^d} \varphi_n d\mu^+ - \lim_{n \rightarrow +\infty} \int_{\mathbb{R}^d} \varphi_n d(-\mu^-) = \int_{\mathbb{R}^d} \lim_{n \rightarrow +\infty} \varphi_n d\mu^+ - \int_{\mathbb{R}^d} \lim_{n \rightarrow +\infty} \varphi_n d(-\mu^-) \\ &= \int_{\mathbb{R}^d} f d\mu^+ - \int_{\mathbb{R}^d} f d(-\mu^-) = \int_{\mathbb{R}^d} f d\mu \lim_{n \rightarrow +\infty} \int_{\mathbb{R}^d} \varphi_n d\mu \Rightarrow \int_{\mathbb{R}^d} f d\mu = \lim_{n \rightarrow +\infty} \int_{\mathbb{R}^d} \varphi_n d\mu \end{aligned}$$

ce qui veut dire que le théorème de convergence monotone est valable pour les mesures signées, **mais à condition que la fonction limite soit bornée.**

en faisant de même pour les mesures ν et $\nu + \mu$ on obtiendra :

$$\lim_{n \rightarrow +\infty} \int_{\mathbb{R}^d} \varphi_n d\nu = \int_{\mathbb{R}^d} \lim_{n \rightarrow +\infty} \varphi_n d\nu = \int_{\mathbb{R}^d} f d\nu$$

et :

$$\lim_{n \rightarrow +\infty} \int_{\mathbb{R}^d} \varphi_n d(\mu + \nu) = \int_{\mathbb{R}^d} \lim_{n \rightarrow +\infty} \varphi_n d(\mu + \nu) = \int_{\mathbb{R}^d} f d(\mu + \nu)$$

donc on obtiendra à travers les résultats obtenus précédemment, et en faisant tendre $n \rightarrow +\infty$ dans \blacklozenge que :

$$\int_{\mathbb{R}^d} f d(\mu + \nu) = \int_{\mathbb{R}^d} f d\mu + \int_{\mathbb{R}^d} f d\nu. **$$

prenons à ce stade f une fonction mesurable **bornée** quelconque (pas forcément positive), on a : $f = f^+ - f^-$ avec $f^+ = \max(f, 0)$ et $f^- = \max(-f, 0)$.

a. **Attention** : cette écriture n'est pas unique, car on peut tout simplement ajouter et retrancher une certaine mesure positive **finie**.

alors il suffit d'utiliser le résultat** pour les fonctions f^+ et f^- , pour déduire que :

pour toute fonction mesurable **bornée** f on a :
$$\int_{\mathbb{R}^d} f d(\mu + \nu) = \int_{\mathbb{R}^d} f d\mu + \int_{\mathbb{R}^d} f d\nu \quad \blacklozenge$$

enfin, considérons une fonction complexe quelconque f , avec **des parties réelles et imaginaires bornées**, on a :

$$\begin{aligned} \int_{\mathbb{R}^d} f d(\mu + \nu) &= \int_{\mathbb{R}^d} \operatorname{Re}(f) + \operatorname{Im}(f)i d(\mu + \nu) = \int_{\mathbb{R}^d} \operatorname{Re}(f) d(\mu + \nu) + i \int_{\mathbb{R}^d} \operatorname{Im}(f) d(\mu + \nu) \\ &= \int_{\mathbb{R}^d} \operatorname{Re}(f) d\mu + \int_{\mathbb{R}^d} \operatorname{Re}(f) d\nu + i \int_{\mathbb{R}^d} \operatorname{Im}(f) d\mu + i \int_{\mathbb{R}^d} \operatorname{Im}(f) d\nu \quad (\text{selon } \blacklozenge) \\ &= \int_{\mathbb{R}^d} f d\mu + \int_{\mathbb{R}^d} f d\nu \end{aligned}$$

donc et en appliquant ce résultat à la fonction $x \mapsto e^{i\langle t, x \rangle} = \underbrace{\cos(\langle t, x \rangle)}_{\text{fonction bornée}} + i \underbrace{\sin(\langle t, x \rangle)}_{\text{fonction bornée}}$ pour n'importe quel $t \in \mathbb{R}$ on déduit que :

$$(\forall \nu, \mu \in M(\mathbb{R}^d)) : \widehat{\mu + \nu} = \widehat{\mu} + \widehat{\nu}$$

en plus, et sans aucune difficulté :

$$\begin{aligned} (\forall \lambda \in \mathbb{R}) (\forall \mu \in M(\mathbb{R}^d)) (\forall t \in \mathbb{R}^d) : \widehat{\lambda\mu}(t) &= \int_{\mathbb{R}^d} e^{i\langle t, x \rangle} d(\lambda\mu(x)) = \int_{\mathbb{R}^d} \lambda e^{i\langle t, x \rangle} d\mu(x) \\ &= \lambda \int_{\mathbb{R}^d} e^{i\langle t, x \rangle} d\mu(x) = \lambda \widehat{\mu}(t). \end{aligned}$$

donc la transformation de Fourier est bien linéaire.

★ **finalement l'injectivité de la transformation de Fourier :**

la preuve est vraiment difficile, et elle se base sur plusieurs lemmes qui nécessitent un pré-requis riche en **analyse harmonique**^a, mais elle est classique et très renommée, on ne rentrera pas complètement dans ses détails, néanmoins on va vous présenter quand même les grandes lignes :

— pour $\sigma > 0$ on définit la fonction suivante sur \mathbb{R}^d , en adoptant la notation $x = (x_1, x_2, \dots, x_d)$:

$$g_{d,\sigma}(x) = \prod_{j=1}^d g_{\sigma}(x_j) = \prod_{j=1}^d \frac{1}{\sqrt{2\pi}} e^{-x_j^2/2\sigma^2} = \frac{1}{\sqrt{2\pi}} e^{-|x|^2/2\sigma^2}.$$

— on peut facilement, et par un simple calcul, établir que la transformée de Fourier de la fonction g est :

$$\widehat{g}_{d,\sigma}(x) = \prod_{j=1}^d \widehat{g}_{\sigma}(x_j) = e^{-2\pi^2\sigma^2|x|^2}.$$

— d'une autre part on peut vérifier en utilisant le théorème de Fubini que :

$$(g_{d,\sigma} \star \mu)(x) := \int_{\mathbb{R}^d} g_{d,\sigma}(x-y) \mu dy = \int_{\mathbb{R}^d} \widehat{\mu}(t) e^{2i\pi\langle t, x \rangle - 2\pi^2\sigma^2|t|^2} dt.$$

et en plus, et c'est ce qui est très intéressant mais un peu dur à établir : toute fonction bornée h sur \mathbb{R}^d , l'intégrale : $\int_{\mathbb{R}^d} h d\mu$ est la limite de : $\int_{\mathbb{R}^d} (g_{d,\sigma} \star \mu) h(x) dx$ lorsque : $\sigma \rightarrow 0$.

— enfin il suffit de relier les trois points précédents pour arriver à l'injectivité de la trans-

a. C'est une branche des mathématiques qui étudie la représentation des fonctions ou des signaux comme superposition d'ondes de base qu'on les appellent les harmoniques, d'où le nom de la discipline.

formation de Fourier, or :

$$\begin{aligned}
 (\forall \nu, \mu \in M(\mathbb{R}^d)) : \hat{\mu} = \hat{\nu} &\implies \int_{\mathbb{R}^d} \hat{\mu}(t) e^{2i\pi\langle t, x \rangle - 2\pi^2 \sigma^2 |t|^2} dt = \int_{\mathbb{R}^d} \hat{\nu}(t) e^{2i\pi\langle t, x \rangle - 2\pi^2 \sigma^2 |t|^2} dt \\
 &\implies (g_{d, \sigma} \star \mu)(x) = (g_{d, \sigma} \star \nu)(x). \\
 &\implies \int_{\mathbb{R}^d} (g_{d, \sigma} \star \mu) h(x) dx = \int_{\mathbb{R}^d} (g_{d, \sigma} \star \nu) h(x) dx \quad (\forall h \in L^\infty(\mathbb{R}^d)) \\
 \text{en faisant tendre } \sigma \rightarrow 0 &\implies \int_{\mathbb{R}^d} h d\mu = \int_{\mathbb{R}^d} h d\nu \quad (\forall h \in L^\infty(\mathbb{R}^d))
 \end{aligned}$$

donc il suffit de prendre : $h = \mathbb{1}_{[a, +\infty[}$ avec $a \in \mathbb{R}$ pour voir que μ et ν se coïncident sur le π -système générant la tribu borélienne $\mathfrak{B}_{\mathbb{R}^d}$ et par suite déduire selon le lemme des classes monotone (appelé également théorème de Dynkin-Sierpinski) l'égalité de ces deux mesures. ■

Exemple 1.1.3 :

Un exemple pratique de la transformée de Fourier d'une mesure est la fameuse **série de Fourier**^a et qui est une série de terme général $a_n e^{2in\pi x}$ indexée par : $n \in \mathbb{Z}$ lorsque les $a_n \in \mathbb{R}$ et $\sum_{n \in \mathbb{Z}} a_n < +\infty$. la somme d'une telle série n'est autre que la transformée de Fourier de la mesure définie sur \mathbb{R} par : $\mu = \sum_{n \in \mathbb{Z}} a_n \delta_{-n}$ où δ_{-n} désigne **la masse de Dirac** centrée, en $-n$ ^b.

Notations :

Pour toute la suite, sauf mention contraire, on va adopter les notations suivantes :

- I_n désigne l'hypercube unité de dimension $n : [0, 1]^n$.
- $C(I_n)$ est l'ensemble des fonctions **numériques** continues sur I_n .
- $(\forall f \in C(I_n))$ on note : $\|f\| = \sup_{x \in I_n} |f(x)|$.
- $M(I_n)$ est l'ensemble des mesures régulières de Borel sur I_n et qui sont **signées et finies**.

Définition 1.1.6 : (fonction discriminatoire)

On dit qu'une fonction : $\sigma : I_n \rightarrow \mathbb{R}$ est **discriminatoire**, dans $M(I_n)$, si et seulement si :

$$(\forall \mu \in M(I_n)) : \left[(\forall \omega \in \mathbb{R}^n) (\forall \theta \in \mathbb{R}) : \int_{I_n} \sigma(\omega^\top x + \theta) d\mu = 0 \right] \implies \mu = 0.$$

Remarque 1.1.4 :

- Afin d'éviter toute confusion, nous proposons le diagramme suivant, qui représente les fonctions qu'on utilise, en montrant les différents paramètres et variables qui apparaissent dans leurs expressions.

a. **les séries de Fourier** qui ont été introduites par **Joseph Fourier** au 19^{ème} siècle, présentent un outil puissant pour **l'étude des fonctions périodiques** et jusqu'à maintenant elles font encore un objet de recherches actives. ces séries ont donné naissance à plusieurs branches nouvelles : analyse harmonique, théorie du signal, ..., etc, **le but de leur introduction est de pouvoir écrire une fonction T-périodique comme la somme de fonctions sinusoïdales** : $f(x) = \sum_{n=-\infty}^{+\infty} \lambda_n(f) e^{i2\pi \frac{n}{T} x}$ avec $\lambda_n(f)$, appelés coefficients de Fourier de f , définis par :

$$\lambda_n(f) = \frac{1}{T} \int_{-T/2}^{T/2} f(t) e^{-i2\pi \frac{n}{T} t} dt.$$

b. la masse de Dirac centrée en un point $x_0 \in \mathbb{R}^d$ (fixé à l'avance) est une mesure positive définie sur la tribu Borélienne de \mathbb{R}^d comme suit : $(\forall A \in \mathfrak{B}_{\mathbb{R}^d}) : \delta_{x_0}(A) = \mathbb{1}_A(x_0)$

$$\begin{array}{ccccc}
 I_n & \xrightarrow{h} & \mathbb{R} & \xrightarrow{\sigma} & \mathbb{R} \\
 \mathbf{x} & \mapsto & \omega^\top \mathbf{x} + \theta & \mapsto & \sigma(\omega^\top \mathbf{x} + \theta)
 \end{array}$$

— pour la définition de la fonction discriminatoire^a et qui était introduite pour la première fois dans l'article [13], le choix de l'espace $M(I_n)$ n'est pas obligatoire, et on peut aisément le changer par un autre espace de mesures.

Exemple 1.1.4 :

1. la fonction **identité** $\sigma_o : x \mapsto x$ est discriminatoire, dans l'ensemble $M^+(I_n)$ des mesures régulières positives de Borel sur I_n car :

$$\begin{aligned}
 (\forall \mu \in M^+(I_n)) : \left[(\forall \omega \in \mathbb{R}^n) (\forall \theta \in \mathbb{R}) : \int_{I_n} \sigma_o(\omega^\top x + \theta) \, d\mu = 0 \right] &\implies \int_{I_n} \sigma_o(\mathbf{0}_{\mathbb{R}^n}^\top x + 1) \, d\mu = 0 \\
 &\implies \int_{I_n} 1 \, d\mu = \mu(I_n) = 0 \\
 \left(\text{et puisque : } (\forall A \in \mathcal{P}(I_n)) : 0 \leq \mu(A) \leq \mu(I_n) \right) &\implies \mu = 0.
 \end{aligned}$$

faites attention, cette fonction n'est pas du tout discriminatoire dans $M(I_n)$.

Définition 1.1.7 : (fonction sigmoïdale :)

On dit qu'une fonction $\sigma : \mathbb{R} \rightarrow \mathbb{R}$ est sigmoïdale si et seulement si : $\lim_{x \rightarrow -\infty} \sigma(x) = 0$ et $\lim_{x \rightarrow +\infty} \sigma(x) = 1$.

Exemples 1.1.5 :

1. la fonction $\varphi(x) : \mathbb{R} \rightarrow \mathbb{R}$ définie comme suit : $\varphi(x) = \frac{1}{2} \times \left(1 + \tanh \left(x + \sin \left(x^2 \right) \right) \right)$. est une fonction sigmoïdale.

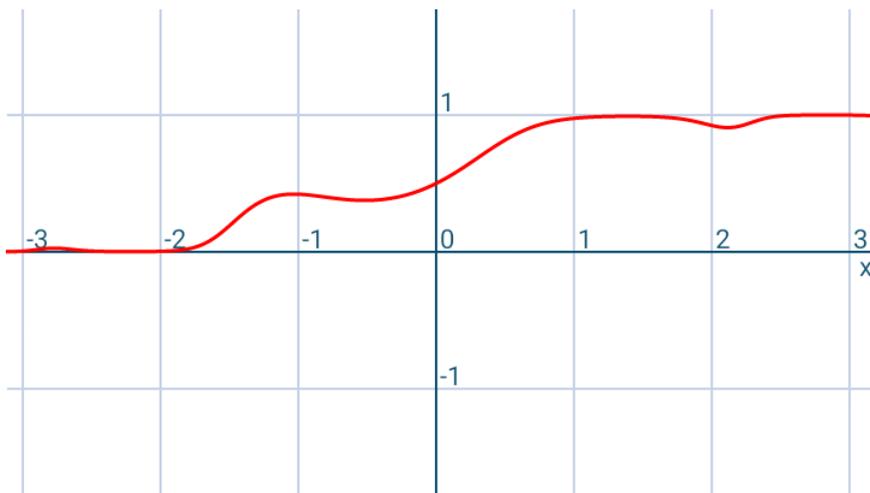


FIGURE III.1: La courbe de la fonction φ .

2. il est facile de remarquer que toute fonction de répartition de probabilité est en particulier sigmoïdale.

a. on peut dire également fonction discriminante, discriminative ou discriminatrice.

1.2 Résultats principaux :

Théorème 1.2.1 :

Soit σ une fonction **continue discriminatoire** dans $M(I_n)$, alors, l'ensemble :

$$\mathcal{N}_\sigma = \left\{ x \mapsto G(x) = \sum_{j=1}^N \alpha_j \sigma(\omega_j^\top x + \theta_j) \mid N \in \mathbb{N}, (\alpha_1, \dots, \alpha_N) \in \mathbb{R}^N, (\theta_1, \dots, \theta_N) \in \mathbb{R}^N, (\omega_1, \dots, \omega_N) \in (\mathbb{R}^n)^N \right\}$$

est dense dans $C(I_n)$ **pour la norme de la convergence uniforme**^a.

et comme il était énoncé dans [13], ceci veut dire que : pour toute fonction $f \in C(I_n)$ on a :

$$(\forall \varepsilon > 0) (\exists N_\varepsilon \in \mathbb{N}) \left(\exists (\alpha^\varepsilon, \theta^\varepsilon) \in (\mathbb{R}^{N_\varepsilon})^2 \right) \left(\exists (\omega_1^\varepsilon, \dots, \omega_{N_\varepsilon}^\varepsilon) \in (\mathbb{R}^n)^{N_\varepsilon} \right) \text{ tq : } \left\| \sum_{j=1}^{N_\varepsilon} \alpha_j^\varepsilon \sigma(\langle \omega_j^\varepsilon, \cdot \rangle + \theta_j^\varepsilon) - f \right\| < \varepsilon$$

Preuve :

D'abord, et avant de commencer, on peut voir que l'ensemble \mathcal{N}_σ est un sous espace vectoriel de $C(I_n)$, (c'est une remarque simple, mais très importante), **notre but dans la suite est d'établir que : $\overline{\mathcal{N}_\sigma} = C(I_n)$** .

on procédera pour atteindre cet objectif par l'absurde : c'est à dire on va supposer que $\overline{\mathcal{N}_\sigma} \neq C(I_n)$, pour tomber enfin de compte sur une contradiction .

Premièrement, on sait que la fermeture d'un sous espace vectoriel est à son rôle un sous espace vectoriel^b ce qui fait que $\overline{\mathcal{N}_\sigma}$, est, **un sous espace vectoriel** de $C(I_n)$, **fermé** pour la norme de la convergence uniforme (car par définition c'est l'adhérence de \mathcal{N}_σ), et en plus **propre**, puisqu'on est parti de la supposition que : $\overline{\mathcal{N}_\sigma} \neq C(I_n)$.

on a : $\overline{\mathcal{N}_\sigma} \neq C(I_n) \implies \overline{\mathcal{N}_\sigma} \subsetneq C(I_n) \implies (\exists x_0 \in C(I_n)) : x_0 \notin \overline{\mathcal{N}_\sigma}$ (bien sûr $x_0 \neq 0$ car : $0 \in \mathcal{N}_\sigma \subset \overline{\mathcal{N}_\sigma}$).

on pose : $F = \overline{\mathcal{N}_\sigma} \oplus \mathbb{R}x_0$ et on considère l'application :

$$\boxed{\begin{array}{l} L : F = \overline{\mathcal{N}_\sigma} \oplus \mathbb{R}x_0 \longrightarrow \mathbb{R} \\ \quad \quad \quad x + \lambda.x_0 \longmapsto L(x) = \lambda \end{array}}$$

il est clair que L est bien définie et linéaire, et on peut affirmer sans problème que : $\text{Ker}(L) = \overline{\mathcal{N}_\sigma}$.

or : $(\forall x \in \overline{\mathcal{N}_\sigma}) : x = x + 0.x_0$, ce qui fait que : $(\forall x \in \overline{\mathcal{N}_\sigma}) : L(x) = 0$, donc : $\text{Ker}(L) \supset \overline{\mathcal{N}_\sigma}$ et :

$(\forall y = x + \lambda x_0 \in \overline{\mathcal{N}_\sigma} \oplus \mathbb{R}x_0) : y \in \text{Ker}(L) \implies L(y) = \lambda = 0 \implies y = x \in \overline{\mathcal{N}_\sigma}$ donc : $\text{Ker}(L) \subset \overline{\mathcal{N}_\sigma}$.

maintenant, à travers le théorème du Noyau fermé A.1.1.1 et car $\text{Ker}(L) = \overline{\mathcal{N}_\sigma}$, on déduit que L est continue sur $F = \overline{\mathcal{N}_\sigma} \oplus \mathbb{R}x_0$, et selon le corollaire A.1.2.1, de théorème de Hahn-Banach on peut en conséquence prolonger L sur l'espace $E = C(I_n)$ tout entier, par une autre forme linéaire continue \tilde{L} et qui vérifie naturellement les deux propriétés suivantes : $\tilde{L} \neq \mathbf{0}$ et $\tilde{L}(\overline{\mathcal{N}_\sigma}) = \tilde{L}(\text{Ker}(L)) = \mathbf{0}$.

si nous prolongeons une deuxième fois, par la même procédure, notre forme \tilde{L} à l'espace $L^2(I_n)$

qui est de Hilbert^c pour le produit scalaire définie par : $\forall (f, g) \in (L^2(I_n))^2 \quad \langle f, g \rangle := \int_{I_n} f.g \, d\lambda$,

on pourra bien par suite s'appuyer sur le théorème de représentation de Riesz A.2.2.1 pour affirmer que :

a. il faut bien déterminer la norme qu'on utilise, car dans $C(I_n)$ qui est de dimension infini, la propriété de l'équivalence de toutes les normes n'est pas vérifiée (voir par exemple $\|\cdot\|_\infty$ et $\|\cdot\|_1$).

b. ceci provient tout simplement de la linéarité de la limite et du fait que dans les espace vectoriels normés, et même métriques, l'adhérence d'une sous partie, est l'ensemble des limites des suites convergentes à valeurs dans cette partie (Attention c'est pas valable dans les espaces topologiques quelconques).

c. ce deuxième prolongement est nécessaire car pour le produit scalaire proposé $C(I_n)$ n'est pas du tout de Hilbert.

$$\left(\exists! g \in L^2(I_n)\right) \text{ tel que : } (\forall h \in C(I_n)) : \tilde{L}(h) = \langle g, h \rangle = \int_{I_n} h(x).g(x) d\lambda(x).$$

si on note : $\int_A g(x) d\lambda = \mu(A) \Leftrightarrow g(x) d\lambda(x) = d\mu(x)$ cette quantité, et selon le premier point de l'exemple III.1.1.2 est une mesure extérieure signée finie^a, et le fait qu'elle est de Borel régulière peut être justifiée comme suit :

$(\forall B \in \mathfrak{B}_{I_n}) : B$ est μ mesurable car :

$$(\forall A \in \mathcal{P}(I_n)) : \mu(B^c \cap A) + \mu(B \cap A) = \int_{B^c \cap A} g(x) d\lambda + \int_{B \cap A} g(x) d\lambda = \int_A g(x) d\lambda = \mu(A).$$

$$\text{et : } (\forall A \subset I_n) (\exists B = \overline{A} \in \mathfrak{B}_{I_n}) \text{ tel que : } A \subset B \text{ et } \mu(B) = \int_A g(x) d\lambda = \int_A g(x) d\lambda = \mu(A).$$

ainsi on obtient :

$$\mu \in M(I_n) \text{ et : } \boxed{(\forall h \in C(I_n)) : \tilde{L}(h) = \int_{I_n} h(x) d\mu(x)} \star.$$

d'après les hypothèses présentées dans l'énoncé, la fonction σ est continue, alors :

$$(\forall (\omega, \theta) \in \mathbb{R}^n \times \mathbb{R}) : \varphi : x \mapsto \sigma(\omega^\top x + \theta) \text{ est continue} \xrightarrow{\text{selon}^*} (\forall (\omega, \theta) \in \mathbb{R}^n \times \mathbb{R}) : \int_{I_n} \sigma(\omega^\top x + \theta) d\mu = \tilde{L}(\varphi)$$

et d'une autre part : $(\forall (\omega, \theta) \in \mathbb{R}^n \times \mathbb{R}) : \varphi : x \mapsto \sigma(\omega^\top x + \theta) \in \mathcal{N}_\sigma$, donc et car : $\tilde{L}(\mathcal{N}_\sigma) = \tilde{L}(\overline{\mathcal{N}_\sigma}) = 0$, on déduit que : $(\forall (\omega, \theta) \in \mathbb{R}^n \times \mathbb{R}) : \int_{I_n} \sigma(\omega^\top x + \theta) d\mu = \tilde{L}(\varphi) = 0$, ce qui implique que $\mu = 0$ car σ est

discriminatoire, alors : $(\forall h \in C(I_n)) : \tilde{L}(h) = \int_{I_n} h(x) d\mu(x) = 0$, cependant : $\tilde{L} \neq 0$, **Absurde!**.

donc $\overline{\mathcal{N}_\sigma} = C(I_n)$, d'où ce qu'il faut démontrer. ■

Ce théorème nous a démontré que les sommes de la forme (III.1) à condition que σ soit **continue** et **discriminatoire**, sont denses dans $C(I_n)$, et son idée de la preuve était assez générale et peut être appliquée dans d'autres cas, comme ça sera indiqué ultérieurement.

Maintenant, nous spécialisons ce résultat en montrant que **toute fonction sigmoïdale continue, est en particulier discriminatoire**, il convient de noter que, dans les applications de réseau neuronal, les fonctions d'activation sigmoïdale sont généralement considérées croissantes, mais **aucune monotonie n'est requise dans nos résultats**.

Lemme 1.2.1 :

Toute fonction **bornée mesurable sigmoïdale**, est discriminatoire dans $M(I_n)$, et en particulier toute fonction continue sigmoïdale l'est aussi,^b (voir pour cela : [13] et [12]).

Preuve :

Pour démontrer ceci, on note que pour n'importe quelle : $x, \omega \in \mathbb{R}^n$, $\theta, \varphi \in \mathbb{R}$ et $\lambda \in \mathbb{R}^+$ on a :

$$\sigma(\lambda(\omega^\top x + \theta) + \varphi) : \begin{cases} \xrightarrow{\lambda \rightarrow +\infty} 1 & \text{si : } \omega^\top x + \theta > 0 \\ \xrightarrow{\lambda \rightarrow +\infty} 0 & \text{si : } \omega^\top x + \theta < 0 \\ = \sigma(\varphi) & \text{si : } \omega^\top x + \theta = 0 \end{cases}$$

donc, la famille des fonctions : $(\sigma_\lambda)_{\lambda \in \mathbb{R}^+}$ converge simplement (ou point par point) vers la fonction : γ définie par :

$$\gamma(x) : \begin{cases} = 1 & \text{si : } \omega^\top x + \theta > 0 \\ = 0 & \text{si : } \omega^\top x + \theta < 0 \\ = \sigma(\varphi) & \text{si : } \omega^\top x + \theta = 0 \end{cases}$$

a. Cette mesure est finie car : $L^2(I_n) \subset L^1(I_n)$.

b. on reviendra à la démonstration de ce fait dans la preuve de deuxième théorème de Cybenko

on note $\Pi_{\omega,\theta}$ l'hyperplan **affine** donné par : $\{x \in \mathbb{R}^n \mid \omega^\top x + \theta = 0\}$ et $H_{\omega,\theta}$, le demi espace affine **ouvert**^a défini par : $\{x \in \mathbb{R}^n \mid \omega^\top x + \theta > 0\}$.

nous avons, et par hypothèse, la bornitude de la fonction σ , donc, les σ_λ sont équi-bornées, et en plus la convergence simple $\sigma_\lambda \xrightarrow{\lambda \rightarrow +\infty} \gamma$ est vérifiée par construction, donc et en appliquant pour n'importe quelle mesure $\mu \in M(I_n)$ **le théorème de la convergence dominée**^b, on obtient simultanément **la μ -intégrabilité de γ et le droit de permuter l'intégrale et la limite**.

$$\left. \begin{array}{l} \bullet (\exists M > 0) (\forall \lambda > 0) : \sigma_\lambda < M. \\ \bullet (\forall x \in \mathbb{R}) : \sigma_\lambda(x) \xrightarrow{\lambda \rightarrow +\infty} \gamma(x). \\ \bullet \int_{I_n} M \, d\mu = M \int_{I_n} 1 \, d\mu < +\infty. \end{array} \right\} \Rightarrow \boxed{\gamma \text{ est intégrable et : } \lim_{\lambda \rightarrow +\infty} \int_{I_n} \sigma_\lambda \, d\mu = \int_{I_n} \lim_{\lambda \rightarrow +\infty} \sigma_\lambda \, d\mu = \int_{I_n} \gamma \, d\mu.}$$

maintenant, passons directement à montrer que σ est discriminatoire, soit $\mu \in M(I_n)$ on a :

$$\begin{aligned} (\forall \omega \in \mathbb{R}^n) (\forall \theta \in \mathbb{R}) : \int_{I_n} \sigma(\omega^\top x + \theta) \, d\mu = 0 &\Rightarrow (\forall \omega \in \mathbb{R}^n) (\forall (\varphi, \lambda, \theta) \in \mathbb{R}^3) : \int_{I_n} \sigma(\lambda(\omega^\top x + \theta) + \varphi) \, d\mu = 0 \\ &\Rightarrow (\forall \omega \in \mathbb{R}^n) (\forall (\varphi, \lambda, \theta) \in \mathbb{R}^3) : \int_{I_n} \sigma_\lambda \, d\mu = 0 \\ &\Rightarrow (\forall \omega \in \mathbb{R}^n) (\forall (\varphi, \theta) \in \mathbb{R}^2) \lim_{\lambda \rightarrow +\infty} \int_{I_n} \sigma_\lambda \, d\mu = \int_{I_n} \lim_{\lambda \rightarrow +\infty} \sigma_\lambda \, d\mu \\ &= \int_{I_n} \gamma \, d\mu = \sigma(\varphi) \cdot \mu(\Pi_{\omega,\theta}) + \mu(H_{\omega,\theta}) = 0. \end{aligned}$$

donc nous obtenons : $\boxed{(\forall \omega \in \mathbb{R}^n) (\forall (\varphi, \theta) \in \mathbb{R}^2) : \sigma(\varphi) \cdot \mu(\Pi_{\omega,\theta}) + \mu(H_{\omega,\theta}) = 0.} \star$

si μ est positive, alors on pourra dire que la mesure de tout les demi plans $(H_{\omega,\theta})_{(\omega,\theta) \in \mathbb{R}^n \times \mathbb{R}}$ de \mathbb{R}^n est nulle, ce qui va nous permettre par la suite de déduire que^c :

$$(\forall (x_1, \dots, x_n) \in \mathbb{R}^n) : \mu \left(\prod_{i=1}^n]-\infty, x_i] \right) = 0.$$

et selon le lemme des classes monotones, ça va nous donner : $\mu = 0$. donc dans $M^+(I_n)$ la fonction σ est aisément discriminatoire.

Mais notre mesure n'est pas supposée positive, et ceci rendra notre tâche un peu plus difficile. afin de régler ce petit problème, fixons $\omega \in \mathbb{R}^n$, et considérons l'application :

$$\boxed{\begin{array}{l} F : L^\infty(\mathbb{R}) \longrightarrow \mathbb{R} \\ h \longmapsto \int_{I_n} h(\omega^\top x) \, d\mu \end{array}}$$

il est facile de voir que F est bien définie et bornée car :

$$(\forall h \in L^\infty(I_n)) : |F(h)| = \left| \int_{I_n} h(\omega^\top x) \, d\mu \right| \leq \int_{I_n} M \, d\mu \leq M \mu(I_n) < +\infty.$$

et si on prend : $h = \mathbb{1}_{[\theta, +\infty[}$, alors :

$$F(h) = \int_{I_n} \mathbb{1}_{[\theta, +\infty[}(\omega^\top x) \, d\mu = \mu(\{\omega^\top x \geq \theta\}) = \mu(\{\omega^\top x > \theta\}) + \mu(\{\omega^\top x = \theta\}) = \underbrace{\mu(H_{\omega,-\theta}) + \mu(\Pi_{\omega,-\theta})}_{\text{selon } \star} = 0.$$

de même, $F(h) = 0$ si h est la fonction indicatrice de l'intervalle ouvert : $]-\infty, \theta[$, donc et par linéarité, $F(h) = 0$ pour n'importe quelle fonction indicatrice d'un intervalle donné de \mathbb{R} , ce qui fait que

a. il est ouvert car c'est l'image réciproque de \mathbb{R}^{*-} par une fonction continue ($x \mapsto \langle \omega, x \rangle + \theta$).

b. ce théorème est valable pour les mesure signées **finies** et vous pouvez le démontrer en imitant la démonstration du troisième point dans la preuve de proposition III.1.1.1.

c. car on peut toujours, plonger les parties de la forme : $\prod_{i=1}^n]-\infty, x_i]$ dans un demi plan de \mathbb{R}^n , pour ceux qui ne sont pas convaincus par l'argument géométrique vous n'avez qu'à utiliser le théorème de séparation de Hahn-Banach.

pour toute fonction simple $\varphi = \sum_{i=1}^n \mathbb{1}_{\mathbb{I}_i}$ qui est sous forme d'une somme finie des fonctions indicatrices des intervalles $\mathbb{I}_i \subset \mathbb{R}$, on a : $F(\varphi) = 0$, et comme les fonctions simples sont denses dans $L^\infty(\mathbb{R})$, on déduit que $F = 0$.

en particulier, les fonctions mesurables bornées $c : x \mapsto \cos(\langle y, x \rangle)$ et $s : x \mapsto \sin(\langle y, x \rangle)$ nous donnent : $F(c + is) = \int_{I_n} \cos(y^\top x) + i \sin(y^\top x) d\mu(x) = \int_{I_n} \exp(i y^\top x) d\mu(x) = 0 \Rightarrow \underbrace{\hat{\mu} = 0}_{\text{voir : III.1.1.1}} \Rightarrow \mu = 0$.

donc, et en guise de conclusion, la fonction σ est discriminatoire, d'où ce qu'il faut prouver. ■

1.3 Application aux réseaux de neurones artificiels :

Dans cette section, nous allons appliquer les résultats précédents dans le cadre le plus intéressant, et qui est celui de la théorie des réseaux neuronaux.

Une combinaison directe du théorème III.1.2.1 et du lemme III.1.2.1 montre que les réseaux avec une seule couche interne, dont **les nœuds disposent d'une fonction d'activation sigmoïdale continue quelconque, peuvent approcher avec une précision arbitraire n'importe quelle fonction continue**, à condition qu'**aucune contrainte n'est imposée sur le nombre de nœuds ou la taille des poids**.

Les conséquences de ce résultat pour le rapprochement des fonctions de décisions seront bien traitées par la suite.

Théorème 1.3.1 :

Soit σ **une fonction sigmoïdale continue**, donc les sommes **finies** de la forme :

$$G(x) = \sum_{j=1}^N \alpha_j \sigma(\omega_j^\top x + \theta_j).$$

sont denses dans $C(I_n)$, et comme ça était formulé dans [13] : pour toute $f \in C(I_n)$ et $\varepsilon > 0$, il existe une somme, $G(x)$, de la forme précédente, pour laquelle :

$$(\forall x \in I_n) : |G(x) - f(x)| < \varepsilon.$$

Preuve :

Soit σ une fonction sigmoïdale continue, alors et sans problème on peut dire qu'elle est bornée, **or s'il n'était pas le cas** l'ensemble : $\mathcal{Jm}(\sigma) := \{\sigma(x) / x \in \mathbb{R}\}$ sera non borné, ce qui fait qu'il va être non majoré ou non minoré, ou les deux à la fois.

sans perte de généralité, on va traiter le cas où cet ensemble est non majoré, et pour les autres, il suffit de procéder d'une façon analogue.

si $\mathcal{Jm}(\sigma)$ est non majorée, alors il existe une suite : $(y_n)_{n \in \mathbb{N}} \in \mathcal{Jm}(\sigma)^{\mathbb{N}}$ tel que : $y_n \xrightarrow[n \rightarrow +\infty]{} +\infty$

ce qui fait qu'il existe une suite : $(x_n)_{n \in \mathbb{N}} \in \mathbb{R}^{\mathbb{N}}$ tel que : $\sigma(x_n) = y_n \xrightarrow[n \rightarrow +\infty]{} +\infty$.

si : $(x_n)_{n \in \mathbb{N}}$ est non bornée, alors il existe une sous suite $(x_{\phi(n)})_{n \in \mathbb{N}}$ tel que : $x_{\phi(n)} \xrightarrow[n \rightarrow +\infty]{} +\infty$ ou $x_{\phi(n)} \xrightarrow[n \rightarrow +\infty]{} -\infty$, et ceci nous donne que : $y_{\phi(n)} = \sigma(x_{\phi(n)}) \xrightarrow[n \rightarrow +\infty]{} \pm 1$, **absurde!**, car :

$$\lim_{n \rightarrow +\infty} y_n = +\infty \Rightarrow \lim_{n \rightarrow +\infty} y_{\phi(n)} = +\infty.$$

d'où : $(x_n)_{n \in \mathbb{N}}$ est bornée, de ce fait, et en vertu du fameux théorème de **Bolzano-Weierstrass** on peut extraire de $(x_n)_{n \in \mathbb{N}}$ une sous suite $(x_{\phi(n)})_{n \in \mathbb{N}}$ qui converge vers un élément $x^* \in \mathbb{R}$.

$\lim_{n \rightarrow +\infty} x_{\phi(n)} = x^* \Rightarrow \sigma \left(\lim_{n \rightarrow +\infty} x_{\phi(n)} \right) = \sigma(x^*) \Rightarrow$ ^a $\lim_{n \rightarrow +\infty} \sigma(x_{\phi(n)}) = \sigma(x^*) \Rightarrow \lim_{n \rightarrow +\infty} y_{\phi(n)} = +\infty = \sigma(x^*)$.

absurde!, donc la fonction σ est bornée.

a. grâce à la continuité de σ .

maintenant on obtient que σ est continue, sigmoïdale et bornée ce qui donne qu'elle est en particulier mesurable (car la continuité implique la mesurabilité), sigmoïdale et bornée, donc, et selon le lemme III.1.2.1, σ sera bien discriminatoire, ce qui donne enfin de compte et à travers le théorème III.1.2.1 la véracité de l'énoncé suggéré. ■

Nous allons présenter et démontrer maintenant le rôle et les implications de ces résultats dans le contexte des **régions de décision**, et qui est vraiment un contexte très important, surtout avec son apparence dans plusieurs domaines d'applications .

soit $\mu := \lambda|_{\mathfrak{B}_{I_n}}$ la mesure de Lebesgue dans I_n , et Soit $:P_1, P_2, \dots, P_k$ une partition de I_n en k sous-ensembles mesurables disjoints. on définit la fonction de décision, f par : $f := \sum_{i=1}^n i \cdot \mathbb{1}_{P_i}$, ce qui veut

dire que : $(\forall i \in \llbracket 1, k \rrbracket) (\forall x \in I_n) : f(x) = i \iff x \in P_i$.

