

DESSINE-MOI UN MOUTON

COMPTE-RENDU DE PROJET LONG

James BARRIOS, Antoine CORDELLE, Benjamin MOUSCADET, Emmanuel TRAN

encadrés par Arpad RIMMEL et Joanna TOMASIK

2017 – 2018



CentraleSupélec

Table des matières

1 Réseaux de neurones	4
1.1 Problème de classification	4
1.2 Apprentissage supervisé	4
1.3 Structure d'un neurone	4
1.4 Réseaux de neurones	5
1.5 Propagation dans un perceptron	6
2 Implémentation des premiers réseaux	7
2.1 Implémentation du perceptron	7
2.1.1 Présentation et diagramme UML	7
2.1.2 NeuronLayer	8
2.1.3 NeuralNetwork	8
2.1.4 DataCollector	8
2.2 Etude du XOR	9
2.2.1 Définition du problème	9
2.2.2 Résolution	9
3 Reconnaissance des chiffres manuscrits	11
3.1 La base de données MNIST	11
3.1.1 Présentation de MNIST	11
3.1.2 Extraction de MNIST	11
3.2 Apprentissage de la base de données	12
3.2.1 Paramétrage du réseau de neurones	12
3.2.2 Résultat	12
4 Implémentation en C++ et problèmes rencontrés	13
4.1 Outils utilisés	13
4.1.1 IDE	13
4.1.2 Eigen	13
5 Generative Adversarial Networks - Partie théorique	15
5.1 Description du fonctionnement des <i>Generative Adversarial Networks</i>	15
5.1.1 GAN - Idée générale	15
5.1.2 Des réseaux en compétition	16
5.2 Discriminateur et générateur	17
5.2.1 Théorie des jeux et équilibre de Nash	17
5.2.2 Les différentes distributions statistiques impliquées	17
5.3 Les fonctions de coût, l'apprentissage	17

5.3.1	Les différentes fonctions de coût	17
5.3.2	Apprentissage par récompense ou par punition?	17
5.3.3	Quand un réseau prend le pas sur l'autre	17
6	Generative Adversarial Networks - Implémentation	18
6.1	Implémentation des <i>Generative Adversarial Networks</i>	18
6.1.1	Présentation de la structure	18
6.1.2	Ratio d'apprentissage	18
6.2	Premiers résultats	18
6.3	Difficultés rencontrées et solutions émises	18
6.3.1	Le problème d'évaluation des performances	19
6.3.2	Mode Collapse du <i>Generateur</i> vers un unique point	20

Introduction

Ce compte-rendu présente le Projet Long Supélec *Dessine-moi un mouton*, mené par 8 élèves de Supélec - séparés en deux équipes de 4 - et encadré par deux enseignants-chercheurs, Monsieur Arpad Rimmel et Madame Joanna Tomasik. Plus particulièrement, c'est le travail de l'équipe **Couleuvre** qui est rapporté ici.

L'objectif de ce projet était d'étudier les réseaux de neurones et de les utiliser pour générer des images de synthèses. L'étude des réseaux de neurones comportait une dimension de recherche afin d'appréhender la théorie de ces objets, et une dimension de développement et de programmation dans l'implémentation de ces réseaux et leur utilisation. Le projet s'est donc naturellement décomposé en plusieurs étapes.

Dans un premier temps, il a fallu étudier le fonctionnement des réseaux de neurones, et établir les équations régissant le comportement des réseaux de type perceptron. Puis nous avons implémenté un perceptron afin de résoudre un problème de décision simple au début (**XOR**), puis plus complexe (application à la base de données de chiffres manuscrits MNIST).

Dans un second temps, nous avons commencé à utiliser la structure qui représentait le coeur du projet : les réseaux semi-supervisés en compétition mutuelle, selon la méthode du **Generative Adversarial Network**. Nous l'avons d'abord appliqué à la génération d'images synthétiques de chiffres, puis nous l'avons utilisé pour faire de la synthèse d'images plus complexes, comme des images de moutons.

Afin de mener à bien ce projet long, il a été nécessaire de mettre en place plusieurs outils de gestion de projet. Nous avons créé une organisation **Github** (<https://github.com/supelec-GAN>) afin de faciliter le partage du code; nous avons utilisé **Trello** pour s'organiser en tâches et optimiser le travail; une bibliographie commune a été mise en place grâce à **Zotero**; la documentation du projet a été faite avec **Doxygen**; enfin, une réunion hebdomadaire réunissant les deux équipes et les enseignants a été imposée, permettant de rendre compte des avancées, de recevoir les directives sur le travail à effectuer, et de manger du chocolat.

Enfin, nous avons choisi le langage C++ pour ce projet. En effet, la structure de réseaux de neurones et de computational graphs (qui sera détaillée plus loin) s'accorde naturellement avec la représentation objet; de plus les réseaux de neurones nécessitent par définition de mener des calculs en grandes quantités. La rapidité de calcul du C++ (héritée du C) ainsi que son approche objet robuste et particulièrement puissante en faisaient le candidat idéal. De plus, un seul membre du groupe n'était pas familier avec le C++ et a dû apprendre en cours de route. **Nous déconseillons vivement aux futurs membres de ce projet long d'utiliser un langage qui ne soit pas maîtrisé par au moins la moitié de l'équipe.**

Chapitre 1

Réseaux de neurones

1.1 Problème de classification

Dans la théorie de l'apprentissage statistique, la classification a pour objectif de déduire d'un nombre fini d'observations indépendantes une partition de l'espace en un ensemble a priori inconnu de domaines de l'espace appelés classes.

1.2 Apprentissage supervisé

1.3 Structure d'un neurone

L'extraordinaire capacité d'apprentissage et d'adaptation des réseaux de neurones biologiques a poussé les scientifiques à tenter de modéliser informatiquement leur fonctionnement afin d'exploiter ces capacités.

Les observations biologiques ont mené au premier modèle du neurone, divisé en 3 parties. Le neurone reçoit des signaux chimio-électriques par ses dendrites; ces signaux sont traités dans le corps cellulaire, où un effet de seuil est appliqué; enfin l'axone permet la transmission du signal chimio-électrique de sortie. De nouvelles découvertes sur les neurones ont permis de complexifier grandement ce modèle, mais c'est celui qui sert de base aux neurones artificiels.

