
RECENT ADVANCES IN MACHINE LEARNING

2023-2024
Théo Lopès-Quintas

Séance 1

Rien ne sert de courir, il faut partir à point

1.1 Partir du bon pied

1.1.1 Explosion et disparition des gradients

Pour mieux comprendre le problème d'explosion et de disparition des gradients, nous devons nous pencher sur l'apprentissage d'un réseau : l'algorithme de backpropagation. Commençons par formaliser un réseau de neurones.

On considère un réseau de neurones composé de L couches avec N_l neurones à la couche $l \leq L$. Chacun de ces neurones est caractérisé par trois objets :

- **Poids** : un vecteur de paramètre qui associe à chacun des coefficients du vecteur d'entrée du neurones un poids
- **Biais** : un scalaire qui rend affine la transformation
- **Activation** : une fonction qui à partir de la transformation affine du vecteur d'entrée renvoie un nombre

La non-linéarité est introduite par la fonction d'activation. Résumons les notations introduites :

L'objectif d'un réseau de neurone, comme tous les algorithmes de Machine Learning, peut s'écrire sous forme d'un problème d'optimisation. Pour cela, on définit \mathcal{C} la fonction de coût et on cherche à la minimiser. Par exemple dans un problème régression une fonction de coût possible est la MSE, et dans le cadre d'une classification la Cross-Entropy.

Puisque les notations sont déjà lourde, nous prendrons quelques liberté de notation pour privilégier des équations simple quand il n'y a pas d'ambiguïté. Ainsi, l'objectif de la back propagation est de calculer les quantités :

- $\frac{\partial \mathcal{C}}{\partial w_{j,k}^l}$: pour calculer la mise à jour du k -ième poids du neurone $j \leq N_l$ à la couche $l \leq L$, avec $k \leq N_{l-1}$
- $\frac{\partial \mathcal{C}}{\partial b_j^l}$: pour calculer la mise à jour du biais du neurone $j \leq N_l$ à la couche $l \leq L$

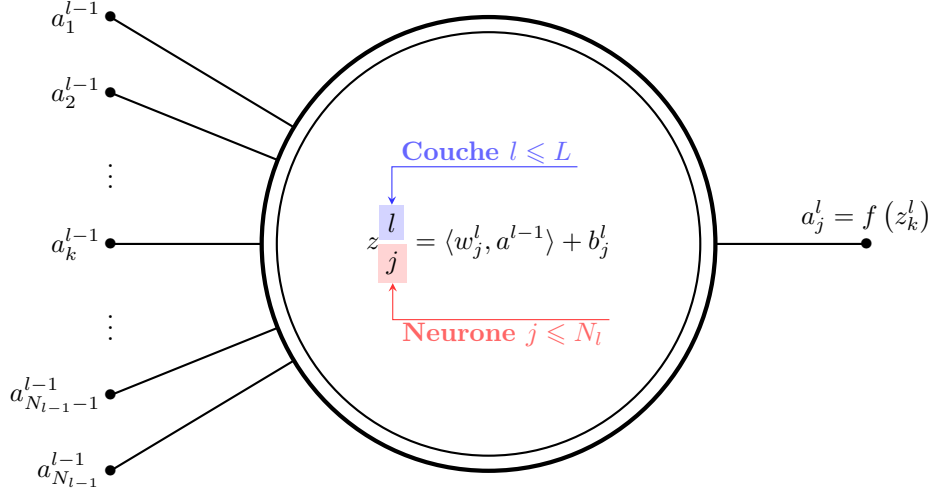


FIGURE 1.1 – Neurone j de la couche l paramétré par les poids $w_i^j \in \mathbb{R}^{N_{l-1}}$ et le biais $b_j^l \in \mathbb{R}$

Pour le faire, résolvons l'exercice suivant.

Exercice 1.1 (Back propagation). *On reprend les notations définies précédemment.*

1. Soit $l \leq L$, $j \leq N_l$ et $k \leq N_{l-1}$. Montrer que :

$$\begin{aligned} \frac{\partial C}{\partial b_j^l} &= \frac{\partial C}{\partial z_j^l} \\ \frac{\partial C}{\partial w_{j,k}^l} &= a_k^{l-1} \frac{\partial C}{\partial z_j^l} \end{aligned}$$

2. On note $\delta_j^l = \frac{\partial C}{\partial z_j^l}$ qui représente l'erreur du neurone $j \leq N_l$ à la couche $l \leq L$.
Montrer que pour la dernière couche L , on a :

$$\forall j \leq N_L, \quad \delta_j^L = f'(z_j^L) \frac{\partial C}{\partial a_j^L}$$

3. On considère à présent la couche $l < L$. Montrer que l'on peut exprimer δ_j^l comme :

$$\delta_j^l = \sum_{k=1}^{n_l} w_{j,k}^{l+1} \delta_k^{l+1} f'(z_j^l)$$

Solution. Nous conservons les notations précédentes et ferons des commentaires supplémentaires à chaque questions pour donner plus de liens à la résolution de l'exercice.