La question qui se pose maintenant est de savoir si une telle fonction de décision peut être mise en œuvre par un réseau neuronal d'une seule couche cachée.

la réponse à cette question a été fournie dans l'article [13], et elle se présente à travers le résultat suivant :

Théorème 1.3.2 :

Soit σ **une fonction sigmoïdale continue**, et f une fonction de décision associée à une certaine partition mesurable finie, de I_n . pour tout $\varepsilon > 0$, il existe une somme finie de la forme :

$$G(x) = \sum_{j=1}^N \alpha_j \sigma(\omega_j^\top x + \theta_j).$$

et un **compact** $D_\varepsilon \subset I_n$ tel que : $\lambda(D_\varepsilon) \geq 1 - \varepsilon$ et : $|G(x) - f(x)| < \varepsilon$ pour tout $x \in D_\varepsilon$, c'est à dire :

$$(\forall \varepsilon > 0) (\exists N_\varepsilon \in \mathbb{N}) \left(\exists (\alpha^\varepsilon, \theta^\varepsilon) \in (\mathbb{R}^{N_\varepsilon})^2 \right) \left(\exists (\omega_1^\varepsilon, \dots, \omega_{N_\varepsilon}^\varepsilon) \in (\mathbb{R}^n)^{N_\varepsilon} \right) \left(\exists D_\varepsilon \in \mathfrak{B}_{I_n} \text{ compact tq : } \mu(D_\varepsilon^c) < \varepsilon \right)$$

tel que : $x \in D_\varepsilon \implies \left| \sum_{j=1}^{N_\varepsilon} \alpha_j^\varepsilon \sigma(\langle \omega_j^\varepsilon, x \rangle + \theta_j^\varepsilon) - f(x) \right| < \varepsilon$.

Preuve :

On a : $f := \sum_{i=1}^k i \cdot \mathbb{1}_{P_i}$ avec : $(P_i)_{1 \leq i \leq k}$ est une partition de I_n en k sous-ensembles mesurables dis-

joint, donc : $\int_{I_n} f \, d\mu = \sum_{i=1}^k i \cdot \mu(P_i) \leq \sum_{i=1}^k i \cdot \mu(I_n) \leq \frac{k \times (k+1)}{2} < +\infty$, ce qui fait que : $f \in L^1(I_n, \mathfrak{B}_{I_n}, \mu)$.

donc, on peut bien appliquer le théorème de Lusin dans sa version faible I.1.3.1, pour déduire que pour tout $\varepsilon > 0$ il existe un sous ensemble **compact** $D \subset I_n$ avec $\mu(D^c) \leq \varepsilon$, et une fonction h continue sur cet ensemble tel que : $f|_D = h$, c'est à dire : $(\forall x \in D) : f(x) = h(x)$.

maintenant h est continue, donc, et selon le théorème III.1.3.1, on peut trouver une somme de la forme G précédente et qui satisfait : $(\forall x \in D) : |G(x) - f(x)| = |G(x) - h(x)| < \varepsilon$.

Remarque : à cause de la continuité, nous sommes toujours exposés au risque de faire quelques décisions incorrectes pour certains points, mais ce résultat nous indique que la mesure totale de ces points mal classés peut être si petite que voulu. à la lumière de cela, le théorème III.1.3.1 **semble être le résultat le plus puissant possible de ce type**.

Nous pouvons et sans perte de généralité, restreindre un peu notre étude en considérant le problème de décision pour un seul ensemble non vide fermé $D \subset I_n$, alors dans ce cas : $f = \mathbb{1}_D$, notre but pour la suite est de trouver une fonction sous forme de la sommation (III.1) et qui approche la fonction de décision f .

soit $\Delta(x, D) = \min \{ |x - y|, y \in D \}$, selon le lemme I.1.2.1, $x \mapsto \Delta(x, D)$ est continue, ce qui fait que :

$$(\forall \varepsilon > 0) : f_\varepsilon(x) = \max \left\{ 0, \frac{\varepsilon - \Delta(x, D)}{\varepsilon} \right\}$$

est une fonction **continue**^a, de sorte que $f_\varepsilon(x) = 0$ pour les points x plus éloignés que ε de D alors que $f_\varepsilon(x) = 1$ pour les points $x \in D$.

par théorème III.1.3.1, nous pouvons trouver et pour n'importe quel $\eta > 0$

une certaine fonction : $G(x) = \sum_{j=1}^N \alpha_j \sigma(\omega_j^\top x + \theta_j)$, tel que : $|G(x) - f_\varepsilon(x)| < \eta$.

alors et en choisissant ce η si petit, de préférence qu'il soit plus petit que ε , on pourra bien ajuster la fonction G à f_ε et par la suite utiliser ce G comme **une fonction de décision approximative** : $G(x) \leq \frac{1}{2}$ permet de deviner que $x \in D^c$ tandis que : $G(x) > \frac{1}{2}$ devine que $x \in D$.

cette procédure de décision est correcte pour tout $x \in D$ et pour tout x à distance au moins ε loin de D , mais Si x est à distance ε de D , sa classification dépend du choix particulier de G . ■

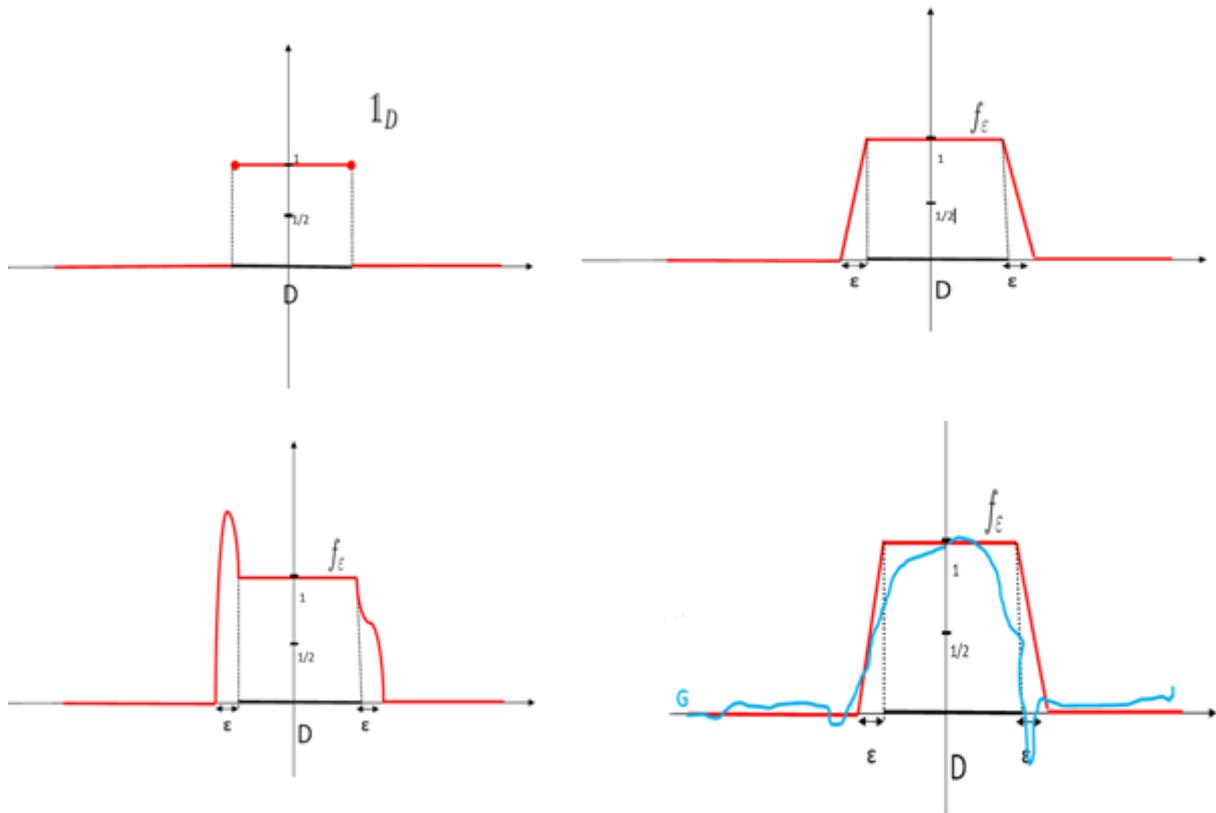


FIGURE III.2: L'idée générale de la preuve.

Remarque 1.3.1 :

En somme ce dernier théorème, montre qu'il existe un réseau qui classe correctement les points suffisamment éloignés de la région de décision, ainsi ceux qui sont à l'intérieur de cette région, en revanche il rend la mesure des points incorrectement classés si petite que voulu **mais ne garantit pas leur emplacement**.

a. le max entre zéro et une certaine fonction continue g , reste toujours une fonction continue, pour justifier ça il suffit de remarquer que : $(\forall x \in \mathbb{R}) : \max(0, g(x)) = \frac{|g(x)| + g(x)}{2}$.

1.4 Résultats pour d'autres fonctions d'activation :

Dans cette section, nous discutons *d'autres classes de fonctions d'activation* et qui ont des propriétés similaires à celles des sigmoïdales continues, mais ces exemples **ont un intérêt un peu moins pratique**, ce qui va justifier le fait de se contenter par des preuves plus résumées, et dont la plupart vont se faire par une simple analogie avec les cas traités précédemment.

On va commencer par l'étude des **fonctions sigmoïdales discontinues** et qui ont vraiment une importance considérable, malgré qu'elles sont moins utilisées que les continues, **en raison du manque d'algorithmes adéquats**.

ces fonctions présentent un intérêt théorique appréciable surtout à travers leur relation avec les perceptrons classiques.

Supposons que σ est une fonction sigmoïdale **bornée et mesurable** (la continuité n'est pas requise dans ce cas), nous avons un analogue du théorème III.1.3.1 et qui était évoqué pour la première fois en 1989 par George Cybenko dans son article [13], et il s'énonce ainsi :

Théorème 1.4.1 :(approximation des fonctions intégrables :)

Soit σ **une fonction sigmoïdale mesurable et bornée**, donc les sommes **finies** de la forme :

$$G(x) = \sum_{j=1}^N \alpha_j \sigma(\omega_j^\top x + \theta_j).$$

sont denses dans $L^1(I_n)$ ^a, En d'autres termes : pour tout $f \in L^1(I_n)$ et $\varepsilon > 0$, il existe une somme, $G(x)$ de la forme précédente, pour laquelle :

$$\|G - f\|_1 = \int_{I_n} |G(x) - f(x)| dx < \varepsilon.$$

Preuve

La preuve va suivre exactement les mêmes étapes que celles des théorèmes III.1.2.1 et III.1.3.1 mais avec des modifications évidentes, telles que remplacer les fonctions continues par des fonctions intégrables, et utiliser le fait que $L^\infty(I_n)$ est le dual de $L^1(I_n)$, la notion de discrimination d'une certaine fonction σ sera donc modifiée en conséquence de la façon suivante :

$$(\forall h \in L^\infty(I_n)) : \left[(\forall \omega \in \mathbb{R}^n) (\forall \theta \in \mathbb{R}) : \int_{I_n} \sigma(\omega^\top x + \theta) h(x) dx = 0 \right] \implies h(x) = 0 \text{ presque partout.}$$

sans surprise et exactement comme on l'a déjà vu au lemme III.1.2.1, les fonctions sigmoïdales mesurables bornées sont en particulier discriminatoires dans ce sens, car les mesures de la forme $h(x) dx$ appartiennent en particulier à $M(I_n)$. ■

Lemme 1.4.1 :(l'inégalité de Markov :)

soit $(\Omega, \mathcal{F}, \mu)$ un espace mesuré, et soit $f \in \mathfrak{M}((\Omega, \mathcal{F}); (\mathbb{R}, \mathcal{B}_{\mathbb{R}}))$ une fonction \mathcal{F} - $\mathcal{B}_{\mathbb{R}}$ mesurables alors on a :

$$(\forall \varepsilon > 0) : \mu(\{|f| \geq \varepsilon\}) \leq \frac{1}{\varepsilon} \int_{\Omega} |f| d\mu$$

Preuve :

Soit $\varepsilon > 0$ on prend $A = \{|f| \geq \varepsilon\}$, et $\varphi = \varepsilon \mathbb{1}_A$, il est clair que : $\varphi \leq |f|$ alors : $\int_{\Omega} \varphi d\mu \leq \int_{\Omega} |f| d\mu$ ce qui fait que : $\varepsilon \mu(\{|f| \geq \varepsilon\}) \leq \int_{\Omega} |f| d\mu$ donc : $(\forall \varepsilon > 0) : \mu(\{|f| \geq \varepsilon\}) \leq \frac{1}{\varepsilon} \int_{\Omega} |f| d\mu.$

a. **Attention!** : cette fois ci on parle de la densité au sens de la norme L^1 **et non plus de la norme de convergence uniforme**.

d'où ce qu'il faut prouver. ■

Proposition 1.4.1 :

Soit $(\Omega, \mathcal{F}, \mu)$ un espace mesuré, et soit f et $(f_n)_{n \in \mathbb{N}^*} \in \mathfrak{M}((\Omega, \mathcal{F}); (\mathbb{R}, \mathcal{B}_{\mathbb{R}}))$ des fonctions $\mathcal{F} - \mathcal{B}_{\mathbb{R}}$ mesurables, alors on a :

$$\left((\exists r \geq 1) : f_n \xrightarrow[n \rightarrow +\infty]{L^r} f \right) \implies f_n \xrightarrow[n \rightarrow +\infty]{\mu} f.$$

c'est à dire que : la convergence dans L^r , implique la convergence en mesure ^a (voir : [2]).

Preuve :

Selon l'inégalité de Markov

$$\text{On a : } (\forall \varepsilon > 0) : \mu(|f_n - f| > \varepsilon) = \mu(|f_n - f|^r > \varepsilon^r) \leq \frac{1}{\varepsilon^r} \int_{\Omega} |f_n - f|^r d\mu$$

donc si : $\int_{\Omega} |f_n - f|^r d\mu \xrightarrow[n \rightarrow \infty]{} 0$ on aura : $(\forall \varepsilon > 0) : \mu(|f_n - f| > \varepsilon) \xrightarrow[n \rightarrow \infty]{} 0$. ce qui achève la preuve. ■

en vertu de cette proposition, la convergence dans L^1 implique la convergence en mesure, et de ce fait, nous avons un analogue du théorème III.1.3.2 et qui était énoncé dans l'article [13] comme suit :

Théorème 1.4.2 :

Soit σ **une fonction sigmoïdale mesurable et bornée**, et f une fonction intégrable sur I_n , alors, pour tout $\varepsilon > 0$, il existe une somme finie de la forme :

$$G(x) = \sum_{j=1}^N \alpha_j \sigma(\omega_j^T x + \theta_j).$$

et un ensemble (**pas forcément compact**) $D \subset I_n$ tel que : $\mu(D) \geq 1 - \varepsilon$ et : $|G(x) - f(x)| < \varepsilon \quad (\forall x \in D)$.

Remarque 1.4.1 :

Si σ est une fonction sigmoïdale **mesurable^b** et bornée, alors elle est en particulier dans $L^1(I_n)$ car :

$$\begin{aligned} (\forall x \in I_n) : |G(x)| &= \left| \sum_{j=1}^N \alpha_j \sigma(\omega_j^T x + \theta_j) \right| \leq \sum_{j=1}^N |\alpha_j| \left| \sigma(\omega_j^T x + \theta_j) \right| \leq \sum_{j=1}^N |\alpha_j| M \\ &\Rightarrow \int_{I_n} |G(x)| d\lambda < \int_{I_n} \sum_{j=1}^N |\alpha_j| M d\lambda = \sum_{j=1}^N |\alpha_j| M \underbrace{\lambda(I_n)}_{=1} < +\infty. \end{aligned}$$

Preuve : (de théorème III.1.4.2)

Selon le théorème III.1.4.1, pour toute fonction **sigmoïdale mesurable et bornée** σ , l'ensemble :

$$\mathcal{N}_{\sigma} = \left\{ x \mapsto G(x) = \sum_{j=1}^N \alpha_j \sigma(\omega_j^T x + \theta_j) \mid N \in \mathbb{N}, (\alpha_1, \dots, \alpha_N) \in \mathbb{R}^N, (\theta_1, \dots, \theta_N) \in \mathbb{R}^N, (\omega_1, \dots, \omega_N) \in (\mathbb{R}^n)^N \right\}$$

a. Attention!, **la réciproque est fautive** en général.

b. le fait que σ soit mesurable est indispensable pour déduire la mesurabilité sur I_n de la fonction : $x \mapsto \sigma(\omega_j^T x + \theta_j)$, ce qui va nous permettre par la suite de parler de l'intégrabilité de cette fonction sur I_n .

est dense dans $L^1(I_n)$ **pour la norme** : $\|\cdot\|_1 : f \mapsto \int_{I_n} f(t) dt$; donc et par la caractérisation de densité dans les espaces métriques ^a, on a :

$$\left(\forall f \in L^1(I_n) \right) \left(\exists (G_n)_{n \in \mathbb{N}} \in \mathcal{M}^{\mathbb{N}} \right) \text{ tq : } \|G_n - f\|_1 \xrightarrow{n \rightarrow +\infty} 0 \left(\iff G_n \xrightarrow[n \rightarrow +\infty]{L^1} f \right).$$

en se basant sur cette équivalence et sur la proposition III.1.4.1, on pourra bien mettre en évidence que : $G_n \xrightarrow[n \rightarrow +\infty]{\lambda} f$, ce qui veut dire, et par définition, que : $(\forall \varepsilon > 0) : \lim_{n \rightarrow +\infty} \lambda(|G_n - f| > \varepsilon) = 0$.

donc : $(\forall \varepsilon > 0) (\forall \varepsilon' > 0) (\exists \eta_{\varepsilon, \varepsilon'} > 0) \text{ tq : } (\forall n \in \mathbb{N}) : n \geq \eta_{\varepsilon, \varepsilon'} \Rightarrow \lambda(|G_n - f| > \varepsilon) \leq \varepsilon'$.

on prend en particulier : $\varepsilon = \varepsilon'$, alors : $(\forall \varepsilon > 0) (\exists \eta_\varepsilon > 0) \text{ tq : } (\forall n \in \mathbb{N}) : n \geq \eta_\varepsilon \Rightarrow \lambda(|G_n - f| > \varepsilon) \leq \varepsilon$.

et ceci implique que :

$$(\forall \varepsilon > 0) \left(\exists N_\varepsilon = \underbrace{[\eta_\varepsilon] + 1}_{\text{partie entière}} \in \mathbb{N}^* \right) \text{ tq : } \lambda(|G_{N_\varepsilon} - f| > \varepsilon) \leq \varepsilon. \quad \star$$

on pose : $D = \{x \in I_n / |G_{N_\varepsilon}(x) - f(x)| \leq \varepsilon\}$, selon \star , $\lambda(|G_{N_\varepsilon} - f| > \varepsilon) \leq \varepsilon \Rightarrow \lambda(D^c) = \lambda(I_n \setminus D) \leq \varepsilon$.
donc : $\lambda(D) \geq 1 - \varepsilon$, et : $(\forall x \in D) : |G_{N_\varepsilon} - f| \leq \varepsilon$ d'où ce qu'il faut démontrer. ■

Remarques 1.4.2 :

1. Pour ce dernier théorème, on peut choisir en particulier la fonction f comme une fonction de décision, et de ce fait on obtiendra un résultat qui ressemble un peu au théorème : III.1.3.2 **mais faites attention**, la compacité de l'ensemble D **n'est plus garanti** pour ce dernier cas.
2. si on ne met pas la compacité de D comme une conséquence dans le théorème III.1.3.2, alors il va être un résultat trivial de III.1.4.2, et l'appel du théorème de Lusin ne sera plus nécessaire.

2 Théorèmes d'approximation au cas d'une fonction d'activation écrasante :

2.1 Définitions et notations :

Notations :

Pour des raisons de cohérence et de simplification, on adoptera dans la suite de ce chapitre, les notations suivantes :

- pour tout $r \in \mathbb{N}^*$, on note : $\mathcal{A}^r = \left\{ A \in (\mathbb{R}^r)^{\mathbb{R}} / (\exists (\omega, \theta) \in \mathbb{R}^r \times \mathbb{R}) : A : x \mapsto \langle \omega, x \rangle + \theta \right\}$.
- pour toute fonction borélienne ^b, $\sigma : \mathbb{R} \rightarrow \mathbb{R}$, on note : $\Sigma^r(\sigma) = \mathcal{N}_\sigma$.

Remarque 2.1.1 :

1. Dans ce contexte, x **correspond à l'entrée d'un réseau de neurones**, ω **correspond aux poids de ce réseau** reliant la couche d'entrée à un nœud appartenant à celle intermédiaire, et b **représente tout simplement le biais** de ce nœud.

a. tenez bien compte que cette caractérisation n'est pas valable en général pour les espaces topologiques non métrisables.

b. Une fonction borélienne n'est autre qu'une fonction mesurable dont l'espace de départ est un espace topologique muni de la tribu borélienne.

2. il est facile de voir que : $\Sigma^r(\sigma)$ n'est autre que l'ensemble des fonctions qui peuvent être réalisées , par un réseau de neurones **ordinaire**^a, dont les nœuds possèdent tous la même fonction d'activation σ .

Définition 2.1.1 :(fonction d'écrasement :)

Une fonction $\Psi : \mathbb{R} \rightarrow [0, 1]$ est dite **écrasante** ou **d'écrasement** lorsqu'elle est **sigmoïdale** et **croissante** (voir :[34]).

Exemples 2.1.1 :

1. la fonction de Heaveside définie par : $\Psi(x) = \mathbb{1}_{\mathbb{R}^+}(x)$ est clairement une fonction d'écrasement.
2. la fonction "rampe" définie par : $\Psi(x) = \mathbb{1}_{[1,+\infty[}(x) + x \cdot \mathbb{1}_{[0,1]}(x)$. est encore une fonction d'écrasement.
3. le cosinus écraseur de Gallant-White défini par : $\Psi(x) = \frac{1}{2} \left(1 + \cos \left(x + \frac{3\pi}{2} \right) \right) \mathbb{1}_{\left[-\frac{\pi}{2}, \frac{\pi}{2}\right]} + \mathbb{1}_{\left[\frac{\pi}{2}, +\infty\right]}$.

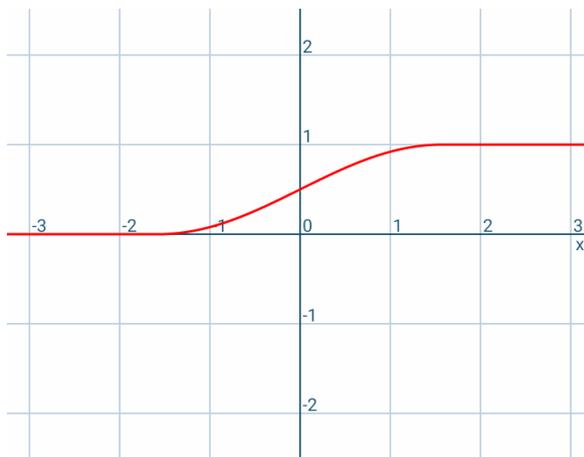


FIGURE III.3: La courbe du cosinus écraseur.

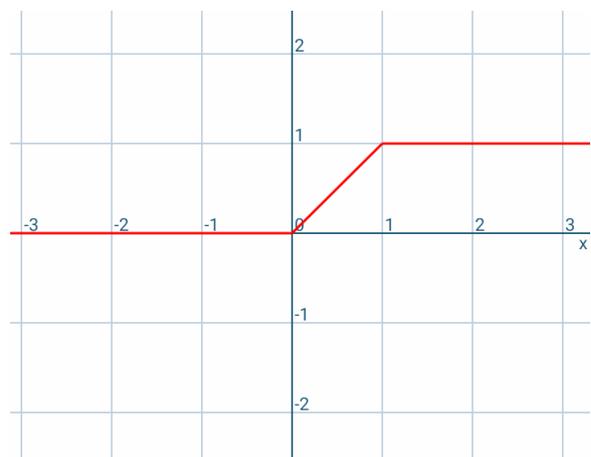


FIGURE III.4: La courbe de la fonction Rampe.

Lemme 2.1.1 :

Soit $f : I \rightarrow \mathbb{R}$ une fonction **croissante** sur un intervalle : $I \subset \mathbb{R}$, alors ,elle admet une limite à gauche et à droite en tout point $a \in I$,et de plus on a :

$$\boxed{\lim_{x \rightarrow a^-} f(x) \leq f(a) \leq \lim_{x \rightarrow a^+} f(x).}$$

Preuve :

La preuve est vraiment très classique,or,et selon [7] il suffit de considérer pour tout $a \in I$ les ensembles :

$$\mathbb{F}_a^+ = \left\{ x \in I / f(x) \geq f(a) \right\} \text{ et : } \mathbb{F}_a^- = \left\{ x \in I / f(x) \leq f(a) \right\}.$$

sans problème,ces deux ensembles sont **non vides** (car I est ouvert) et par la croissance de notre fonction,on obtient respectivement que le premier est minoré tandis que le deuxième est majoré

a. on désigne par le terme ordinaire ici que tous les nœuds constituant notre réseau sont doté d'une fonction d'entrée totale qui est la somme pondérée,et d'une fonction de sortie confondue avec l'identité.

donc les quantités : $\inf \mathbb{F}_a^+, \sup \mathbb{F}_a^-$ auront bien un sens, et en plus :

$$(\forall x \in \mathbb{F}_a^-) : f(x) \leq f(a) \Rightarrow \sup \mathbb{F}_a^- \leq f(a) \text{ et } (\forall x \in \mathbb{F}_a^+) : f(x) \geq f(a) \Rightarrow \inf \mathbb{F}_a^+ \geq f(a).$$

donc et à travers la monotonie de la fonction f et la définition de la borne supérieure ainsi que l'inférieure : on peut déduire que : $\lim_{x \rightarrow a^-} f(x) = \sup \mathbb{F}_a^- \leq f(a)$ et $\lim_{x \rightarrow a^-} f(x) = \inf \mathbb{F}_a^+ \geq f(a)$

d'où : $\lim_{x \rightarrow a^-} f(x) \leq f(a) \leq \lim_{x \rightarrow a^+} f(x)$. ■

Théorème 2.1.1 : (voir : [7])

Si une fonction f est croissante sur un intervalle ouvert I , alors l'ensemble des points où elle n'est pas continue est au plus dénombrable ^a.

Preuve :

Soient a et b deux réels de I tels que $a < b$, supposons que f est non continue en ces deux points. alors : $\lim_{x \rightarrow a^-} f(x) < \lim_{x \rightarrow a^+} f(x) \leq \lim_{x \rightarrow b^-} f(x) < \lim_{x \rightarrow b^+} f(x)$ ce qui fait que les intervalles ouverts :

$$\left] \lim_{x \rightarrow a^-} f(x), \lim_{x \rightarrow a^+} f(x) \right[\text{ et } \left] \lim_{x \rightarrow b^-} f(x), \lim_{x \rightarrow b^+} f(x) \right[$$

sont disjoints et non vides, par la densité de \mathbb{Q} dans \mathbb{R} , chacun d'eux contient au moins un rationnel. on peut donc et **par l'axiome de choix** construire une application injective, associant à tout point de discontinuité de f , un rationnel. comme l'ensemble des rationnels est dénombrable, le résultat s'ensuit. ■

Remarque 2.1.2 :

On peut étendre les deux résultats précédents pour les fonctions décroissantes à travers un simple remplacement de : f par $-f$.

Corollaire 2.1.1 :

1. Toute fonction $f : I \rightarrow \mathbb{R}$ croissante, sur un intervalle ouvert I de \mathbb{R} est $(\mathcal{B}_{\mathbb{R}} \cap I, \mathcal{B}_{\mathbb{R}})$ mesurable.
2. en particulier **les fonctions d'écrasement sont toujours boréliennes.**

Preuve :

1. Afin de démontrer ce corollaire, on note : \mathcal{D} l'ensemble des points de discontinuité de f et $\mathcal{C} = \mathcal{D}^c$ l'ensemble de ses points de continuité, selon le théorème III.2.1.1 : \mathcal{D} est au plus dénombrable, ce qui fait qu'il est mesurable ^b donc $\mathcal{C} = \mathcal{D}^c$ sera à son tour mesurable comme étant le complémentaire d'un ensemble mesurable, et par suite la fonction indicatrice : $\mathbb{1}_{\mathcal{C}}$ sera immédiatement une fonction mesurable ce qui implique que : $f \cdot \mathbb{1}_{\mathcal{C}}$ le sera aussi, car c'est le produit de deux fonctions mesurables (la fonction f est continue sur \mathcal{C}).

et d'une autre part : l'application : $\varphi : x \mapsto f \cdot \mathbb{1}_{\mathcal{D}}$ est mesurable car :

$$(\forall A \in \mathcal{B}_{\mathbb{R}}) : \varphi^{-1}(A) \subset \mathcal{D} \implies \text{card}(\varphi^{-1}(A)) \leq \aleph_0 \implies \varphi^{-1}(A) \in \mathcal{B}_{\mathbb{R}} \cap I.$$

donc et sans aucun problème la fonction : $f = f \cdot \mathbb{1}_{\mathcal{C}} + f \cdot \mathbb{1}_{\mathcal{D}}$ sera mesurable.

2. par définition les fonctions d'écrasement sont croissantes donc et selon le premier point ça va impliquer qu'elles sont mesurables. ■

a. c'est à dire de cardinalité $\leq \aleph_0$.

b. On laisse au lecteur le soin de vérifier que tout sous ensemble au plus dénombrable de \mathbb{R} est mesurable en particulier.

Définition 2.1.2 : (les réseaux de neurones $\Sigma\Pi$)

Pour toute fonction $\sigma : \mathbb{R} \rightarrow \mathbb{R}$, et $r \in \mathbb{N}^*$ on définit la classe des fonctions $\Sigma\Pi^r(\sigma)$ par :

$$\Sigma\Pi^r(\sigma) = \left\{ x \mapsto G(x) = \sum_{j=1}^q \beta_j \cdot \prod_{k=1}^{l_j} \sigma(A_{jk}(x)) \mid q \in \mathbb{N}^*, l_j \in \mathbb{N}, (\beta_1, \dots, \beta_q) \in \mathbb{R}^q, A_{jk}(x) \in \mathcal{A}^r \right\}.$$

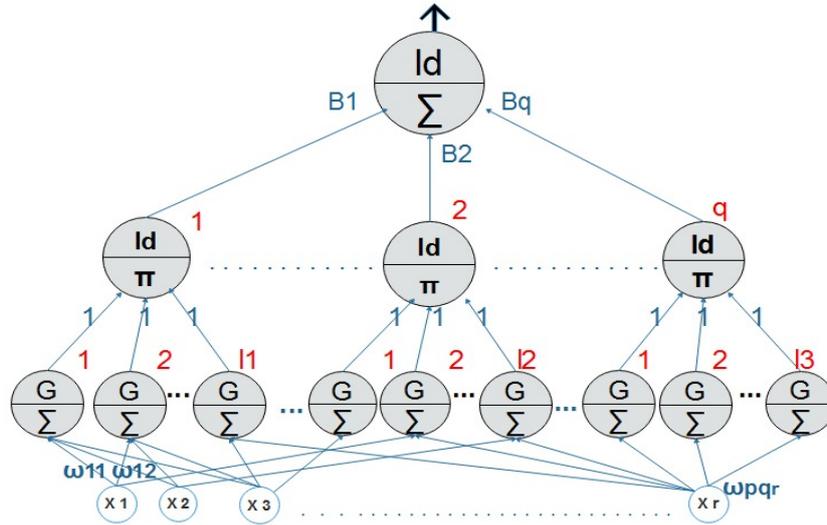


FIGURE III.5: L'architecture d'un réseaux de neurones $\Sigma\Pi^r(G)$.

Remarque 2.1.3 :

À l'instar de [34], on prouvera dans la suite, nos résultats généraux pour les réseaux $\Sigma\Pi$, et ensuite on va les étendre aux réseaux Σ , qui ne sont qu'un simple cas particulier des $\Sigma\Pi$, pour lesquels : $l_j = 1$ pour tout $j \in \{1, \dots, q\}$.

Définition 2.1.3 : (la densité uniforme compacte :)

1. Un sous-ensemble S de $C(\mathbb{R}^r)$ est dit **uniformément compactement dense** dans $C(\mathbb{R}^r)$, si **pour tout compact** $K \subset \mathbb{R}^r$, $S_K = \{f|_K / f \in S\}$ est ρ_K -dense^a dans $C(K)$ avec : $(\forall (f, g) \in C(K)) : \rho_K(f, g) = \sup_{x \in K} |f(x) - g(x)|$.
2. Une suite de fonctions $(f_n)_{n \in \mathbb{N}} \in C(\mathbb{R}^r)^{\mathbb{N}}$ **converge** vers une fonction $f \in C(\mathbb{R}^r)$ **uniformément compactement** sur $C(\mathbb{R}^r)$ si pour tout compact $K \subset \mathbb{R}^r$: $\rho_K(f_n, f) \xrightarrow{n \rightarrow +\infty} 0$.

Remarque 2.1.4 :

On peut voir facilement que la densité uniforme d'un ensemble $S \subset C(\mathbb{R}^r)$, implique directement sa densité uniforme compacte et ceci provient du fait que lorsque : $\sup_{x \in \mathbb{R}^r} |f(x) - g(x)| < \varepsilon$ est vérifiée alors : $\sup_{x \in K} |f(x) - g(x)| < \varepsilon$ le sera aussi sans doute. cette petite remarque va faire l'objet d'un petit passage dans la preuve du théorème III.2.2.3.

a. il est très important de spécifier la norme, car dans les espaces de dimension infinie **l'équivalence de toute les normes n'est pas vérifiée.**

Définition 2.1.4 : (la μ -équivalence :)

l'espace de fonctions $\mathcal{B}_{\mathbb{R}^r} - \mathcal{B}_{\mathbb{R}}$ mesurables

Soit μ une mesure **quelconque** sur $(\mathbb{R}^r, \mathcal{B}_{\mathbb{R}^r})$, si f et g appartiennent à $\overbrace{\mathfrak{M}((\mathbb{R}^r, \mathcal{B}_{\mathbb{R}^r}); (\mathbb{R}, \mathcal{B}_{\mathbb{R}}))}^{\text{l'espace de fonctions } \mathcal{B}_{\mathbb{R}^r} - \mathcal{B}_{\mathbb{R}} \text{ mesurables}}$, alors on dit qu'elles sont μ -équivalentes, ssi : $\mu \{x \in \mathbb{R}^r / f(x) \neq g(x)\} = 0$, et on note dans ce cas : $f \mathcal{R}_\mu g$.

Remarque 2.1.5 :

1. Si μ est finie avec $\mu(\mathbb{R}^r) = m$, (pensez aux mesures de probabilité par exemple), alors la définition précédente peut être exprimée ainsi :

$$f \mathcal{R}_\mu g \iff \mu \{x \in \mathbb{R}^r / f(x) = g(x)\} = m.$$

2. comme son nom l'indique la relation " \mathcal{R}_μ " **est bien une relation d'équivalence** sur : $\mathfrak{M}((\mathbb{R}^r, \mathcal{B}_{\mathbb{R}^r}); (\mathbb{R}, \mathcal{B}_{\mathbb{R}}))$, ce qui fait qu'on peut bien donner un sens à l'espace quotient :

$$M((\mathbb{R}^r, \mathcal{B}_{\mathbb{R}^r}); (\mathbb{R}, \mathcal{B}_{\mathbb{R}})) = \mathfrak{M}((\mathbb{R}^r, \mathcal{B}_{\mathbb{R}^r}); (\mathbb{R}, \mathcal{B}_{\mathbb{R}})) / \mathcal{R}_\mu$$

et ce dernier fera bien l'objet de plusieurs résultats qui seront présentés par la suite.

Notation :

On opte pour la suite les notations suivantes :

$$M^r = M((\mathbb{R}^r, \mathcal{B}_{\mathbb{R}^r}); (\mathbb{R}, \mathcal{B}_{\mathbb{R}})) \text{ et } : \mathfrak{M}^r = \mathfrak{M}((\mathbb{R}^r, \mathcal{B}_{\mathbb{R}^r}); (\mathbb{R}, \mathcal{B}_{\mathbb{R}})).$$

Définition 2.1.5 : (distance associée à une mesure :)

De toute mesure **finie** μ sur $(\mathbb{R}^r, \mathcal{B}_{\mathbb{R}^r})$, provient une métrique ρ_μ sur M^r définie par :

$$\begin{aligned} \rho_\mu & : M^r \times M^r \longrightarrow \mathbb{R}^+ \\ (f, g) & \longmapsto \inf \left\{ \varepsilon > 0 / \mu \{x \in \mathbb{R}^r / |g(x) - f(x)| > \varepsilon\} < \varepsilon \right\}. \end{aligned}$$

Remarque 2.1.6 :

1. On peut justifier facilement le fait que ρ_μ est une distance sur M^r , mais faites attention!, elle ne l'est pas du tout sur : \mathfrak{M}^r , car la propriété de séparation ne devient plus valable ^a.
2. l'hypothèse de la finitude de μ est très importante, car elle nous permet de justifier l'existence de l'image :

or si on note : $(\forall (f, g) \in M^r) : \mathcal{A}_{f,g} = \left\{ \varepsilon > 0 / \mu \{x \in \mathbb{R}^r / |g(x) - f(x)| > \varepsilon\} < \varepsilon \right\} \neq \emptyset$

on aura : $\mu \{x \in \mathbb{R}^r / |g(x) - f(x)| > \varepsilon\} \leq \mu(\mathbb{R}^r) \in \mathbb{R}^+ \implies \mu(\mathbb{R}^r) + 1 \in \mathcal{A}_{f,g} \implies \mathcal{A}_{f,g} \neq \emptyset$
 $\implies \rho_\mu(f, g) = \inf \mathcal{A}_{f,g}$ existe.

2.2 Résultats fondamentaux :

Théorème 2.2.1 :

Soit σ une fonction **continue non constante quelconque**, alors et comme ça était mentionné dans [34], l'ensemble $\Sigma \Pi^r(\sigma)$ est uniformément compactement dense dans $C(\mathbb{R}^r)$.

a. Dans ce cas : $\rho_\mu(f, g) = 0 \iff f = g$ μ -presque partout.