Définition 1 (Neurone). Un neurone à n entrées permet de modéliser une fonction de \mathbb{R}^n dans \mathbb{R} . Un neurone est défini par la données de trois paramètres :

- un vecteur de poids $\omega \in \mathbb{R}^n$
- un biais $b \in \mathbb{R}$
- une fonction d'activation $g \in \mathbb{R}^{\mathbb{R}}$

La fonction f modélisée par le neurone s'écrit alors :

$$\forall x \in \mathbb{R}^n, f(x) = g(\omega^T x - b) = g\left(\sum_{i=0}^{n-1} \omega_i x_i - b\right) \quad (1.1)$$

Remarque 1 (Biais). Il est possible de rajouter une composante $x_n = -1$ à tous les vecteurs d'entrées afin de pouvoir considérer le biais b comme la composante ω_n du vecteur de poids. L'équation (1.1) devient alors

$$\forall x \in \mathbb{R}^n, f(x) = g(\omega^T x) = g\left(\sum_{i=0}^n \omega_i x_i\right) \quad (1.2)$$

Cependant, considérer le biais comme un poids pose des problèmes lorsqu'on utilise l'approche des graphes de calculs, qui sera explicitée dans la suite de ce rapport. Nous n'avons donc pas appliqué cette simplification pour nos propres réseaux.

1.4 Réseaux de neurones

Un neurone permet, en prenant n entrées, de renvoyer en sortie une valeur réelle entre 0 (ou -1) et 1. Au vu des fonctions d'activation utilisées, les réponses renvoyées peuvent être assimilées à des booléens (0 : 1). La fonction de transfert du réseau de neurone étant une combinaison linéaire à qui on applique le créneau (fonction d'activation), un neurone seul peut séparer un plan en deux. Par exemple, en prenant un réseau à deux entrées, une sortie, dans un problème de classification classique, ce réseau permet de séparer le plan (x,y) en deux parties séparées par une droite, comme montré ci-dessous :

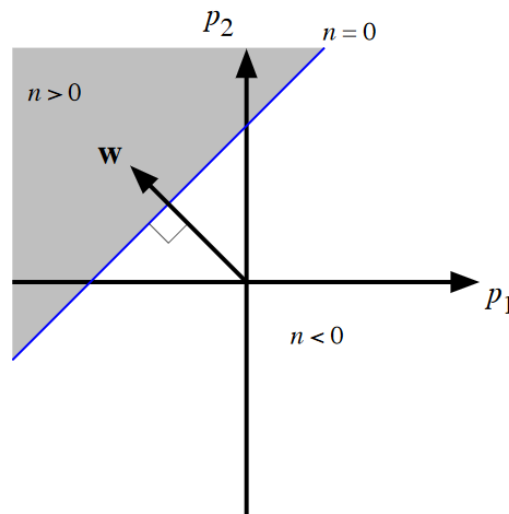


FIG. 1.1 – *Classification pour un neurone seul à deux entrées, une sortie*

Un neurone seul ne peut donc que séparer un plan en deux parties. Les possibilités limitées d'un neurone seul ont donc poussé l'informatique à agencer les neurones en réseaux afin de créer des modèles plus complexes. Ainsi, un réseau à deux couches peut (à condition d'avoir suffisamment de neurones dans ces couches) identifier toute zone convexe d'un espace vectoriel, alors qu'un réseau à 3 couches peut séparer toute partie de cet espace. On agence alors les neurones en réseaux afin de créer des modèles plus complexes.

Définition 2 (Réseau de neurone - Perceptrons). Un réseau de neurones est défini par un graphe orienté $\mathcal{G}(V,A)$ où les nœuds sont des neurones, les arêtes des liens entre les neurones et par un ensemble de neurones d'entrée $V_{in} \subset V$ et de neurones de sorties $V_{out} \subset V$. Une arête partant d'un neurone i vers un neurone j signifie que la sortie du neurone i est une entrée pour le neurone j . Notons f_j la fonction représentant le neurone j . Finalement, nous noterons un réseau $\mathcal{N}(V,A,V_{in},V_{out},f)$ où f est l'ensemble des fonctions des neurones.

Définition 3 (Perceptron). Un réseau de neurones est de type perceptron lorsque le graphe qui le représente est acyclique (réseau feedforward) et qu'il peut être représenté

comme une suite finie de couches de neurones telle que les sorties d'une couche sont exactement les entrées de la couche suivante.

1.5 Propagation dans un perceptron

On cherche à obtenir les équations qui régissent la propagation des entrées dans un réseau perceptron. Comme les neurones sont agencées par couches, on peut indiquer chaque neurone par son numéro de couche et sa position dans la couche. De plus, on parlera dans la suite indistinctement du neurone ou de la fonction qu'il modélise.

Alors, soit un réseau perceptron à N couches. Pour $i \in \llbracket 1, N \rrbracket$ soit n_i le nombre de neurones de la couche i . Pour tout i, j tel que $i \in \llbracket 1, N \rrbracket, j \in \llbracket 1, n_i \rrbracket$, le neurone (i, j) est défini par

- sa fonction d'activation $f_{i,j}$
- son vecteur de poids $\omega_{i,j}$
- son biais $b_{i,j}$

Comme on est dans un réseau perceptron, tous les neurones de la couche $i - 1$ sont connectés au neurone (i, j) . Par conséquent le vecteur d'entrée de ce neurone est un vecteur X_n de dimension $n_{i-1} \times 1$. Alors la sortie du neurone (i, j) est égale à $f_{i,j}(\omega_{i,j} X_{i-1}^T - b_{i,j})$. En écrivant $X_n = \begin{pmatrix} X_0^n & \dots & X_{n_i-1}^n \end{pmatrix}^T$, on a alors la relation de récurrence

$$X_i^{j+1} = f_{j,i}(\omega_{j,i} X_j^T - b_{j,i}) \quad (1.3)$$

Dans un réseau perceptron, on suppose de plus que toutes les fonctions d'activations des neurones d'une même couche sont identiques. On a alors $\forall i \in \llbracket 1, N \rrbracket, \forall j \in \llbracket 1, n_i \rrbracket, f_{i,j} = f_i$. On peut alors réécrire l'équation (1.3) sous forme matricielle. On pose

$$\begin{aligned} \forall i \in \llbracket 1, N \rrbracket, F_i &= \begin{pmatrix} f_i & \dots & f_i \end{pmatrix} \quad (\text{Fonction d'activation de la couche}) \\ W_i &= \begin{pmatrix} \omega_{i,0} & \dots & \omega_{i,n_i-1} \end{pmatrix}^T \quad (\text{Matrice des poids de la couche}) \\ B_i &= \begin{pmatrix} b_{i,n_i-1} & \dots & b_{i,n_i-1} \end{pmatrix}^T \end{aligned}$$