1. On considère le neurone $j \leq N_l$ à la couche $l \leq L$. Alors par le théorème de dérivation des

fonctions composées, on a :

$$\begin{aligned}\frac{\partial C}{\partial b_j^l} &= \frac{\partial C}{\partial b_j^l} \frac{\partial b_j^l}{\partial z_j^l} \quad \text{mais} \quad \frac{\partial b_j^l}{\partial z_j^l} = 1 \\ &= \frac{\partial C}{\partial z_j^l}\end{aligned}$$

On considère maintenant le k -ième poids du neurones. On a à nouveau avec le théorème de dérivation des fonctions composées :

$$\begin{aligned}\frac{\partial C}{\partial w_{j,k}^l} &= \frac{\partial C}{\partial z_j^l} \frac{\partial z_j^l}{\partial w_{j,k}^l} \quad \text{mais} \quad \frac{\partial z_j^l}{\partial w_{j,k}^l} = a_k^{l-1} \\ \frac{\partial C}{\partial w_{j,k}^l} &= a_k^{l-1} \frac{\partial C}{\partial z_j^l}\end{aligned}$$

On observe que pour le poids et le biais on a la même quantité $\frac{\partial C}{\partial z_j^l}$ qui apparait. Elle représente l'erreur du neurone $j \leq N_l$ à la couche $l \leq L$. L'exercice nous propose de la noter δ_j^l pour la suite.

2. On considère à présent spécifiquement la dernière couche. Soit $j \leq N_L$ un neurone de cette couche. Encore avec le théorème de dérivation des fonctions composées :

$$\begin{aligned}\delta_j^L &= \frac{\partial C}{\partial a_j^L} \frac{a_j^L}{\partial z_j^L} \quad \text{mais} \quad a_j^L = f(z_j^L) \\ &= f'(z_j^L) \frac{\partial C}{\partial a_j^L}\end{aligned}$$

Calculer la valeur de $\frac{\partial C}{\partial a_j^L}$ est aisé puisqu'il suffit de comparer la valeur prédite par le réseau à la valeur attendue. En combinant ce résultat à la question précédente, nous savons calculer la mise à jour des poids et des biais de la dernière couche.

Pour être capable de calculer la mise à jour des poids et des biais des couches précédentes, il faudrait être capable d'exprimer les δ_j^l en fonction de δ_j^L : probablement de proche en proche.

3. On considère à présent le j -ième neurone de la couche $l < L$. Par définition des notations, on a :

$$\begin{aligned}z_j^l &= \sum_{k=1}^{N_l} w_{j,k}^l a_k^{l-1} + b_j^l \\ z_k^{l+1} &= \sum_{i=1}^{N_{l+1}} w_{k,i}^{l+1} f(z_i^l) + b_k^{l+1}\end{aligned}$$

Ainsi, en appliquant plusieurs fois le théorème de dérivation des fonctions composées :

$$\begin{aligned}
\delta_j^l &= \frac{\partial C}{\partial z_j^l} \\
&= \sum_{k=1}^{N_l} \frac{\partial C}{\partial z_k^{l+1}} \frac{\partial z_k^{l+1}}{\partial z_j^l} \quad \text{et} \quad \frac{\partial z_k^{l+1}}{\partial z_j^l} = w_{k,j}^{l+1} f'(z_j^l) \\
&= \sum_{k=1}^{N_l} w_{k,j}^{l+1} f'(z_j^l) \frac{\partial C}{\partial z_k^{l+1}} \\
&= \sum_{k=1}^{N_l} w_{k,j}^{l+1} \delta_k^{l+1} f'(z_j^l) \quad \text{par définition des } \delta_j^l
\end{aligned}$$

Nous calculons donc l'actualisation des poids et des biais de proche en proche : d'où le terme back propagation !

□

En décrivant l'algorithme de back propagation, on comprend que chaque poids est calculés de proche en proche, en partant de la couche la plus proche à la plus profonde (celle juste après l'input). Puisque les calculs se font de proche en proche, alors il y a une possibilité qu'une suite géométrique apparaisse et fasse disparaître la magnitude du signal.