Preuve :

Nous appliquons pour démontrer ceci le théorème de Stone-Weierstrass I.2.2.1. soit $K \subset \mathbb{R}^r$ un compact^a, pour toute fonction $\sigma, \Sigma\Pi^r(\sigma)$ est clairement une sous algèbre dans $C^r(K)$, et si on prend $(x, y) \in K$ avec : $x \neq y$, alors on trouvera forcément un $A \in \mathcal{A}^r$ tel que : $\sigma(A(x)) \neq \sigma(A(y))$, pour voir ceci, choisissez $a, b \in \mathbb{R}$ avec $a \neq b$ et : $\sigma(a) \neq \sigma(b)$. (de tels éléments existent car **la fonction σ est supposée non constante**) et prenez un $A(\cdot) \in \mathcal{A}^r$ satisfaisant : $A(x) = a, A(y) = b$. alors : $\sigma(A(x)) \neq \sigma(A(y))$, ce qui garantit que : $\Sigma\Pi^r(\sigma)$ sépare les points de K . d'une autre part, il y en a des fonctions $\sigma(A(\cdot))$ qui sont constantes et non égales à zéro, pour voir cela choisissez $b \in \mathbb{R}$ tel que : $\sigma(b) \neq 0$ et définir $\sigma(A(x)) = \langle 0, x \rangle + b$, donc : $(\forall x \in K) : \sigma(A(x)) = \sigma(b) \neq 0$. alors toutes les hypothèses du théorème de Stone-Weierstrass I.2.2.1 sont vérifiées, d'où ce qu'il faut prouver. ■

Remarques 2.2.1 :

1. En d'autres termes, ce théorème énonce que les réseaux de neurones $\Sigma\Pi$ sont capables d'approximer avec une précision arbitraire toute fonction continue à valeurs réelles sur un ensemble compact.
2. l'exigence de la compacité apparait chaque fois où les valeurs possibles des entrées x sont bornées.
3. une caractéristique intéressante de ce résultat est que la fonction d'activation σ peut être n'importe quelle fonction continue non constante.
4. nos résultats ultérieurs découleront tous de ce théorème.

Lemme 2.2.1 : (caractérisation de la ρ_μ -convergence :

Soit μ **une mesure de probabilité** sur $(\mathbb{R}^r, \mathcal{B}_{\mathbb{R}^r})$, et soit ρ_μ sa distance associée, (voir la définition : III.2.1.5, alors pour toute suite : $(f_n)_{n \in \mathbb{N}} \in (M^r)^\mathbb{N}$ et pour toute fonction : $f \in M^r$, les assertions suivantes sont équivalentes (voir : [34]) :

- a) $\rho_\mu(f_n, f) \xrightarrow{n \rightarrow +\infty} 0$.
- b) $(\forall \varepsilon > 0) : \mu \left\{ x \in \mathbb{R}^r / |f_n(x) - f(x)| > \varepsilon \right\} \xrightarrow{n \rightarrow +\infty} 0$.
- c) $\int_{\mathbb{R}^r} \min \left\{ |f_n(x) - f(x)|, 1 \right\} d\mu \xrightarrow{n \rightarrow +\infty} 0$.

Preuve :

a) \iff b) : Évidente, il suffit d'écrire la définition de la limite pour les deux assertions et ça ira de soi.

b) \implies c) : si : $(\forall \varepsilon > 0) : \mu \left\{ x \in \mathbb{R}^r / |f_n(x) - f(x)| > \varepsilon \right\} \xrightarrow{n \rightarrow +\infty} 0$ alors : à partir d'un certain rang on aura : $\mu \left\{ x \in \mathbb{R}^r / |f_n(x) - f(x)| > \frac{\varepsilon}{2} \right\} < \frac{\varepsilon}{2}$ donc :

$$\int_{\mathbb{R}^r} \min \left\{ |f_n - f|, 1 \right\} d\mu < \int_{\mathbb{R}^r} |f_n - f| d\mu < \overbrace{\int_{\mathbb{R}^r} |f_n - f| \times \mathbb{1}_{\left\{ |f_n - f| \leq \frac{\varepsilon}{2} \right\}} d\mu}^{< \frac{\varepsilon}{2}} + \overbrace{\int_{\mathbb{R}^r} |f_n - f| \times \mathbb{1}_{\left\{ |f_n - f| > \frac{\varepsilon}{2} \right\}} d\mu}^{< \frac{\varepsilon}{2}} < \varepsilon.$$

c) \implies b) : on obtient cette implication directement de la proposition : III.1.4.1. ■

a. selon le théorème de Borel-Lebesgue, les compacts dans \mathbb{R}^r sont exactement les fermés bornés.

Remarque 2.2.2 :

On peut remplacer dans le point "b)" du lemme précédent la différence en valeur absolue par n'importe quelle distance générant la topologie Euclidienne dans \mathbb{R}^r , ainsi, la quantité intégrée au point "c)" peut être substitué par une métrique bornée sur \mathbb{R}^r .

Lemme 2.2.2 : (La convergence uniforme compacte \Rightarrow la ρ_μ -convergence :

Soit μ une mesure de probabilité sur $(\mathbb{R}^r, \mathfrak{B}_{\mathbb{R}^r})$ et soit $(f_n)_{n \in \mathbb{N}} \in (M^r)^\mathbb{N}$ une suite de fonctions mesurables qui converge uniformément compactement vers une fonction $f \in M^r$, alors : $\rho_\mu(f_n, f) \xrightarrow{n \rightarrow +\infty} 0$.

Preuve :

Choisissons un $\varepsilon > 0$, d'après le lemme III.2.2.1, et pour qu'on puisse déduire que : $\rho_\mu(f_n, f) \xrightarrow{n \rightarrow \infty} 0$ il suffit de trouver un $N_\varepsilon \in \mathbb{N}$ tel que : pour tout $n \geq N_\varepsilon$ on a : $\int_{\mathbb{R}^r} \min\{|f_n(x) - f(x)|, 1\} d\mu < \varepsilon$.

soit $n \in \mathbb{N}$, on note : $B_n = B(0, n] = \left\{ x \in \mathbb{R}^r / \sqrt{\sum_{i=1}^r x_i^2} \leq n \right\} \subset \mathbb{R}^r$ la boule **fermée** de centre 0 et de rayon n pour la norme Euclidienne dans \mathbb{R}^r .

il est facile de voir que : $\mathbb{R}^r = \bigcup_{i \in \mathbb{N}} B_n$, donc : $\mu(\mathbb{R}^r) = \mu\left(\bigcup_{i \in \mathbb{N}} B_n\right)$, et d'une autre part, la suite $(B_n)_{n \in \mathbb{N}}$ est clairement croissante car : $(\forall n \in \mathbb{N}) : B_n \subset B_{n+1}$, aussi, et sans aucune difficulté, les $(B_n)_{n \in \mathbb{N}}$ sont mesurables car : B_n fermé $\Rightarrow B_n^c \in \Theta_{\mathbb{R}^r} \Rightarrow B_n^c \in \tau(\Theta_{\mathbb{R}^r}) := \mathfrak{B}_{\mathbb{R}^r} \Rightarrow B_n \in \mathfrak{B}_{\mathbb{R}^r}$, donc de tout cela **et selon la continuité croissante de la mesure**, on obtient : $\mu(\mathbb{R}^r) = \mu\left(\bigcup_{i \in \mathbb{N}} B_n\right) = \lim_{n \rightarrow +\infty} \mu(B_n)$.

ce qui se traduit comme suit :

$$(\forall \varepsilon > 0) (\exists N_\varepsilon \in \mathbb{R}^+) \text{ tq : } (\forall n \in \mathbb{N}) : n \geq N_\varepsilon \Rightarrow |\mu(B_n) - 1| = 1 - \mu(B_n) \leq \varepsilon.$$

et ceci implique que : $(\forall \varepsilon > 0) (\exists \underbrace{M_\varepsilon}_{=[N_{\varepsilon/2}] + 1} \in \mathbb{N}) \text{ tq : } |\mu(B_{M_\varepsilon}) - 1| = 1 - \mu(B_{M_\varepsilon}) = \mu(\mathbb{R}^r \setminus B_{M_\varepsilon}) \leq \frac{\varepsilon}{2}$. ★

or, $(f_n)_{n \in \mathbb{N}}$ converge uniformément compactement vers $f \in C(\mathbb{R}^r)$, alors : pour tout compact $K \subset \mathbb{R}^r$ on a : $\sup_{x \in K} |f_n(x) - f(x)| \xrightarrow{n \rightarrow \infty} 0$ donc, et en particulier ^a : $\sup_{x \in B_{M_\varepsilon}} |f_n(x) - f(x)| \xrightarrow{n \rightarrow \infty} 0$

ce qui fait que : $(\forall \varepsilon' > 0) (\exists \eta_{\varepsilon'} \in \mathbb{R}^+) \text{ tq : } (\forall n \in \mathbb{N}) : n \geq \eta_{\varepsilon'} \Rightarrow \sup_{x \in B_{M_\varepsilon}} |f_n(x) - f(x)| \leq \varepsilon'$.

on prend en particulier $\varepsilon' = \frac{\varepsilon}{2}$, alors :

$$(\exists \eta_{\varepsilon'} \in \mathbb{R}^+) \text{ tq : } (\forall n \in \mathbb{N}) : n \geq \eta_{\varepsilon'} \Rightarrow \sup_{x \in B_{M_\varepsilon}} |f_n(x) - f(x)| \leq \varepsilon' = \frac{\varepsilon}{2}.$$

d'où :

$$\begin{aligned} (\forall n \in \mathbb{N}) : \int_{\mathbb{R}^r} \min\{|f_n(x) - f(x)|, 1\} d\mu &= \int_{\mathbb{R}^r \setminus B_{M_\varepsilon}} \min\{|f_n(x) - f(x)|, 1\} d\mu + \int_{B_{M_\varepsilon}} \min\{|f_n(x) - f(x)|, 1\} d\mu \\ &\leq \int_{\mathbb{R}^r \setminus B_{M_\varepsilon}} 1 d\mu + \int_{B_{M_\varepsilon}} |f_n(x) - f(x)| d\mu \\ &\leq \mu(\mathbb{R}^r \setminus B_{M_\varepsilon}) + \int_{B_{M_\varepsilon}} \sup_{x \in B_{M_\varepsilon}} |f_n(x) - f(x)| d\mu \\ &\leq \underbrace{\mu(\mathbb{R}^r \setminus B_{M_\varepsilon})}_{\leq \frac{\varepsilon}{2} \text{ selon : } \star} + \int_{\mathbb{R}^r} \sup_{x \in B_{M_\varepsilon}} |f_n(x) - f(x)| d\mu \leq \frac{\varepsilon}{2} + \frac{\varepsilon}{2} = \varepsilon \end{aligned}$$

a. la boule $B_{M_\varepsilon} = B(0, M_\varepsilon]$ est un fermé borné de \mathbb{R}^r donc, elle est compact selon le théorème de **Borel-Lebesgue**.

donc : $\int_{\mathbb{R}^r} \min \left\{ |f_n(x) - f(x)|, 1 \right\} d\mu \leq \varepsilon$, d'où ce qu'il faut prouver. ■

Lemme 2.2.3 : (la ρ_μ densité de $C(\mathbb{R}^r)$ dans M^r , [34])

Si μ est une mesure **finie** sur : $(\mathbb{R}^r, \mathfrak{B}_{\mathbb{R}^r})$, alors l'ensemble $C(\mathbb{R}^r)$ sera ρ_μ -dense dans M^r .

Preuve :

Pour montrer que l'ensemble $C(\mathbb{R}^r)$ est ρ_μ -dense dans **l'espace métrique** (M^r, ρ_μ) , **il faut et il suffit** de trouver pour n'importe quel fonction : $f \in M^r$, une certaine suite : $\Phi_{n \in \mathbb{N}} \in C(\mathbb{R}^r)$ tel que : $\Phi_n \xrightarrow[n \rightarrow +\infty]{\rho_\mu} f$, on rappelle que selon le lemme III.2.2.1, cette convergence revient à dire que :

$$\int_{\mathbb{R}^r} \min \left\{ |\Phi_n(x) - f(x)|, 1 \right\} d\mu \xrightarrow[n \rightarrow +\infty]{} 0$$

soit : $f \in M^r$, on pose : $(\forall n \in \mathbb{N}) : \varphi_n := f \times \mathbb{1}_{\{|f| \leq n\}} \in L^1(\mathbb{R}^r, \mathfrak{B}_{\mathbb{R}^r}, \mu)$, il est bien claire que :

$$\lim_{n \rightarrow +\infty} \int_{\mathbb{R}^r} \min \left\{ |\varphi_n - f|, 1 \right\} d\mu = 0^a$$

donc, pour tout : $\varepsilon > 0$ on peut trouver un N_ε suffisamment grand tel que :

$$(\forall n \geq N_\varepsilon) : \int_{\mathbb{R}^r} \min \left\{ |\varphi_n - f|, 1 \right\} d\mu < \frac{\varepsilon}{2}. \quad \star$$

D'une autre part, et selon la proposition I.3.1.1 μ sera régulière (voir la définition I.3.1.1) ce qui fait que $(\mathbb{R}^r, \mathfrak{B}_{\mathbb{R}^r}, \mu)$ sera en particulier **un espace de Radon**, alors et par le théorème de Lusin généralisé I.3.2.1 qui nous garantit la possibilité de restreindre toutes les fonctions mesurables sur un espace de Radon à des autres qui sont continues sur une partie **fermée**, aussi grande que voulue, on peut trouver forcément une fonction continue Φ_n tel que : $\varphi_n = \Phi_n$ sur un fermé F qui vérifie : $\mu(F) < 1 - \frac{\varepsilon}{2}$. et dans ce cas on obtient pour tout $n \geq N_\varepsilon$:

$$\begin{aligned} \int_{\mathbb{R}^r} \min \left\{ |\varphi_n - \Phi_n|, 1 \right\} d\mu &= \int_F \min \left\{ |\varphi_n - \Phi_n|, 1 \right\} d\mu + \int_{F^c} \min \left\{ |\varphi_n - \Phi_n|, 1 \right\} d\mu \\ &< \int_{F^c} 1 d\mu < \mu(F^c) < 1 - \mu(F) = 1 - \left(1 - \frac{\varepsilon}{2}\right) = \frac{\varepsilon}{2} \quad \star\star\star \end{aligned}$$

alors et par une simple application de l'inégalité triangulaire on obtient pour tout $n \geq N_\varepsilon$:

$$\int_{\mathbb{R}^r} \min \left\{ |f - \Phi_n|, 1 \right\} d\mu \leq \int_{\mathbb{R}^r} \min \left\{ |f - \varphi_n|, 1 \right\} d\mu + \int_{\mathbb{R}^r} \min \left\{ |\varphi_n - \Phi_n|, 1 \right\} d\mu < \frac{\varepsilon}{2} + \frac{\varepsilon}{2} = \varepsilon$$

d'où ce qu'il faut établir. ■

Théorème 2.2.2 :

Pour toute fonction **continue non constante** σ , tout $r \in \mathbb{N}^*$, et toute **mesure de probabilité** μ sur (X, \mathfrak{B}_X) , l'ensemble $\Sigma\Pi^r(\sigma)$ est ρ_μ -dense dans M^r .

Autrement dit, et en adoptant la même formulation que [34], les réseaux de neurones $\Sigma\Pi$ à une seule couche cachée et à propagation en avant, peuvent bien approcher n'importe quelle fonction mesurable, indépendamment de la fonction continue non constante σ utilisée, de la dimension de l'espace d'entrée r , et d'environnement spatial d'entrée^b μ , en l'occurrence les réseaux $\Sigma\Pi$ sont des approximateurs universels.

a. Cette convergence peut être justifiée géométriquement par le fait que lorsque n devient plus grand les courbes de φ_n et f deviennent plus confondues.

b. La mesure adoptée pour l'espace d'entrée s'appelle l'environnement spatial d'entrée, cette terminologie était introduite pour la première fois en 1988 par Halbert White.

Preuve :

Prenons une fonction continue non constante σ , selon le théorème III.2.2.1 l'ensemble $\Sigma\Pi^r(\sigma)$ est uniformément compactement dense dans $C(\mathbb{R}^r)$, en outre, notre mesure μ est de probabilité donc on peut déduire en vertu du lemme III.2.2.2 que $\Sigma\Pi^r(\sigma)$ est ρ_μ -dense dans M^r et en conséquence $\Sigma\Pi^r(\sigma)$ sera ρ_μ -dense dans M^r , ce qui achève bien notre preuve. ■

Remarque 2.2.3 :

On doit signaler quand même, que la dernière partie de la preuve précédente a été basée sur un résultat simple mais très important de la topologie classique et qui s'énonce ainsi :
Soit $A \subset B \subset C$ trois parties emboîtées d'un espace métrique (X, d) si A est dense dans B et B est dense dans C . alors A sera dense dans C .

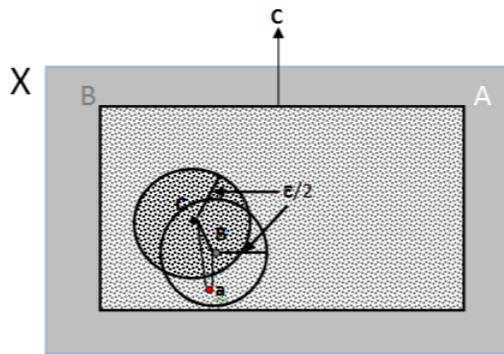


FIGURE III.6: Indication de la démonstration .

La preuve de ce résultat est très facile et elle découle essentiellement de l'inégalité triangulaire, on laissera au lecteur le plaisir de la démontrer avec détails en lui donnant comme indication la figure : III.6 .

Lemme 2.2.4 : (l'approximation des fonctions d'écrasement par des Σ réseaux de neurones :)

Soit $F : \mathbb{R} \rightarrow [0, 1]$ une fonction d'écrasement continue, et $\Psi : \mathbb{R} \rightarrow [0, 1]$ une fonction d'écrasement quelconque alors on a :

$$(\forall \varepsilon > 0) \left(\exists H_\varepsilon \in \Sigma^1(\Psi) \right) \text{ tel que : } \sup_{\lambda \in \mathbb{R}} |F(\lambda) - H_\varepsilon(\lambda)| < \varepsilon .$$

Preuve :

Soit $\varepsilon > 0$, sans perte de généralité on peut considérer que : $\varepsilon < 1$, notre but est de trouver un entier $Q \in \mathbb{N}$, une collection finie des réels, $(\beta_1, \dots, \beta_Q) \in \mathbb{R}^Q$, et des fonctions affines : $(A_1, \dots, A_Q) \in (\mathcal{A}^1)^Q$ tel que :

$$\sup_{\lambda \in \mathbb{R}} \left| F(\lambda) - \sum_{j=1}^Q \beta_j \Psi(A_j(\lambda)) \right| < \varepsilon .$$

prenons pour cela un entier $Q \in \mathbb{N}$ de sorte que : $\frac{1}{1+Q} < \frac{\varepsilon}{2}$, et pour tout $j \in \{1, \dots, Q\}$ posons $\beta_j = \frac{1}{Q+1}$, on choisit un $M > 0$ tel que : $\Psi(-M) < \frac{\varepsilon}{2Q}$, et $\Psi(M) > 1 - \frac{\varepsilon}{2Q}$, (un tel M existe toujours car : $\lim_{\lambda \rightarrow +\infty} \Psi(\lambda) = 1$ et $\lim_{\lambda \rightarrow -\infty} \Psi(\lambda) = 0$), et on considère ensuite le $Q+1$ -uplet : (r_1, \dots, r_{Q+1})

défini par :

$$\begin{cases} r_j = \sup \left\{ \lambda \in \mathbb{R} / F(\lambda) = \frac{j}{Q+1} \right\} & \text{si : } j \in \{1, \dots, Q\}. \\ r_j = \sup \left\{ \lambda \in \mathbb{R} / F(\lambda) = 1 - \frac{1}{2(Q+1)} \right\} & \text{si : } j = Q+1. \end{cases}$$

► **Note explicative** : il est indispensable de mentionner que les deux bornes supérieures précédentes ont bien un sens et ceci se justifie comme suit :

F est continue avec : $\lim_{\lambda \rightarrow -\infty} F(\lambda) = 0$ et $\lim_{\lambda \rightarrow +\infty} F(\lambda) = 1$, donc et **selon le théorème des valeurs intermédiaires** :

$$(\forall \delta \in]0, 1[) (\exists \lambda_\delta \in \mathbb{R}) \text{ tq : } F(\lambda_\delta) = \delta.$$

en posant en particulier $\delta = \frac{j}{Q+1}$ pour tout $j \in \llbracket 1, Q \rrbracket$, on pourra affirmer facilement que :

$$\left\{ \lambda \in \mathbb{R} / F(\lambda) = \frac{j}{Q+1} \right\} \neq \emptyset.$$

d'une autre part : $(\forall j \in \llbracket 1, Q \rrbracket) : \left\{ \lambda \in \mathbb{R} / F(\lambda) = \frac{j}{Q+1} \right\}$ est borné. car s'il ne l'est pas on aura :

$$(\forall n \in \mathbb{N}) (\exists \alpha_n \geq n) \text{ tel que : } F(\alpha_n) = \frac{j}{Q+1}.$$

et ceci va impliquer que : $\lim_{n \rightarrow +\infty} \alpha_n = +\infty$ et $\lim_{n \rightarrow +\infty} F(\alpha_n) = \frac{j}{Q+1} < 1$, ce qui est absurde car $\lim_{\lambda \rightarrow +\infty} F(\lambda) = 1$.

pour l'ensemble $\left\{ \lambda \in \mathbb{R} / F(\lambda) = 1 - \frac{1}{2(Q+1)} \right\}$ la bornitude et la non videsse^a peuvent être établies d'une façon similaire. □

pour tout $\theta < \theta'$ notons : $A_{\theta, \theta'} \in \mathcal{A}^1$ l'unique forme affine vérifiant : $A_{\theta, \theta'}(\theta) = M$ et $A_{\theta, \theta'}(\theta') = -M$.

l'approximation désirée est donc : $H_\varepsilon(\lambda) = \sum_{j=1}^Q \beta_j \Psi(A_{\theta, \theta'}(\lambda))$ puisqu'on peut voir facilement que

sur les intervalles : $] -\infty, r_1],] r_1, r_2], \dots,] r_Q, +\infty [$ on a : $|F(\lambda) - H_\varepsilon| < \varepsilon$ ^b ■

Théorème 2.2.3 :

Pour toute fonction **d'écrasement** Ψ , tout $r \in \mathbb{N}^*$, et toute **mesure de probabilité** μ sur (X, \mathfrak{B}_X) , l'ensemble $\Sigma \Pi^r(\Psi)$ est uniformément compactement dense dans $C(\mathbb{R}^r)$ et au même temps il est ρ_μ -dense dans M^r .

c'est à dire que les théorèmes :III.2.2.1 et III.2.2.2 **restent valables pour les fonctions d'écrasement**,^c et c'est exactement ce qui était établi dans l'article : [34].

Preuve :

Nous allons démontrer en premier lieu que : l'ensemble $\Sigma \Pi^r(\Psi)$ est uniformément compactement dense dans $C(\mathbb{R}^r)$, et le fait de la ρ_μ densité dans M^r sera déduite immédiatement via les lemmes :III.2.2.2, III.2.2.3 et la remarque III.2.2.3.

rappelons que selon le théorème III.2.2.1, on a : l'ensemble $\Sigma \Pi^r(\sigma)$ est uniformément compactement dense dans $C(\mathbb{R}^r)$ pour toute fonction **continue non constante** σ , alors si on établit que : $\Sigma \Pi^r(\Psi)$ est uniformément compactement dense dans $\Sigma \Pi^r(\sigma)$ pour une certaine fonction continue non constante σ , le résultat désiré sera obtenu immédiatement grâce à la remarque III.2.2.3.

prenons $\varepsilon > 0$ et $\ell \in \mathbb{N}^*$, car la multiplication est continue^d et $[0, 1]^\ell$ est compact, on pourra dé-

a. Ce terme a été introduit en 2015 pour remplacer le mot vuidesse qui est devenu moins utilisé et presque oublié dans le vocabulaire et le lexique français voir pour cela le **dictionnaire DMF** dans le site officiel du centre national de ressources textuelles et lexicales : www.cnrtl.fr.

b. Ça découle des inégalité précédentes, de la croissance des fonctions F et Ψ et la décroissance de la forme affine $A_{\theta, \theta'}$.

c. On rappelle que les fonctions d'écrasement ne sont pas forcément continues.

d. C'est à dire que : $(x_1, \dots, x_\ell) \mapsto \prod_{k=1}^{\ell} x_k$ est une fonction continue.

duire selon le fameux théorème de Heine ^a, qu'elle est uniformément continue sur l'intervalle $[0, 1]^\ell$, ce qui fait qu'il existe un $\delta_\varepsilon > 0$ tel que :

$$\left(\forall \left((a_1, \dots, a_\ell), (b_1, \dots, b_\ell) \right) \in \left([0, 1]^\ell \right)^2 \right) : \left[\left(\forall k \in \{1, \dots, \ell\} \right) : |a_k - b_k| < \delta_\varepsilon \right] \Rightarrow \left| \prod_{k=1}^{\ell} a_k - \prod_{k=1}^{\ell} b_k \right| < \varepsilon.$$

par le lemme : III.1.3.4, on peut trouver une fonction : $H_{\delta_\varepsilon}(\cdot) = \sum_{t=1}^T \beta_t \Psi(A_t(\cdot))$ tel que : $\sup_{\lambda \in \mathbb{R}^r} |F(\lambda) - H_{\delta_\varepsilon}(\lambda)| < \delta_\varepsilon$

ce qui fait que :

$$\sup_{x \in \mathbb{R}} \left| \prod_{k=1}^{\ell} \sigma(A_k(x)) - \prod_{k=1}^{\ell} H_{\delta_\varepsilon}(A_k(x)) \right| = \sup_{x \in \mathbb{R}} \left| \prod_{k=1}^{\ell} \sigma(A_k(x)) - \prod_{k=1}^{\ell} \sum_{t=1}^T \beta_t \Psi(A_t(A_k(x))) \right| < \varepsilon.$$

et car : $A_t(A_k(x)) \in \mathcal{A}^r$, on peut déduire que : $\prod_{k=1}^{\ell} H_{\delta_\varepsilon}(A_k(\cdot)) \in \Sigma\Pi^r(\Psi)$. ^b

donc la fonction : $\prod_{k=1}^{\ell} \sigma(A_k(x))$ peut être uniformément approchée par des éléments de $\Sigma\Pi^r(\Psi)$, et selon la remarque III.2.1.4 ceci implique qu'elle peut être uniformément compactement approchée par des éléments de $\Sigma\Pi^r(\Psi)$ donc, et sans aucun problème, $\Sigma\Pi^r(\Psi)$ sera uniformément compactement dense dans $\Sigma\Pi^r(\sigma)$ d'où ce qu'il faut établir. ■

Lemme 2.2.5 :

Pour toute fonction d'écrasement Ψ , tout $\varepsilon > 0$, et tout $M > 0$, il existe une fonction $\varphi_{M,\varepsilon} \in \Sigma^1(\Psi)$ tel que : $\sup_{\lambda \in [-M, M]} |\varphi_{M,\varepsilon}(\lambda) - \cos(\lambda)| < \varepsilon$.

Preuve :

Soit Γ la fonction du cosinus écraseur de Gallant-White ^c, on peut remarquer que :

$$2\Gamma\left(x - \frac{3\pi}{2}\right) - 1 = \cos(x) \times \mathbb{1}_{[\pi, 2\pi]} - \mathbb{1}_{]-\infty, \pi[} + \mathbb{1}_{]2\pi, +\infty[}.$$

Donc en ajoutant, soustrayant et multipliant par un scalaire un nombre fini de versions translatées

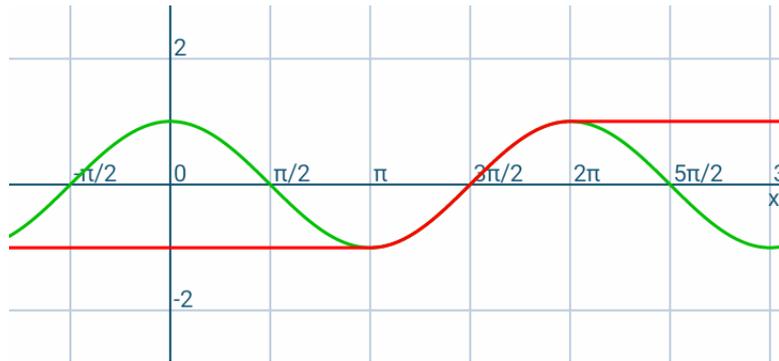


FIGURE III.7: La courbe de la fonction : $x \mapsto 2\Gamma\left(x - \frac{3\pi}{2}\right) - 1$.

de Γ on **peut obtenir une fonction** $\varphi_{M,\varepsilon} \in \Sigma^1(\Psi)$ qui se restreint à la fonction cosinus sur n'importe quel intervalle $[k\pi, k'\pi]$, et par la suite sur n'importe quel intervalle $[-M, M]$ ^d.

Le résultat découle maintenant du lemme III.2.2.3 et de **l'inégalité triangulaire**. ■

a. Un théorème démontré par Eduard Heine en 1872 et qui énonce que toute application continue d'un **espace métrique compact** à valeurs dans un espace métrique quelconque est uniformément continue.

b. Il suffit de développer l'expression de notre fonction pour rendre ce fait plus clair.

c. voir le troisième point de l'exemple : III.2.1.1.

d. il suffit de plonger cet intervalle dans un autre de la forme : $[k\pi, k'\pi]$.

Lemme 2.2.6 :

Soit $g : \mathbb{R}^r \rightarrow \mathbb{R}$ une certaine fonction numérique définie par $g : x \mapsto \sum_{j=1}^N \beta_j \cos(A_j(x))$, avec : $N \in \mathbb{N}^*$, $(\beta_1, \dots, \beta_N) \in \mathbb{R}^N$ et $A_j \in \mathcal{A}^r$, alors, pour toute fonction d'écrasement Ψ , tout compact $K \subset \mathbb{R}^r$, et tout $\varepsilon > 0$, il existe un $f \in \Sigma^r(\Psi)$ tel que : $\sup_{x \in K} |g(x) - f(x)| < \varepsilon$.

Preuve :

Prenons un $M > 0$ tel que pour tout $j \in \{1, \dots, N\}$ $A_j(K) \subset [-M, M]$, (car : $N < +\infty$, K est compact, et les A_j sont continues un tel M existe toujours, et ceci revient au fait qu'une fonction continue sur un compact est toujours bornée et atteint ses bornes).

on pose : $\varepsilon' = \frac{\varepsilon}{\sum_{j=1}^N |\beta_j|}$, par le lemme III.2.2.5, il existe une fonction $\varphi_{M, \varepsilon'} \in \Sigma^1(\Psi)$ tel que :

$$\sup_{\lambda \in [-M, M]} |\varphi_{M, \varepsilon'}(\lambda) - \cos(\lambda)| < \varepsilon'.$$

ce qui fait que : $(\forall j \in \{1, \dots, N\}) (\forall x \in K) : |\varphi_{M, \varepsilon'}(A_j(x)) - \cos(A_j(x))| < \frac{\varepsilon}{N}$.

$$\begin{aligned} \text{donc : } (\forall x \in K) : \left| \sum_{j=1}^N \beta_j \varphi_{M, \varepsilon'}(A_j(x)) - g(x) \right| &= \left| \sum_{j=1}^N \beta_j \varphi_{M, \varepsilon'}(A_j(x)) - \sum_{j=1}^N \beta_j \cos(A_j(x)) \right| \\ &= \left| \sum_{j=1}^N \beta_j \left(\varphi_{M, \varepsilon'}(A_j(x)) - \cos(A_j(x)) \right) \right| \\ &\leq \sum_{j=1}^N |\beta_j| \times \left| \varphi_{M, \varepsilon'}(A_j(x)) - \cos(A_j(x)) \right| < \sum_{j=1}^N |\beta_j| \times \varepsilon' = \varepsilon. \end{aligned}$$

ce qui achève notre preuve. ■

Théorème 2.2.4 :

Pour toute fonction **d'écrasement** Ψ , tout $r \in \mathbb{N}^*$, et toute **mesure de probabilité** μ sur (X, \mathfrak{B}_X) , l'ensemble $\Sigma^r(\Psi)$ est uniformément compactement dense dans $C(\mathbb{R}^r)$ et au même temps il est ρ_μ -dense dans M^r .

c'est à dire que les théorèmes :III.2.2.1 et III.2.2.2 **restent valables pour les réseaux de neurones.**(voir :[34] et [33])

Preuve :

On note \mathfrak{T} l'ensemble des polynômes trigonométriques défini par :

$$\mathfrak{T} = \Sigma \Pi^r(\cos(\cdot)) = \left\{ x \mapsto G(x) = \sum_{j=1}^q \beta_j \cdot \prod_{k=1}^{l_j} \cos(A_{jk}(x)) \mid q \in \mathbb{N}^*, l_j \in \mathbb{N}, (\beta_1, \dots, \beta_q) \in \mathbb{R}^q, A_{jk}(x) \in \mathcal{A}^r \right\}.$$

Par le théorème III.2.2.1 cet ensemble est uniformément compactement dense dans : $C(\mathbb{R}^r)$ ^a donc et en appliquant plusieurs fois l'identité trigonométrique ^b :

$$(\forall (a, b) \in \mathbb{R}^2) : \cos(a) \times \cos(b) = \frac{1}{2} (\cos(a+b) + \cos(a-b)).$$

nous pouvons réécrire les éléments de \mathfrak{T} sous la forme : $\sum_{i=1}^T \alpha_i \cos(A_i(\cdot))$ où : $\alpha_i \in \mathbb{R}$ et $A_i \in \mathcal{A}^r$.

maintenant, la densité uniforme compacte de \mathfrak{T} dans $C(\mathbb{R}^r)$ provient directement du lemme :III.2.2.6, et par la suite la ρ_μ densité de \mathfrak{T} dans M^r découlera immédiatement des lemmes :III.2.2.1, III.2.2.2 et la remarque III.2.2.2 ■

a. Car la fonction $x \mapsto \cos(x)$ est continue et non constante sur \mathbb{R} .

b. Cette formule symbolise un peu la nostalgie de l'époque de lycée, n'est-ce pas cher lecteur?!

3 Théorèmes d'approximation de Funahashi :

3.1 Mise en situation :

On complétera notre liste de théorèmes d'approximation par cette section dans laquelle on va présenter, et d'une façon un peu résumée, les résultats principaux qui ont été proposés par le mathématicien japonais Ken-Ichi Funahashi (1952, ...), la plupart des énoncés de ces résultats ne vont pas surprendre notre lecteur, car elles ne sont que des cas particuliers et des applications directes des résultats cités précédemment, mais l'ajout de cette section, vient comme une sorte d'appréciation, et de gratitude envers ce mathématicien qui se considère parmi les pionniers et les fondateurs de la théorie d'approximation, et à qui revient l'originalité des idées de la plupart des théorèmes traités auparavant.

3.2 Résultats principaux :

Théorème 3.2.1 : (Premier théorème de Funahashi, [21]) :

Soit $\phi(x)$ une fonction non constante, bornée, continue, et croissante. Soit $K \subset \mathbb{R}^n$ un compact et $f(x_1, \dots, x_n)$ une fonction réelle continue sur K , alors pour tout $\varepsilon > 0$, il existe un entier $N_{\varepsilon, K}$ des vecteurs réels : $(c_1, \dots, c_n) \in \mathbb{R}^n$ et $(\theta_1, \dots, \theta_n) \in \mathbb{R}^n$, et une matrice $(\omega_{i,j})_{\substack{1 \leq i \leq n \\ 1 \leq j \leq N}}$ tel que :

$$\sup_{x \in K} \left| \sum_{j=1}^N c_j \phi \left(\sum_{i=1}^n \omega_{i,j} x_i - \theta_i \right) - f \right| < \varepsilon .^a$$

Preuve :

La fonction ϕ est bornée et croissante donc on peut déduire que :

$$(\forall x \in \mathbb{R}) : \lim_{n \rightarrow -\infty} \phi(x) := m \leq \phi(x) \leq M := \lim_{n \rightarrow +\infty} \phi(x)$$

alors si on reprend les théorèmes précédents pour ce type de fonctions à la place des fonctions d'écrasement on obtiendra un résultat exactement analogue au théorème :III.2.2.4 et de cette façon le résultat sera aisément prouvé. ■

Remarque 3.2.1 :

On peut remarquer que l'hypothèse de la continuité de ϕ n'a pas été évoquée dans la dernière démonstration, et ceci revient au fait que **le résultat de Hornik-White-Stinchcombe pour les Σ réseaux de neurones est très fort, et il généralise celui énoncé dans ce dernier théorème**, mais faites attention!, avant les travaux de ces premiers, cette Hypothèse a joué un grand rôle dans la preuve proposée par Funahashi et qui était basée essentiellement sur le théorème de Paley-Wiener et la formule d'Irie-Miyake.

Théorème 3.2.2 : (Deuxième théorème de Funahashi, [21]) :

Pour n'importe quelle fonction $g \in M^n$, et pour tout $\varepsilon > 0$ il existe un compact $K \subset \mathbb{R}^n$ et une fonction $f \in \Sigma(\Psi)$ tel que : $\mu(K) = 1 - \varepsilon$ et : $\sup_{x \in K} |f(x) - g(x)| < \varepsilon$, et ceci indépendamment du choix de la fonction d'écrasement Ψ , de l'entier n , et de la mesure de probabilité μ .

a. bien sûr les c_i , les θ_j et les $\omega_{i,j}$ dépendent du choix de ε et K .

Preuve :

Fixons un $\varepsilon > 0$, selon le théorème de Lusin I.1.3.1, il existe un certain compact $K^1 \subset \mathbb{R}$ tel que : $\mu(\mathbb{K}_1) > 1 - \frac{\varepsilon}{2}$ et : $g|_{\mathbb{K}_1}$ est continue sur \mathbb{K}_1 , et selon le théorème de Tietze I.4.2.1 on peut prolonger $g|_{\mathbb{K}_1}$ par une certaine fonction continue $g' \in C(\mathbb{R}^r)$, ce qui fait qu'on aura :

$$\sup_{x \in \mathbb{K}_1} g'(x) = \sup_{x \in \mathbb{K}_1} g(x).$$

par le théorème : III.2.2.4 $\Sigma^r(\Psi)$ est uniformément compactement dense dans $C(\mathbb{R}^r)$, donc si on choisi un compact $\mathbb{K}_2 \subset \mathbb{R}$ tel que : $\mu(\mathbb{K}) = 1 - \frac{\varepsilon}{2}$, on trouvera forcément une fonction $f \in \Sigma(\Psi)$ qui vérifie :

$$\sup_{x \in \mathbb{K}_2} |g(x) - f(x)| < \varepsilon.$$

donc : $\sup_{x \in \mathbb{K} \cap \mathbb{K}_1} |g(x) - f(x)| \leq \sup_{x \in \mathbb{K}_2} |g(x) - f(x)| < \varepsilon$, et d'une autre part :

$$\mu(\mathbb{K}_1) > 1 - \frac{\varepsilon}{2} \text{ et } \mu(\mathbb{K}_2) > 1 - \frac{\varepsilon}{2} \Rightarrow \mu(\mathbb{K}_1^c) < \frac{\varepsilon}{2} \text{ et } \mu(\mathbb{K}_2^c) < \frac{\varepsilon}{2} \Rightarrow \mu(\mathbb{K}_1^c \cup \mathbb{K}_2^c) \leq \mu(\mathbb{K}_2^c) + \mu(\mathbb{K}_1^c) < \frac{\varepsilon}{2} + \frac{\varepsilon}{2}$$

$$\text{alors : } \mu(\mathbb{K}_1^c \cup \mathbb{K}_2^c) = \mu((\mathbb{K}_1 \cap \mathbb{K}_2)^c) < \varepsilon \Rightarrow \mu(\mathbb{K}_1 \cap \mathbb{K}_2) > 1 - \varepsilon.$$

d'où le résultat. ■

Remarque 3.2.2 :

Ce dernier théorème nous a présenté un résultat très important et qui reflète la puissance des réseaux de neurones d'approcher n'importe quelle application mesurable en restant dans un cadre d'erreur tolérable, mais il a exigé d'adopter pour les nœuds constituant notre réseau une fonction d'activation d'écrasement, une hypothèse que Cybenko a fait du bon travail pour l'affaiblir et de considérer juste des fonctions sigmoïdales mais et à son tour il était obligé d'ajouter la condition de continuité et d'affaiblir le type de la fonction à approcher et la rendre juste une fonction de décision et c'est ce qu'on voit au théorème : III.1.3.2 .