En nommant X_0 le vecteur en entrée du perceptron, on a la nouvelle équation matricielle :

$$\forall i \in \llbracket 1, N \rrbracket, X_i = F_i(W_i X_{i-1} - B_i) \quad (1.4)$$

Chapitre 2

Implémentation des premiers réseaux

Ce chapitre décrit l'implémentation de nos réseaux de neurones sous forme de perceptrons, que nous utiliserons par la suite dans les différentes applications. Nous verrons ici cette implémentation et le premier exemple avec le XOR.

2.1 Implémentation du perceptron

2.1.1 Présentation et diagramme UML

Comme décrit dans la partie théorique sur les perceptrons (1.5), nous avons choisi de représenter nos réseaux de neurones sous forme matricielle. Les neurones ne seront donc pas représentés individuellement, mais par couche de neurones (`neuronLayer`). Cela permet d'effectuer les calculs sous forme matricielle par rapport à des propagations et rétro-propagations neurone par neurone, ce qui fait gagner beaucoup de temps de calcul. Cela est possible en gardant les mêmes fonctions d'activation pour chaque neurone d'une même couche, ainsi que les mêmes entrées et sorties. Dans le cadre du perceptron, cela convient.

Le diagramme UML comporte une classe principale, `Application`, qui va gérer les expériences et réseaux de neurones. Cette classe contient un réseau de neurones (`NeuralNetwork`) composé de plusieurs couches de neurones (`NeuronLayer`). `Application` possède aussi un collecteur qui est utilisé pour récupérer les données des expériences, les traiter et les exporter dans un fichier `.csv`. Enfin, la classe `Teacher` a pour rôle de gérer la rétropropagation et l'apprentissage du réseau de neurones.

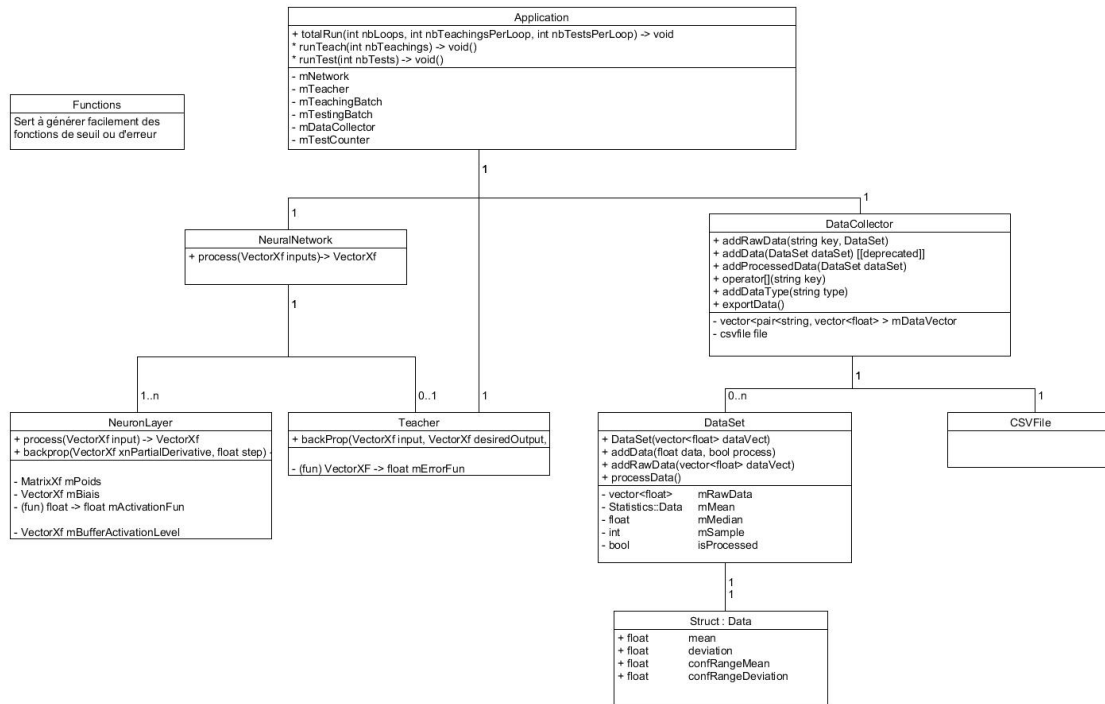


FIG. 2.1 – Diagramme UML

2.1.2 NeuronLayer

Une instance de la classe **NeuronLayer** représente une couche de neurones. Comme nous utilisons la représentation matricielle des réseaux perceptrons, une couche de neurones n'est constituée que d'une matrice de poids, d'un vecteur de biais et d'une fonction d'activation. Si *input* et *output* représentent respectivement la dimension du vecteur d'entrée et de sortie de la couche, alors on a les dimensions suivantes (on travaille en ligne) :

- La matrice de poids est de dimension $input \times output$
- Le vecteur de biais est un vecteur ligne de dimension *output*

2.1.3 NeuralNetwork

Une instance de la classe **NeuralNetwork** représente un réseau perceptron, c'est à dire plusieurs couches de neurones (**NeuronLayer**) mises les unes à la suite des autres. Chaque couche de neurones reçoit en entrée la sortie de la couche précédente et transmet sa sortie à la couche suivante. Pour représenter cet objet, nous avons utilisé la structure de liste simplement chaînée. La class **NeuralNetwork** hérite donc de la classe **std::list<NeuronLayer>**. Cette structure légère et robuste permet d'itérer facilement sur toutes les couches du réseau.

2.1.4 DataCollector

Afin de collecter les données du réseau et de pouvoir effectuer un traitement statistique dessus (calcul de moyenne, ecart-type, intervalle de confiance...) nous avons écrit la classe

DataCollector. Elle se charge de récupérer les sorties du réseau, d'effectuer des calculs statistiques standards et de les exporter au format **.csv** afin de permettre une visualisation facile de toutes les données.