En effet, le signal de correction du j-ième neurones de la couche $l \leq L$ est :

$$\delta_j^l = \sum_{k=1}^{n_l} w_{j,k}^{l+1} \delta_k^{l+1} f'(z_j^l)$$

Dépendant clairement du précédent, on peut l'écrire sous la forme :

$$\begin{aligned}
\delta_j^l &= \sum_{k=1}^{n_l} w_{j,k}^{l+1} \delta_k^{l+1} f'(z_j^l) \\
&= \sum_{k=1}^{n_l} w_{j,k}^{l+1} f'(z_j^l) \times \left(\sum_{k=1}^{n_{l+1}} w_{j,k}^{l+2} \delta_k^{l+2} f'(z_j^{l+1}) \right) \\
&= \sum_{k=1}^{n_l} w_{j,k}^{l+1} f'(z_j^l) \times \left(\sum_{k=1}^{n_{l+1}} w_{j,k}^{l+2} f'(z_j^{l+1}) \times \left(\dots \left(\sum_{k=1}^{n_{L-1}} w_{j,k}^L f'(z_j^L) \delta_j^L \right) \dots \right) \right)
\end{aligned}$$

Si l'on note $\gamma = \max_{x \in \mathbb{R}} f'(x)$ alors :

$$|\delta_j^l| \leq \gamma^{L-l+1} \left| \sum_{k=1}^{n_l} w_{j,k}^{l+1} \times \left(\sum_{k=1}^{n_{l+1}} w_{j,k}^{l+2} \times \left(\dots \left(\sum_{k=1}^{n_{L-1}} w_{j,k}^L \delta_j^L \right) \dots \right) \right) \right|$$

Or, dans les premiers réseaux de neurones, la fonction sigmoid était utilisée. Dans ce cas : $\gamma = \frac{1}{4}$, d'où une probabilité forte que la force du signal des dernières couches soit faible et donc que le réseau n'apprennent pas correctement. Notons que la fonction tanh était également utilisée et que l'argument que nous avons construit ici ne suffit pas à conclure car $\gamma = 1$. Cependant, la même logique explique le problème : une composition répétée de nombreuse fois d'une quantité

de valeur inférieure à 1 conduit à une disparition des gradients. Le phénomène inverse s'explique de la même manière.

Ces phénomènes ont fait perdre un temps considérable aux chercheurs : l'entraînement d'un réseau de neurones profond est lent, encore plus avec une puissance machine réduite. De plus, si l'on arrive pas à montrer que nos recherches peuvent fonctionner mieux que l'existant, on n'obtient pas de fonds supplémentaires pour continuer de mener les recherches. Face à ces deux problèmes, on pourrait s'en remettre à la chance.

Je ne crois pas qu'il y ait de bonne ou de mauvaise situation.

— Édouard Baer, jouant Otis dans Astérix et Obélix Mission Cléopâtre (2002)

En revanche, il existe un *bon* et un *mauvais* hasard. Voyons pourquoi.

1.1.2 Réponses Glorot et He

Une des manières de lutter contre ces deux phénomènes qui interviennent au cours de l'apprentissage du réseau est d'adopter une bonne initialisation des poids. Cela peut paraître surprenant : comment une initialisation peut-elle améliorer durablement un réseau de neurones ? Pour le comprendre, on considère un réseau de neurones à $L \in \mathbb{N}^*$ couches. On reprend la majorité des notations précédentes, mais on note z_k^l le résultat du neurone k à la couche l défini par :

$$\forall l \leq L, \forall k \in N_l, z_k^l = \langle w_k^l, x^l \rangle$$

Nombre de couche du réseau
Vecteur d'entrée de la couche l

Nombre de neurones de la couche l
Matrice des poids du neurone k à la couche l

Nous avons omis ici le biais car notre analyse se situe à l'initialisation du réseau : il n'y a pas de raison d'introduire déjà un biais. Nous faisons les hypothèses supplémentaires :

- **Indépendance et distribution** : $(w^l)_{l \leq L}$, $(x^l)_{l \leq L}$ et $(z^l)_{l \leq L}$ sont indépendants et identiquement distribués (respectivement)
- **Indépendance** : $\forall l \leq L, w^l \perp x^l$
- **Poids centrés** : $\mathbb{E}[w^l] = 0$ et la distribution est symétrique autour de 0
- **Résultats centrés** : $\mathbb{E}[z^l] = 0$ et la distribution est symétrique autour de 0 : il n'y a pas de raison que les signaux qui se propagent soient biaisés