4 Propriété de parcimonie chez les réseaux de neurones :

4.1 Mise en contexte :

Dans la pratique, le nombre de paramètres nécessaires pour réaliser une approximation est un critère très important, ce qui fait que dans la recherche d'un approximateur, le concepteur du modèle doit toujours prendre ce critère en considération, et ceci dans l'objectif de **construire ce modèle d'une sorte que le nombre de paramètres ajustables soit le plus faible possible**, on dit dans ce cas qu'on cherche à faire l'approximation la plus parcimonieuse .

Après la justification de leur capacité approximative, dans cette partie nous allons montrer, que certains types de réseaux neuronaux peuvent vraiment se présenter comme des **approximateurs convenables et parcimonieux**, ainsi on insistera une autre fois sur l'importance de la fonction d'activation adoptée, et qui sera dans ce cas aussi, un élément clé et décisif pour la réalisation de la parcimonie.

4.2 Résultats principaux :

Le mathématicien américain **Andrew.R.Barron** a prouvé en 1993 dans son article [4] que, l'approximation par des modèles qui dépendent, non linéairement, des paramètres ajustables, est beaucoup plus parcimonieuse que celle qui se fait par des modèles dépendants linéairement de ces paramètres et plus précisément on a l'énoncé suivant :

Théorème 4.2.1 :(Barron,[4] :)

Pour une précision donnée, le nombre de paramètres ajustables, croît exponentiellement par rapport au nombre de variables dans le cas des **approximateurs linéaires par rapport à leurs paramètres**, alors qu'il croît juste linéairement par rapport à ce nombre pour des approximateurs non linéaires. ce qui reflète la grande puissance de ces derniers.

Remarque 4.2.1 :

1. À la lumière du théorème précédent on peut remarquer que la parcimonie devient donc plus importante lorsque le nombre d'entrées du modèle est grand, autrement dit, si on a juste une ou deux entrées, on peut utiliser indifféremment un modèle linéaire comme les polynômes par exemple ou également un autre non linéaire (par rapport à ses paramètres) comme les réseaux de neurones artificiels.
2. exactement comme la propriété d'approximation universelle, la parcimonie n'est pas vérifiées pour tous les types des réseaux neuronaux, or, l'importance du choix de la fonction d'activation s'impose une autre fois dans ce cadre comme un critère critique qui peut améliorer la capacité du réseau si elle est bien choisie, ou l'entraver et le rendre par conséquent non parcimonieux, dans le cas contraire.
voir à titre d'exemple le réseau Madaline ^a à plusieurs entrées ci dessous, et qui est linéaire par rapport à ces paramètres.

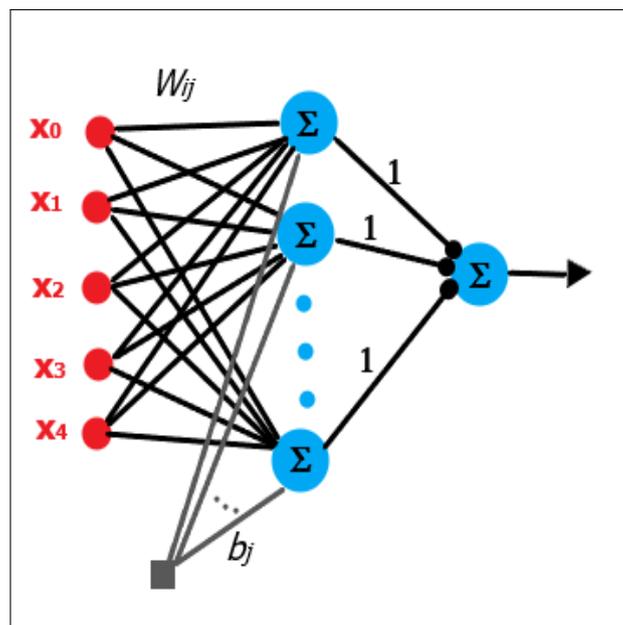


FIGURE III.8: Exemple d'un réseau non parcimonieux.

Théorème 4.2.2 :(Hornik,[35])

Pour les réseaux de neurones à **fonction d'activation sigmoïdale**, l'erreur commise dans l'approximation varie comme l'inverse du nombre de neurones cachés, et elle est donc **indépendante du nombre de variables de la fonction à approcher**.

Par conséquent, pour une précision fixée, et donc pour un nombre de neurones cachés donné, le nombre de paramètres du réseau est proportionnel au nombre de variables de la fonction à approcher.

a. Un Madaline (Many adaline) est tout simplement un réseau à trois couches constitué de plusieurs Adaline

Remarque 4.2.2 :

Ce résultat s'applique aux réseaux de neurones à fonction d'activation sigmoïdale puisque la sortie de ces neurones n'est pas linéaire par rapport aux poids synaptiques.

Cette propriété montre l'intérêt des réseaux de neurones sigmoïdes par rapport à d'autres approximateurs comme les polynômes dont la sortie est une fonction linéaire des paramètres ajustables, et pour mieux éclaircir les choses on peut dire que :

pour un même nombre d'entrées, le nombre de paramètres ajustables à déterminer est plus faible pour un réseau de neurones que pour un polynôme.

Cette propriété est indispensable dans le cas du filtrage des textes car le nombre d'entrées est typiquement très élevé.

Conclusion :

Les résultats de ce chapitre ont établi que les réseaux de neurones artificiels feed forward à une seule couche cachée, sont capables d'approcher et, au degré souhaité de précision, plusieurs types de fonctions, y compris les fonctions continues et mesurables, et on a vu que la capacité de ces réseaux, et d'un point de vue théorique, dépend essentiellement du sens désiré d'approximation et du type de la fonction d'activation adoptée, ce qui va nous permettre enfin de compte de déduire que ces réseaux sont des *approximateurs universels*, et d'affirmer que *tout manque de succès dans le cadre applicatif doit survenir forcément d'un apprentissage inadéquat, un nombre insuffisant d'unités cachées ou d'une absence complète de la relation déterministe entre l'entrée et la sortie désirée (cas aléatoire)*.

Après toutes ces affirmations qu'on a présenté jusqu'à l'instant dans l'espoir de justifier théoriquement la grande aptitude des réseaux de neurones artificiels, on mettra dans le chapitre suivant cette machine à l'épreuve et ceci dans le cadre précis de la résolution des équations différentiels fractionnaires.

Application : Résolution numérique des équations différentielles fractionnaires par les Réseaux de Neurones Artificiels :

Résumé :

Dans ce chapitre, nous allons commencer par un aperçu sur la dérivation fractionnaire, et à travers lequel on découvrira comment cette notion et qui était définie habituellement pour des ordres entiers, s'est étendue pour englober des ordres quelconques, ensuite on verra que cette généralisation sera à l'origine d'apparition d'un nouveau type d'équations différentielles, appelées les équations différentielles fractionnaires.

L'invalidité des schémas numériques classiques pour ce types d'équations va nous obliger de se mettre dans le cadre d'apprentissage automatique, dans lequel on montrera que *les réseaux de neurones artificiels* adoptant une fonction sigmoïdale en tant que fonction d'activation, et une architecture feed-forward à trois couches, sont *capables de générer des schémas itératifs pour résoudre les équations différentielles* linéaires d'ordre *fractionnaire*.

L'algorithme classique d'apprentissage du *rétro-propagation* et qui était basé sur la méthode de descente de gradient comme on a déjà vu au deuxième chapitre, sera, et après l'introduction de quelques modifications, capable d'*approcher, et au degré souhaité de précision, la solution* de notre problème et même sur n'importe quel intervalle d'étude donné.

Enfin de compte, et pour être plus précis, certains problèmes de test seront également présentés lors d'une étude illustrative des résultats numériques fournis par la méthode qu'on va proposer.

Mots clés : Réseaux de neurones, dérivation fractionnaire, équations différentielles d'ordre fractionnaire, algorithme de rétro-propagation.

Introduction :

Les équations différentielles sont des outils permettant d'étudier des phénomènes naturels dans un large éventail d'applications, et au fil des ans, des recherches scientifiques fondamentales ont conclu que celles qui se posent dans le cadre des problèmes du monde réel ont rarement des solutions analytiques exactes, par conséquent, nous étions obligés de les analyser numériquement via des techniques et des simulations sur l'ordinateur, ce qui va nous offrir un moyen puissant pour les résoudre plus facilement et rapidement.

En raison de leur complexité particulière, la classe des équations différentielles d'ordre fractionnaire s'est marquée par le comportement de ses solutions qui sont très difficile à trouver, à décrire et à comprendre, et pour cela on a proposé d'utiliser des modèles d'apprentissage pour estimer cette solution sur un domaine d'étude donné.

Dans ce sens le modèle le plus classique qui peut être proposé, sera les réseaux de neurones artificiels bien évidemment, la particularité de ces réseaux et comme on déjà vu au chapitre III réside dans le fait que chaque structure régulière peut être considérée comme un approximateur universel parcimonieux, ce qui fait, que ce modèle est approprié au problème mentionné précédemment et qu'il peut estimer avec une précision approximative très élevée la fonction inconnue sur le domaine dont on dispose.

L'objectif de ce chapitre est d'utiliser un réseau neuronal feed-forward à trois couches pour résoudre le problème de l'équation différentielle ordinaire fractionnaire (E.D.F.O), et afin d'arriver à ce but nous allons procéder comme suit :

dans un premier temps, on va discrétiser l'intervalle d'étude, bien sûr, il existe plusieurs façons pour le faire mais ici, nous utiliserons juste le schéma de discrétisation standard, ensuite l'équation différentielle d'origine sera reformulée dans le cas discret sous la forme d'un système d'équations dont les inconnus sont les images des points nodaux, et enfin ce système va être convertie en un problème d'optimisation qui sera résolu par l'auto-apprentissage des réseaux de neurones artificiels et qu'on le fera dans ce cadre par l'algorithme de rétropropagation du gradient qu'on a vu au deuxième chapitre.

1 Initiation à l'analyse fractionnaire :

L'*analyse fractionnaire* est une discipline des mathématiques qui s'intéresse à l'étude des différentes possibilités permettant de définir des puissances numériques réelles ^a pour l'opérateur de différenciation, et en d'autres termes, c'est une branche qui cherche à attribuer un sens à la dérivation selon un ordre qui n'est pas forcément entier.

La première apparition du concept de la dérivée fractionnaire était en 1694 dans une lettre écrite par le mathématicien allemand **Wilhelm Leibniz** (1646 – 1716) à son collègue **Guillaume de l'Hôpital** (français, 1661, 1704), qui n'a pas exclu la possibilité de définir une telle notion. le fondement effective de cette théorie s'est fait deux siècles après et il revient à **Joseph Liouville** (1809 – 1882) et à **Oliver Heaviside** (1850 – 1925) qui a introduit l'utilisation pratique de ces opérateurs différentiels fractionnaires dans l'analyse des lignes de transport d'électricité .

1.1 Quelques fonctions spéciales :

1.1.1 La fonction Gamma :

Définition 1.1.1 : (la fonction Gamma,[23]) :

On appelle "*fonction Gamma*" la fonction définie par :

$$\Gamma : \mathbb{R}^{+*} \longrightarrow \mathbb{R}^+ \\ x \longmapsto \int_0^{+\infty} t^{x-1} \times e^{-t} dt$$

a. Attention!, ne vous laissez pas tromper par le terme "fractionnaire", l'ordre de dérivation qu'on adoptera dans cette théorie sera réel.

Proposition 1.1.1 :

La fonction Gamma Γ présentée précédemment est bien définie, continue et même de classe C^∞ sur \mathbb{R}^{+*} , en outre, elle est convexe et vérifie les caractérisations suivantes :

1. $\Gamma(1) = 1$, et $\Gamma\left(\frac{1}{2}\right) = \sqrt{\pi}$.
2. $(\forall x > 0) : \Gamma(x+1) = x \times \Gamma(x)$.
3. la fonction : $x \mapsto \ln(\Gamma(x))$ est convexe.

Preuve :

On note :

$$f : \mathbb{R}^+ \times \mathbb{R}^{+*} \longrightarrow \mathbb{R}^+ \quad \text{et} \quad (\forall x \in \mathbb{R}^{+*}) : f_x : \mathbb{R}^{+*} \longrightarrow \mathbb{R}^+ \\ (x, t) \longmapsto t^{x-1} \times e^{-t} \quad t \longmapsto t^{x-1} \times e^{-t}$$

quel que soit $x \in \mathbb{R}$, f_x est continue, positive sur $]0, +\infty[$ donc $\Gamma(x)$ est définie lorsque f_x est intégrable sur $]0, +\infty[$, et bien sûr il est bien le cas car :

- $f_x(t) \underset{0}{\sim} t^{1-x} \implies f_x$ est intégrable sur $]0, 1[$.
- $f_x(t) \underset{+\infty}{=} o\left(\frac{1}{t^2}\right) \implies f_x$ est intégrable sur $]0, +\infty[$.

en plus on a : $(\forall [a, b] \subset \mathbb{R}) (\forall x \in [a, b]) (\forall t > 0) : f(x, t) \leq \varphi(t)$ avec : $\varphi(t) = \begin{cases} t^{a-1} & \text{si } t \in]0, 1[\\ e^{-t} \times t^{b-1} & \text{si } t \in]1, +\infty[\end{cases}$

donc et car φ est positive, continue par morceaux et intégrable sur $]0, +\infty[$, on déduit et selon le théorème de continuité des intégrales paramétriques que : $x \mapsto \Gamma(x)$ est continue.

pour le fait que $\Gamma \in C^\infty(\mathbb{R}^{+*})$, on va établir que : $(\forall k \in \mathbb{N}^*) : \Gamma \in C^k(\mathbb{R}^{+*})$.

d'abord : $(\forall k \in \mathbb{N}^*) : f \in C^k((\mathbb{R}^{+*})^2)$ et :

$$(\forall (x, t) \in (\mathbb{R}^{+*})^2) : \frac{\partial^k f}{\partial x^k}(x, t) = (\ln(t))^k \times e^{-t} \times t^{x-1}$$

Considérons de nouveau un intervalle compact $[a, b] \subset \mathbb{R}^{+*}$ et soit $\Psi_h :]0, +\infty[\rightarrow \mathbb{R}$ telle que :

$$(\forall t \in]0, +\infty[) : \Psi(t) = |\ln(t)|^k \times \varphi(t)$$

on a : $\Psi_k \underset{0}{=} o\left(t^{\frac{a}{2}-1}\right)$ et $\Psi_k \underset{+\infty}{=} o\left(t^{-2}\right)$ donc Ψ_k est positive, continue par morceaux, intégrable sur $]0, +\infty[$ et vérifie :

$$(\forall (x, t) \in [a, b] \times \mathbb{R}^{+*}) : \left| \frac{\partial^k f}{\partial x^k}(x, t) \right| \leq \Psi_k(t)$$

d'où et selon le théorème de dérivabilité des intégrales paramétriques on déduit que :

$$(\forall k \in \mathbb{N}^*) (\forall x \in \mathbb{R}^{+*}) : \Gamma^{(k)}(x) = \frac{d^k}{dx^k} \int_0^{+\infty} f(x, t) dt = \int_0^{+\infty} \frac{\partial^k f}{\partial x^k}(x, t) dt = \int_0^{+\infty} (\ln(t))^k \times e^{-t} \times t^{x-1} dt$$

de cette dernière formule on peut constater que la dérivée seconde de Γ sera strictement positive, ce qui implique immédiatement que cette fonction est strictement convexe.

► pour les caractérisation :

- $\Gamma(1) = \int_0^{+\infty} e^{-t} dt = \left[-e^{-t}\right]_0^{+\infty} = 1$ et $\Gamma\left(\frac{1}{2}\right) = \int_0^{+\infty} \frac{e^{-t}}{\sqrt{t}} dt$ en posant : $u = \sqrt{t}$ on aura :
 $\Gamma\left(\frac{1}{2}\right) = 2 \times \int_0^{+\infty} e^{-t^2} dt = \sqrt{\pi}$ (intégrale de Gauss).

- $(\forall x > 0) : \Gamma(x+1) = \int_0^{+\infty} t^x \times e^{-t} dt = \underbrace{\left[-e^{-t} \times t^x\right]_0^{+\infty}}_{=0} - \int_0^{+\infty} -e^{-t} \times x t^{x-1} dt = x \times \Gamma(x)$

- $x \mapsto \ln(\Gamma(x))$ est convexe comme étant le composé d'une fonction convexe et d'un autre croissante $\ln(\cdot)$. ■

Remarques 1.1.1 :

1. Le théorème de **Liouville** en algèbre différentielle montre en particulier que de nombreuses primitives de fonctions usuelles, telle que la fonction Γ **ne peuvent jamais s'exprimer comme combinaison de fonctions élémentaires**.
2. le théorème de **Bohr-Mollerup**, (voir : [24]) caractérise la fonction gamma, comme la seule fonction définie sur $x > 0$ et qui vérifie simultanément les trois caractérisations mentionnées précédemment.
3. la fonction Γ est appelée parfois la fonction factorielle, car elle vérifie :

$$(\forall n \in \mathbb{N}) : \Gamma(n+1) = n \times \Gamma(n) = n \times n-1 \times \Gamma(n-1) = \dots = n! \times \Gamma(1) = n!$$

et ceci nous permet de la considérer comme une généralisation de cette notion, ce qui va nous permettre de l'étendre au cas des nombres réels, et on peut écrire dans ce cas :

$$(\forall x > 0) : x! = \Gamma(x+1).$$

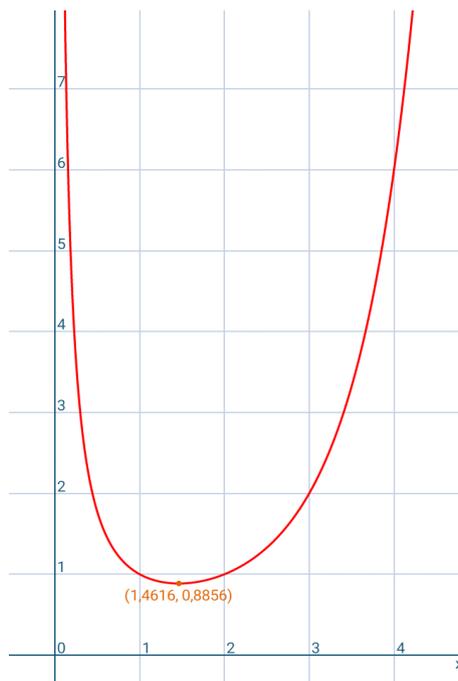


FIGURE IV.1: La courbe de la fonction Gamma.

1.1.2 La fonction Beta :

Définition 1.1.2 : (la fonction Bêta, [56] :) :

On appelle "*fonction Bêta*", (ou fonction d'Euler) la fonction définie par :

$$B : \mathbb{R}^{+*} \times \mathbb{R}^{+*} \longrightarrow \mathbb{R}^+$$

$$(x, y) \longmapsto \int_0^1 t^{x-1} \times (1-t)^{y-1} dt$$

Proposition 1.1.2 :

La fonction Bêta présentée précédemment est bien définie, symétrique, et elle vérifie l'égalité fondamentale suivante :

$$\left(\forall (x, y) \in (\mathbb{R}^{+*})^2 \right) : B(x, y) = \frac{\Gamma(x) \times \Gamma(y)}{\Gamma(x+y)}$$

Preuve :

- **la bonne définition de la fonction Bêta** : pour démontrer ce fait, on va vérifier en premier lieu que l'application φ donnée ainsi :

$$\varphi : \mathbb{R}^{+*} \times \mathbb{R}^{+*} \longrightarrow \mathbb{R}^+$$

$$(x, y) \longmapsto \int_0^{+\infty} \frac{t^{x-1}}{(1+t)^{x+y}} dt$$

a bien un sens, puis, on établira qu'elle coïncide avec notre fonction B .

on a :

$$\frac{t^{x-1}}{(1+t)^{x+y}} \underset{0^+}{\sim} t^{x-1} \quad \text{et} \quad \frac{t^{x-1}}{(1+t)^{x+y}} \underset{+\infty}{\sim} \frac{1}{t^{1+y}}$$

donc et selon les critères de Riemann, et d'équivalence, l'intégrale généralisée $\int_0^{+\infty} \frac{t^{x-1}}{(1+t)^{x+y}} dt$ sera convergente pour tout $x, y \in \mathbb{R}^{+*}$, ce qui fait que φ est bien définie.

-en effectuant le changement de variable : $\rho = \frac{t}{1+t}$ on aura : $t = \frac{\rho}{1-\rho}$ et $dt = \frac{d\rho}{(1-\rho)^2}$.

donc pour tout $x > 0$ et $y > 0$ on a :

$$\begin{aligned} \varphi(x, y) &= \int_0^{+\infty} \frac{t^{x-1}}{(1+t)^{x+y}} dt = \int_0^1 \left(\frac{\rho}{1-\rho} \right)^{x-1} \left(1 + \frac{\rho}{1-\rho} \right)^{-x-y} \frac{1}{(1-\rho)^2} d\rho \\ &= \int_0^1 \left(\frac{\rho}{1-\rho} \right)^{x-1} (1-\rho)^{x+y} \frac{1}{(1-\rho)^2} d\rho \\ &= \int_0^1 \rho^{x-1} \times (1-\rho)^{y-1} d\rho = B(x, y) \end{aligned}$$

et ceci nous donne en conséquence que la fonction bêta est bien définie sur $(\mathbb{R}^{+*})^2$, car φ l'est.

- **la symétrie** : Il suffit d'effectuer dans l'expression de notre fonction le changement de variable $\rho = 1 - t$ pour obtenir l'égalité voulue :

$$\left(\forall (x, y) \in (\mathbb{R}^{+*})^2 \right) : B(x, y) = B(y, x).$$

► **la relation avec la fonction Gamma** : soit a, b deux nombres strictement positifs, on a :

$$\Gamma(a) \times \Gamma(b) = \int_0^{+\infty} \int_0^{+\infty} e^{-x-y} \times x^{a-1} \times y^{b-1} dx dy .$$

on fait le changement de variables : $\rho = x + y$ et $x = \rho \cdot w$, donc $\rho \in \mathbb{R}^+$, $w \in [0, 1]$ et :

$$d\rho = dx + dy, dx = w d\rho + \rho dw, dy = (1 - w) d\rho - \rho dw$$

donc $dx dy = \rho dw d\rho$, ce qui fait que :

$$\Gamma(a) \times \Gamma(b) = \int_0^1 w^{a-1} (1-w)^{b-1} dw \times \int_0^{+\infty} e^{-\rho} \cdot \rho^{a+b-1} d\rho = B(a, b) \times \Gamma(a+b) .$$

d'où ce qu'il faut établir. ■

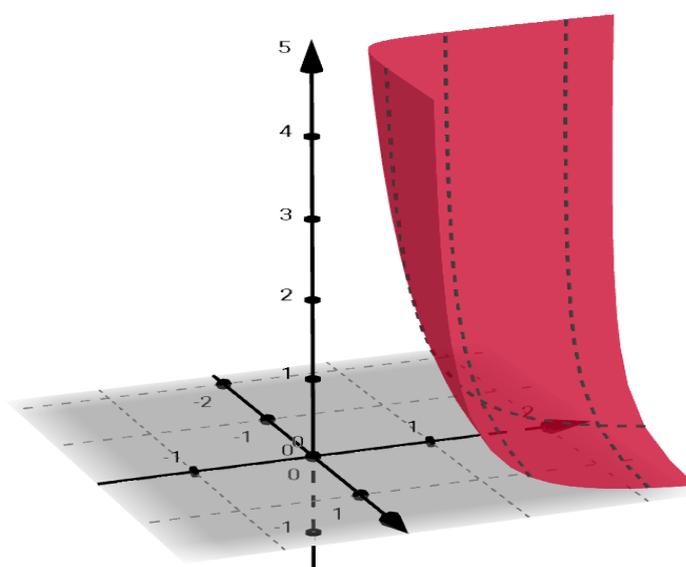


FIGURE IV.2: La courbe de la fonction bêta.

1.2 Intégrale fractionnaire de Riemann- Liouville :

Proposition 1.2.1 : (Formule de Cauchy pour l'intégration itérée)

Soit f une fonction continue sur \mathbb{R} , alors on a pour tout $a \in \mathbb{R}$ fixé, pour tout $n \in \mathbb{N}^*$, et pour tout $(x, y_1, \dots, y_{n-1}) \in [a, +\infty[)^n$:

$$\int_a^x \int_a^{y_1} \dots \int_a^{y_{n-1}} f(y_n) dy_n \dots dy_2 dy_1 = \frac{1}{(n-1)!} \int_a^x (x-t)^{n-1} f(t) dt$$

Preuve :

On va procéder par récurrence sur n :

- pour $n = 1$: l'égalité est triviale ce qui fait que la raison de récurrence est vérifiée.
- soit $n \in \mathbb{N}^*$, on suppose que notre proposition est vérifiée pour n , et on montre qu'elle l'est

aussi pour $n + 1$:

$$\begin{aligned}
 \int_a^x \int_a^{y_1} \cdots \int_a^{y_{n-1}} f(y_n) dy_{n+1} \cdots dy_2 dy_1 &= \int_a^x \frac{1}{(n-1)!} \times \int_a^{y_1} (y_1 - t)^{n-1} f(t) dt \\
 &= \int_a^x \left[\int_a^{y_1} \frac{1}{n!} \times \frac{d}{dy_1} (y_1 - t)^n f(t) \right] dt \\
 &= \int_a^x \frac{1}{n!} \times \frac{d}{dy_1} \left[\int_a^{y_1} (y_1 - t)^n f(t) \right] dt \\
 &= \frac{1}{n!} \times \int_a^x (y_1 - t)^n f(t) dt
 \end{aligned}$$

$$\text{d'où : } (\forall n \in \mathbb{N}^*) : \int_a^x \int_a^{y_1} \cdots \int_a^{y_{n-1}} f(y_n) dy_n \cdots dy_2 dy_1 = \frac{1}{(n-1)!} \int_a^x (x-t)^{n-1} f(t) dt \quad \blacksquare$$

▷ La dérivée usuelle qu'on a eu l'habitude de l'utiliser dans le cadre de l'analyse réelle classique, peut être décrite plus correctement comme un opérateur défini de l'ensemble $\mathcal{D}(\mathbb{R})$ des fonctions dérivable vers l'ensemble des fonctions numériques $\mathcal{F}(\mathbb{R}, \mathbb{R})$, de la façon suivante :

$$\begin{array}{ccc}
 \Delta : \mathcal{D}(\mathbb{R}) & \longrightarrow & \mathcal{F}(\mathbb{R}, \mathbb{R}) \\
 f & \longmapsto & f'
 \end{array}$$

maintenant et à ce stade une question assez naturelle qui peut se poser, peut-on trouver un opérateur linéaire H , ou une **semi-dérivée**, tel que :

$$H^2 := H \circ H = \Delta ?$$

pour garder en principe la cohérence de cette interrogation, et surtout pour des ordres supérieurs à deux, on sera obligé de rendre l'espace de départ et d'arrivée de l'opérateur Δ identiques, pour cela on va le prendre et sans perte de généralité $C^\infty(\mathbb{R})$, l'espace des fonctions continument dérivable.

▷ Considérant une fonction $x \mapsto f(x)$ qui est **définie et continue** sur \mathbb{R}^+ , nous pouvons et sans aucun problème exprimer sa primitive (qui s'annule en 0), par l'intégrale impropre définie de a à x par :

$$(\forall x > 0) : [P(f)](x) = \int_0^x f(t) dt .$$

selon la proposition 2.2.1, une application itérée de cet opérateur pour la fonction f nous donne :

$$(\forall n \in \mathbb{N}^*) (\forall x > 0) : [P^n(f)](x) := [(P \circ \cdots \circ P)(f)](x) = \frac{1}{(n-1)!} \times \int_0^x \underbrace{(x-t)^{n-1}}_{\geq 0} \times f(t) dt .$$

l'utilisation de la fonction gamma $x \mapsto \Gamma(x)$ pour franchir la nature entière de la factorielle sera un choix naturel et évident pour donner un sens aux applications fractionnaires de l'opérateur intégral, et ainsi on peut écrire :

$$\boxed{(\forall f \in C^0(\mathbb{R}^+)) (\forall \alpha > 0) (\forall x > 0) : [P^\alpha(f)](x) := \frac{1}{\Gamma(\alpha)} \times \int_0^x (x-t)^{\alpha-1} \times f(t) dt .} \quad (\blacklozenge)$$

Définition 1.2.1 : (intégrale fractionnaire de Riemann- Liouville :)

On appelle (\blacklozenge) décrite ci-dessus la formule fondamentale du calcul fractionnaire, ou encore l'intégrale fractionnaire de Riemann-Liouville.

Proposition 1.2.2 : (la propriété du semi-groupe pour l'opérateur « P ») :

Soit $f \in C^0(\mathbb{R}^+)$, l'opérateur d'intégration P satisfait la propriété suivante pour tout α et β **strictement positifs** :

$$(\forall x \in (\mathbb{R}^{+*})) : \left[P^\alpha \left[P^\beta (f) \right] \right] (x) = \left[P^\beta \left[P^\alpha (f) \right] \right] (x) = \left[P^{\alpha+\beta} (f) \right] (x) = \frac{1}{\Gamma(\alpha+\beta)} \int_0^x (x-t)^{\alpha+\beta-1} f(t) dt.$$

Preuve :

Soit α et β deux nombres positifs strictement, on a et pour tout $x > 0$:

$$\begin{aligned} \left[P^\alpha \left[P^\beta (f) \right] \right] (x) &= \frac{1}{\Gamma(\alpha)} \int_0^x (x-t)^{\alpha-1} \left[P^\beta (f) \right] (t) dt \\ &= \frac{1}{\Gamma(\alpha)\Gamma(\beta)} \int_0^x \int_0^t (x-t)^{\alpha-1} (t-s)^{\beta-1} f(s) ds dt \\ &= \frac{1}{\Gamma(\alpha)\Gamma(\beta)} \int_0^x f(s) \left(\int_s^x (x-t)^{\alpha-1} (t-s)^{\beta-1} dt \right) ds \end{aligned}$$

et en effectuant le changement de variables : $\rho = \frac{t-s}{x-s}$, on obtient :

$$\left[P^\alpha \left[P^\beta (f) \right] \right] (x) = \frac{1}{\Gamma(\alpha)\Gamma(\beta)} \int_0^x (x-s)^{\alpha+\beta-1} f(s) \left(\int_0^1 (1-r)^{\alpha-1} r^{\beta-1} dr \right) ds \quad (\blacktriangle)$$

on remarque que l'intégrande dans la dernière expression n'est autre que la fameuse fonction bêta et qui vérifie la propriété suivante :

$$\int_0^1 (1-\rho)^{\alpha-1} \rho^{\beta-1} d\rho = B(\alpha, \beta) = \frac{\Gamma(\alpha)\Gamma(\beta)}{\Gamma(\alpha+\beta)}.$$

en remplaçant cela dans l'équation (\blacktriangle) on obtient :

$$\left[P^\alpha \left[P^\beta (f) \right] \right] (x) = \frac{1}{\Gamma(\alpha+\beta)} \int_0^x (x-s)^{\alpha+\beta-1} f(s) ds = \left[P^{\alpha+\beta} (f) \right] (x).$$

d'où ce qu'il faut prouver. ■

1.3 Dérivée fractionnaire de Caputo :

Le mathématicien italien **Michele Caputo** (1927, ...) s'est inspiré des travaux de Riemann et Liouville pour améliorer et généraliser les notions d'intégration et de dérivation dans le cadre fractionnaire, et en 1967 il a proposé dans son article [8] une idée intelligente basée sur la propriété du semi groupe mentionnée à la proposition IV.2.2.2 pour introduire une généralisation fractionnaire de la dérivation usuelle et qui porte jusqu'à l'instant son nom : "**la dérivation de Caputo**".

Notations :

On note pour la suite : $\hat{\mathbb{R}} = \mathbb{R} \setminus \mathbb{Z}$ et $\hat{\mathbb{R}}^+ = \mathbb{R}^+ \setminus \mathbb{N}$.

Définitions 1.3.1 : (intégrale fractionnaire de Caputo :)

- Soit $f \in C^0(\mathbb{R})$ **une fonction continue** sur \mathbb{R} , et $a \in \mathbb{R}$ un nombre réel fixé, l'intégrale fractionnaire de Caputo **pour un ordre positif** est une généralisation de celui de Riemann-Liouville qui s'exprime sous les deux formes suivantes :

$$\left\{ \begin{array}{l} \text{-supérieure : } (\forall \alpha > 0) (\forall x > a) : \left[{}_a P_x^\alpha (f) \right] (x) = \frac{1}{\Gamma(\alpha)} \int_a^x (x-t)^{\alpha-1} f(t) dt. \\ \text{et :} \\ \text{-inférieure : } (\forall \alpha > 0) (\forall x < a) : \left[{}_x P_a^\alpha (f) \right] (x) = \frac{1}{\Gamma(\alpha)} \int_x^a (t-x)^{\alpha-1} f(t) dt. \end{array} \right.$$

- afin d'étendre la formulation précédente aux ordres réels négatifs^a, Caputo s'est inspiré de la propriété du semi groupe chez l'intégrale de Riemann-Liouville pour nous fournir la définition suivante :

$$\text{— d'une part : } (\forall n \in \mathbb{N}^*) : \left[{}_x P_a^{-n} (f) \right] (\cdot) = \frac{d^n f}{dx^n} (\cdot) = f^{(n)}(\cdot)$$

— et d'une autre :

$$\begin{aligned} (\forall \alpha \in \widehat{\mathbb{R}}^+) (\forall f \in C^{[\alpha]}(\mathbb{R})) (\forall x > a) : \left[{}_a P_x^{-\alpha} (f) \right] (x) &:= \left[{}_a P_x^{[\alpha]-\alpha} \left[{}_a P_x^{-[\alpha]} (f) \right] \right] (x) \\ &= \frac{1}{\Gamma([\alpha] - \alpha)} \int_a^x (x-t)^{[\alpha]-\alpha-1} f^{([\alpha])}(t) dt \\ &= \frac{1}{\Gamma([\alpha] - \alpha)} \int_a^x \frac{f^{([\alpha])}(t)}{(x-t)^{\alpha-[\alpha]+1}} dt. \end{aligned}$$

et d'une façon analogue à la forme supérieure présentée en dessus, on obtient pour la forme inférieure :

$$(\forall \alpha \in \widehat{\mathbb{R}}^+) (\forall f \in C^{[\alpha]}(\mathbb{R})) (\forall x > a) : \left[{}_x P_a^{-\alpha} (f) \right] (x) := \frac{1}{\Gamma([\alpha] - \alpha)} \int_x^a \frac{f^{([\alpha])}(t)}{(t-x)^{\alpha-[\alpha]+1}} dt.$$

Remarques 1.3.1 :

1. Dans la définition précédente, on a adopté la notation anglo-saxonne pour la partie entière et selon laquelle on écrit :

$$\underbrace{[\alpha]}_{\text{partie entière de } \alpha} \leq \alpha < [\alpha] + 1.$$

2. il est bien facile de justifier que : $(\forall n \in \mathbb{N}^*) : \left[{}_x P_a^{-n} (f) \right] (\cdot) = \frac{d^n f}{dx^n} (\cdot) = f^{(n)}(\cdot)$ et ceci par le simple fait que la dérivation n'est autre que l'opérateur réciproque de l'intégration.
3. **faites attention!**, dans le cas où $\alpha \in \widehat{\mathbb{R}} = (\mathbb{R}^+ \setminus \mathbb{N})$, la définition d'intégration fractionnaire proposée par Caputo ne s'adapte plus avec la propriété du commutativité (voir[3] page :22), c'est à dire que :

$$(\forall \alpha \in \widehat{\mathbb{R}}^+) (\forall x > a) : \left[{}_a P_x^{-\alpha} (f) \right] (x) := \left[{}_a P_x^{[\alpha]-\alpha} \left[{}_a P_x^{-[\alpha]} (f) \right] \right] (x) \neq \underbrace{\left[{}_a P_x^{-[\alpha]} \left[{}_a P_x^{[\alpha]-\alpha} (f) \right] \right]}_{\text{La dérivée fractionnaire de Riemann-Liouville}} (x)$$

La dérivée fractionnaire de Riemann-Liouville.

a. On rappelle que le problème dans ce cas se pose au niveau de la fonction Gamma qui est définie juste pour des valeurs strictement positifs.

Définition 1.3.2 : (dérivation fractionnaire de Caputo :)

Soit $a \in \mathbb{R}$ et $\alpha \in \mathbb{R}$ des réels fixés, et soit $f \in C^{[\alpha]}(\mathbb{R})$ une fonction de classe $C^{[\alpha]}$ sur \mathbb{R} , on définit la dérivée de Caputo (supérieure) de f pour un l'ordre réel α , comme suit :

$$(\forall x > a) : \left[{}_a D_x^\alpha (f) \right] (x) = \begin{cases} f^{(\alpha)}(x) & \text{si } \alpha \in \mathbb{Z}. \\ \left[{}_a P_x^{-\alpha} (f) \right] (x) & \text{si } \alpha \in \widehat{\mathbb{R}}. \end{cases}$$

ce qui fait que :

$$(\forall x > a) : \left[{}_a D_t^\alpha (f) \right] (x) = \begin{cases} f^{(\alpha)}(x) & \text{si } \alpha \in \mathbb{N}. \\ \frac{1}{\Gamma([\alpha] - \alpha)} \int_a^x \frac{f^{([\alpha])}(t)}{(x-t)^{\alpha-([\alpha]+1)}} dt & \text{si } \alpha \in \widehat{\mathbb{R}}^+. \\ \frac{1}{\Gamma(-\alpha)} \int_a^x (x-t)^{-\alpha-1} f(t) dt & \text{si } \alpha < 0. \end{cases}$$

et pour ce qui concerne la dérivée de Caputo inférieure, elle sera et sans surprise définie d'une façon analogue, juste il faut inverser les bornes de l'intégrale et remplacer $(x - t)$ par $(t - x)$ dans l'intégrande.