2.2 Etude du XOR

Le problème du XOR est un problème de classification standard, auquel toutes les méthodes de classification sont confrontées. Nous avons donc naturellement commencé par le traitement de ce problème pour nous exercer à la manipulation des perceptrons.

2.2.1 Définition du problème

Le problème du XOR est un problème en deux dimensions : les entrées sont des paires (x_1, x_2) et la sortie est un booléen. Deux approches différentes existent, selon le choix du domaine de définition :

- 1^{er} choix : on veut obtenir un XOR purement booléen, dont les 4 entrées possibles sont les paires $(0,0)$, $(0,1)$, $(1,0)$, $(1,1)$. Les résultats attendus sont respectivement 0, 1, 1, 0.
- 2nd choix : on veut obtenir un XOR défini sur le carré unité du plan : $(x,y) \in [-1,1]^2$. Dans ce cas, on attend du XOR qu'il renvoie 1 lorsque $x \times y < 0$ et 0 sinon. On fait donc un XOR sur le signe des coordonnées.

Selon l'approche, le réseau requis sera différent. Dans le premier cas, il est possible de classer les 4 points selon leur appartenance à un ensemble convexe (qui prend la forme d'une bande). Par conséquent, le problème peut se résoudre avec deux couches de neurones.

Dans le second cas, il faut séparer le plan en 4 quadrants, et regrouper les quadrants par deux. On obtient des classes non convexes, et le problème requiert 3 couches de neurones pour être résolu.

Nous avons choisi de travailler sur la deuxième situation, le but est donc d'obtenir un réseau qui découpe le carré unité en quatre quadrants bien distincts, qui correspondent aux quatre quadrants délimités par les axes des abscisses et des ordonnées.

2.2.2 Résolution

Il s'agit d'un problème à deux dimensions, renvoyant une seule valeur, et nécessitant 3 couches cachées. Nous optons donc pour un réseau 2-2-1.

Qualitativement, on peut dire que

- les deux neurones de la première couche vont se charger de tracer les deux droites correspondant aux abscisses et aux ordonnées
- les deux neurones de la seconde couche vont se charger de réaliser les fonctions ouët ẽtšur les quadrants ainsi tracés, afin de d'obtenir deux quadrants correspondants au résultat 1 et deux quadrants correspondants au résultat 0.
- la troisième couche se charge de réaliser la fonction oušur les résultats renvoyés par la couche 2, afin que le résultat final soit 1 si le point se trouve dans l'un des deux quadrants correspondant au résultat 1.

De plus, nous avons choisi d'utiliser des fonctions d'activations sigmoïdes en $\frac{1}{1+\exp(-\lambda x)}$. Une telle fonction d'activation est nécessaire pour pouvoir permettre un apprentissage par rétro-propagation (continue et dérivable), cependant elle empêche d'avoir un résultat booléen à la fin, puisque la valeur de sortie varie continuellement entre 0 et 1. Le résultat est donc une densité de probabilité.

Enfin, nous utilisons un pas d'apprentissage de $\eta \leq 0.01$. Bien que la littérature suggère souvent le pas empirique $\eta = 0.2$, la valeur de 0.01 est la valeur maximale que nous ayons trouvée qui permette au XOR de converger. Au delà, le réseau se bloque presque systématiquement dans un état non désiré.

Avec une dizaine de milliers d'apprentissages, on obtient le résultat suivant :

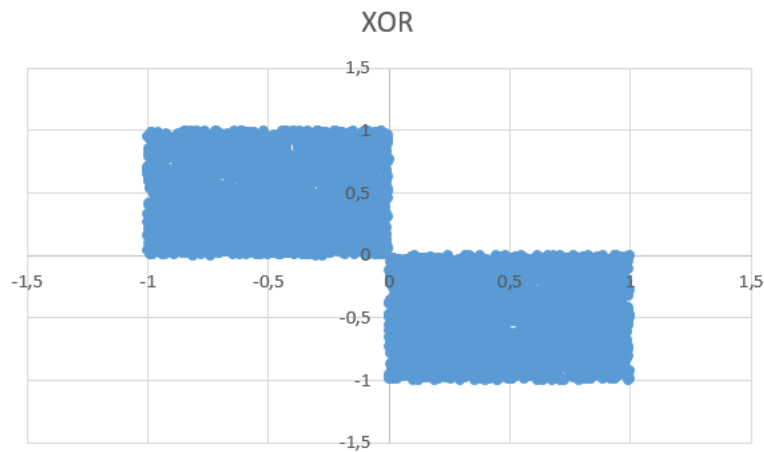


FIG. 2.2 – *Résultat d'un XOR sur 1000 tests après 10000 apprentissages pour un réseau 2-2-1*

Chapitre 3

Reconnaissance des chiffres manuscrits

3.1 La base de données MNIST

3.1.1 Présentation de MNIST

La base de données MNIST (ou Mixed National Institute of Standards and Technology) est un regroupement de 70 000 chiffres manuscrits. 60 000 d'entre eux servent à l'apprentissage et les 10 000 autres permettent de tester le réseau de neurones après l'apprentissage. Ce sont des images normalisées en noir et blanc, de 28 pixels de côté chacun codé sur un octet. Nous nous appuyerons sur cette base de données pour faire apprendre à notre réseau de neurones la reconnaissance de chiffres.

3.1.2 Extraction de MNIST

Les données sont dans le format idx1 et idx3 tel que :

offset	type	valeur	description
0000	32 bit integer	0x00000803(2051)	magic number
0004	32 bit integer	60000	number of items
0008	32 bit integer	28	number of rows
0012	32 bit integer	28	number of columns
0016	unsigned byte	??	pixel
0017	unsigned byte	??	pixel
...
xxxx	unsigned byte	??	pixel

On peut distinguer :

- Le *magic number* qui permet d'identifier le format de la base de données
- Les *number of items*, *number of rows*, *number of columns* donne des informations sur les données, ce qui nous permettra d'extraire les images
- Les *pixels* qui sont les octets en niveau de gris des images

L'extraction se fait donc en lisant successivement les octets et en les stockant dans des *vector<Eigen:Matrix>* grâce aux données récupérées. En faisant de même avec les labels dont le format est très sensiblement le même, on obtient une liste de vecteur colonne

d'*Eigen::Matrix* de taille 784 (28*28) ainsi que leur label sur un autre *Eigen::Matrix*. On fait de même avec les échantillons de test et on est fin prêt pour l'apprentissage.