Proposition 1.3.1 :

Soit $\alpha \in \widehat{\mathbb{R}}$ et $f \in C^{[\alpha]}(\mathbb{R})$ tel que $n - 1 < \alpha < n = [\alpha] \in \mathbb{N}$, alors et dans ce cas, l'opérateur de dérivation fractionnaire de Caputo $[{}_a D_x^\alpha (\cdot)]$ vérifie les propriétés suivantes :

$$\lim_{\alpha \rightarrow n} \left[{}_a D_x^\alpha (f) \right] (x) = f^{(n)}(x) \text{ et } \lim_{\alpha \rightarrow n-1} \left[{}_a D_x^\alpha (f) \right] (x) = f^{(n-1)}(x) - f^{(n-1)}(a)$$

preuve :

On va juste utiliser une intégration par partie :

$$\begin{aligned} \left[{}_a D_x^\alpha (f) \right] (x) &= \frac{1}{\Gamma([\alpha] - \alpha)} \int_a^x \frac{f^{(n)}(t)}{(x-t)^{\alpha-([\alpha]+1)}} dt \\ &= \frac{1}{\Gamma([\alpha] - \alpha)} \times \left(\left[-f^{(n)}(t) \cdot \frac{(x-t)^{[\alpha]-\alpha}}{[\alpha] - \alpha} \right]_{t=a}^{t=x} - \int_a^x -f^{(n+1)}(t) \times \frac{(x-t)^{[\alpha]-\alpha}}{n-\alpha} dt \right) \\ &= \frac{1}{\Gamma([\alpha] - \alpha)} \times \left(-f^{(n)}(a) \cdot \frac{(x-a)^{[\alpha]-\alpha}}{[\alpha] - \alpha} + \int_a^x f^{(n+1)}(t) \times \frac{(x-t)^{[\alpha]-\alpha}}{n-\alpha} dt \right) \\ &= \frac{1}{\Gamma([\alpha] - \alpha + 1)} \times \left(-f^{(n)}(a) \cdot (x-a)^{[\alpha]-\alpha} + \int_a^x f^{(n+1)}(t) \times \frac{(x-t)^{[\alpha]-\alpha}}{n-\alpha} dt \right) \end{aligned}$$

en faisant tendre respectivement $\alpha \rightarrow n$ et $\alpha \rightarrow n - 1$, on obtiendra :

$$\lim_{\alpha \rightarrow n} \left[{}_a D_x^\alpha (f) \right] (x) = f^{(n)}(a) + \left[f^{(n)}(x) \right]_{t=a}^{t=x} = f^{(n)}(x)$$

et :

$$\begin{aligned} \lim_{\alpha \rightarrow n} \left[{}_a D_x^\alpha (f) \right] (x) &= f^{(n)}(a) \cdot (x-a) + \left[f^{(n)}(x) \cdot (x-t) \right]_{t=a}^{t=x} + \int_a^x f^{(n)}(t) dt \\ &= \left[f^{(n-1)}(t) \right]_{t=a}^{t=x} \\ &= f^{(n-1)}(x) - f^{(n-1)}(a) \end{aligned}$$

d'où ce qu'il faut établir. ■

Exemple 1.3.1 :

Supposons que f est un monôme de la forme $f(x) = x^n$ définie sur \mathbb{R}^+ , c'est à dire :

$$\boxed{\begin{array}{ccc} f & : & \mathbb{R}^+ \longrightarrow \mathbb{R}^+ \\ & & t \longmapsto t^n \end{array}}$$

on a clairement : $f \in C^\infty(\mathbb{R}^+)$, et si on pose $a = 0$ et on prend $\alpha \in \widehat{\mathbb{R}}^+$ on obtiendra :

$$(\forall x \in \mathbb{R}^+) : \left[{}_a D_t^\alpha (f) \right] (x) = \frac{1}{\Gamma([\alpha] - \alpha)} \int_a^x \frac{f^{([\alpha])}(t)}{(x-t)^{\alpha-[\alpha]+1}} dt.$$

- si $[\alpha] > n$ alors : $f^{([\alpha])}(t) = 0$, ce qui fait que :

$$(\forall x \in \mathbb{R}^+) : \left[{}_0 D_t^\alpha (t^n) \right] (x) = 0.$$

- si $[\alpha] < n$ alors : $f^{([\alpha])}(t) = n \times (n-1) \times \dots \times (n-[\alpha]+1) \times t^{n-[\alpha]}$, ce qui fait que :

$$\begin{aligned} (\forall x \in \mathbb{R}^+) : \left[{}_0 D_t^\alpha (t^n) \right] (x) &= \frac{1}{\Gamma([\alpha] - \alpha)} \times \int_0^x \frac{f^{([\alpha])}(t)}{(x-t)^{\alpha-[\alpha]+1}} dt \\ &= \frac{1}{\Gamma([\alpha] - \alpha)} \times \int_0^x \frac{n \times (n-1) \times \dots \times (n-[\alpha]+1) \times t^{n-[\alpha]}}{(x-t)^{\alpha-[\alpha]+1}} dt \\ &= \frac{1}{\Gamma([\alpha] - \alpha)} \times \int_0^x \frac{n! \times t^{n-[\alpha]}}{(n-[\alpha])! \times (x-t)^{\alpha-[\alpha]+1}} dt \\ &= \frac{\Gamma(n+1)}{(n-[\alpha])! \times \Gamma([\alpha] - \alpha)} \times \int_0^x t^{n-[\alpha]} \times (x-t)^{[\alpha]-\alpha-1} dt \\ \left(\text{on pose : } \rho = \frac{t}{x} \right) &= \frac{\Gamma(n+1)}{(n-[\alpha])! \times \Gamma([\alpha] - \alpha)} \times \int_0^1 (x \cdot \rho)^{n-[\alpha]} \times (x - \rho \cdot x)^{[\alpha]-\alpha-1} x d\rho \\ &= \frac{\Gamma(n+1) \times x^{n-\alpha}}{(n-[\alpha])! \times \Gamma([\alpha] - \alpha)} \times \int_0^1 \rho^{n-[\alpha]} \times (1-\rho)^{[\alpha]-\alpha-1} d\rho \\ &= \frac{\Gamma(n+1) \times x^{n-\alpha}}{(n-[\alpha])! \times \Gamma([\alpha] - \alpha)} \times B(n-[\alpha]+1, [\alpha]-\alpha) \\ &= \frac{\Gamma(n+1) \times x^{n-\alpha}}{\cancel{(n-[\alpha])!} \times \cancel{\Gamma([\alpha] - \alpha)}} \times \frac{\Gamma(\cancel{(n-[\alpha]+1)} + 1) \times \Gamma(\cancel{[\alpha]-\alpha})}{\Gamma(n-[\alpha]+1 + [\alpha]-\alpha)} \\ &= \boxed{\frac{\Gamma(n+1)}{\Gamma(n+1-\alpha)} \times x^{n-\alpha}} \end{aligned}$$

donc : $\boxed{(\forall \alpha \in \widehat{\mathbb{R}}^+) (\forall x \in \mathbb{R}^+) : [\alpha] < n \Rightarrow \left[{}_0 D_t^\alpha (t^n) \right] (x) = \frac{\Gamma(n+1)}{\Gamma(n+1-\alpha)} \times x^{n-\alpha}. \quad (\star)}$

Dans (\star) , si on affecte à α une valeur entière, alors on tombera sur la formule classique de la dérivation répétée d'un monôme, et qui est donnée par :

$$(\forall \alpha \in \widehat{\mathbb{R}}^+) (\forall x \in \mathbb{R}^+) : \alpha \leq n \Rightarrow \left[{}_0 D_t^\alpha (t^n) \right] (x) = \frac{n!}{(n-\alpha)!} \times x^{n-\alpha} = \frac{d^\alpha}{dx^\alpha} x^n.$$

▷ **La semi-dérivée de la fonction identité :** pour expliciter son expression on va poser dans (\star) $k = 1$ et $\alpha = 0.5 \in \widehat{\mathbb{R}}^+$, et ainsi on obtiendra :

$$\frac{d^{\frac{1}{2}}}{dx^{\frac{1}{2}}} x := \left[{}_0 D_t^\alpha (t) \right] (x) = \frac{\Gamma(1+1)}{\Gamma\left(1+1-\frac{1}{2}\right)} \times x^{1-\frac{1}{2}} = \frac{\Gamma(2)}{\Gamma\left(\frac{3}{2}\right)} \times x^{\frac{1}{2}} = \frac{\Gamma(2)}{\Gamma\left(1+\frac{1}{2}\right)} \times \sqrt{x} = \frac{\Gamma(2)}{\frac{1}{2} \cdot \Gamma\left(\frac{1}{2}\right)} \times \sqrt{x} = \frac{1}{\frac{\sqrt{\pi}}{2}} \times \sqrt{x}.$$

Pour garantir qu'elle s'agit bien de la "*demi-dérivée*" nous allons la réappliquer une seconde fois pour obtenir :

$$\frac{d^{\frac{1}{2}}}{dx^{\frac{1}{2}}} \left(\frac{d^{\frac{1}{2}}}{dx^{\frac{1}{2}}} x \right) = \frac{d^{\frac{1}{2}}}{dx^{\frac{1}{2}}} \frac{2x^{\frac{1}{2}}}{\sqrt{\pi}} = \frac{2}{\sqrt{\pi}} \times \frac{\Gamma(1+\frac{1}{2})}{\Gamma(\frac{1}{2}-\frac{1}{2}+1)} x^{\frac{1}{2}-\frac{1}{2}} = \frac{2}{\sqrt{\pi}} \times \frac{\Gamma(\frac{3}{2})}{\Gamma(1)} x^0 = \frac{2\frac{\sqrt{\pi}}{2} \times x^0}{\sqrt{\pi}} = 1$$

$$\Rightarrow \frac{d^{\frac{1}{2}}}{dx^{\frac{1}{2}}} \left(\frac{d^{\frac{1}{2}}}{dx^{\frac{1}{2}}} x \right) = 1 = \frac{d}{dx} x.$$

et qui est bien le résultat attendu.

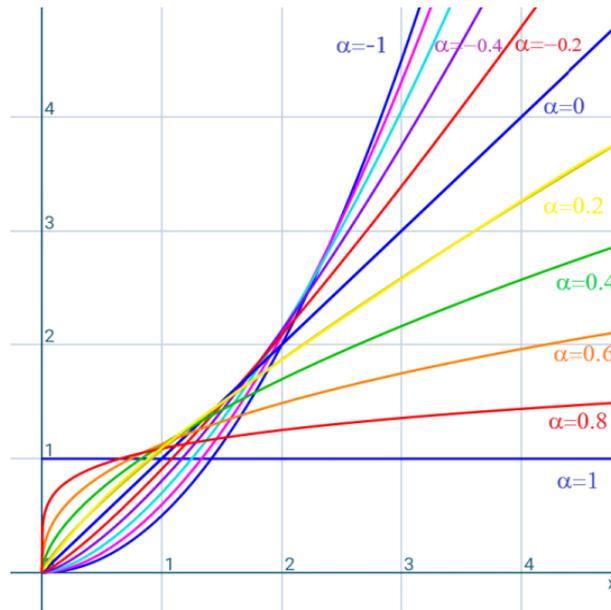


FIGURE IV.3: La dérivée fractionnaire de $x \mapsto x$ pour certains ordres $\alpha \in [-1, 1]$.

Remarques 1.3.2 :

1. La $k^{\text{ème}}$ dérivée d'une fonction f **est une propriété locale uniquement lorsque k est un nombre entier**, mais pour des ordres qui ne le sont pas ce n'est plus le cas, autrement dit, la dérivée fractionnaire d'une fonction f en un point x_0 pour des ordres non entières ne dépend pas de son comportement au voisinage de ce point, contrairement aux dérivées de puissance entière qui ont été fondées sur la notion du limite, et ceci rend bien évidemment **l'approximation de ce type de dérivée par des formulations numériques une tâche compliquée et très difficile**^a.
2. la dérivée d'ordre fractionnaire d'une fonction nécessite sa connaissance sur tout l'intervalle d'étude, alors que dans le cas entier, seule la connaissance locale de f autour d'un point est nécessaire, cette propriété s'interprète par le fait que les systèmes d'ordres fractionnaires sont des systèmes à mémoire longue, tandis que les systèmes entiers sont vus comme des systèmes à mémoire courte.
3. dans ce qui précède, on a présenté la définition la plus facile ou la plus naturelle pour la dérivation fractionnaire et qui est tout simplement celle de Caputo, mais dans la littérature beaucoup d'autre choix sont proposés (voir [17]), par exemple : la dérivée fractionnaire de Riesz, de Liouville, de Miller, ... etc.

a. Malgré sa difficulté une telle tâche n'est pas complètement impossible, et vous pouvez trouver dans la littérature quelques schémas qui le font, voir par exemple les pages : 27-28 de [17] (Schéma de Grünwald-Letnikov et schéma G^α).

4. la proposition IV.1.3.1 se traduit par le fait que l'opérateur de dérivation de Caputo est toujours continue à gauche des ordres entiers, en revanche, sa continuité à droite d'un $n \in \mathbb{N}$ est conditionnée par la réalisation de la propriété $f^{(n-1)}(a) = 0$, ce qui rend cet opérateur un peu moins avantageux que celui de Riemann-Liouville, (voir Remarque IV.1.3.1), et qui est continue dans les deux côtés de n'importe quel élément de \mathbb{N} sans aucune exigence supplémentaire. ^a.
5. par contre à l'opérateur de différentiation fractionnaire de Riemann-Liouville et qui généralise automatiquement la dérivation usuelle pour les ordres entiers, l'opérateur de Caputo nécessite toujours une définition par parties pour le renforcer à s'adapter avec cette dérivation, et ceci provient essentiellement du fait que le premier est basé sur une différentiation après intégration, ce qui lui permet de franchir les problèmes des constantes supplémentaires, tandis que le deuxième est basé sur une intégration après différentiation et qui ne lui permet pas de surmonter ce problème.

2 Description et discrétisation du problème :

2.1 Description du problème :

Notre but principal dans la suite est d'utiliser le réseau neuronal multicouches feed-forward qu'on a étudié au deuxième chapitre, (voir :II.3.2.1 page : 45), pour trouver une approximation convenable, convaincante, ou au moins acceptable de la solution de l'équation différentielle fractionnaire suivante :

$$(\mathcal{E}_\alpha) \begin{cases} \left[{}_a D_t^\alpha (V(t)) \right] + h(t) \times V(t) = s(t). & \text{si : } a < t < b. \\ V(t) = V(a) = \beta \in \mathbb{R}. & \text{si : } t = a. \end{cases} \quad (\text{IV.1})$$

dans cette dernière équation, $V(\cdot) \in C^1([a, b])$ joue le rôle de **la fonction inconnue** qu'on cherche à trouver, tandis que $h(\cdot)$ et $s(\cdot)$ sont des fonctions continues à valeurs réelles, fixées et connues préalablement, c'est à dire qu'elles **sont considérées comme des paramètres** pour notre équation, en plus, et afin de simplifier le processus de calcul, on choisira $\alpha \in]0, 1]$, et on considère l'intervalle d'étude comme $[a, b] = [0, T]$, et même s'il ne l'était pas on peut le rendre ainsi, car tout intervalle de \mathbb{R} peut être transformé à cette forme via une correspondance affine simple.

Maintenant, **l'idée qu'on propose consiste à remplacer $V(y)$ sur le domaine $\Omega = [0, T]$ par une architecture neuronale bien conçue**, et puis la modifier au fur et à mesure pour aboutir une bonne approximation de notre solution, et on signale que cette façon de procéder peut être facilement utilisée pour n'importe quelles relations fractionnaires, même si elles sont d'une nature un peu plus compliquée que (IV.1).

2.2 Discrétisation du problème :

La première étape dans l'approche itérative suggérée pour résoudre ce problème, est de le rendre d'abord plus adaptable à des reformulations discrètes, et ceci en combinant sa condition initiale avec la solution éventuelle qu'on lui propose ^b, ce qui fait que dans le cadre de l'équa-

a. On laisse au lecteur le soin de prouver ce fait d'une façon détaillée, et pour quelques indications, on peut l'orienter vers [3] page 21.

tion (\mathcal{E}_α), cette fonction qu'on présente comme un bon candidat pour être la solution effective de notre problème, doit concilier son inspiration de l'architecture neuronale feed-forward avec la vérification de la condition initiale, et naturellement, ceci va nous conduire à la considération de la fonction : $\tilde{V}(t)$ définie par :

$$\left(\forall t \in [0, T] \right) : \tilde{V}(t) = \beta + t \times N(t) = \beta + t \times \sum_{i=1}^I \omega_i^{(2)} \times \sigma \left(\omega_i^{(1)} \times t + b_i \right) \quad (*)$$

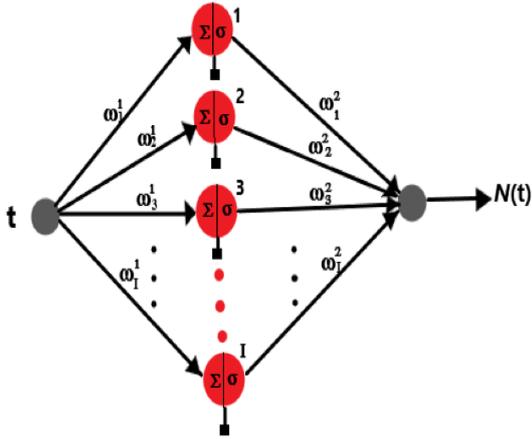


FIGURE IV.4: Le réseau feed-forward réalisant $N(\cdot)$.

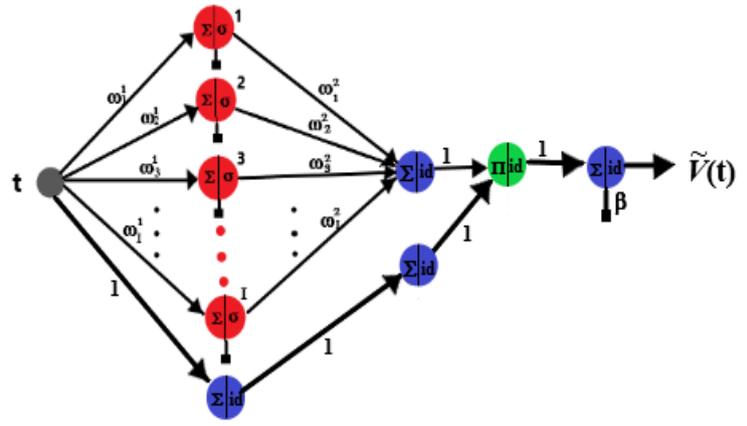


FIGURE IV.5: Le $\Sigma\Pi$ réseau réalisant $\tilde{V}(\cdot)$.

Le choix de cette formulation peut être justifiée par le fait que $\tilde{V}(\cdot)$ contient l'architecture d'un réseau de neurones feed-forward, et qui se reflète par $N(\cdot)$, et au même temps, elle satisfait la condition de départ : $\tilde{V}(0) = \beta$.

En principe la fonction $\tilde{V}(\cdot)$ présente une solution possible, et non forcément effective, qui fait intervenir les trois vecteurs de paramètres ajustables $\omega^1 \in \mathbb{R}^I$, $\omega^2 \in \mathbb{R}^I$ et $b \in \mathbb{R}^I$.

une substitution directe de $V(t)$ par (*) dans l'équation (IV.1) transformera notre problème de sa forme d'origine à une autre plus simple ,plus pratique, non soumise à des contraintes ou des conditions de départ, et qui se présente ainsi :

$$\left(\forall t \in [0, T] \right) : \left[{}_0D_t^\alpha \left(t \times \sum_{i=1}^I \omega_i^{(2)} \times \sigma \left(\omega_i^{(1)} \times t + b_i \right) \right) \right] + h(t) \times \left(\beta + t \times \sum_{i=1}^I \omega_i^{(2)} \times \sigma \left(\omega_i^{(1)} \times t + b_i \right) \right) = s(t) \quad (\blacktriangledown)$$

Généralement, un réseau neuronal ne peut pas prendre facilement la place de l'inconnue dans une équation différentielle fractionnaire, mais avec le type particulier de la fonction d'activation choisie ,et qui était sigmoïdale dans ce cas, ce fait devient justifiable par la propriété d'approximation universelle qu'on a traité au chapitre III, et selon laquelle l'espace des réseau neuronaux sigmoïdes est dense dans celui des fonctions continues, et dans ce sens, on profite de l'occasion pour insister sur l'importance pratique tant que théorique du type de la fonction d'activation choisie.

b. Le recours à cette idée revient essentiellement à l'absence de la notion de la dérivée fractionnaire à gauche ou à droite d'un point.

À ce stade, le problème principal qui se pose encore réside dans la grande difficulté de déterminer explicitement l'expression de la dérivée fractionnaire $D_t^\alpha (\tilde{V}(t))$, et spécialement lorsque cette fonction $\tilde{V}(\cdot)$ est non polynomiale, ceci va nous pousser à chercher un moyen qui peut nous aider à franchir le calcul directe de cette quantité, et surtout à cause de la complexité pratique qu'il présente.

L'approche la plus classique qu'on peut proposer pour résoudre le problème mentionné ci-dessus, sera de substituer la fonction d'activation σ qu'on la suppose de classe C^∞ dorénavant par **son développement en séries entières**, et qui peut être obtenu de la formule de **Maclaurin** par exemple.

► Avant de continuer, on suppose d'abord que le rayon de convergence de la série entière ^a

$$\Phi = \sum_{n \in \mathbb{N}} \frac{\sigma^{(n)}(0)}{n!} \times x^n$$

est égale à $+\infty$, c'est à dire que :

$$R_\Phi := \sup \left\{ |x| / x \in \mathbb{R} \text{ et : } \sum_{n \geq 0} \frac{\sigma^{(n)}(0)}{n!} \times x^n \text{ converge} \right\} = \sup \left\{ |x| / x \in \mathbb{R} \text{ et : } \frac{\sigma^{(n)}(0)}{n!} \times x^n \text{ bornée} \right\} = +\infty.$$

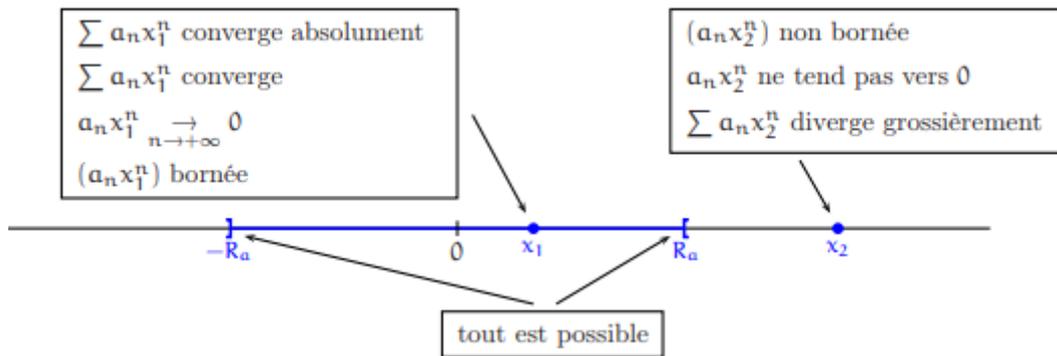


FIGURE IV.6: Le comportement d'une série entière de rayon de convergence R .

pratiquement, et pour garantir la réalisation de cette condition, il suffit de s'assurer que l'une des conditions suivantes est vérifiée :

- $(\forall x \in \mathbb{R}) \left(\frac{\sigma^{(n)}(0)}{n!} \times x^n \right)_{n \in \mathbb{N}}$ est bornée. (critère d'Abel.)
- $\limsup_{n \rightarrow +\infty} \sqrt[n]{\left| \frac{\sigma^{(n)}(0)}{n!} \right|} = 0$. (critère de Cauchy-Hadamard,[26].)
- $\lim_{n \rightarrow +\infty} \left| \frac{\sigma^{(n+1)}(0)}{(n+1)!} \right| \times \left| \frac{n!}{\sigma^{(n)}(0)} \right| = \frac{1}{n+1} \times \left| \frac{\sigma^{(n+1)}(0)}{\sigma^{(n)}(0)} \right| = 0$. (critère d'Alembert.)

a. C'est une série de fonctions réelles (ou complexes) φ_n de définies par : $(\forall z \in \mathbb{R}) : \varphi_n(z) = a_n \times z^n$.
pour un z donné $\sum_{n \geq 0} a_n \times z^n$, **n'est pas nécessairement convergente**, et s'il était le cas on pose : $f(z) = \sum_{n=0}^{+\infty} a_n \times z^n$

le fait que $R_\Phi = +\infty$ va nous permettre d'attribuer un sens, à la fonction somme ϕ sur tout l'ensemble \mathbb{R} , et qui est effectivement le disque de convergence de Φ dans ce cas, et on écrit donc :

$$\boxed{\begin{array}{l} \phi :]-\infty, \overbrace{+\infty}^{=R_\Phi}[\longrightarrow \mathbb{R} \\ x \longmapsto \sum_{k=0}^{+\infty} \frac{\sigma^{(k)}(0)}{k!} \times x^k \end{array}} \quad (\diamond)$$

Définition 2.2.1 : (fonction développable en séries entières) :

Soit $x_0 \in \mathbb{R}$ un réel, et $\rho \in \mathbb{R}^+ \cup \{+\infty\}$ un rayon strictement positif, on dit qu'une fonction g est développable en série entière en x_0 sur $]x_0 - \rho, x_0 + \rho[$, si et seulement si :

$$\boxed{\rho \leq R_S \quad \text{et} \quad (\forall x \in]x_0 - \rho, x_0 + \rho[) : g(x) = \sum_{n=1}^{+\infty} a_n \times (x - x_0)^n = g(x_0) + \sum_{n=1}^{+\infty} a_n \times (x - x_0)^n .}$$

et dans ce cas là cette identité s'appelle le développement en série entière de g en x_0 sur $]x_0 - \rho, x_0 + \rho[$.

Remarques 2.2.1 :

1. L'hypothèse $\rho \leq R_S$ où S est la série entière définie par :

$$S = \sum_{n \geq 0} a_n \times (x - x_0)^n$$

est nécessaire pour attribuer un sens à sa somme pour tout $x \in]x_0 - \rho, x_0 + \rho[$.

2. Une fonction g sera développable en série entière en x_0 sur $]x_0 - \rho, x_0 + \rho[$, si et seulement si sa translatée $\tau_{x_0}(g) : x \mapsto g(x_0 + x)$ l'est en 0 sur $] - \rho, \rho[$, ce qui fait qu'on pourra toujours se ramener à des développements en série entière en 0.

Lemme 2.2.1 :

Soit $\sum_{n \geq 0} a_n \cdot x^n$ une série entière, de rayon de convergence R , et soit ρ un réel tel que $0 < \rho < R$, alors :

$$(\forall \varepsilon > 0) (\exists n_0 \in \mathbb{N}) \text{ tel que : } (\forall n \geq n_0) \sup_{|x| \leq \rho} \left| \sum_{k=0}^n a_n \cdot x^n - \sum_{k=0}^{+\infty} a_n \cdot x^n \right| < \varepsilon .$$

c'est à dire qu'une série entière converge uniformément sur n'importe quel intervalle borné inclus dans son disque de convergence.

Preuve :

Fixons un ρ' dans $] \rho, R[$, pour tout $n \in \mathbb{N}$ on a :

$$|a_n \cdot x^n| \leq \underbrace{|a_n| (\rho')^n}_{\text{bornée}} \times \frac{\rho^n}{(\rho')^n} \leq M \left(\frac{\rho^n}{\rho'} \right)^n .$$

alors, pour tout réel $x \in] - \rho, \rho[$:

$$\left| \sum_{k=0}^n a_n \cdot z^n - \sum_{k=0}^{+\infty} a_n \cdot z^n \right| \leq \sum_{k=0}^{+\infty} M \left(\frac{\rho^n}{\rho'} \right)^n = \left(\frac{Mr}{r' - r} \right) \left(\frac{\rho}{\rho'} \right)^{n+1}$$

cette dernière majoration est indépendante de x , donc la convergence est bien uniforme, ce qui achève notre preuve ■

Théorème 2.2.1 :

Soit f une fonction développable en série entière au voisinage de 0, telle que :

$$(\forall x \in]-\rho, \rho[) : f(x) = \sum_{n=0}^{+\infty} a_n \cdot x^n \quad (\text{avec : } a_{n \in \mathbb{N}} \in \mathbb{R}^{\mathbb{N}})$$

alors :

- la fonction f est indéfiniment dérivable sur $]-\rho, \rho[$, et sa dérivée sur cet intervalle n'est autre que la somme de la série des dérivées terme à terme, et on écrit :

$$(\forall x \in]-\rho, \rho[) : f'(x) = \frac{d}{dx} \left(\sum_{n=0}^{+\infty} a_n \cdot x^n \right) = \sum_{n=0}^{+\infty} \frac{d}{dx} (a_n \cdot x^n) = \sum_{n=1}^{\infty} n \cdot a_n \cdot x^{n-1}$$

- la primitive de f sur $]-\rho, \rho[$ et qui s'annule en 0 n'est autre que la somme de la série intégrée terme à terme, et on écrit :

$$(\forall x \in]-\rho, \rho[) : \int_0^x f(t) dt = \int_0^x \sum_{n=0}^{+\infty} a_n \cdot t^n dt = \sum_{n=0}^{+\infty} \int_0^x a_n \cdot t^n dt = \sum_{n=0}^{\infty} \frac{a_n}{n+1} \cdot x^{n+1}$$

Preuve :

- Pour montrer le premier résultat, on commence d'abord par remarquer que la série $\sum_{n \geq 0} a_n \cdot z^n$ et la série dérivée $\sum_{n \geq 0} n \cdot a_n \cdot z^{n-1}$ ont des rayons de convergence identiques, en effet :

$$\lim_{n \rightarrow \infty} \frac{n+1}{n} = 1 \text{ et } \lim_{n \rightarrow \infty} \frac{a_{n+1}}{a_n} = R \implies \lim_{n \rightarrow \infty} \frac{(n+1) \cdot a_{n+1}}{n \cdot a_n} = "1 \times R" = R$$

le lemme 3.2.1 entraîne que la convergence est uniforme sur tout intervalle $[-\rho, \rho]$ inclus dans $]-R, R[$, et bien sûr si une suite de fonctions dérivables converge uniformément sur un intervalle, ainsi que la suite de ses dérivées, alors la limite de la suite sera dérivable à l'intérieur de cet intervalle, et sa dérivée ne va être que la limite des dérivées.

ceci implique que la fonction f est dérivable sur tout intervalle $]-\rho, \rho[\subset]-R, R[$ donc sur $]-R, R[$ tout entier, et par récurrence, f sera donc indéfiniment dérivable sur $]-R, R[$.

- pour le second point du théorème il suffit d'appliquer le résultat de dérivabilité à la série primitive.

d'où ce qu'il faut prouver. ■

Définition 2.2.2 : (série de Maclaurin)

Si f est une fonction indéfiniment dérivable sur un intervalle $]-R, R[$, alors on appelle **série de Maclaurin** associée à f en 0 la série entière définie par :

$$S_f = \sum_{n \geq 0} \frac{f^{(n)}(0)}{n!} \times x^n.$$

Corollaire 2.2.1 :

Si f est **une fonction développable en série entière** sur un intervalle $]-\rho, \rho[$, alors :

1. elle sera indéfiniment dérivable au voisinage de 0, et son développement en 0 sur cet intervalle est :

$$f(x) = \sum_{n=0}^{+\infty} \frac{f^{(n)}(0)}{n!} x^n = f(0) + f'(0)x + \dots + \frac{f^{(n)}(0)}{n!} x^n + \dots$$

2. elle admettra un développement limité en 0 à tout ordre, et dans ce cas là, et pour un certain ordre fixé : $n \in \mathbb{N}$ il sera la somme partielle d'ordre n de la série de Maclaurin de f :

$$f(x) = f(0) + f'(0) \cdot x + \dots + \frac{f^{(n)}(0)}{n!} \cdot x^n + o(x^n).$$

Preuve :

1. Posons :

$$(\forall x \in]-\rho, \rho[) : f(x) = \sum_{n=0}^{+\infty} a_n \cdot x^n$$

par une application répétée du premier point de théorème IV.3.2.1, on déduit que pour tout $k \in \mathbb{N}$ la $k^{\text{ème}}$ dérivée de f est donnée par :

$$f^{(k)}(z) = \sum_{n=k}^{+\infty} n \times (n-1) \cdots \times (n-k+1) \times a_n \times z^{n-k}.$$

et dans ce cas, le premier terme est : $k! \times a_k = f^{(k)}(0)$, ce qui fait que $f \in C^\infty(\mathbb{R})$ et $a_k = \frac{f^{(k)}(0)}{k!}$.

2. Il suffit d'écrire le développement en séries entières correctement et ça ira de soi par la suite :

$$(\forall x \in]-\rho, \rho[) : f(x) = \sum_{k=0}^{+\infty} \frac{f^{(k)}(0)}{k!} \times x^k = f(0) + f'(0) \cdot x + \dots + \frac{f^{(n)}(0)}{n!} \cdot x^n + \underbrace{\left(x^n \times \sum_{k=n+1}^{+\infty} \frac{f^{(k)}(0)}{k!} \cdot x^{k-n} \right)}_{=o(x^n)}.$$

d'où le résultat désiré. ■

Remarque 2.2.2 :

la réciproque du deuxième point dans le dernier théorème est fautive en général, or, ***il se peut qu'une fonction f admette un développement limité à tous ordres au voisinage de 0, sans qu'elle soit somme de sa série de Maclaurin***, un exemple classique de ce fait est la fonction définie par :

$$(\forall x \in \mathbb{R}) : \varphi(x) = \begin{cases} \text{sign}(x) \cdot \exp\left(-\frac{1}{x^2}\right) & \text{si } x \neq 0 \\ 0 & \text{sinon.} \end{cases}$$

pour tout : $n \geq 0$, on a :

$$|x^{-n} \times \varphi(x)| \leq |x|^{-n} \times \exp\left(-\frac{1}{x^2}\right) \xrightarrow{x \rightarrow 0} 0.$$

on en déduit donc que φ ainsi que toutes ses dérivées sont nulles en 0, et que le développement limité de f à tout ordre est nul aussi, pourtant f n'est pas du tout identiquement nulle ^a.

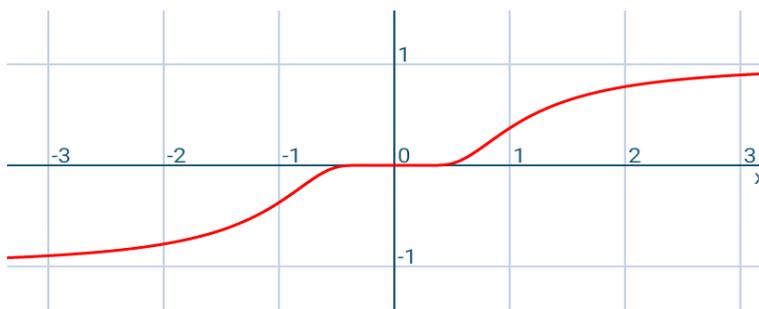


FIGURE IV.7: La courbe de la fonction : $x \mapsto \varphi(x)$.

a. Dans le cas où vous n'êtes pas bien convaincus, une deuxième lecture de cette remarque sera bien souhaitable.

Maintenant, et pour garder la véracité de la réciproque nous allons essayer de donner une condition suffisante pour laquelle une fonction indéfiniment dérivable sera forcément développable en série entière.
pour cela on fera appel à la formule classique de Taylor avec reste intégral.

Théorème 2.2.2 : (formule de Taylor avec reste intégrale :)

Soit f une fonction de classe C^{n+1} sur un intervalle $] -R, R[$ alors on a :

$$(\forall x \in] -R, R[) : f(x) = f(0) + f'(0) \cdot x + \dots + \frac{f^{(n)}(0)}{n!} \cdot x^n + \int_0^x \frac{(x-t)^n}{n!} f^{(n+1)}(t) dt .$$

Preuve :

On procédera par récurrence sur n , pour $n = 0$, la formule est vraie car :

$$(\forall x \in] -R, R[) : f(x) = f(0) + (f(x) - f(0)) = f(0) + \int_0^x f'(t) dt .$$

soit maintenant $n \geq 0$, on suppose que :

$$(\forall x \in] -R, R[) : f(x) = f(0) + f'(0) \cdot x + \dots + \frac{f^{(n)}(0)}{n!} \cdot x^n + \int_0^x \frac{(x-t)^n}{n!} f^{(n+1)}(t) dt$$

alors et pour tout $x \in] -R, R[$:

$$\begin{aligned} \int_0^x \frac{(x-t)^{n+1}}{(n+1)!} \times f^{(n+2)}(t) dt &= \left[\frac{(x-t)^{n+1}}{(n+1)!} \times f^{(n+1)}(t) \right]_0^x - \int_0^x \frac{-(x-t)^n}{n!} \times f^{(n+1)}(t) dt \\ &= -\frac{f^{(n+1)}(0)}{(n+1)!} \times x^{n+1} + \int_0^x \frac{(x-t)^n}{n!} \times f^{(n+1)}(t) dt \end{aligned}$$

(et selon l'hypothèse de récurrence)
$$= -\frac{f^{(n+1)}(0)}{(n+1)!} \times x^{n+1} + f(x) - \left(f(0) + f'(0) \cdot x + \dots + \frac{f^{(n)}(0)}{n!} \cdot x^n \right)$$

donc :

$$(\forall x \in] -R, R[) : f(x) = f(0) + f'(0) \cdot x + \dots + \frac{f^{(n+1)}(0)}{(n+1)!} \cdot x^{n+1} + \int_0^x \frac{(x-t)^{n+1}}{(n+1)!} f^{(n+2)}(t) dt$$

d'où ce qu'il faut démontrer. ■

Corollaire 2.2.2 :

Si f est indéfiniment dérivable sur $] -R, R[$ et s'il existe $M > 0$ et $a \in] 0, +\infty[$ positifs tels que pour tout n et pour tout $x \in] -R, R[$,

$$|f^{(n)}(x)| \leq M \times a^n .$$

alors f est développable en série entière en 0 sur $] -R, R[$, et son développement sera bien sûr la somme de sa série de Maclaurin.