3.2 Apprentissage de la base de données

3.2.1 Paramétrage du réseau de neurones

Bien que le format de l'entrée soit le même (vecteur colonne si ce n'est la taille qui change), un réseau à 3 ou 4 neurones ne suffit pas. Il a fallu adapter de façon empirique les différents paramètres. On peut jouer sur :

- Le nombre de neurones
- Le nombre de couches
- Le pas d'apprentissage
- Les fonctions d'activations par couches
- Le nombre d'apprentissages

Le nombre de couches et de neurones

Le pas d'apprentissage

Les fonctions d'activations par couches

Le nombre d'apprentissages

3.2.2 Résultat

Chapitre 4

Implémentation en C++ et problèmes rencontrés

Ce chapitre a pour but d'aider les prochains qui souhaiteraient coder en C++ pour qu'ils ne se heurtent pas aux mêmes écueils que nous, et qu'ils puissent passer plus de temps sur la conception du projet plutôt que sur le debug et la recherche d'outils efficaces.

4.1 Outils utilisés

4.1.1 IDE

Il est important que chaque membre développe dans un IDE qui lui soit familier, et qui offre un minimum d'ergonomie et d'efficacité. Dans le projet, les IDE utilisés ont été :

- XCode, pour développer sous MAC.
- QtCreator, disponible sous Linux et Windows (bien qu'il soit un peu plus pénible à installer sous windows).
- Vim, déconseillé à ceux qui ne maîtrisent pas l'outil, mais qui peut se révéler très puissant avec les bons plugins.

Chaque IDE venant avec ses propres fichiers de configuration, nous conseillons de nommer ceux-ci du nom de leur utilisateur (typiquement `john_doe.pro.user`) ou de les ajouter au `.gitignore` afin d'éviter les conflits.

4.1.2 Eigen

Les réseaux de neurones de type perceptron reposent largement sur les produits matriciels. Afin d'effectuer ces calculs le plus efficacement possible, notre code utilise le framework C++ Eigen. Ce framework permet de manipuler l'algèbre linéaire et les matrices de manière la plus optimale possible. De plus, Eigen utilise `openmp` ce qui permet de paralléliser les calculs sur l'ensemble des cœurs.

De ce que nous avons pu voir, il n'existe aucune alternative à Eigen qui offre autant de possibilités et d'optimisation. Cependant Eigen n'est pas facile à manipuler au premier abord. Sa construction entièrement à base de template et organisée de façon à pouvoir effectuer le plus de calculs possible en compile-time rend les types des objets et les messages d'erreur difficilement déchiffrables. Debugguer un problème lié à Eigen nécessite un bon débogueur et de la patience. De plus, la méthode de parallélisation d'Eigen reste assez obscure. Le framework semble paralléliser les calculs en colonnes, c'est à dire qu'il calcule

simultanément les colonnes de la matrices résultante d'un calcul. Autrement dit, si le résultat est une ligne, le calcul sera parallélisé, mais si c'est une colonne il ne tournera que dans un seul thread. Ceci implique de n'utiliser que des matrices lignes dans le code, bien que tous les calculs menés théoriquement soient faits avec des colonnes. Attention donc à transposer correctement.

Chapitre 5

Generative Adversarial Networks - Partie théorique

Après avoir appris à manipuler correctement les réseaux perceptrons simples, nous nous sommes intéressés à la structure de **Generative Adversarial Networks**, qui a constitué le coeur du projet, et son enjeu majeur.

La naissance du GAN se place dans un contexte de recherche de moyens innovants et efficaces de génération de données arbitraires. Dans l'ère de l'information, les données constituent une ressource fondamentale, et la capacité d'en générer facilement de nouvelles représente un pan entier de la recherche, si ce n'est plusieurs. Le GAN sont une solution à base de réseaux à apprentissage semi-supervisés, dont l'objectif est de générer des données arbitraires en extrapolant à partir de données initialement fournies.

Dans ce chapitre, nous décrirons qualitativement le fonctionnement attendu des GANs et leur théorie, puis nous détaillerons les aspects plus techniques et mathématiques du problème.

5.1 Description du fonctionnement des *Generative Adversarial Networks*

Le GAN est une structure qui fait intervenir deux réseaux placés en compétition, qui vont s'émuler mutuellement afin d'atteindre un optimum, ceci en partant de rien. Les GAN représentent un domaine de recherche qui s'est ouvert en 2014 avec la publication d'un article de Ian Goodfellow.

En 2014, l'équipe de Ian Goodfellow publie dans NIPS un article intitulé **Generative Adversarial Nets**, dans lequel il décrit une nouvelle structure censée répondre au problème de génération de donnée. Cette structure est fondée sur l'entraînement simultané de deux réseaux neuronaux en compétition, respectivement nommés **Générateur** et **Discriminateur**.

5.1.1 GAN - Idée générale

Pour commencer, on souhaite obtenir un réseau qui soit capable de générer sur commande des données

- Pertinentes par rapport aux contraintes qu'on lui impose (exemple : on souhaite pouvoir obtenir un réseau qui dessine spécifiquement des moutons, et non un réseau

qui dessine n'importe quoi comme bon lui semble)

- Variées, c'est à dire que les données générées sont différentes deux à deux
- Réalistes, voire indiscernables de données réelles aux yeux d'un humain

En résumé on souhaite construire un réseau (nommé dans toute la suite **Générateur**) qui modélise une projection d'un ensemble d'entrée vers l'ensemble des données réalistes et pertinentes vis-à-vis des contraintes imposées, et dont l'image ne soit pas réduite à un point de cet ensemble.

Résumée de cette façon, la tâche semble extraordinairement ardue, c'est pourquoi la structure fait appel à un second réseau (nommé **Discriminateur** dans toute la suite). En effet, la difficulté principale liée aux réseaux de neurones est l'apprentissage. Dans ce cas en particulier, la difficulté majeure est de caractériser ce qu'est une donnée de synthèse pertinente et réaliste, voire indiscernable d'une donnée réelle; car s'il est très facile de distinguer à l'oeil une image de mouton de n'importe quoi d'autre, il est extrêmement difficile de construire un critère mathématique permettant de distinguer les images de mouton des autres. Or c'est ce critère dont nous avons besoin pour faire apprendre correctement à notre **Générateur**.

Cependant, nous avons vus avec les perceptrons et la base de données MNIST qu'il était parfaitement possible d'obtenir un réseau de neurones modélisant avec une grande précision des critères visuels, que l'on serait bien incapables d'exprimer mathématiquement. Il est par exemple impossible d'écrire analytiquement ce qu'est un chiffre manuscrit, mais il est possible de générer un réseau de neurones qui sache discriminer les chiffres manuscrits des autres images. Ainsi, en déléguant à un second réseau la responsabilité de modéliser le critère selon lequel nous voulons que nos données soient générées (exemple: on ne veut générer que des images de mouton, on va faire apprendre au **Discriminateur** à différencier un mouton de toutes les autres images possibles), on résout le (premier) problème majeur qui se pose.

5.1.2 Des réseaux en compétition

Le GAN va donc faire intervenir deux réseaux, dont l'un sert seulement à rendre possible l'apprentissage de l'autre. Le **Discriminateur** a pour unique rôle de modéliser le critère que l'on souhaite utiliser; son apprentissage se fait donc de la même manière que tout ce qui a été fait jusqu'ici, en apprentissage supervisé. C'est un réseau dont l'entrée est de la dimension des données que l'on souhaite construire, et dont la sortie est un réel entre 0 et 1, caractérisant la certitude du réseau que l'entrée qui lui été donnée est une donnée réelle ou non. Plus précisément, pour une entrée donnée, plus la sortie est proche de 1, plus le **Discriminateur** est convaincu que l'entrée est une donnée réelle, et non de synthèse. Son objectif est donc de sortir 1 pour toute entrée réelle, et 0 pour toute donnée de sythèse.

Le **Générateur** de son côté est le réseau qui doit générer des données synthétiques suffisamment convaincantes pour que le **Discriminateur** ne puisse les différencier des données réelles. Sa sortie est donc naturellement de la dimension des données qu'on souhaite construire. L'entrée du **Générateur** est un bruit (blanc) à partir duquel le réseau doit construire une donnée convaincante. Le **Générateur** est donc un réseau qui réalise une projection de l'ensemble d'entrée vers l'ensemble des données suffisamment convaincantes pour tromper un humain: Passer un bruit en entrée permet d'obtenir des sorties variées. On comprend cependant que, le réseau réalisant un calcul déterministe, la dimension de la

sorties est nécessairement inférieure à la dimension de l'entrée. Pour être certain d'obtenir un ensemble de sortie aussi grand que possible, on s'assurera toujours d'avoir un ensemble d'entrée de dimension suffisamment grande.

Exemple 1. Cas des images de mouton

On souhaite obtenir un réseau qui génère des images 28×28 de moutons, suffisamment réalistes pour ne pouvoir être distinguées de réelles photos de moutons. On souhaite également observer de la variation dans ces images (avoir plein de moutons différents).

L'ensemble de sortie du **Générateur** est donc un ensemble inclu dans l'ensemble des images 28×28 et qui contient toutes les images que le **Générateur** considère comme étant des moutons réalistes. Pour être certain que l'ensemble d'entrée soit de dimension au moins aussi importante que l'ensemble de sortie, on mettra en entrée du réseau un bruit blanc de dimension 28×28 . On est certain de cette façon que la taille de l'ensemble d'entrée n'impose pas de restriction sur les sorties.

5.2 Discriminateur et générateur

5.2.1 Théorie des jeux et équilibre de Nash

$D(G(z)) = D(x) = \frac{1}{2}$: un premier critère d'optimalité

5.2.2 Les différentes distributions statistiques impliquées

5.3 es fonctions de coût, l'apprentissage

5.3.1 Les différentes fonctions de coût

5.3.2 Apprentissage par récompense ou par punition ?

5.3.3 Quand un réseau prend le pas sur l'autre

Chapitre 6

Generative Adversarial Networks - Implémentation

La structure théorique des **Generative Adversarial Networks** étant maintenant définie, nous devons implémenter ces réseaux. Pour cela, nous nous sommes lancés dans la continuité de MNIST, avec la génération de chiffres manuscrits par le **Générateur**, qui doit alors convaincre le **Discriminateur**. Ce cas de figure est idéal pour une première implémentation des GAN, puisque nous restons avec des données relativement simple (Dimension 28×28 , en niveau de gris), avec un format normalisé et une importante base de données à disposition (La base de données MNIST : 60000 images d'apprentissages, 10000 données de test).

Le **Discriminateur** s'entraîne donc à reconnaître les chiffres provenant de la base de données de MNIST en les différenciant des chiffres factices créés par le **Générateur**. Dans un premier temps, nous travaillons sur des chiffres uniques, en apprenant au Discriminateur à reconnaître uniquement de 3 par exemple.

Nous verrons dans ce chapitre tout d'abord l'implémentation de notre GAN, puis les premiers résultats obtenus. Ces résultats nous ont confronté à divers problèmes, que nous aborderons alors avec les différentes solutions envisagées.

6.1 Implémentation des *Generative Adversarial Networks*

6.1.1 Présentation de la structure

6.1.2 Ratio d'apprentissage

6.2 Premiers résultats

6.3 Difficultés rencontrées et solutions émises

Au cours du développement des GANs, nous avons été confrontés à de très nombreuses difficultés, que nous avons tenté de répertorier avec les solutions envisagées. À l'heure actuelle, certains problèmes sont toujours sans solution probante et restent des domaines de recherches ouverts : c'est le cas en particulier du **mode collapse** que nous détaillons ci-après.

6.3.1 Le problème d'évaluation des performances

Le premier et plus grand problème que nous avons rencontré est **l'absence de méthode systématique d'évaluation des performances d'un réseau donné**.

En effet, pour qu'un protocole expérimental puisse être mené à bien, il faut comparer un résultat témoin avec les résultats pour lesquels des paramètres ont varié, et déduire de la différence des performances si le paramètre testé est signifiant, et dans quelle mesure.