Preuve :

L'idée de la preuve consiste à utiliser la formule de Taylor avec reste intégral et de majorer ce reste comme suit :

$$(\forall x > 0) : \left| \int_0^x \frac{(x-t)^n}{n!} f^{(n+1)}(t) dt \right| \leq \int_0^x (x-t)^n \frac{M \cdot a^n}{n!} dt = \frac{x^{n+1} M}{(n+1)!} \times a^n \leq \underbrace{\frac{(R \cdot a)^{n+1}}{(n+1)!}}_{\xrightarrow{n \rightarrow +\infty} 0} \times \frac{M}{a} \xrightarrow{n \rightarrow +\infty} 0.$$

et pour le cas de $x \leq 0$ un raisonnement analogue peut garantir que cet intégrale reste encore convergeant vers 0.

et ceci entraînera le résultat désiré ^a. ■

► Maintenant, et en restant encore dans le cadre de la reformulation de notre problème d'origine (\mathcal{E}_a), on supposera pour la suite que :

$$\boxed{(\exists M > 0) (\exists a \in]0, +\infty[) (\forall x \in]-\infty, +\infty[) : \left| \sigma^{(n)}(x) \right| \leq M \times a^n. \quad (\diamond\diamond)}$$

on a et selon (\diamond) la fonction : $x \mapsto \sum_{k=0}^{+\infty} \frac{\sigma^{(k)}(0)}{k!} \times x^k$ est bien définie sur \mathbb{R} toute entière.

et selon ($\diamond\diamond$) et le corollaire : IV.2.2.2, σ sera développable en série entière en 0 sur $\mathbb{R} =]-\infty, \overbrace{+\infty}^{=R_\Phi}[$, et son développement ne va être que la somme de sa série de Maclaurin, donc on peut écrire dans ce cas :

$$\boxed{(\forall x \in \mathbb{R}) : \sigma(x) = \sum_{k=0}^{+\infty} \frac{\sigma^{(k)}(0)}{k!} \times x^k.}$$

alors, et en remplaçant σ par ce développement dans (\blacktriangledown) on aura, et pour tout $t \in [0, T]$:

$$\boxed{\left[{}_0D_t^\alpha \left(t \times \sum_{i=1}^I \omega_i^{(2)} \times \sum_{k=0}^{+\infty} \frac{\sigma^{(k)}(0)}{k!} \times \left(\omega_i^{(1)} \times t + b_i \right)^k \right) \right] + h(t) \times \left(\beta + t \times \sum_{i=1}^I \omega_i^{(2)} \times \sum_{k=0}^{+\infty} \frac{\sigma^{(k)}(0)}{k!} \times \left(\omega_i^{(1)} \times t + b_i \right)^k \right) = s(t)}$$

or, et selon la formule de binôme de newton on a ^b :

$$\left(\omega_i^{(1)} \times t + b_i \right)^k = \sum_{j=0}^k C_k^j \cdot \left(\omega_i^{(1)} \right)^j \cdot t^j \cdot (b_i)^{n-j}.$$

d'où on obtient d'une part que :

$$\begin{aligned} t \times \sum_{i=1}^I \omega_i^{(2)} \times \sum_{k=0}^{+\infty} \frac{\sigma^{(k)}(0)}{k!} \times \left(\omega_i^{(1)} \times t + b_i \right)^k &= t \times \sum_{i=1}^I \omega_i^{(2)} \times \sum_{k=0}^{+\infty} \overbrace{\frac{\sigma^{(k)}(0)}{k!}}{:=\lambda_k} \times \sum_{j=0}^k C_k^j \cdot \left(\omega_i^{(1)} \right)^j \cdot t^j \cdot (b_i)^{n-j} \\ &= \sum_{i=1}^I \sum_{k=0}^{+\infty} \sum_{j=0}^k \lambda_k \cdot \omega_i^{(2)} \cdot C_k^j \cdot \left(\omega_i^{(1)} \right)^j \cdot t^{j+1} \cdot (b_i)^{n-j} \end{aligned}$$

a. On rappelle que : $(\forall a \in \mathbb{R}^{+*}) : \frac{a^n}{n!} \xrightarrow{n \rightarrow +\infty} 0$ (il suffit d'utiliser le critère d'Alembert pour les suites.).

b. Ici C_k^j désigne le nombre de combinaisons de k éléments parmi n et : $C_k^j = \frac{k!}{j! \times (k-j)!}$.

et d'une autre :

$$\left[{}_0D_t^\alpha \left(t \cdot \sum_{i=1}^I \omega_i^{(2)} \cdot \sum_{k=0}^{+\infty} \lambda_k \cdot (\omega_i^{(1)} \cdot t + b_i)^k \right) \right] = \left[{}_0D_t^\alpha \left(\sum_{i=1}^I \sum_{k=0}^{+\infty} \sum_{j=0}^k \lambda_k \cdot \omega_i^{(2)} \cdot C_k^j \cdot (\omega_i^{(1)})^j \cdot t^{j+1} \cdot (b_i)^{n-j} \right) \right]$$

$$= \sum_{i=1}^I \sum_{k=0}^{+\infty} \sum_{j=0}^k \lambda_k \cdot \omega_i^{(2)} \cdot C_k^j \cdot (\omega_i^{(1)})^j \cdot (b_i)^{n-j} \left[{}_0D_t^\alpha (t^{j+1}) \right]$$

$$\left(\text{de l'exemple : IV.1.3.1 et car : } \overbrace{[\alpha] = 1}^{\alpha \in [0,1[} \leq n \right) = \sum_{i=1}^I \sum_{k=0}^{+\infty} \sum_{j=0}^k \lambda_k \cdot \omega_i^{(2)} \cdot C_k^j \cdot (\omega_i^{(1)})^j \cdot (b_i)^{n-j} \cdot \left(\frac{\Gamma(j+2)}{\Gamma(j+2-\alpha)} \cdot t^{j+1-\alpha} \right)$$

donc d'après (▼), et en vertu de tout ce qui était présenté auparavant, on peut déduire que quel que soit $t \in [0, T]$ on a :

$$\sum_{i=1}^I \sum_{k=0}^{+\infty} \sum_{j=0}^k \lambda_k \cdot \omega_i^{(2)} \cdot C_k^j \cdot (\omega_i^{(1)})^j \cdot (b_i)^{n-j} \left(\frac{\Gamma(j+2)}{\Gamma(j+2-\alpha)} t^{j+1-\alpha} \right) + h(t) \cdot \left(\beta + \sum_{i=1}^I \sum_{k=0}^{+\infty} \sum_{j=0}^k \lambda_k \cdot \omega_i^{(2)} \cdot C_k^j \cdot (\omega_i^{(1)})^j \cdot t^{j+1} \cdot (b_i)^{n-j} \right) = s(t)$$

et par une simple factorisation on peut réécrire cette équation ainsi :

$$\sum_{i=1}^I \sum_{k=0}^{+\infty} \sum_{j=0}^k \lambda_k \cdot \omega_i^{(2)} \cdot C_k^j \cdot (\omega_i^{(1)})^j \cdot (b_i)^{n-j} \left(\frac{\Gamma(j+2)}{\Gamma(j+2-\alpha)} t^{j+1-\alpha} + h(t) \cdot t^{j+1} \right) + \beta \cdot h(t) = s(t)$$

où $\lambda_k \in \mathbb{R}$ désigne ici : $\frac{\sigma^{(k)}(0)}{k!}$.

en notant : $\zeta_{i,j}^k = \lambda_k \cdot \omega_i^{(2)} \cdot C_k^j \cdot (\omega_i^{(1)})^j \cdot (b_i)^{n-j}$ on obtient :

$$\sum_{i=1}^I \sum_{k=0}^{+\infty} \sum_{j=0}^k \zeta_{i,j}^k \times \left(\frac{\Gamma(j+2)}{\Gamma(j+2-\alpha)} t^{j+1-\alpha} + h(t) \cdot t^{j+1} \right) + \beta \cdot h(t) = s(t) \quad (\blacktriangledown\blacktriangledown)$$

▷ Afin d'achever notre procédure de discrétisation, on va essayer bien évidemment de construire un maillage convenable pour notre intervalle d'étude $\Omega = [0, T]$, et pour cela, on prendra un entier positif $P \in \mathbb{N}^*$, et on adopte une partition régulière^a et qui va se déterminer tout simplement via les points nodaux suivants :

$$\left(\forall m \in \llbracket 0, P \rrbracket \right) : t_m = m \times \frac{T}{P}.$$

et dans ce cas, notre maillage sera clairement donné par la famille : $(\Omega_m)_{0 \leq m \leq P-1}$ avec :

$$\left(\forall m \in \{0, \dots, P-1\} \right) : \Omega_m = [t_m, t_{m+1}].$$

on rappelle que les éléments Ω_m , et qui portent habituellement le nom des mailles, sont moins nombreux d'une seule unité que les point nodaux, voir pour cela la figure IV.8.

a. C'est à dire que toutes les mailles constituant cette partition sont de même amplitude.

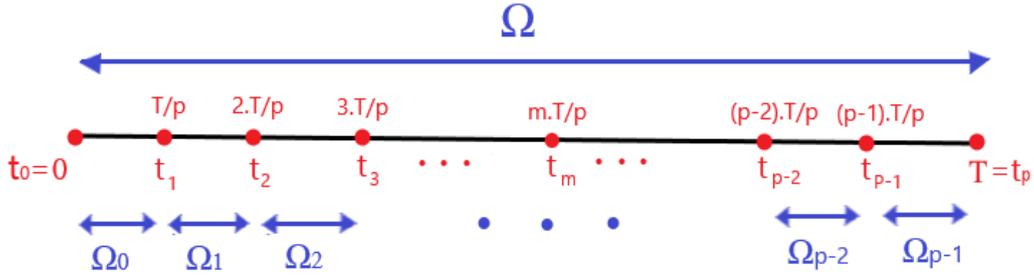


FIGURE IV.8: Maillage proposé du domaine d'étude $\Omega = [0, T]$.

à l'instar des procédures numériques usuelles, la substitution de $+\infty$ par un ordre de limitation^a acceptablement grand, et l'affectation des valeurs nodales t_m à la variable t dans ($\blacktriangledown\blacktriangledown$) va ramener notre problème d'origine (\mathcal{E}_α) à sa forme finale, et qui se présente dans notre cas par le système des $P + 1$ équations suivant :

$$\left(\forall m \in \{0, \dots, P\} \right) : \sum_{i=1}^I \sum_{k=0}^N \sum_{j=0}^k \zeta_{i,j}^k \times \left(\frac{\Gamma(j+2)}{\Gamma(j+2-\alpha)} (t_m)^{j+1-\alpha} + h(t_m) \cdot (t_m)^{j+1} \right) + \beta \cdot h(t_m) = s(t_m) \quad (\mathfrak{S})$$

avec : $\zeta_{i,j}^k = \lambda_k \cdot \omega_i^{(2)} C_k^j \left(\omega_i^{(1)} \right)^j (b_i)^{n-j}$ et $\lambda_k = \frac{\sigma^{(k)}(0)}{k!}$.

en principe, notre but dans la suite sera de trouver les pondérations : $\hat{\omega}_i^{(1)} \in \mathbb{R}^I$ et $\hat{\omega}_i^{(2)} \in \mathbb{R}^I$ qui résolvent ce système^b.

3 Principe de résolution par un processus d'apprentissage :

3.1 Passage à un problème d'optimisation :

Dans ce qui précède, nous avons vu comment on a pris notre problème d'origine (\mathcal{E}_α), et on l'a ramené, après une longue procédure de discrétisation, à une forme plus simplifiée relativement, et qui s'exprime par le système d'équations (\mathfrak{S}), mais malgré cela, notre but de trouver une bonne approximation à la solution présumée de notre problème reste encore non achevé, et pour le faire il faut trouver comme on a déjà dit les vecteurs résolvant ce système, or et même au cas de leur existence, cette dernière tâche reste très difficile, si elle n'est pas impossible, à réaliser analytiquement, ce qui fait qu'on sera obligés de procéder numériquement à travers la générations d'une suite itérative de pondérations qui convergent vers les solutions : $\hat{\omega}_i^{(1)} \in \mathbb{R}^I$ et $\hat{\omega}_i^{(2)} \in \mathbb{R}^I$ qu'on les suppose existantes dans ce contexte.

maintenant, la question qui se pose encore, est comment déterminer la façon qu'on doit adopter pour effectuer cette génération, et surtout que le système (\mathfrak{S}) qu'on traite n'est pas linéaire, ou d'un autre type classique dont lequel on dispose déjà d'une gamme d'algorithmes de résolution prête à l'avance^c, et consacrée pour ce fait ?.

la réponse qu'on propose sera de réécrire notre système d'équations sous la forme d'un problème d'optimisation convenable, puis le résoudre par un algorithme classique comme celui de descente de gradient par exemple.

Dans la suite on va justifier d'abord l'existence d'une solution exacte pour notre problème d'origine (\mathcal{E}_α), et après on rentrera dans les détails du passage de (\mathfrak{S}) à un problème d'optimisation (\mathfrak{P}) qu'on le déterminera ultérieurement, et on verra que l'application d'un algorithme de

a. On désigne ici l'ordre de limitation des séries de Maclaurin.

b. On signale que le problème d'existence de ces pondérations, est bien présent dans ce cas.

c. Comme l'algorithme de Jacobi ou de Gauss-Seidel pour les systèmes linéaires usuels (voir :[7]).

minimisation itérative dans le cadre de ce problème *va nous conduire vers une bonne approximation de la solution exacte de (\mathcal{E}_α) et ceci même si ce problème (\mathbb{P}) n'admet pas de solutions.*

3.1.1 Existence d'une solution exacte :

Théorème 3.1.1 : (d'existence d'une solution pour un problème différentiel fractionnaire, [28]) :

Considérons le problème différentiel défini par :

$$(\mathcal{D}_\alpha) \begin{cases} \left[{}_0D_t^\alpha (V(t) + \Psi(t, V(t))) \right] = \langle A(t, V(t)), V(t) \rangle + f(t, V(t)) & \text{si : } t > 0. \\ V(0) = V_0 \in \mathbb{R}. & \text{si : } t = 0. \end{cases}$$

avec :

- $(X, \|\cdot\|)$ est un espace de Banach .
- $[{}_0D_t^\alpha (\cdot)]$ désigne la dérivée fractionnaire supérieure d'ordre $\alpha \in]0, 1]$ au sens de Caputo.
- $f : J \times X \rightarrow X$, $\Psi : J \times X \rightarrow X$ sont des fonctions données où $J = [0, +\infty[$.
- $A(\cdot, \cdot) : J \times X \rightarrow \mathcal{B}(X)$ est une application qui associe à chaque couple $(t, x) \in J \times X$ un opérateur linéaire borné^a $A(t, x) \in \mathcal{B}(X)$ définie de X dans X .

si les hypothèses suivantes sont vérifiées :

(H₁) : la fonction $(t, x) \mapsto A(t, x) \in \mathcal{B}(X)$ est continue et il existe deux fonctions continues, bornées et non négatives $\varphi(\cdot)$ et $\psi(\cdot)$ définies sur $J = \mathbb{R}^+$ telles que :

$$\begin{cases} (\forall x \geq 0) : 0 \leq \psi(x) \leq x + 1. \\ \text{et :} \\ (\forall t \geq 0) (\forall x \in X) : \|A(t, x)\|_{\mathcal{B}(X)} \leq \frac{\varphi(t)}{1 + t^{\alpha+1}} \times \psi\left(\frac{\|x\|}{1 + t^{\alpha+1}}\right). \end{cases}$$

(H₂) : la fonction $f : J \times X \rightarrow X$ est continue, et il existe deux fonctions continues, bornées et non négatives $a(\cdot)$, $b(\cdot)$ définies sur J telles que :

$$(\forall t \geq 0) (\forall x \in X) : \|f(t, x)\| \leq \frac{a(t)}{1 + t^{\alpha+1}} \|x\| + b(t).$$

(H₃) : la fonction $\Psi : J \times X \rightarrow X$ est continue et il existe une constante $\zeta \geq 0$ telle que :

$$\begin{cases} (\forall t \geq 0) (\forall (u, v) \in X^2) : \|\Psi(t, u) - \Psi(t, v)\| \leq \zeta \times \|u - v\|. \\ \text{et :} \\ \delta = \sup_{t \geq 0} \|\Psi(t, 0)\| < +\infty. \end{cases}$$

(H₄) : il existe $\rho > 0$ tel que :

$$\begin{cases} C_0(\rho) + C_1(\rho) < \rho \quad \text{où :} \\ C_0(\rho) = \|V_0\| + \|\Psi(0, V(0))\| + \delta + \zeta \times \rho. \\ C_1(\rho) = \frac{1}{\Gamma(\alpha+1)} (\|\varphi\|_\infty \rho^2 + \|a\|_\infty \rho + \|b\|_\infty). \end{cases}$$

alors le problème (\mathcal{D}_α) admet au moins une solution.

a. C'est à dire que c'est une application linéaire entre deux espaces vectoriels normés X et Y (dans ce cas $X = Y$) telle que l'image de la boule unité de X est une partie bornée de Y , on peut montrer facilement qu'ils s'identifient aux applications linéaires continues de X dans Y .

Preuve :

La démonstration de ce résultat se base essentiellement sur deux idées principales, la première est de mettre le système (\mathcal{D}_α) sous la forme d'un problème du point fixe, et la deuxième consiste à appliquer le fameux **théorème du point fixe de Schauder**^a pour garantir l'existence d'une solution à ce problème.

► **se ramener à un problème du point fixe**: supposons que $V(\cdot)$ est une solution du problème (\mathcal{D}_α) alors on a :

$$(\forall t > 0) : \left[{}_0D_t^\alpha \left(V(t) + \Psi(t, V(t)) \right) \right] = \left\langle A(t, V(t)), V(t) \right\rangle + f(t, V(t)).$$

en appliquant l'opérateur d'intégration fractionnaire $[{}_0P_x^\alpha(\cdot)]$ aux deux membres de l'équation ci-dessus nous obtenons :

$$\begin{aligned} (\forall x > 0) : \quad V(x) + \Psi(x, V(x)) + \widehat{\lambda}^{\text{constante}} &= \left[{}_0P_x^\alpha \left(\left\langle A(t, V(t)), V(t) \right\rangle + f(t, V(t)) \right) \right] \\ &= \frac{1}{\Gamma(\alpha)} \int_0^x (x-t)^{\alpha-1} \times \left(\left\langle A(t, V(t)), V(t) \right\rangle + f(t, V(t)) \right) dt. \end{aligned}$$

maintenant, et en utilisant la condition initiale : $V(0) = V_0$ on trouve :

$$\lambda = -V(0) - \Psi(0, V(0)) = -V_0 - \Psi(0, V(0))$$

ce qui fait que :

$$\begin{aligned} (\forall x > 0) : \quad V(x) + \Psi(x, V(x)) + \lambda &= V(x) + \Psi(x, V(x)) - V_0 - \Psi(0, V(0)) \\ &= \frac{1}{\Gamma(\alpha)} \int_0^x (x-t)^{\alpha-1} \times \left(\left\langle A(t, V(t)), V(t) \right\rangle + f(t, V(t)) \right) dt. \end{aligned}$$

donc nous obtenons enfin de compte, et pour tout $x > 0$ l'égalité suivante :

$$V(x) = -\Psi(x, V(x)) + V_0 + \Psi(0, V(0)) + \frac{1}{\Gamma(\alpha)} \int_0^x (x-t)^{\alpha-1} \cdot \left(\left\langle A(t, V(t)), V(t) \right\rangle + f(t, V(t)) \right) dt$$

et qui s'appelle **l'équation intégrale de Volterra**, (Vito Volterra, italien, 1860 1940).

Réciproquement, si $V(\cdot)$ est une solution de l'équation de Volterra, alors on doit juste procéder inversement en appliquant l'opérateur différentiel $[{}_0D_t^\alpha(\cdot)]$ à ses deux membres, pour déduire que notre fonction $V(\cdot)$ sera forcément une solution de l'équation différentielle (\mathcal{D}_α) donc on peut dire qu'il y a une équivalence entre ce problème et l'équation de Volterra.

► **utiliser le théorème du point fixe de Schauder**: dans le but de démontrer l'existence d'une solution de l'équation intégrale de Volterra on introduit le sous-ensemble :

$$B_\rho := \left\{ V \in C^1(\mathbb{R}^+) / \|V\|_\infty \leq \rho \right\}.$$

avec $\rho > 0$ est la constante positive introduite dans l'hypothèse : (\mathbf{H}_4) .

on considère sur B_ρ l'opérateur \mathcal{F} défini par :

$$\begin{aligned} \mathcal{F} : B_\rho &\longrightarrow C^1(\mathbb{R}^+) \\ g &\longmapsto \mathcal{F}(g) = V_0 + \Psi(0, g(0)) - \Psi(\cdot, g(\cdot)) + \frac{1}{\Gamma(\alpha)} \int_0^\bullet (\cdot-t)^{\alpha-1} \times \Phi(t, g(t)) dt \end{aligned}$$

a. Ce théorème a été présenté pour la première fois en 1930 par le mathématicien polonais **Juliusz Pawel Schauder** (1899-1943) dans son article[55], et il énonce que si C est un **convexe fermé** non vide d'un \mathbb{R} - espace vectoriel topologique séparé $(E, +, \cdot, \tau)$, alors toute application L continue de C dans C telle que $L(C)$ est relativement compact, c'est à dire $L(C)$ est compact, admettra forcément un point fixe.

l'importance de ce théorème (dans cette version) réside dans le fait qu'**il ne nécessite pas que l'espace E soit de Banach**, (voir [9] page : 79).

avec :

$$\Phi(t, g(t)) = \langle A(t, g(t)), g(t) \rangle + f(t, g(t))$$

sur cet ensemble, l'équation intégrale de Volterra peut être écrite sous la forme de problème du point fixe suivant :

$$\boxed{\text{trouver un } V \in B_\rho \text{ tel que : } \mathcal{F}(V) = V.}$$

et qu'on va montrer le fait qu'il admet au moins une solution en s'appuyant sur le théorème du point fixe de Schauder dont on a parlé précédemment.

maintenant et pour appliquer ce théorème, et déduire l'existence d'une solution, il faut vérifier la réalisation des trois hypothèses suivantes :

- i. B_ρ est un convexe fermé non vide.
- ii. B_ρ est stable par \mathcal{F} , ou autrement dit : $\mathcal{F}(B_\rho) \subset B_\rho$.
- iii. \mathcal{F} est continue et $\mathcal{F}(B_\rho)$ est relativement compact.

— **La première condition** : il est clair que l'ensemble B_ρ est une boule fermée de centre 0 et de rayon ρ pour la norme de la convergence uniforme, donc elle sera clairement une partie non vide (car elle contiendra au moins la fonction nulle), convexe et fermée de $C^1(\mathbb{R}^+)$ muni de cette norme.

— **La deuxième condition** : d'abord notons pour tout $g \in B_\rho$ $\mathcal{F}(g)$ par $(\mathcal{F}g)$, alors dans ce cas on a :

$$\begin{aligned} \left\| \frac{(\mathcal{F}g)(x)}{1+x^{\alpha+1}} \right\| &= \left\| \frac{V_0 + \Psi(0, g(0)) - \Psi(x, g(x))}{1+x^{\alpha+1}} + \frac{1}{\Gamma(\alpha)} \int_0^x \frac{(x-t)^{\alpha-1}}{1+x^{\alpha+1}} \times \Phi(t, g(t)) dt \right\| \\ &\leq \frac{1}{1+x^{\alpha+1}} \left(\|V_0\| + \|\Psi(0, g(0))\| + \|\Psi(x, g(x)) - \Psi(x, 0)\| + \|\Psi(x, 0)\| \right) \\ &\quad + \frac{1}{\Gamma(\alpha)} \int_0^x \frac{(x-t)^{\alpha-1}}{1+x^{\alpha+1}} \times \|\Phi(t, g(t))\| dt \\ \text{selon } (\mathbf{H}_3) &\leq \frac{1}{1+x^{\alpha+1}} \left(\|V_0\| + \|\Psi(0, g(0))\| + \varsigma \times \|g(x) - 0\| + \delta \right) \\ &\quad + \frac{1}{\Gamma(\alpha)} \int_0^x \frac{(x-t)^{\alpha-1}}{1+x^{\alpha+1}} \times \left\| \langle A(t, g(t)), g(t) \rangle + f(t, g(t)) \right\| dt \\ &\leq \frac{1}{1+x^{\alpha+1}} \left(\|V_0\| + \|\Psi(0, g(0))\| + \varsigma \times \rho + \delta \right) \\ &\quad + \frac{1}{\Gamma(\alpha)} \int_0^x \frac{(x-t)^{\alpha-1}}{1+x^{\alpha+1}} \times \left(\|A(t, g(t))\|_{\mathcal{B}(X)} \cdot \|g(t)\| + \|f(t, g(t))\| \right) dt \end{aligned}$$

Nous obtenons grâce aux hypothèses (\mathbf{H}_1) , (\mathbf{H}_2) et (\mathbf{H}_3) l'estimation suivante :

$$\left\| \frac{(\mathcal{F}g)(x)}{1+x^{\alpha+1}} \right\| \leq \frac{C_0(\rho)}{1+x^{\alpha+1}} + C_1(\rho) \cdot \frac{t^\alpha}{1+t^{\alpha+1}}.$$

d'où :

$$\lim_{t \rightarrow +\infty} \frac{(\mathcal{F}g)(x)}{1+x^{\alpha+1}} = 0.$$

et selon l'hypothèse (\mathbf{H}_4) :

$$\|(\mathcal{F}g)\|_\infty \leq C_0(\rho) + C_1(\rho) \leq \rho.$$

ce qui prouve que : $g \in B_\rho \Rightarrow (\mathcal{F}g) \in B_\rho$.

- **La troisième condition** : cette dernière condition est un peu difficile à démontrer, et surtout qu'elle se base à son tour sur un lemme un peu compliqué et qui donne des caractérisation de la compacité relative dans l'espace des fonction continues , pour plus de détails sur la démonstration de cette hypothèse veuillez voir [28] pages :27 – 32.

et ceci achève bien notre preuve . ■

Corollaire 3.1.1 :

Dans le cadre de l'équation différentielle fractionnaire suivante :

$$(\mathcal{E}_\alpha) \begin{cases} \left[{}_0D_t^\alpha (V(t)) \right] + h(t) \times V(t) = s(t). & \text{si : } 0 < t < T. \\ V(t) = V(0) = \beta \in \mathbb{R}. & \text{si : } t = 0. \end{cases}$$

et qui fait effectivement l'objet de notre étude dans ce chapitre, l'existence d'une solution est toujours garantie si la condition suivante est satisfaite :

$$\left(\|\varphi\|_\infty \cdot \left(\|s\|_\infty + \Gamma(\alpha + 1) \cdot |\beta| \right) \right) < \frac{\Gamma(\alpha + 1)}{4}$$

où φ ici désigne la fonctions définie sur $[0, T]$ par : $\varphi(t) = (1 + t^{\alpha+1}) \times |h(t)|$

Preuve :

L'idée pour démontrer ce corollaire est de vérifier en premier lieu la réalisation des quatre hypothèses $(\mathbf{H}_1) - (\mathbf{H}_4)$ mentionnées au théorème IV.3.1.1, et puis appliquer ce dernier pour aboutir immédiatement le résultat désiré.

d'abord commençons par remarquer que le système fractionnaire (\mathcal{E}_α) n'est autre qu'un cas particulier de (\mathcal{D}_α) et pour lequel on a :

- $(X, \|\cdot\|) = (\mathbb{R}, |\cdot|)$ et qui est bien un espace de Banach .
- $V_0 = \beta \in \mathbb{R}$.
- $f : [0, T] \times \mathbb{R} \rightarrow \mathbb{R}$, $\Psi : [0, T] \times \mathbb{R} \rightarrow \mathbb{R}$ sont les fonctions définies par : $f(t, x) = s(t)$ et $\Psi(t, x) = 0$.
- $A(\cdot, \cdot) : [0, T] \times \mathbb{R} \rightarrow \mathcal{B}(\mathbb{R})$ est l'application qui associe à chaque couple $(t, x) \in [0, T] \times \mathbb{R}$ l'application linéaire $A(t, x) := \langle A(t, x), \bullet \rangle = (-h(t) \times \bullet) \in \mathcal{B}(\mathbb{R}) \cong \mathbb{R}$.

alors et dans ce cas on a :

(\mathbf{H}_1) : la fonction $(t, x) \mapsto A(t, x) = (-h(t) \times \bullet) \in \mathcal{B}(\mathbb{R})$ est continue car, et par hypothèse $h(\cdot)$ est continue, en plus et en considérant les fonctions définies sur $[0, T]$ par : $\varphi(t) = (1 + t^{\alpha+1}) \times |h(t)|$ et $\psi(t) = \mathbb{1}$ on aura :

$$\begin{cases} (\forall x \in [0, T]) : 0 \leq \psi(x) = 1 \leq x + 1. \\ \text{et :} \\ (\forall t \in [0, T]) (\forall x \in \mathbb{R}) : \overbrace{\|A(t, x)\|_{\mathcal{B}(\mathbb{R})}}^{=|h(t)|} \leq \frac{\varphi(t)}{1 + t^{\alpha+1}} \times \psi\left(\frac{\|x\|}{1 + t^{\alpha+1}}\right) = \overbrace{\frac{|h(t)| \cdot (1 + t^{\alpha+1})}{1 + t^{\alpha+1}}}^{=|h(t)|} \times 1. \end{cases}$$

(\mathbf{H}_2) : puisque : $s(\cdot)$ est continue la fonction $f : [0, T] \times \mathbb{R} \rightarrow \mathbb{R}$ le sera aussi, et en considérant les fonctions définies sur $[0, T]$ par : $a(\cdot) = 0$ et $b(\cdot) = |s(\cdot)|$ on aura :

$$(\forall t \in [0, T]) (\forall x \in \mathbb{R}) : \underbrace{\|f(t, x)\|}_{=|s(t)|} \leq \frac{a(t)}{1 + t^{\alpha+1}} \|x\| + b(t) = |s(t)|.$$

(H₃) : puisqu'elle est identiquement nulle, la fonction $\Psi : [0, T] \times \mathbb{R} \rightarrow \mathbb{R}$ est continue, et en posant $\varsigma = 0$ on aura facilement :

$$\left\{ \begin{array}{l} (\forall t \in [0, T]) (\forall (u, v) \in \mathbb{R}^2) : \|\Psi(t, u) - \Psi(t, v)\| \leq \varsigma \times \|u - v\|. \\ \text{et :} \\ \delta = \sup_{0 \leq t \leq T} \|\Psi(t, 0)\| = 0 < +\infty. \end{array} \right.$$

(H₄) : Dans ce cas particulier nous avons :

$$\left\{ \begin{array}{l} C_0(\rho) = \overbrace{\|V_0\|}^{=|\beta|} + \overbrace{\|\Psi(0, V(0))\|}^{=0} + \overbrace{\delta}^{=0} + \overbrace{\varsigma \times \rho}^{=0} = |\beta|. \\ C_1(\rho) = \frac{1}{\Gamma(\alpha + 1)} \left(\|\varphi\|_\infty \rho^2 + \underbrace{\|a\|_\infty \rho + \|b\|_\infty}_{=0} \right). \end{array} \right.$$

or, et d'une autre part on a :

$$\begin{aligned} \left(\|\varphi\|_\infty \cdot (\|s\|_\infty + \Gamma(\alpha + 1) \cdot |\beta|) \right) &< \frac{\Gamma(\alpha + 1)}{4} \Rightarrow 1 - 4 \times \left(\frac{1}{\Gamma(\alpha + 1)} \|\varphi\|_\infty \cdot \left(\frac{1}{\Gamma(\alpha + 1)} \times \|s\|_\infty + |\beta| \right) \right) < 0 \\ &\Rightarrow (\exists x_0 > 0) : \frac{1}{\Gamma(\alpha + 1)} \|\varphi\|_\infty x_0^2 - x_0 + \frac{1}{\Gamma(\alpha + 1)} \times \|s\|_\infty + |\beta| = 0 \end{aligned}$$

ce qui fait qu'il existera forcément un $\rho > 0$ au voisinage de x_0 tel que :

$$C_0(\rho) + C_1(\rho) = |\beta| + \frac{1}{\Gamma(\alpha + 1)} \left(\|\varphi\|_\infty \rho^2 + \|s\|_\infty \right) < \rho$$

donc et selon le théorème IV.3.1.1, l'équation (\mathcal{E}_α) admettra forcément une solution. ■

Remarques 3.1.1 :

1. Le corollaire précédent donne une condition suffisante pour l'existence d'une solution au problème (\mathcal{E}_α) , mais faites attention elle n'est pas du tout nécessaire, ce qui fait qu'on peut bien trouver des cas où la solution existe sans que cette condition soit vérifiée.
2. malgré sa nécessité théorique, la question d'unicité de cette solution, n'est pas très importante dans notre cadre d'étude, et surtout qu'on cherche juste à fournir une approximation numérique d'une fonction qui résout notre problème, mais si le lecteur est intéressé par cette question, on peut lui suggérer d'aller voir le théorème 2.3.1 page 36 de la source [28] et qui se base sur six hypothèses pour garantir à la fois l'existence et l'unicité d'une solution au problème différentiel quasi-linéaire (\mathcal{D}_α) et même dans un cas plus général où la condition initiale est non locale ^a.

3.1.2 Détails du passage de système (\mathfrak{S}) à un problème d'optimisation :

Dans la section précédente on a remis notre problème d'origine (\mathcal{E}_α) sous sa forme discrétisée, et qui se présente comme on a déjà signalé par le système des $P + 1$ équations suivant :

$$\left(\forall m \in \{0, \dots, P\} \right) : \sum_{i=1}^I \sum_{k=0}^N \sum_{j=0}^k \zeta_{i,j}^k \times \left(\frac{\Gamma(j+2)}{\Gamma(j+2-\alpha)} (t_m)^{j+1-\alpha} + h(t_m) \cdot (t_m)^{j+1} \right) + \beta \cdot h(t_m) = s(t_m)$$

a. La condition non locale est une condition donnée sur la solution en un nombre fini ou infini d'instants, celle ci est présentée sous forme de la relation $V(0) = V_0 + g(V)$ liant $V(0)$ et d'autres valeurs de la solution, par exemple : $g(V) = \sum_{k=0}^P c_i V(\tau_i)$ où $(c_i)_{1 \leq i \leq P}$ sont des constantes données et $0 < \tau_1 < \dots < \tau_P < T$ sont des instants connus.

avec :

$$\zeta_{i,j}^k = \lambda_k \cdot \omega_i^{(2)} C_k^j \left(\omega_i^{(1)} \right)^j (b_i)^{n-j} \text{ et } \lambda_k = \frac{\sigma^{(k)}(0)}{k!}.$$

la première chose à connaître sur ce système, c'est ses inconnus, et qui sont dans ce cas les pondérations $\omega^{(1)} \in \mathbb{R}^I, \omega^{(2)} \in \mathbb{R}^I$ et les biais $b \in \mathbb{R}^I$, la deuxième chose, est de garder à l'esprit que ce système n'admet pas toujours des solutions et même s'il est le cas l'unicité ne sera pas forcément garantie, la troisième et la dernière chose est de remarquer qu'il n'est pas linéaire, ce qui rend sa résolution d'un point de vue numérique une tâche un peu plus difficile et surtout avec l'absence d'un arsenal algorithmique consacré pour cela, ce fait va nous mettre sous l'obligation de ramener ce système à un problème d'optimisation de la façon suivante :

d'abord et afin d'alléger nos prochaines écritures, on adoptera et pour tout : $m \in \llbracket 0, P \rrbracket$ la notation suivante :

$$\left(\forall \left(\omega^{(1)}, \omega^{(2)}, b \right) \in \left(\mathbb{R}^I \right)^3 \right) : \Lambda_m \left(\omega^{(1)}, \omega^{(2)}, b \right) = \sum_{i=1}^I \sum_{k=0}^N \sum_{j=0}^k \zeta_{i,j}^k \times \left(\frac{\Gamma(j+2)}{\Gamma(j+2-\alpha)} (t_m)^{j+1-\alpha} + h(t_m) \cdot (t_m)^{j+1} \right) + \beta \cdot h(t_m)$$

nous avons :

$$\begin{aligned} \left(\omega^{(1)}, \omega^{(2)}, b \right) \in \left(\mathbb{R}^I \right)^3 \text{ est une solution de } (\mathfrak{S}) &\iff \left(\forall m \in \{0, \dots, P\} \right) : \Lambda_m \left(\omega^{(1)}, \omega^{(2)}, b \right) = s(t_m) \\ &\iff \left(\forall m \in \{0, \dots, P\} \right) : \Lambda_m \left(\omega^{(1)}, \omega^{(2)}, b \right) - s(t_m) = 0 \\ &\iff \left(\forall m \in \{0, \dots, P\} \right) : \left(\Lambda_m \left(\omega^{(1)}, \omega^{(2)}, b \right) - s(t_m) \right)^2 = 0 \\ &\iff \sum_{m=0}^P \left(\Lambda_m \left(\omega^{(1)}, \omega^{(2)}, b \right) - s(t_m) \right)^2 = 0 \\ &\implies \left(\omega^{(1)}, \omega^{(2)}, b \right) \in \left(\mathbb{R}^I \right)^3 \text{ est une solution du problème} \\ &\hspace{15em} \text{d'optimisation : } (\mathfrak{P}) \end{aligned}$$

avec (\mathfrak{P}) est donné par :

$$(\mathfrak{P}) \left\{ \begin{array}{l} \min E \left(\omega^{(1)}, \omega^{(2)}, b \right) = \frac{1}{2} \times \sum_{m=0}^P \underbrace{\left(\Lambda_m \left(\omega^{(1)}, \omega^{(2)}, b \right) - s(t_m) \right)^2}_{\text{l'erreur locale : } E_m(\omega)} \\ \left(\omega^{(1)}, \omega^{(2)}, b \right) \in \left(\mathbb{R}^I \right)^3 \end{array} \right.$$

à ce stade, la résolution de ce problème, fera notre objectif pour la suite, et ceci car la détermination de son ensemble des points optimaux fournit une très bonne localisation des pondérations résolvants (\mathfrak{S}) dans $(\mathbb{R}^I)^3$, or et comme on peut le remarquer facilement, l'inclusion, $\mathfrak{S} \subset \mathfrak{S}_{\mathfrak{P}}$ est toujours vérifiée^a.