Ce protocole est **fondamental** en recherche expérimentale, et il est **inapplicable** ici. En effet, nous n'avons trouvé aucun moyen de mesurer la qualité d'une image, à l'exception de la mesure à l'oeil par les membres de l'équipe (ce qui ne peut être qualifié de méthode rigoureuse et systématique d'évaluation).

Les données à disposition pour mesurer les performances

Afin de bien cerner le problème, résumons ce que nous avons à disposition à la fin d'une expérience. Si dans certains cas nous avons fait des essais nous permettant de mesurer d'autres grandeurs que les suivantes (comme la valeur des poids etc), seules les deux ci-après sont accessibles systématiquement, pour chaque expérience. Nous disposons donc de :

- Un ensemble d'images de synthèse générées à différents instants d'apprentissage; par exemple nous pouvons demander à produire 100 images tous les 10 apprentissages, ce qui permet de suivre assez finement l'évolution du GAN au fur et à mesure de son apprentissage
- Les courbes de score des deux réseaux, générées à la fin de l'expérience. La courbe de score d'un réseau représente la proportion d'erreurs dans toutes les réponses qu'il donne. Ainsi la courbe de score du discriminateur compte le nombre de bonnes réponses du réseau à des entrées aussi bien de synthèse (réponse attendue : 0) que réelles (réponse attendue : 1); la courbe du générateur est calculée en fonction des réponses du discriminateur aux images de synthèse. Ces courbes sont graduées en abscisse par rapport à nombre d'apprentissages effectués.

Sur l'utilisation des courbes de score À première vue, ces courbes paraissent être un indicateur pertinent de la qualité des réseaux. Elles permettent **en théorie** :

- d'observer l'équilibre de Nash (qui, en pratique, n'a jamais été observé)
- d'observer les déséquilibres entre les deux réseaux afin d'affiner les paramètres d'équilibre sur les expériences suivantes
- de comparer les images produites avec les scores des réseaux au moment de la génération de l'image

En pratique, nous n'avons jamais réussi à interpréter efficacement ces courbes. Leur étude représente une perte pure de temps, temps qui aurait dû être alloué à la recherche d'une méthode d'évaluation efficace. Les courbes paraissent extrêmement instables, on peut observer des formes récurrentes mais très difficiles à expliquer (chute brutale de l'une des deux malgré une bonne qualité d'image, etc), et par dessus tout nous n'avons pas réussi à mettre en lumière des corrélations entre les paramètres choisis et les courbes résultantes.

Enfin, nous avons été incapables de trouver des corrélations pertinentes entre la qualité des images et les courbes de score. Nous avons au mieux été capables de trouver des explications à posteriori sur les liens qui peuvent exister; c'est à dire que nous n'avons

jamais pu prévoir (de manière reproductible et indépendante du hasard) la qualité des images à partir de la seule analyse des courbes.

Sur l'analyse des échantillons d'images Deux problèmes principaux s'opposent à l'analyse

Évaluation par réseau(x) tiers

Une première idée pour l'évaluation, serait de passer les images de synthèses dans un 3ème réseau (appelons le évaluateur, mais il se comporte exactement comme le discriminateur) qui aurait appris à part, en dehors de la structure de GAN et qui serait un juge objectif sur les résultats du GAN. Ce réseau aurait été entraîné sur les données réelles, et resterait immuable une fois son apprentissage terminé.

Cette méthode est en pratique inapplicable (et inappliquée, pour ce que nous en savons), notamment à cause de la difficulté à équilibrer ce 3ème réseau avec les deux autres; équilibrer un GAN est déjà une tâche particulièrement ardue (nous y reviendrons), ajouter cet évaluateur la rend insurmontable. Car si l'évaluateur est très performant, il risque de donner une performance nulle à tous les générateurs qu'il testera; à l'inverse si on l'a délibérément fait peu apprendre, il ne permettra pas de discriminer efficacement les générateurs qu'il testera.

On pourrait envisager de créer une échelle d'évaluateurs, constituée de plusieurs évaluateurs ayant chacun arrêté leur apprentissage à un moment différent, afin de pouvoir graduer les résultats obtenus par les générateurs sur une échelle de performance. Nous n'avons pas nous même implémenté cette technique et n'avons pas trouvé d'articles faisant mention de l'application de cette méthode, mais elle a au moins l'avantage majeur d'être rigoureuse, systématique et indépendante du jugement humain.

6.3.2 Mode Collapse du *Generateur* vers un unique point

Présentation du problème

Un problème récurrent que l'on peut observer dans les GAN est le *mode collapse*, ou *missing modes*. Prenons un GAN, le générateur doit tenter de créer des images les plus réalistes possible par rapport à une distribution de données réelles, par l'intermédiaire du discriminateur. Cette distribution peut présenter plusieurs modes, plusieurs zones en quelque sorte (En reprenant l'exemple des chiffres de MNIST, on peut vouloir générer des 3, ou des 4).

Cependant, on observe que le générateur a souvent tendance à apprendre à ne générer qu'un seul mode. On appelle cela le mode collapse : le générateur apprendra à aller vers une des zones de la distribution voulue, se rendra compte que cette zone marche, et y restera. Sur l'exemple simplifié ci-dessous, la distribution d'origine présente 4 zones, et le GAN n'en couvre qu'une seule. Le discriminateur apprendra alors à se méfier de cette zone, puisque, même si elle est très proche de la distribution d'origine, il apprend au fur et à mesure à reconnaître les images du générateur. Ainsi, cette zone est peu à peu considérée comme artificielle pour le discriminateur, et le générateur va collapse sur une autre zone considérée comme valide. On aura alors un jeu de rebonds, le générateur collapse sur des zones qui seront alors désapprouvées par le discriminateur, et se déplacera vers une autre zone. On a donc un cercle qui fera que le générateur sera et restera dans un état

de collapse. Il faut alors trouver des méthodes pour empêcher le générateur de ne générer qu'une distribution très restreinte (un mode).