Remarque 3.1.2 :

On peut avoir la réciproque de la dernière étape de notre passage de (\mathfrak{S}) à (\mathfrak{P}) , mais en ajoutant l'hypothèse d'existence d'une solution à (\mathfrak{S}) , c'est à dire que :

$$\left\{ \begin{array}{l} \left(\omega^{(1)}, \omega^{(2)}, b \right) \in \left(\mathbb{R}^I \right)^3 \text{ est une solution du problème } (\mathfrak{P}) \\ (\mathfrak{S}) \text{ admet une solution.} \end{array} \right. \implies \left(\omega^{(1)}, \omega^{(2)}, b \right) \text{ est une solution de } (\mathfrak{S})$$

a. **Attention!**, l'autre inclusion est fautive en général, et on remarque que la dernière étape dans le passage de (\mathfrak{S}) à (\mathfrak{P}) était juste en implication et non plus en équivalence comme les précédentes, et si vous avez encore des doutes, pensez par exemple à l'équation $(E) : x^2 + 1 = 0$.

Maintenant, si notre problème d'origine (\mathcal{E}_α) vérifie certaines conditions qui le rendent admettant une solution exacte $V(\cdot)$, par exemple, ceux qu'on a traité au corollaire IV.3.1.1, alors, et à la lumière de ce qu'on a vu au troisième chapitre, il existera forcément un réseaux neuronal qui approche cette solution, ce qui fait qu'on trouvera sans doute des pondérations $(w^{(1)}, w^{(2)}) \in (\mathbb{R}^I)^2$ et des biais $b \in \mathbb{R}^I$ qui résolvent le système (\mathfrak{S}) mais bien sûr en commettant une certaine erreur $\varepsilon(t)$ dont la norme uniforme $\|\varepsilon\|_\infty$ est négligeable, ceci rendra logiquement ces vecteurs très proches de la valeur minimale de $E(w)$, et tout cela va nous permettre d'approcher ces pondération, en faisant l'appel à des algorithmes d'apprentissage convenables qui vont, et en vertu de ce qui précède, nous conduire vers une très bonne approximation de la solution exacte de (\mathcal{E}_α) , et cela même si ce problème (\mathbb{P}) n'admet pas de solutions.

3.2 Algorithme d'apprentissage proposé :

Le processus d'apprentissage dans les réseaux de neurones artificiels consiste à modifier les paramètres du réseau avec une règle de correction d'erreur appropriée, et ceci dans le but de trouver une collection de poids et de biais convenable pour la résolution du problème qu'on traite.

Jusqu'à présent, plusieurs algorithmes d'apprentissage ont été construits pour réaliser ce processus, et ils se basent tous sur la règle du delta classique, et qui provient en principe de la méthode de descente de gradient.

Dans notre cas, et afin d'avoir plus de performance pratique, on adoptera la vision habituelle d'apprentissage en ligne et qui s'appuie en général sur l'optimisation de l'erreur après la présentation de chaque exemple de notre ensemble de données, ou plus clairement elle se base sur une minimisation progressive de l'erreur local $E_m(w)$ avec $m \in \llbracket 0, P \rrbracket$, pour approcher les pondérations optimales de notre système^a.

Afin d'atteindre cet objectif l'algorithme de rétro-propagation de gradient avec un élan donné est choisi dans ce cas pour s'occuper de la modification quantitative et qualitative des paramètres $\omega_i^{(1)}$, $\omega_i^{(2)}$ et b_i qui apparaissent clairement dans la formulation du problème (\mathbb{P}) , ces paramètres seront sélectionnés dans un premier temps d'une façon aléatoire ensuite, le signal d'entrée traversera notre réseau pour générer une sortie nette à valeur réelle, et qui sera comparée à celle souhaitée pour calculer en fin de compte l'erreur en sortie.

Maintenant, la fonction du coût $E_m(\cdot)$, et qui est heureusement différentiable dans ce cas, doit être minimisée dans la direction opposée à sa dérivée, or, et comme on l'a déjà dit, l'algorithme d'apprentissage supervisé choisi en l'occurrence est celui de rétro-propagation avec élan et qui modifie le paramètre $\omega_i^{(2)}$ comme suit :

$$\left(\forall i \in \{1, \dots, I\} \right) : \begin{cases} \omega_i^{(2)}(\tau + 1) = \omega_i^{(2)}(\tau) + \Delta \omega_i^{(2)}(\tau) \\ \Delta \omega_i^{(2)}(\tau) = -\eta \times \frac{\partial E_m}{\partial \omega_i^{(2)}} + \gamma \times \Delta \omega_i^{(2)}(\tau - 1) \end{cases} \quad (\star)$$

dans ce qui précède, l'indice τ désigne l'ordre de l'itération courante, tandis que l'indice i est l'étiquette de la connexion du poids, de plus, $\omega_i^{(2)}(\tau + 1)$ et $\omega_i^{(2)}(\tau)$ décrivent respectivement les paramètres de poids ajusté et actuel.

pour les deux paramètres η et γ , ils sont respectivement le taux d'apprentissage et le terme d'élan, et qui contrôlent en réalité la vitesse de convergence de cet algorithme. Afin d'accomplir la procé-

a. De l'autre part il y a ce qu'on appelle un apprentissage en batch ou aussi off ligne et qui consiste à minimiser l'erreur après la présentation de tout les exemples de la base, ou plus clairement il consiste à minimiser l'erreur global $E(w) = \sum_{m=0}^P E_m(w)$.

de l'apprentissage, la dérivée partielle ci-dessus s'exprime ainsi :

$$\begin{aligned}
\frac{\partial E_m}{\partial \omega_i^{(2)}}(\omega^{(1)}, \omega^{(2)}, b) &= \frac{\partial E_m}{\partial \Lambda_m}(\omega^{(1)}, \omega^{(2)}, b) \times \frac{\partial \Lambda_m}{\partial \omega_i^{(2)}}(\omega^{(1)}, \omega^{(2)}, b) \\
&= \left(\Lambda_m(\omega^{(1)}, \omega^{(2)}, b) - s(t_m) \right) \\
&\quad \times \left[\sum_{k=0}^N \sum_{j=0}^k \lambda_k \cdot C_k^j (\omega_i^{(1)})^j (b_i)^{n-j} \times \left(\frac{\Gamma(j+2)}{\Gamma(j+2-\alpha)} (t_m)^{j+1-\alpha} + h(t_m) \cdot (t_m)^{j+1} \right) \right] \\
&= \left(\sum_{i=1}^I \sum_{k=0}^N \sum_{j=0}^k \zeta_{i,j}^k \times \left(\frac{\Gamma(j+2)}{\Gamma(j+2-\alpha)} (t_m)^{j+1-\alpha} + h(t_m) \cdot (t_m)^{j+1} \right) + \beta \cdot h(t_m) - s(t_m) \right) \\
&\quad \times \left[\sum_{k=0}^N \sum_{j=0}^k \lambda_k \cdot C_k^j (\omega_i^{(1)})^j (b_i)^{n-j} \times \left(\frac{\Gamma(j+2)}{\Gamma(j+2-\alpha)} (t_m)^{j+1-\alpha} + h(t_m) \cdot (t_m)^{j+1} \right) \right]
\end{aligned}$$

le réglage des paramètres $\omega_i^{(1)}$ et b_i se fait par une procédure similaire à (★), et le vecteur gradient dans ce cas se présente comme suit :

pour les $\omega_i^{(1)}$:

$$\begin{aligned}
\frac{\partial E_m}{\partial \omega_i^{(1)}}(\omega^{(1)}, \omega^{(2)}, b) &= \frac{\partial E_m}{\partial \Lambda_m}(\omega^{(1)}, \omega^{(2)}, b) \times \frac{\partial \Lambda_m}{\partial \omega_i^{(1)}}(\omega^{(1)}, \omega^{(2)}, b) \\
&= \left(\Lambda_m(\omega^{(1)}, \omega^{(2)}, b) - s(t_m) \right) \\
&\quad \times \left[\sum_{k=0}^N \sum_{j=0}^k \lambda_k \cdot C_k^j \cdot \omega_i^{(2)} \cdot j \cdot (\omega_i^{(1)})^{j-1} (b_i)^{n-j} \times \left(\frac{\Gamma(j+2)}{\Gamma(j+2-\alpha)} (t_m)^{j+1-\alpha} + h(t_m) \cdot (t_m)^{j+1} \right) \right] \\
&= \left(\sum_{i=1}^I \sum_{k=0}^N \sum_{j=0}^k \zeta_{i,j}^k \times \left(\frac{\Gamma(j+2)}{\Gamma(j+2-\alpha)} (t_m)^{j+1-\alpha} + h(t_m) \cdot (t_m)^{j+1} \right) + \beta \cdot h(t_m) - s(t_m) \right) \\
&\quad \times \left[\sum_{k=0}^N \sum_{j=0}^k \lambda_k \cdot C_k^j \cdot \omega_i^{(2)} \cdot j \cdot (\omega_i^{(1)})^{j-1} (b_i)^{n-j} \times \left(\frac{\Gamma(j+2)}{\Gamma(j+2-\alpha)} (t_m)^{j+1-\alpha} + h(t_m) \cdot (t_m)^{j+1} \right) \right]
\end{aligned}$$

et pour les b_i :

$$\begin{aligned}
\frac{\partial E_m}{\partial b_i}(\omega^{(1)}, \omega^{(2)}, b) &= \frac{\partial E_m}{\partial \Lambda_m}(\omega^{(1)}, \omega^{(2)}, b) \times \frac{\partial \Lambda_m}{\partial b_i}(\omega^{(1)}, \omega^{(2)}, b) \\
&= \left(\Lambda_m(\omega^{(1)}, \omega^{(2)}, b) - s(t_m) \right) \\
&\quad \times \left[\sum_{k=0}^N \sum_{j=0}^k \lambda_k \cdot C_k^j \cdot \omega_i^{(2)} \cdot (\omega_i^{(1)})^j (n-j) \cdot (b_i)^{n-j-1} \times \left(\frac{\Gamma(j+2)}{\Gamma(j+2-\alpha)} (t_m)^{j+1-\alpha} + h(t_m) \cdot (t_m)^{j+1} \right) \right] \\
&= \left(\sum_{i=1}^I \sum_{k=0}^N \sum_{j=0}^k \zeta_{i,j}^k \times \left(\frac{\Gamma(j+2)}{\Gamma(j+2-\alpha)} (t_m)^{j+1-\alpha} + h(t_m) \cdot (t_m)^{j+1} \right) + \beta \cdot h(t_m) - s(t_m) \right) \\
&\quad \times \left[\sum_{k=0}^N \sum_{j=0}^k \lambda_k \cdot C_k^j \cdot \omega_i^{(2)} \cdot (\omega_i^{(1)})^j (n-j) \cdot (b_i)^{n-j-1} \times \left(\frac{\Gamma(j+2)}{\Gamma(j+2-\alpha)} (t_m)^{j+1-\alpha} + h(t_m) \cdot (t_m)^{j+1} \right) \right]
\end{aligned}$$

3.3 Exemple illustratif :

Dans cette section, l'efficacité de la méthode itérative proposée précédemment sera examinée sur un problème de test concret, de plus, les résultats numériques obtenus seront présentés graphiquement à travers des courbes qui rendront l'interprétation des résultats une opération plus facile.

Au cours des prochaines simulations, et pour plus de facilité, nous utilisons les paramètres suivants et qui peuvent être modifiés ultérieurement pour des raisons comparatives, ou pour montrer l'effet d'un certain paramètre sur la qualité d'apprentissage :

1. Nombre de neurones cachés : $I = 6$.
2. La fonction d'activation (logistique dans ce cas^a) : $\sigma(\bullet) = \text{Logsig}(\bullet) = \frac{\exp(\bullet)}{1 + \exp(\bullet)}$.

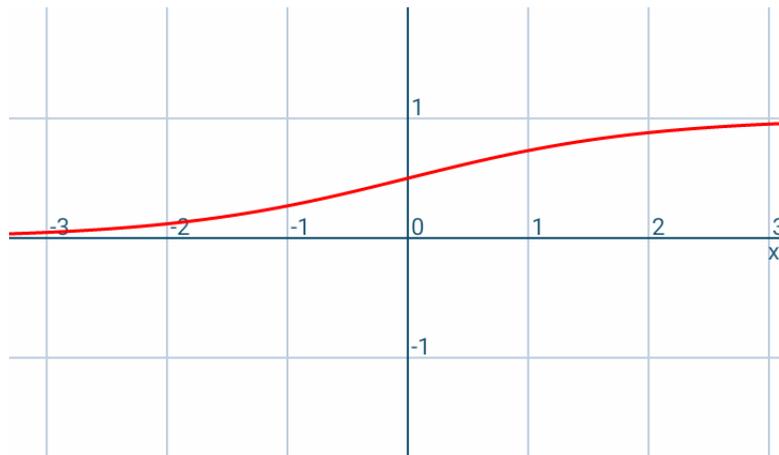


FIGURE IV.9: La courbe de la fonction "Logsig".

3. Le taux d'apprentissage (constant dans ce cas) : $\eta = 0.01$.
4. L'élan (constant dans ce cas) : $\gamma = 0.05$.
5. L'ordre de limitation pour les séries de Maclaurin : $N = 5$.
6. Nombre des points nodaux : $P = 11$.

► **Description du problème** : soit $\alpha \in]0, 1]$, considérons l'équation différentielle fractionnaire suivante :

$$(\mathcal{P}_\alpha) \begin{cases} \left[{}_0 D_t^\alpha (V(t)) \right] + V(t) = t^2 + \frac{\Gamma(3)}{\Gamma(3-\alpha)} \times t^{2-\alpha} + 1. & \text{si : } 0 < t < 1. \\ V(t) = V(0) = 1. & \text{si : } t = 0. \end{cases}$$

d'abord nous pouvons remarquer que ce système n'est autre qu'un cas particulier, de (\mathcal{E}_α) et pour lequel on a :

a. Les fonctions logistiques sont en général de la forme : $f(t) = \frac{K}{1 + a.e^{-kt}}$ où K, a et k sont des réels positifs quelconques.

ces fonctions sont des solutions en temps continu du modèle de **Verhulst** (belge, 1804-1849) et qui se présente par :

$$y' = k.y \left(1 - \frac{y}{K} \right) \text{ avec : } k > 0 \text{ et : } K > 0$$

ce modèle et qui s'exprime en temps discret par : $u_{n+1} - u_n = a.u_n \left(1 - \frac{u_n}{K} \right)$, imagine que le taux de natalité et le taux de mortalité d'une population donné sont des fonctions affines de la taille de la population, respectivement décroissante et croissante, autrement dit, plus la taille de la population augmente, plus son taux de natalité diminue et son taux de mortalité augmente, ce modèle est très utile et il a été proposé dans le cadre de la dynamique des populations pour remplacer celui de **Malthus** (britannique, 1766-1834), et qui proposait un taux d'accroissement constant conduisant à une croissance exponentielle de la population.

- $\Omega = [0, 1]$.
- $h(t) = \mathbf{1}$.
- $s(t) = t^2 + \frac{\Gamma(3)}{\Gamma(3-\alpha)} \times t^{2-\alpha} + 1$.
- $\beta = 1$.

et aussi il est bien facile de prédire sa solution exacte, et qui est définie dans ce cas par :

$$\left(\forall t \in [0, 1] \right) : V(t) = t^2 + 1$$

notre but maintenant c'est de vérifier si cette approche de résolution numérique par une procédure d'apprentissage peut nous fournir une bonne approximation de cette solution exacte.

► **Discrétisation du problème** : Dans le cadre de ce problème, la procédure de discrétisation adoptée est celle qui était décrite dans la section IV.2.2 et dans laquelle la technique du maillage utilisée pour le domaine d'étude : $\Omega = [0, 1]$ était basée sur la subdivision régulière décrite dans la figure IV.7.

► **Étapes de résolution par un processus d'apprentissage** : Dans un premier temps, le processus d'apprentissage incrémentiel, et qui est basé dans ce cas sur un algorithme de rétro-propagation avec élan, commence par une affectation de petites constantes réelles aléatoires aux paramètres $\omega_i^{(1)}$, $\omega_i^{(2)}$ et b_i du réseau, ensuite, la règle de correction (★) est appliquée pour ajuster successivement les poids de connexion et les termes de biais jusqu'à l'obtention d'une solution appropriée.

À ce stade et pour justifier la précision de cette technique, on exposera dans un premier temps les valeurs de l'erreur commise pour l'ordre $\alpha = 0.25$, et qui seront présentés dans le tableau IV.1, après, et vers la fin de cette section on fera du même pour d'autres ordres fractionnaires et qui seront également présentés pour un nombre de neurones cachés $N = 6$ et un nombre maximal d'itérations $\tau_{\max} = 500$ dans un autre tableau IV.3, puis, et pour une bonne visualisation de ces résultats, la fonction d'erreur correspondante à ces tableaux sera tracée, et pour des différents nombres d'itérations à la figure IV.10, en plus, les solutions exactes et approximatives seront à leurs tour illustrées sur la figure IV.11.

Les erreurs absolues entre les solutions exactes et approximatives sur les points nodaux vont être lues à l'aide de la fonction de l'erreur locale absolue :

$$E_{\text{abs}}(m) = \left| V(t_m) - \tilde{V}(t_m) \right|$$

et qu'on la présentera pour $\tau_{\max} = 100$ dans la figure IV.12.

Enfin, l'efficacité de la structure neuronale proposée sera étudiée, et pour les différents éléments de contrôle (L'ordre de limitation N , et le nombre des neurones I .) à travers la fonction de l'erreur moyenne absolue E_{mid} définie par :

$$E_{\text{mid}} = \sum_{m=0}^P \left| V(t_m) - \tilde{V}(t_m) \right| = \sum_{m=0}^P E_{\text{abs}}(m)$$

et qui va être présentée par le diagramme :IV.13 (pour $\tau_{\max} = 100$).

$\alpha=0.25$			
t_m	Solution exacte	Solution approchée	Erreur locale
0.1	1.0100	1.040000000090579	9.0579×10^{-11}
0.2	1.0400	1.010000000081472	8.1472×10^{-11}
0.3	1.0900	1.090000000012699	1.2698×10^{-11}
0.4	1.1600	1.160000000091338	9.1337×10^{-11}
0.5	1.2500	1.250000000063236	6.3235×10^{-11}
0.6	1.3600	1.36000000009754	9.7540×10^{-11}
0.7	1.4900	1.490000000027850	2.7849×10^{-11}
0.8	1.6400	1.640000000054688	5.4688×10^{-11}
0.9	1.8100	1.810000000095751	9.5750×10^{-11}

TABLE IV.1: Les solutions exactes et approchées dans le cas où : $\alpha = 0.25, I = 6$ et $\tau_{\max} = 500$.

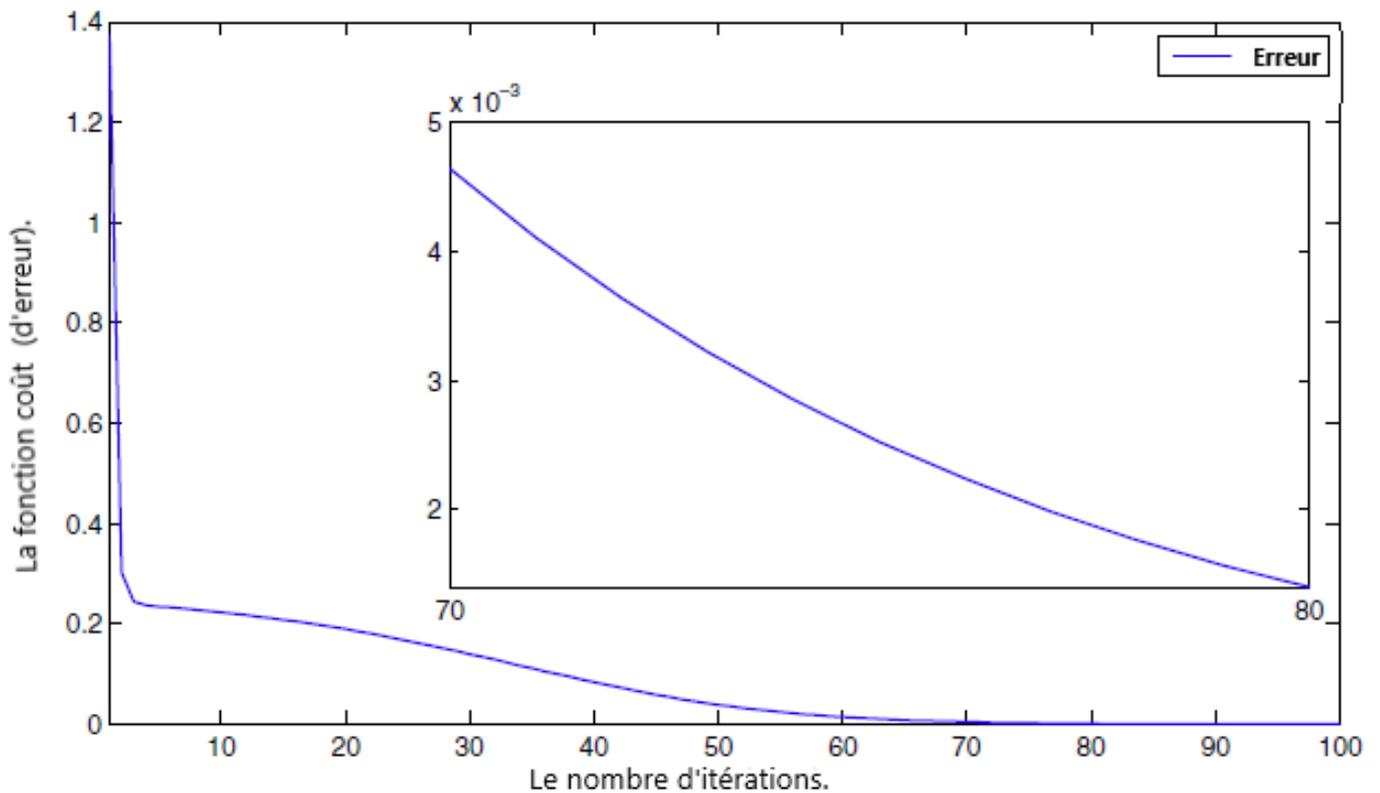


FIGURE IV.10: L'erreur moyenne en fonction du nombre d'itérations τ pour : $\tau \in \{0, \dots, 100\}$.

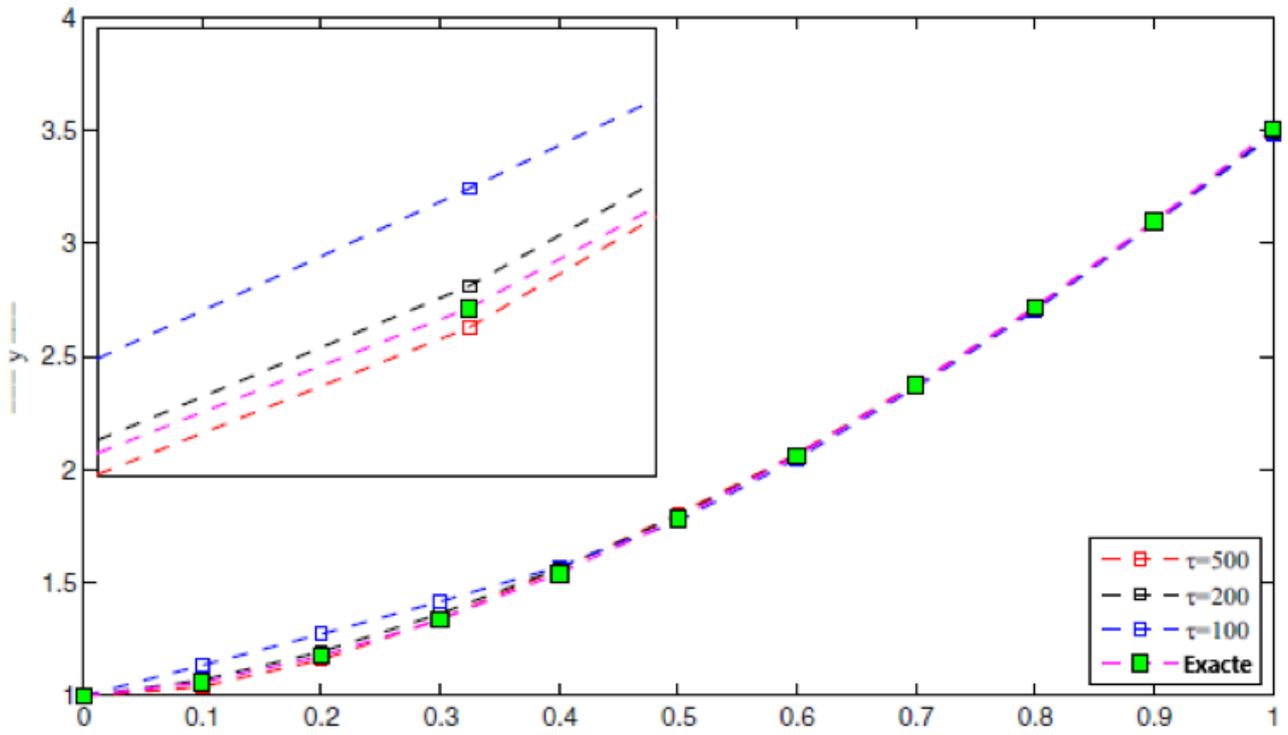


FIGURE IV.11: L'erreur moyenne en fonction du nombre d'itérations τ pour : $\tau = 100, 200$ et 500 .

Le nombre des neurones cachés.				
t_m	$I = 6$	$I = 8$	$I = 10$	$I = 12$
0.1	6.9483×10^{-11}	4.8976×10^{-14}	9.1900×10^{-15}	6.9908×10^{-17}
0.2	9.1710×10^{-12}	8.4559×10^{-13}	4.9836×10^{-16}	8.9090×10^{-17}
0.3	9.5022×10^{-12}	8.4631×10^{-13}	1.5974×10^{-16}	9.5929×10^{-16}
0.4	3.3445×10^{-11}	7.0936×10^{-14}	3.4039×10^{-16}	5.4722×10^{-17}
0.5	4.3874×10^{-11}	5.8527×10^{-16}	5.8527×10^{-16}	1.3862×10^{-17}
0.6	3.8156×10^{-11}	2.7603×10^{-14}	2.2381×10^{-16}	1.4929×10^{-17}
0.7	7.5952×10^{-11}	6.7970×10^{-14}	7.5127×10^{-16}	2.5751×10^{-17}
0.8	7.5920×10^{-11}	6.5510×10^{-14}	2.5510×10^{-16}	8.4072×10^{-17}
0.9	1.8687×10^{-11}	1.6261×10^{-14}	5.0596×10^{-16}	2.5428×10^{-17}

TABLE IV.2: Les valeurs de l'erreur locale absolue pour des différents nombres de neurones cachés.

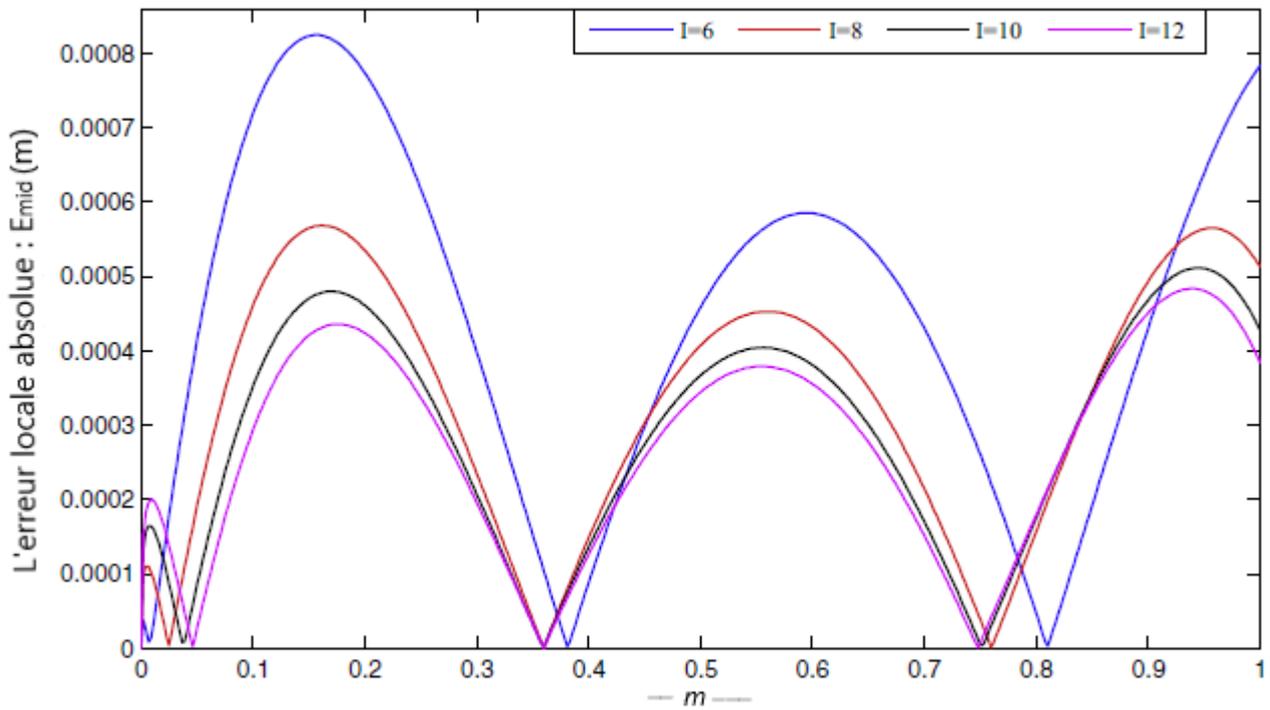


FIGURE IV.12: La courbe de l'erreur locale absolue pour des différents nombres de neurones cachés.

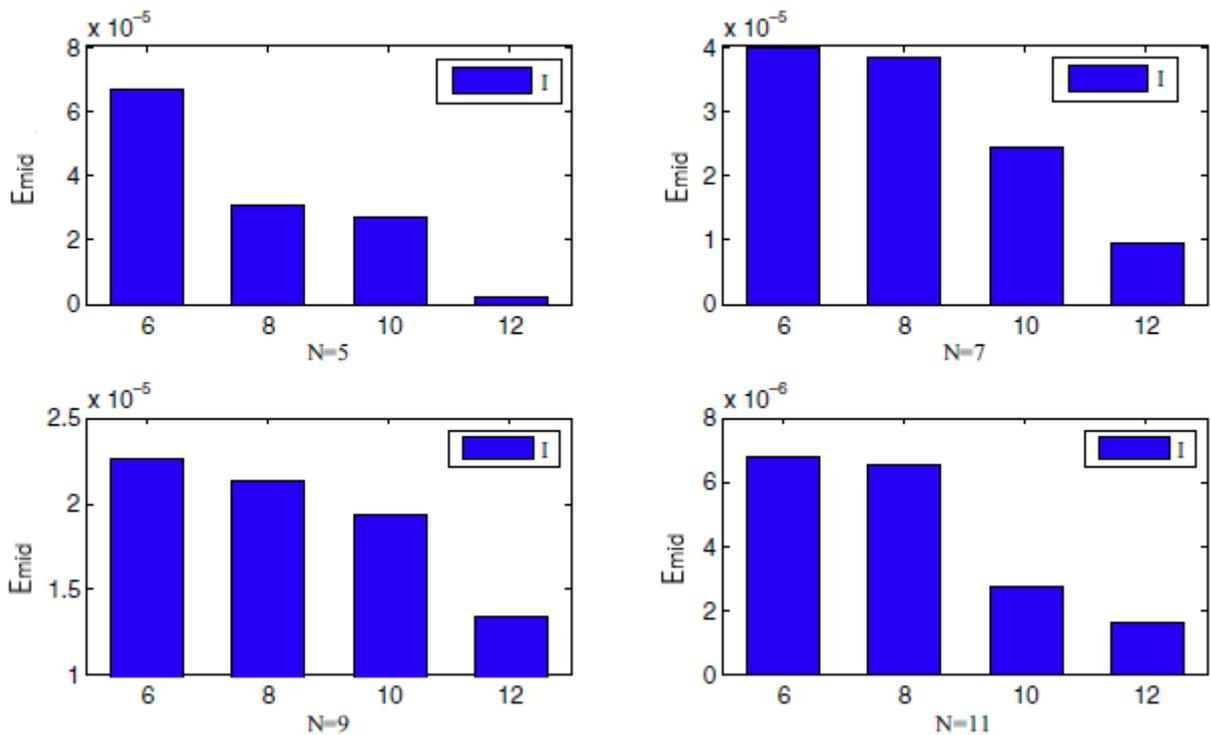


FIGURE IV.13: L'effet des paramètres N et I sur la qualité d'apprentissage.

	Les valeurs de α .		
t_m	$\alpha = 0.5$	$\alpha = 0.75$	$\alpha = 1$
0.1	9.6488×10^{-11}	4.5283×10^{-11}	2.5602×10^{-11}
0.2	9.5716×10^{-11}	3.5472×10^{-11}	1.4518×10^{-11}
0.3	8.0028×10^{-11}	3.0180×10^{-11}	9.2834×10^{-12}
0.4	8.8530×10^{-11}	9.2143×10^{-12}	2.7308×10^{-11}
0.5	5.4188×10^{-11}	5.7390×10^{-11}	8.9179×10^{-12}
0.6	4.2176×10^{-11}	5.8903×10^{-11}	1.7321×10^{-11}
0.7	7.3576×10^{-11}	6.4140×10^{-11}	3.8641×10^{-11}
0.8	7.1573×10^{-11}	6.3517×10^{-11}	2.4577×10^{-11}
0.9	8.7059×10^{-11}	7.3820×10^{-11}	3.8429×10^{-11}
L'erreur locale absolue.			

TABLE IV.3: Les valeurs de l'erreur locale absolue pour des différents choix de α .

4 Quelques domaines d'applications des systèmes fractionnaires :

Les systèmes fractionnaires apparaissent de plus en plus dans les différents domaines de la recherche, toutefois, l'intérêt progressif que l'on porte à ces systèmes et surtout grâce à leurs applications qui apparaissent en sciences fondamentales et appliquées à la fois.

On peut noter que pour la majeure partie des domaines présentés ci dessous, les opérateurs fractionnaires sont utilisés pour prendre en compte des effets de mémoire .

Mentionnons à cet égard les ouvrages [51, 52] et qui regroupent diverses applications du calcul fractionnaire.

4.1 En Physique :

Une des applications les plus remarquables du calcul fractionnaire en physique était dans le contexte de la mécanique classique.

Le physicien Fred Riewe a montré en 1996 dans son article [49] que le Lagrangien contenant des dérivées temporelles d'ordres fractionnaires **conduit à une modélisation pertinente du mouvement avec des forces non-conservatives** comme les frottements à titre d'exemple.

Ce résultat est remarquable du fait que les forces de frottement ou même les forces non-conservatives^a en général sont essentielles dans le traitement variationnel macroscopique habituel, et par conséquent, dans les méthodes les plus avancées de la mécanique classique.

Riewe a généralisé le calcul des variations habituel au Lagrangien qui dépend des dérivées fractionnaires afin de s'adapter avec les forces non-conservatives habituelles.

4.2 En automatique :

En automatique, ce n'est qu'au début des années 1990 que le régulateur CRONE (commande robuste d'ordre non entier) a été proposé par Oustaloup [45].

En profitant des propriétés avantageuses des systèmes d'ordres fractionnaires, ce régulateur permet d'assurer la robustesse de la commande dans une bande de fréquence donnée, la réussite de

a. Par contre aux forces conservatives, le travail produit par ce type de forces dépend du chemin suivi par son point d'action.

cette approche fut énorme, et plusieurs variantes de cette commande ont vu le jour (1^{ère}, 2^{ème} et 3^{ème} générations).

4.3 En Acoustique :

Pour certains instruments de musique à vent, les pertes visco-thermiques peuvent être modélisées efficacement à l'aide de dérivées fractionnaires temporelles [30].

Conclusion :

Dans ce chapitre on a utilisé les réseaux de neurones artificiels, et qui se présentent comme des approximateurs universels, pour approcher la solution d'une équation différentielle fractionnaire à valeur initiale.

Pour arriver à cet objectif, une structure neuronale à trois couches avec une fonction d'activation sigmoïde a été utilisée lors de la modélisation de ce problème.

Le placement de cette structure neuronale dans l'équation différentielle concernée, le recours à la discrétisation, l'emploi du développement de Maclaurin pour franchir la complexité pratique du calcul de la dérivée fractionnaire pour la fonction d'activation, et l'utilisation de l'erreur quadratique moyenne, nous ont conduit à un problème d'optimisation relativement classique et qui sera traité numériquement par l'algorithme d'apprentissage de rétro-propagation, et qui est basé sur la règle de correction delta.

Enfin, un problème de test a été fait pour mieux illustrer la méthodologie proposée, les résultats numériques et les simulations informatiques obtenus mettent en évidence le fait que cette technique itérative a introduit vraiment une puissante amélioration sur la résolution des problèmes fractionnaires à valeurs initiales.

Conclusion générale et perspectives

Au cours de ce travail, nous avons traité le problème de la résolution numérique des équations différentielles fractionnaires ordinaires par des réseaux de neurones formels, et plus précisément par des perceptrons multicouches.

En premier lieu on a commencé par un prérequis essentiel d'analyse fonctionnelle, puis on est passé à une étude générale de l'apprentissage automatique et du modèle des réseaux de neurones artificiels, ensuite on a utilisé ce prérequis pour justifier l'aptitude de ces derniers, et qui se considèrent pour un choix particulier de la fonction d'activation comme des approximateurs universels parcimonieux, enfin on a investi cette propriété fondamentale dans le cadre de l'étude numérique des équations différentielles fractionnaires et qui seront résolus dans ce contexte à l'aide d'un processus d'apprentissage, fondé du côté algorithmique sur la méthode de rétro-propagation de gradient.

Comme perspectives de notre travail, nous souhaitons d'abord avoir plus de choix et de possibilités pour la fonction d'activation ^a, et qui peuvent garder la vérification de la propriété d'approximation parcimonieuse de notre réseau neuronal, ainsi on espère bien de généraliser la forme du problème différentiel qu'on traite, et ceci en considérant par exemple un ordre de dérivation $\alpha \geq 1$, ou un domaine d'étude Ω non bornée, ou même des modèles non linéaires, ainsi la recherche d'une gamme algorithmique plus puissante sera vraiment désirable, enfin on signale que l'extension de cette méthode pour équations aux dérivées fractionnaire partielles reste encore une question très importante, et qui se pose jusqu'à l'instant dans le cadre de la recherche.

a. Récemment, quelques travaux de recherches ont fait vraiment du bon effort dans ce sens, et on cite à titre exemple l'article [25] de Vigor Ismailov et Namig Guliyev qui ont proposé en 2018 une construction algorithmique d'une fonction d'activation qui peut garantir la vérification de la propriété d'approximation, et ceci en faisant appel à la notion des fonctions λ -croissantes.