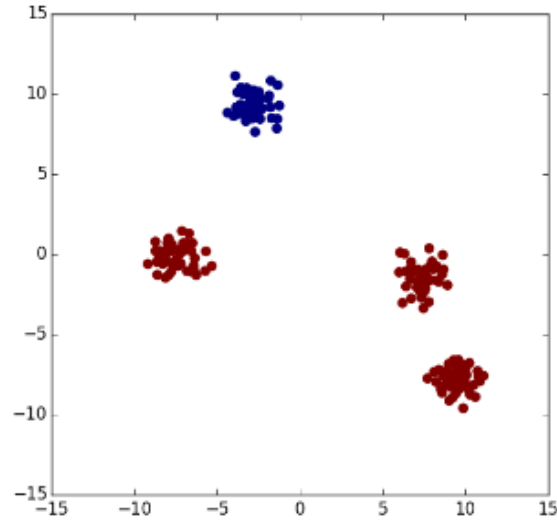


FIG. 6.1 – *Missing modes* : les points représentent les éléments de la distribution de données réelles. En bleu l'ensemble de sortie du GAN appartenant à la distribution réelle, et en rouge les zones de la distribution initiale non couverte par le GAN.

Ce phénomène récurrent est observable à différents degrés : on peut avoir des collapses partiels, les données générées seront donc variables mais toute la distribution ne sera pas représentée par le générateur, ou des collapses totaux, où le générateur ne sera qu'une projection vers un point. Les données générées seront alors très proches les unes des autres. Nous avons été confrontés à ce problème de collapse total en générant des chiffres à partir de la base de données MNIST. Comme montré ci-dessous, en essayant de générer des 3, nous obtenons pour un réseau donné, des images générées qui sont toutes très proches, avec la même forme de 3, les mêmes pixels singuliers.

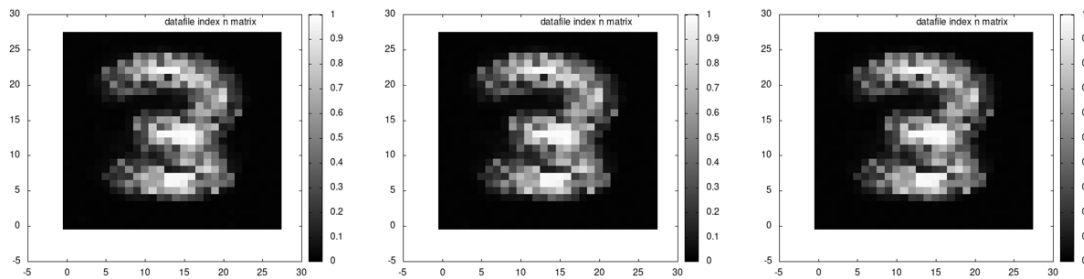


FIG. 6.2 – *Missing modes* : On obtient pour un même réseau, 3 échantillons générés qui sont sensiblement les mêmes.

Solution envisagée : Retrait du biais dans les réseaux (non fructueuse)

Une première hypothèse concernant ce problème de Mode Collapse était la suivante : l'unique image produite par le générateur n'était en réalité qu'une image du masque de biais qui figeait l'image, par exemple en plaçant certains biais trop haut ou trop bas, annulant les entrées bruitées en début de réseau. On aurait alors une image déterminée uniquement par les biais, ce qui expliquerait un collapse vers une unique image. De plus, en observant précisément le masque de biais sous forme matricielle, on remarque que ses poids forment effectivement l'image d'un nombre, comme on peut le voir ci-dessous :

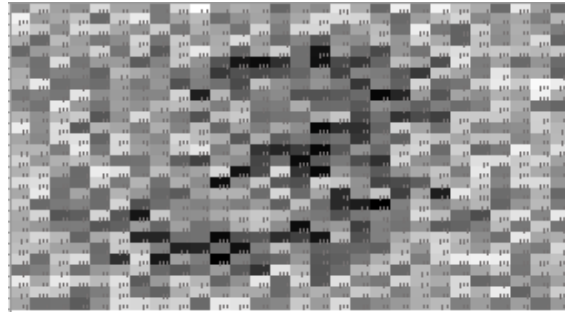


FIG. 6.3 – *Masque du biais : on peut observer que le biais forme lui même les bases d'une image, un 3 ici.*

On peut vérifier cette supposition en supprimant les biais de nos réseaux de neurones. On aurait alors plus de valeurs constantes qui figeraient l'image. Cependant, en implémentant ces réseaux, on n'observe qu'un ralentissement de l'apprentissage vers un état de Collapse similaire : les résultats sont moins intéressants (chiffres moins beaux), l'apprentissage est plus lent, et on converge toujours vers un Mode Collapse.

La suppression du biais ne permet donc pas de résoudre ce problème, c'est bien la fonction de transfert globale du réseau de neurones qui converge vers cet unique point.

Solution envisagée : Injection de bruit dans chaque couche du Générateur (fructueuse)

Une solution est proposée par Ian Goodfellow dans son article dans NIPS 2016 intitulé **Tutorial: Generative Adversarial Networks**. La structure du générateur est assez libre dans l'architecture des GAN, ainsi, on peut placer des entrées (du bruit) au début du réseau, mais aussi sur n'importe quelle couche de celui-ci notamment, sans réellement compromettre les performances. En effet, injecter des entrées de bruit dans des couches intermédiaires revient à faire des entrées classiques qui traverseront moins de couches de neurones. Cependant, le bruit au niveau de la dernière couche semble perturber le réseau et l'empêche de tomber dans un mode collapse.

Cela peut s'expliquer par le fait que les entrées bruitées perturbent la dernière couche du réseau, l'empêchant de se stabiliser sur une unique sortie trompant le discriminateur. Cette perturbation va forcer le générateur à diversifier ses sorties. Une autre hypothèse serait que la dernière couche ne peut pas différencier les entrées provenant des couches antérieures des entrées de bruit si celles-ci sont suffisamment nombreuses. Ainsi, il sera obligé de toutes les prendre en compte (il ne peut pas juste annuler les poids correspondant) ce qui générera de la variabilité dans la sortie, puisqu'une partie non négligeable des entrées

n'est que du bruit.

Cependant on peut observer dans des réseaux bruités, pour un même nombre d'apprentissages, des patterns généraux dans chaque image, ce qui nous montre un mode collapse partiel mais très léger comparé au précédent. De tels réseaux ont de plus besoin de beaucoup plus d'apprentissages pour donner des résultats convaincants. Cette solution nous permet donc de limiter partiellement ce problème de mode collapse.