Prérequis essentiel pour la lecture du troisième chapitre :

Dans cette annexe, on met à la disposition du lecteur un rappel résumé de quelques théorèmes de topologie et d'analyse fonctionnelle, ainsi que leurs preuves et leurs corollaires, et ceci dans le but de stimuler et de renforcer son bagage et son pré-acquis mathématique.

La lecture de cette partie n'est pas du tout obligatoire ou indispensable, mais elle sera en revanche souhaitable, et surtout pour une bonne compréhension de quelques preuves dans lesquelles l'importance de ces théorèmes apparaît clairement, (voir par exemple la démonstration du résultat : [III.1.2.1](#)).

1 Théorèmes du Noyau fermé et de Hahn-Banach :

1.1 Théorème du noyau fermé :

Théorème 1.1.1 :

Une forme linéaire sur un \mathbb{R} -espace vectoriel E est continue **si et seulement si** son noyau est fermé.

Preuve :

Il est clair que si $f \in E^*$ est continue $H = f^{-1}(\{0\})$ est fermé comme image réciproque d'un singleton.

Pour la réciproque considérons $a \notin H$ tel que : $f(a) = 1$, comme H est fermé il en est de même de $\{a\} \times H$ et par suite de $a + H$, qui est l'image réciproque du fermé $\{a\} \times H$ par l'application continue $(x, y) \mapsto x + y$, et comme $0 \notin a + H$, il existe une boule $B(0, r[$ qui ne rencontre pas $a + H$, par suite, $x \in B(0, r[\Rightarrow f(x) \neq 1$.

montrons maintenant que : $x \in B(0, r[\Rightarrow |f(x)| \leq 1$, on suppose le contraire et on note : $f(x) = \alpha$ avec $x \in B(0, r[$ et $|\alpha| > 1$; alors $\left\| \frac{x}{\alpha} \right\| = \left(\frac{1}{|\alpha|} \right) \|x\| < r$ et $f\left(\frac{x}{\alpha}\right) = 1$ ce qui contredit la définition de $B(0, r[$. f étant bornée sur la boule $B(0, r[$ est bornée sur la boule unité et donc continue. ■

Remarque 1.1.1 :

Ce théorème est classique et très célèbre, et il nous servira dans la suite pour légitimer quelques passages au niveau de la démonstration du théorème universel d'approximation [II.1.2.1](#).

1.2 Théorème de Hahn-Banach :

Lemme 1.2.1 : (de Zorn, [6])

Tout ensemble ordonné, *inductif*^a et non vide admet un élément maximal.

Théorème 1.2.1 : (Hahn-Banach Analytique, [7])

Soit $(E, +, \cdot)$ un \mathbb{R} -espace vectoriel, et $p : E \rightarrow \mathbb{R}^+$ une semi norme sur cet espace, c'est à dire une application vérifiant :

$$\bullet p(\lambda \cdot x) = \lambda \times p(x) \quad (\forall x \in E), (\forall \lambda \geq 0).$$

$$\bullet p(x + y) \leq p(x) + p(y) \quad (\forall (x, y) \in E^2).$$

-Si G est un sous-espace vectoriel de E , et $g : G \rightarrow \mathbb{R}$ est une forme linéaire telle que :

$$(\forall x \in G) : g(x) \leq p(x)$$

alors il existe une forme linéaire $f : E \rightarrow \mathbb{R}$ qui prolonge g et qui vérifie :

$$(\forall x \in E) : f(x) \leq p(x).$$

Preuve :

On considère :

$$\mathfrak{P} := \left\{ (S, h) / \text{avec } S \text{ sev de } E, h \in E^*, G \subset S, h \text{ prolonge } g \text{ et } \forall x \in D(h), h(x) \leq p(x) \right\}^b.$$

On munit \mathfrak{P} de la relation d'ordre :

$$(S_1, h_1) \leq (S_2, h_2) \iff S_1 \subset S_2 \text{ et } h_2 \text{ prolonge } h_1.$$

► $\mathfrak{P} \neq \emptyset$ car $g \in \mathfrak{P}$.

► \mathfrak{P} est *inductif* : soit $Q \subset \mathfrak{P}$ un sous-ensemble totalement ordonné, on note $Q = (S_i, h_i)_{i \in I}$, on définit $D := \bigcup_{i \in I} S_i$ et $h(x) := h_i(x)$ si $x \in S_i$.

il est facile de voir que (D, h) est un majorant de Q . Donc d'après le lemme de Zorn, \mathfrak{P} possède un élément maximal (H, f) . Prouvons que $H = E$:

Supposons que $H \neq E$ et soit $x_0 \notin H$, on pose $F = H + \mathbb{R}x_0$ et :

$$\begin{aligned} h & : F = H \oplus \mathbb{R}x_0 & \longrightarrow & \mathbb{R} \\ & x + \lambda \cdot x_0 & \longmapsto & h(x) = f(x) + \lambda \cdot \alpha \end{aligned}$$

pour un certain α , le but étant que $h \in P$. On doit donc avoir :

$$f(x) + t\alpha \leq p(x + tx_0) \quad (\forall x \in H), (\forall t \in \mathbb{R}).$$

c'est à dire :

$$\begin{cases} f(x) + \alpha \leq p(x + x_0) & (\forall x \in H) \\ f(x) - \alpha \leq p(x - x_0) \end{cases}$$

par (i). Il suffit donc de choisir α tel que :

$$\sup_{y \in H} (f(y) - p(y - x_0)) \leq \alpha \leq \inf_{x \in H} (p(x + x_0) - f(x))$$

Ceci est possible car :

$$\begin{aligned} f(x) + f(y) & \leq p(x + y) \\ & \leq p(x + x_0) + p(y - x_0) \end{aligned}$$

Donc $f(y) - p(y - x_0) \leq p(x + x_0) - f(x) \quad \forall (x, y) \in H$.

Ainsi, $h \in P$ et $f \neq h$, on obtient une contradiction. ■

a. un ensemble ordonné, (E, \leq) est dit *inductif* si toutes ses *chaînes* (sous ensembles totalement ordonnées) admettent un élément maximal.

b. E^* désigne le *dual algébrique* de E , c'est à dire l'espace des formes linéaires sur E .

Corollaire 1.2.1 :

Soit $(E, +, \cdot, \|\cdot\|)$ un espace vectoriel normé, G un sous-espace vectoriel de E . Soit $g : G \rightarrow \mathbb{R}$ linéaire et **continue**. Alors il existe $f \in E'$ **prolongeant** g telle que $\|f\| = \|g\|$.

Preuve :

On pose $p(x) := g\|x\|$, les hypothèses du théorème sont bien vérifiées et on a $|f(x)| \leq \|g\| \cdot \|x\|$, d'où $\|f\| \leq \|g\|$, donc $\|f\| = \|g\|$ car f prolonge g . ■

Corollaire 1.2.2 :

Soit $(E, +, \cdot, \|\cdot\|)$ un espace vectoriel normé, Pour tout $x_0 \in E$, il existe $f_0 \in E'$ ^a tel que $\|f_0\| = \|x_0\|$ et $\langle f_0, x_0 \rangle = \|x_0\|^2$.

Preuve :

Appliquer le corollaire avec $G := \mathbb{R}x_0$ et $g(tx_0) := t\|x_0\|^2$. ■

Remarques 1.2.1 :

1. Pour l'application qu'on a choisi dans la preuve, on a :

$$\|g\| := \sup_{x \in E \setminus \{0\}} \frac{\|g(x)\|}{\|x\|} = \sup_{t \in \mathbb{R}^*} \frac{\|t\|x_0\|^2\|}{\|t\|x_0\|} = \|x_0\|$$

2. ce corollaire va jouer un rôle indispensable, dans la démonstration du théorème III.1.2.1 .

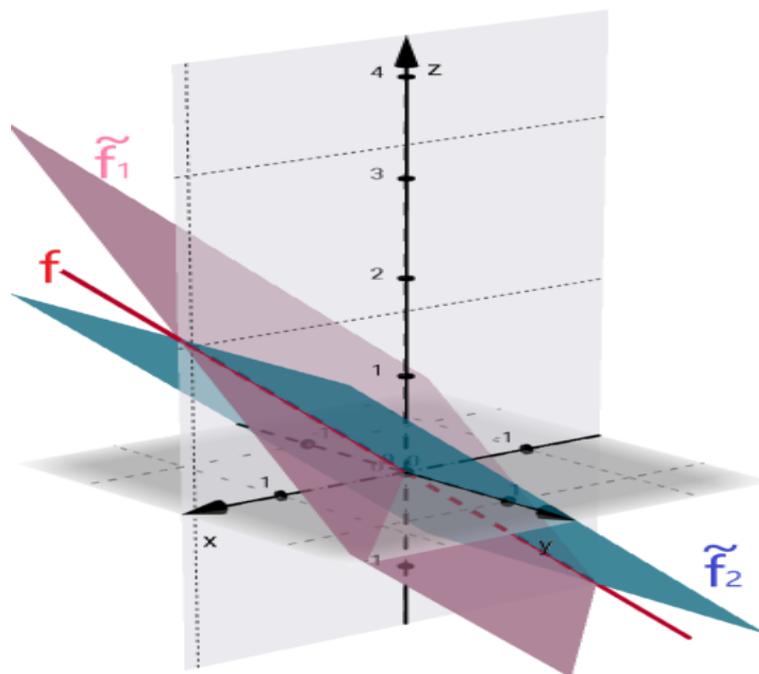


FIGURE A.1: L'idée du théorème de Hahn-Banach pour $E = \mathbb{R}^2$, $F = \mathbb{R} \times \{0\}$ et $f(x) = 2x$.

a. E' désigne le **dual topologique** de E , c'est à dire l'espace des formes linéaires **continues** sur E , juste comme remarque :il est facile de voir que $E' \subseteq E^*$.

2 Espaces Hilbertiens et théorème de représentation de Riesz :

2.1 Espaces Hilbertiens :

Définition 2.1.1 :(Espace de Hilbert) :

Soit $(E, +, \cdot, \langle \cdot, \cdot \rangle)$ **un espace préhilbertien** réel, c'est à dire, un \mathbb{R} -espace vectoriel munit d'un **produit scalaire**,^a on dit que cet espace est de Hilbert ou Hilbertien s'il est complet pour la norme $\|\cdot\|$ définie par : $(\forall x \in E) : \|x\| = \sqrt{\langle x, x \rangle}$ ^b.

Exemples 2.1.1 :

1. Sans surprise, l'espace $E = \mathbb{R}^n$ munit du produit scalaire usuel définit par :

$$\begin{array}{l} \langle \cdot, \cdot \rangle : E \times E \longrightarrow \mathbb{R} \\ (x, y) \longmapsto \sum_{i=0}^n x_i \cdot y_i \end{array}$$

est un espace de Hilbert, et en général : **tout espace euclidien^c est de Hilbert.**

2. soit $(\Omega, \mathcal{F}, \mu)$ un espace mesuré, on note : $\mathcal{L}^2(\Omega, \mathcal{F}, \mu)$ l'espace des fonctions **réelles** de carré μ -intégrable, c'est à dire :

$$\mathcal{L}^2(\Omega, \mathcal{F}, \mu) := \left\{ \underbrace{f \in \mathcal{M}((\Omega, \mathcal{F}); (\mathbb{R}, \mathcal{B}_{\mathbb{R}}))}_{\text{l'espace de fonctions } \mathcal{F}\text{-}\mathcal{B}_{\mathbb{R}} \text{ mesurables}} \mid \int_{\Omega} |f|^2 d\mu < +\infty \right\}$$

-on définit ensuite l'espace $L^2(\Omega, \mathcal{F}, \mu) := \mathcal{L}^2(\Omega, \mathcal{F}, \mu) / \mathfrak{R}$, où \mathfrak{R} est la relation d'équivalence sur $\mathcal{L}^2(\Omega, \mathcal{F}, \mu)$ définie par l'égalité μ -presque partout, cet espace quotient va nous permettre de corriger l'application :

$$\begin{array}{l} \|\cdot\|_2 : \mathcal{L}^2(\Omega, \mathcal{F}, \mu) \longrightarrow \mathbb{R}^+ \\ f \longmapsto \left(\int_{\Omega} |f|^2 d\mu \right)^{\frac{1}{2}} \end{array}$$

et qui est juste une **semi-norme^d** sur $\mathcal{L}^2(\Omega, \mathcal{F}, \mu)$, et la rendre une norme sur $L^2(\Omega, \mathcal{F}, \mu)$. Maintenant, et sur l'espace $L^2(\Omega, \mathcal{F}, \mu)$, on va définir un produit scalaire comme suit :

$$\forall (f, g) \in \left(L^2(\Omega, \mathcal{F}, \mu) \right)^2 \quad \langle f, g \rangle := \int_{\Omega} f \cdot g d\mu$$

il est facile de remarquer que : $(\forall f \in L^2(\Omega, \mathcal{F}, \mu)) : \|f\|_2 = \sqrt{\langle f, f \rangle}$ donc et en se basant sur le fameux **théorème de Riesz-Fischer^e** on peut déduire facilement que : $(L^2(\Omega, \mathcal{F}, \mu), +, \cdot, \langle \cdot, \cdot \rangle)$ est un espace de Hilbert.

- a. **un produit scalaire**(réel) est une forme bilinéaire réelle, symétrique, positive et définie positive.
- b. cette norme s'appelle la **norme Euclidienne** ou la **norme induite** du produit scalaire de notre espace.
- c. c'est à dire un \mathbb{R} -espace vectoriel préhilbertien **de dimension fini**.
- d. l'axiome de séparation n'est pas vérifiée.
- e. ce théorème affirme que pour tout $p \geq 1$, $L^p(\Omega, \mathcal{F}, \mu)$ est un espace de Banach pour la norme :

$$\|f\|_p := \left(\int_{\Omega} |f|^p d\mu \right)^{\frac{1}{p}}$$

Définition 2.1.2 : (l'orthogonalité) :

Soit $(E, +, \cdot, \langle \cdot, \cdot \rangle)$ **un espace préhilbertien réel** :

1. On dit que deux vecteurs x et y de E sont **orthogonaux** (pour le produit scalaire $\langle \cdot, \cdot \rangle$), et on note $x \perp y$, si et seulement si : $\langle x, y \rangle = 0$.
2. **L'orthogonal d'une partie** $X \subset E$, noté X^\perp , est le sous ensemble de E constituée des vecteurs orthogonaux à tous les vecteurs de X , autrement dit :

$$X^\perp = \{y \in E \mid \forall x \in X, \langle x, y \rangle = 0\}.$$

Proposition 2.1.1 :

Soit $(E, +, \cdot, \langle \cdot, \cdot \rangle)$ **un espace préhilbertien réel**, et soit $X \subset E$ une partie de E^a alors on a :

1. X^\perp est toujours **un sous espace vectoriel fermé** de E .
2. l'égalité $X^\perp = (\text{vect}(X))^\perp$ est toujours vérifiée.

Preuve :

Soit $(E, +, \cdot, \langle \cdot, \cdot \rangle)$ un espace préhilbertien réel, et soit $X \subset E$ une partie de E

1. Montrons que X^\perp est toujours **un sous espace vectoriel fermé** de E :
-d'abord, on remarque que :

$$\begin{aligned} X^\perp &= \{y \in E \mid \forall x \in X, \langle x, y \rangle = 0\} \\ &= \bigcap_{x \in X} \text{Ker}(\langle x, \cdot \rangle) \end{aligned}$$

avec $\langle x, \cdot \rangle$ est la forme linéaire ^b définie par :

$$\begin{aligned} \langle x, \cdot \rangle &: E \longrightarrow \mathbb{R} \\ y &\longmapsto \langle x, y \rangle \end{aligned}$$

et selon **l'inégalité de Cauchy-Schwartz** on a : $\boxed{(\forall (x, y) \in X^2) : \langle x, y \rangle \leq \|x\| \|y\|}$, donc :

$(\forall x \in X) : y \mapsto \langle x, y \rangle$ est **continue**, ce qui fait que :

$(\forall x \in X) : \text{Ker}(\langle x, \cdot \rangle)$ est un sous espace vectoriel fermé de E .

Rappel : le noyau d'une application linéaire, est toujours un sous espace vectoriel et si de plus cette application est continue, alors il sera sans doute fermé car c'est l'image réciproque de $\{0\}$ qui est fermé par une application continue.

donc et en guise de conclusion :

$$X^\perp = \bigcap_{x \in X} \text{Ker}(\langle x, \cdot \rangle) \text{ est un sous espace vectoriel fermé de } E^c.$$

2. montrons que $X^\perp = (\text{vect}(X))^\perp$ est vraie :
d'abord on a $X \subset (\text{vect}(X))$ donc :

$$(\text{vect}(X))^\perp \subset X^\perp \quad \star.$$

et d'une autre part : si on prend $u \in X^\perp$:

-on a : $(\forall v \in \text{vect}(X)) (\exists (\lambda_1, \lambda_2, \dots, \lambda_n) \in \mathbb{R}^n) (\exists (x_1, x_2, \dots, x_n) \in X^n) : v = \sum_{i=1}^n \lambda_i x_i$.

a. **on n'exige pas de cette partie qu'elle soit un sous espace vectoriel de E**
b. elle est clairement linéaire car le produit scalaire, et par définition, est bilinéaire.
c. on signale que l'intersection quelconque des sous espaces vectoriels est un sous espace vectoriel et cette propriété reste valable pour les ensembles fermés aussi

donc : $(\forall v \in \text{vect}(X)) : \langle u, v \rangle = \left\langle u, \sum_{i=1}^n \lambda_i x_i \right\rangle = \sum_{i=1}^n \lambda_i \langle u, x_i \rangle = 0$ (car $u \in X^\perp$). donc : $u \in \text{vect}(E)^\perp$, ce qui fait que :

$$X^\perp \subset (\text{vect}(X))^\perp \quad \star \star.$$

donc de \star et $\star \star$ on déduit que : $X^\perp = (\text{vect}(X))^\perp$. d'où le résultat ■

2.2 Théorème de représentation de Riesz :

Théorème 2.2.1 : (Riesz, [6])

Soit $(E, +, \cdot, \langle, \rangle)$ un espace de Hilbert réel, pour toute forme linéaire f continue sur E , il existe un unique $y \in E$ tel que :

$$(\forall x \in E) : f(x) = \langle y, x \rangle.$$

Preuve :

$F := \ker f$ est fermé, car f est continue.

— **Existence :** Si $F^\perp = \{0\}$, alors : $\underbrace{(F^\perp)^\perp}_{=\bar{F}} = \{0\}^\perp = E$, or F est un sous espace vectoriel fermé donc

$f = 0$ et on prend $y = 0$. Sinon, soit $w \in F^\perp$ tel que $\|w\| = 1$.

on a $f(w) \neq 0$ et pour : $x \in E$: $\left(x - \frac{f(x)}{f(w)} \cdot w\right) \in \ker f = F$ donc : $\left\langle w, x - \frac{f(x)}{f(w)} \cdot w \right\rangle = 0$ d'où, pour $x \in E$,

$$\underbrace{f(x) = f(x) \langle w, w \rangle}_{\text{car } w \text{ est un vecteur unitaire}} = f(w) \langle w, x \rangle = \langle y, x \rangle$$

en posant $y = f(w) \cdot w$.

— **Unicité :** Soit $y_1, y_2 \in E$ tel que pour tout $x \in E$, $\langle y_1, x \rangle = f(x) = \langle y_2, x \rangle$. Alors, en prenant $x = y_1 - y_2$, on obtient :

$$\langle y_1 - y_2, y_1 - y_2 \rangle = \|y_1 - y_2\|^2 = 0$$

d'où $y_1 = y_2$. ■

Remarque 2.2.1 :

1. Selon la propriété de représentation matricielle des applications linéaires définies sur des espaces de dimension fini^a vers d'autres qui sont à leurs tours de dimension fini, on peut voir que ce théorème devient presque trivial.
2. ce résultat joue un rôle important, surtout en calcul différentiel, car il nous permet de **généraliser** quelque notions, comme le **vecteur gradient** en **dimension infini**^b.
3. signalons que ce théorème n'est pas là pour faire joli !, mais il sera d'une très grande importance pour la justification et la preuve du fameux théorème d'approximation par les réseaux neuronaux.

a. on rappelle que si E est de dimension fini, alors le dual algébrique E^* et le dual topologique E' deviennent confondus.

b. en fait on peut voir le vecteur gradient comme un représentant de la différentielle.

Index

A

Adaline, 28, 42, 44
Algèbre différentielle, 89
Algèbre sur un corps., 11, 75
Algorithme d'apprentissage, 20, 33, 34, 37
Algorithme d'apprentissage d'Adaline, 39, 42
Algorithme d'apprentissage du Perceptron, 39, 40
Algorithme d'apprentissage du perceptron simple, 36
Algorithme incrémentiel, 117
Analyse fonctionnelle, 3
Apprentissage, 5, 10, 19, 20, 30, 31, 33, 35
Apprentissage automatique, 20, 39, 86
Apprentissage en batch., 114
Apprentissage en ligne., 114
Apprentissage non supervisé, 31
Apprentissage off ligne., 114
Apprentissage par renforcement, 32
Apprentissage par transfert, 32
Apprentissage statistique, 33
Apprentissage supervisé, 31, 33, 35, 36, 48
Architecture, 73, 86, 98, 99
Axone, 22, 23

B

Base d'apprentissage, 31, 35
Biais, 24, 27, 35, 50–52, 70
Binôme de Newton, 105
Bio-informatique, 21

C

Carte, 47
Carte auto adaptative, 34, 47
Classification, 27, 31, 33, 52, 67
Clustering, 52
Compact, 7–11, 14, 66, 69, 73, 75, 76, 79, 81–83, 88

Complet, 7, 10, 46, 127

Condition non locale., 112

Convexe, 7, 88, 89

Corps cellulaire, 22

Cosinus écraseur, 80

Couche cachée, 1, 7, 38, 48, 50, 51, 53, 66, 77, 85

Couche d'entrée, 47, 70

Couche de sortie, 48–50

Couche de sortie, 47

Couche(s) interne(s), 51, 53, 54

Critère d'équivalence, 90

Critère d'Abel, 100

Critère d'Alembert, 100

Critère d'arrêt, 37, 52

Critère de Cauchy-Hadamard, 100

Critère de Riemann, 90

Critère de Smirnov, 18

Critère des moindres carrés, 35

D

Dérivée fractionnaire, 87, 97

Dérivée fractionnaire de Caputo, 93, 95

Dérivée fractionnaire de Riemann Liouville, 95

Dérivée fractionnaire inférieure de Caputo, 95

Dérivée fractionnaire supérieure de Caputo, 95

Densité, 3, 10, 53, 68, 72

Densité uniforme, 73

Densité uniforme compacte, 73, 81

Descente de Gradient, 29, 36–38, 40, 42, 48, 86

Discretisation, 87, 98, 107

E

Echantillon, 43

Equation différentielle, 3, 46, 86
Equation différentielle fractionnaire, 1, 86, 87, 98, 99, 111, 116
Equation intégrale de Volterra., 109
Erreur, 35, 37, 50, 83
Erreur globale, 48
Erreur locale, 35–37
Erreur moyenne, 35
Erreur quadratique moyenne, 52
Erreur totale, 35
Erreur totale moyenne, 35
Espace L^p , 8, 10
Espace Compact, 13
Espace compact, 80, 81
Espace de Banach., 11, 108, 111, 127
Espace fonctionnel, 3
Espace métrique, 6–8
Espace mesuré, 4, 9, 55, 57, 68, 69, 127
Espace mesuré de Radon, 15, 17
Espace mesurable, 55, 56
Espace paracompact, 18
Espace précompact, 18
Espace préhilbertien., 127
Espace topologique, 11, 15, 17, 18, 70
Espace topologique à base dénombrable, 17
Espace topologique métrisable, 18
Espace topologique paracompact, 18
Espace topologique séparé, 15, 18
Espace vectoriels topologiques, 3

F
Feedforward, 86
Fonction Bêta, 90
Fonction convexe, 89
Fonction d'écrasement, 71, 78, 80–82
Fonction d'activation, 23, 25, 27–30, 35, 37–40, 42, 52, 54, 99
Fonction d'entrée totale, 23, 25, 27, 34, 71
Fonction développable en séries entières, 101, 102
Fonction de décision, 10, 20, 25, 29, 66, 70
Fonction de Heavside, 27, 39, 40
Fonction de répartition de probabilité, 61
Fonction de sortie, 23, 25, 27, 71
Fonction discriminatoire, 61
Fonction discriminatoire, 60
Fonction Gamma, 87
Fonction logistique, 116
Fonction sigmoïde, 29
Forme bilinéaire., 127
forme discrétisée, 112

Forme(s) linéaire(s)., 3, 62, 124, 125, 128, 129
Formule de Taylor avec reste intégrale, 104

G
Gradient, 50, 51

I
Intégrale, 95
Intégrale de Gauss, 89
Intégrale fractionnaire de Caputo, 94
Intégrale fractionnaire de Riemann-Liouville, 92
Intégrale généralisée, 90
Intégrale impropre, 92
Intégrale paramétrique, 57
Intégrande, 95
Intégration répétée de Cauchy-Riemann, 91
Intelligence artificielle, vi, 20, 21

L
La ρ_μ densité, 77, 79, 81
Lemme de Zorn, 125
Lemme des classes monotone, 60
Linéairement séparable., 26, 28, 40
Loi de composition interne., 11

M
Machine d'apprentissage, 20, 28, 33
Madaline, 34, 84
Maillage, 106, 117
Maillage., 106
mailles, 106
Masse de Dirac., 60
Maximum de vraisemblance., 34
Mesure à densité, 57
Myéline, 21

N
Nervoglies, 21, 22
Neurone(s) biologique(s), 20–24
Neurone(s) formel(s), 20, 22–24, 36

O
Opérateur de différenciation, 87
Opérateur inverse, 94
Opérateur linéaire borné., 108
Opérateur(s), 87, 92, 93

P
Parcimonie, 83, 84
Partition de l'unité., 17, 18
Partition., 66, 106
Perceptron, 20, 27–30, 35, 40, 42–44, 48
Perceptron de Rosenblatt, 27, 29

Perceptron multicouche, 34
 Perceptron Simple, 34
 Perceptron(s) multicouche(s), 48
 Perceptron(s) multicouches, 45, 48
 Perceptron(s) sigmoïdes, 29
 Poids, 49–52, 70
 Points nodaux., 87, 106
 Pondération(s), 27, 35, 37, 40, 42, 43
 Population, 116

Q
 Qualité d'apprentissage., 33, 116
 Quantification vectorielle., 47

R
 Régression statistique., 31, 33, 34, 52
 Réseau bayésien., 33
 Réseau de Kohonen, 47
 Réseaux de Hopfield, 20, 34, 46, 47
 Réseaux de neurones $\Sigma\Pi$, 75
 Réseaux de neurones artificiels, 4, 7, 19, 20
 Réseaux de vecteurs de support., 33
 Réseaux dynamiques, 45, 46
 Réseaux feedforward, 45, 53
 Réseaux non récurrents, 45
 Réseaux récurrents, 45, 46
 Réseaux statiques, 45
 Rétro-propagation avec élan., 114
 Rétropropagation de Gradient, 38, 87
 Règle de delta., 114
 Règle de Hebb, 29, 39, 40
 Relativement compact, 109

S
 Séries de Maclaurin, 100, 102–104
 Séries entières, 100, 101, 103
 Semi dérivée, 96
 Seuil, 20, 26–29, 40, 52
 Seuillage, 25, 29
 Soma, 22
 Sous algèbre unitaire., 12, 14
 Sous algèbre., 12
 Support d'une fonction, 8

Synapses, 22, 23, 40

T
 Théorème d'approximation, 58
 Théorème d'approximation polynomiale de Weierstrass, 14
 Théorème d'Egorov, 4–6, 10
 Théorème de Bohr-Mollerup, 89
 Théorème de Bolzano-Weierstrass, 9, 65
 Théorème de Borel-Lebesgue, 75, 76
 Théorème de continuité des intégrales paramétriques, 57, 88
 Théorème de convergence bornée, 5
 Théorème de convergence dominée, 10
 Théorème de convergence monotone, 8
 Théorème de décomposition de Hahn., 56
 Théorème de dérivabilité des intégrales paramétriques, 88
 Théorème de densité des fonctions continues dans les espaces L^p , 6
 Théorème de Dini, 11, 13
 Théorème de Dungundji, 18, 19
 Théorème de Dynkin-Sierpinski, 56, 60
 Théorème de Fubini, 59
 Théorème de Funahashi, 82
 Théorème de Hahn-Banach, 62, 64, 125, 126
 Théorème de Heine, 80
 Théorème de la convergence dominée, 64
 Théorème de la convergence monotone, 58
 Théorème de Levy, 56
 Théorème de Liouville, 89
 Théorème de Lusin, 3, 5, 9, 17, 66, 70, 77, 83
 Théorème de Novikov, 40
 Théorème de projection, 7
 Théorème de représentation de Riesz, 62, 129
 Théorème de Riesz-Fischer, 127
 Théorème de Stone-Weierstrass, 11, 12, 14, 75
 Théorème de Tietze, 18, 83
 Théorème des valeurs intermédiaires, 79
 Théorème du Noyau fermé, 62, 124
 Théorème du point fixe de Schauder, 109

Bibliographie

- [1] BJ Albert. Novikov. on convergence proofs for perceptrons. Technical report, DTIC Document, 1963.
- [2] Robert B Ash. *Real Analysis and Probability : Probability and Mathematical Statistics : a Series of Monographs and Textbooks*. Academic press, 2014.
- [3] Benaïssa Assia. Quelques propriétés et applications de l'opérateur fractionnaire de Caputo. Master's thesis, Université Dr Tahar Moulay - Saïda, Algérie, 2016/2017.
- [4] A Barron. Universal approximation bounds for superposition of a sigmoidal function. *IEEE Transaction on Information Theory*, 19 :930–944, 1993.
- [5] Pierre Borne, Mohamed Benrejeb, and Joseph Haggège. *Les réseaux de neurones : présentation et applications*, volume 15. Editions OPHRYS, 2007.
- [6] Nicolas Bourbaki. les structures fondamentales de l'analyse : livre iii : topologie générale. *Éléments de mathématique : première partie* :, 1949.
- [7] Franck Boyer. Agrégation externe de mathématiques analyse numérique. *ellipses*, 2013.
- [8] Michele Caputo. Linear models of dissipation whose q is almost frequency independent—ii. *Geophysical Journal International*, 13(5) :529–539, 1967.
- [9] Antoine Chambert-Loir, Stéphane Fermigier, and Vincent Maillot. *Exercices de mathématiques pour l'agrégation : analyse 1 : Antoine Chambert-Loir, Stéphane Fermigier, Vincent Maillot*. masson, 1995.
- [10] David E Chauvin, Yves et Rumelhart. *Backpropagation : théorie, architectures et applications*. Psychology Press, 2013.
- [11] E.W. Cheney. *Introduction to approximation theory*. International series in pure and applied mathematics. McGraw-Hill Book Co., 1966.
- [12] G. Cybenko. Correction : approximation by superpositions of a sigmoidal function. *Mathematics of Control, Signals and Systems*, 5(4) :455–455, Dec 1992.
- [13] George Cybenko. Approximation by superpositions of a sigmoidal function. *Mathematics of control, signals and systems*, 2(4) :303–314, 1989.
- [14] Claude Dellacherie. Un complément au théorème de Weierstrass-Stone. *Séminaire de probabilités de Strasbourg*, 1 :52–53, 1967.
- [15] Jean-Pierre Demailly. *Analyse numérique et équations différentielles-4ème Ed*. EDP sciences, 2016.
- [16] Gérard Dreyfus. *Réseaux de neurones : méthodologie et applications*. Eyrolles, 2004.
- [17] François Dubois, Ana Cristina Galucio, and Nelly Point. Introduction à la dérivation fractionnaire-théorie et applications. *Techniques de l'Ingénieur*, 2008.

-
- [18] James Dugundji et al. An extension of tietze's theorem. *Pacific Journal of Mathematics*, 1(3) :353–367, 1951.
- [19] Mohamed Ettaouil. Cours master :optimisation non linéaire et modélisation., 2017-2018.
- [20] Jerome Friedman, Trevor Hastie, and Robert Tibshirani. *The elements of statistical learning*, volume 1. Springer series in statistics New York, 2001.
- [21] Ken-Ichi Funahashi. On the approximate realization of continuous mappings by neural networks. *Neural networks*, 2(3) :183–192, 1989.
- [22] Federico Girosi and Tomaso Poggio. Networks and the best approximation property. *Biological cybernetics*, 63(3) :169–176, 1990.
- [23] Maurice Godefroy. *La fonction gamma : théorie, histoire, bibliographie*. Gauthier-Villars, 1901.
- [24] Detlef Gronau and Janusz Matkowski. Geometrical convexity and generalizations of the bohr-mollerup theorem on the gamma function. *Mathematica Pannonica*, 4(2) :153–160, 1993.
- [25] Namig J Guliyev and Vugar E Ismailov. Approximation capability of two hidden layer feedforward neural networks with fixed weights. *Neurocomputing*, 316 :262–269, 2018.
- [26] Jacques Hadamard. *Sur le rayon de convergence des series ordonnees suivant les puissances d'une variable*. Ecole royale polytechnique, 1888.
- [27] Paul R Halmos. *Measure theory*, volume 18. Springer, 2013.
- [28] Medjekal Hamza. *Existence et unicité de la solution d'une équation différentielle fractionnaire impulsive de temps infini dans un espace de Banach*. PhD thesis, Université Badji Mokhtar Annaba,Algérie, 2015.
- [29] D Hebb. The organization of behavior john wiley & sons. *New York*, 1949.
- [30] Th Hélié and Denis Matignon. Diffusive representations for the analysis and simulation of flared acoustic pipes with visco-thermal losses. *Mathematical Models and Methods in Applied Sciences*, 16(04) :503–536, 2006.
- [31] Francis Hirsch and Gilles Lacombe. *Eléments d'analyse fonctionnelle : cours et exercices*. Masson, 1997.
- [32] John J Hopfield. Neural networks and physical systems with emergent collective computational abilities. *Proceedings of the national academy of sciences*, 79(8) :2554–2558, 1982.
- [33] Kurt Hornik. Approximation capabilities of multilayer feedforward networks. *Neural networks*, 4(2) :251–257, 1991.
- [34] Kurt Hornik, Maxwell Stinchcombe, and Halbert White. Multilayer feedforward networks are universal approximators. *Neural networks*, 2(5) :359–366, 1989.
- [35] Kurt Hornik, Maxwell Stinchcombe, Halbert White, and Peter Auer. Degree of approximation results for feedforward networks approximating unknown mappings and their derivatives. *Neural Computation*, 6(6) :1262–1275, 1994.
- [36] Ahmad Jafarian, Safa Measoomy Nia, Alireza Khalili Golmankhaneh, and Dumitru Baleanu. On artificial neural networks approach with new cost functions. *Applied Mathematics and Computation*, 339 :546–555, 2018.
- [37] Jean-Pierre Kahane and Pierre Gilles Lemarié. *Séries de Fourier et ondelettes*. Presse Universitaire, 1998.
- [38] Teuvo Kohonen. Self-organized formation of topologically correct feature maps. *Biological cybernetics*, 43(1) :59–69, 1982.
- [39] Ph Leray. *Réseaux bayésiens : Apprentissage et diagnostic de systemes complexes*. PhD thesis, Université de Rouen, 2006.
-

-
- [40] K Ming Leung. Adaline for pattern classification. Technical report, Polytechnic institute of new york, 2008.
- [41] N Lusin. Sur les propriétés des fonctions mesurables. *CR Acad. Sci. Paris (Projet Numdam)*, 154(25) :1688–1690, 1912.
- [42] Warren S McCulloch and Walter Pitts. A logical calculus of the ideas immanent in nervous activity. *The bulletin of mathematical biophysics*, 5(4) :115–133, 1943.
- [43] Marvin Minsky and Seymour Papert. Perceptrons : An Introduction to computational geometry. MIT Press, Cambridge, Massachusetts, 1969.
- [44] Patrick Naïm, Pierre-Henri Wuillemin, Philippe Leray, Olivier Pourret, and Anna Becker. *Réseaux bayésiens*. Editions Eyrolles, 2011.
- [45] Alain Oustaloup. *La dérivation non entière*. Number 2. Hermes, 1995.
- [46] K Parthasarathy. Probability measure on metric spaces. *Journal of the American Statistical Association*, 63, 09 1968.
- [47] George Pólya and Gabor Szegő. *Problems and Theorems in Analysis II : Theory of Functions. Zeros. Polynomials. Determinants. Number Theory. Geometry*. Springer Science & Business Media, 1997.
- [48] M.J.D. Powell. *Approximation Theory and Methods*. Cambridge University Press, 1981.
- [49] Fred Riewe. Nonconservative lagrangian and hamiltonian mechanics. *Physical Review E*, 53(2) :1890, 1996.
- [50] Frank Rosenblatt. The perceptron : a probabilistic model for information storage and organization in the brain. *Psychological review*, 65(6) :386, 1958.
- [51] Hilfer Rudolf. *Applications of fractional calculus in physics*. World Scientific, 2000.
- [52] JATMJ Sabatier, Ohm Parkash Agrawal, and JA Tenreiro Machado. *Advances in fractional calculus*, volume 4. Springer, 2007.
- [53] Arthur Samuel. Some studies in machine learning using the game of checkers. *IBM Journal of Research and Development*, 3(3) :210–219, July 1959.
- [54] Arthur L Samuel. Some studies in machine learning using the game of checkers. *IBM Journal of research and development*, 44(1.2) :206–226, 2000.
- [55] Juliusz Schauder. Der fixpunktsatz in funktionalräumen. *Studia Mathematica*, 2(1) :171–180, 1930.
- [56] Vadim Schechtman. Introduction aux fonctions spéciales. *cours de calcul fractionnaire*, 2010.
- [57] Basel Solaiman and Lepage Richard. *Les réseaux de neurones artificiels et leurs applications en imagerie et en vision par ordinateur*. Presse Universitaire du Québec, 2003.
- [58] Marshall Harvey Stone. Applications of the theory of boolean rings to general topology. *Transactions of the American Mathematical Society*, 41(3) :375–481, 1937.
- [59] Claude Touzet. *les réseaux de neurones artificiels, introduction au connexionnisme*. EC2, 1992.
- [60] Onno van Gaans. Probability measures on metric spaces. Technical report, Delft University of Technology, 2002.
- [61] Vladimir N Vapnik. An overview of statistical learning theory. *IEEE transactions on neural networks*, 10(5) :988–999, 1999.
- [62] Bernard Widrow and Marcian E Hoff. Adaptive switching circuits. Technical report, Stanford Univ Ca Stanford Electronics Labs, 1960.
-