Puisque l'on souhaite une bonne propagation des signaux dans les deux sens afin d'éviter une explosion ou une disparition des gradients, ce que l'on souhaite est que :

$$\forall l \leq L, \mathbb{V}[z^l] = \mathbb{V}[z^1] \quad (1.1)$$

Si c'est le cas, alors nous aurons lutté efficacement contre ces deux phénomènes. Calculons pour $l \leq L$ et $k \in N_l$:

$$\begin{aligned}
\mathbb{V}[z_k^l] &= \mathbb{V}[\langle w_k^l, x^l \rangle] \\
&= \sum_{j=1}^{n_{\text{input}}} \mathbb{V}[w_{k,j}^l x_j^l] \text{ par indépendance} \\
&= n_{\text{input}} \mathbb{V}[w_{k,j}^l x_j^l] \text{ car identiquement distribué, pour } j \leq n_{\text{input}} \\
&= n_{\text{input}} \left(\mathbb{V}[w_{k,j}^l] \mathbb{V}[x_j^l] + \mathbb{E}[w_{k,j}^l]^2 \mathbb{V}[x_j^l] + \mathbb{E}[x_j^l]^2 \mathbb{V}[w_{k,j}^l] \right) \\
&= n_{\text{input}} \mathbb{V}[w_{k,j}^l] \left(\mathbb{V}[x_j^l] + \mathbb{E}[x_j^l]^2 \right) \text{ car } \mathbb{E}[w_{k,j}^l] = 0 \\
&= n_{\text{input}} \mathbb{V}[w_{k,j}^l] \mathbb{E}[(x_j^l)^2]
\end{aligned}$$

Rappelons que $x^l = f(z^{l-1})$ avec f la fonction d'activation entre la couche $l-1$ et la couche l . Si f est impaire, alors $\mathbb{E}[x^l] = \mathbb{E}[f(z^{l-1})] = 0$ donc $\mathbb{E}[(x_j^l)^2] = \mathbb{V}[z_k^{l-1}]$. Ainsi,

$$\begin{aligned}
\forall l \leq L, \forall k \in N_l, \quad \mathbb{V}[z_k^l] &= n_{\text{input}} \mathbb{V}[w_{k,j}^l] \mathbb{V}[z_k^{l-1}] \\
&= \prod_{i=1}^{l-1} (n_{\text{input}} \mathbb{V}[w_{k,j}^i]) \mathbb{V}[z_k^1]
\end{aligned}$$

Une condition suffisante pour que l'on respecte l'équation (1.1) est que :

$$\mathbb{V}[w_{k,j}^l] = \frac{1}{n_{\text{input}}}$$

La condition précédente a été obtenue en considérant une passe forward avec une fonction d'activation impaire. On peut se convaincre que la passe backward abouti à la même forme d'équation sauf que cette fois n_{input} représente le nombre de neurones de la couche inférieure. Nous venons de reproduire le raisonnement de l'article *Understanding the difficulty of training deep feedforward neural networks* [Glorot and Bengio, 2010]. Pour résoudre le problème final, choisir la condition qui conviendra pour les deux passes, les auteurs proposent de prendre la moyenne entre les deux conditions. Plus tard, nous ferons références au nombre de neurones de la couche précédente (respectivement suivante) par *fan in* (respectivement *fan out*).

Notons que nous n'avons supposé aucune loi pour l'initialisation des poids à part des conditions d'indépendance, distribution et que la loi doit être symétrique autour de 0 et d'espérance 0.

Exercice 1.2 (Forme uniforme et normale). *On considère un réseau de neurones avec une fonction d'activation impaire. Donner l'expression de la variance en respectant la condition précédente, en supposant que les poids suivent une loi normale puis une loi uniforme.*

Solution. Si les poids sont initialisé selon une loi normale $\mathcal{N}(0, \sigma^2)$: $\sigma = \sqrt{\frac{2}{\text{fan in} + \text{fan out}}}$

On rappelle que pour $U \sim \mathcal{U}([-a, a])$ on a $\mathbb{V}[U] = \frac{a^2}{3}$. Ainsi, si les poids sont initialisé selon une loi uniforme $\mathcal{U}([-a, a])$: $a = \sqrt{\frac{6}{\text{fan in} + \text{fan out}}}$ □

Le résultat précédent s'applique pour des réseaux qui utiliserons des fonctions d'activation impaire comme la fonction sigmoid ou tangente hyperbolique comme présenté à la figure (1.2a).

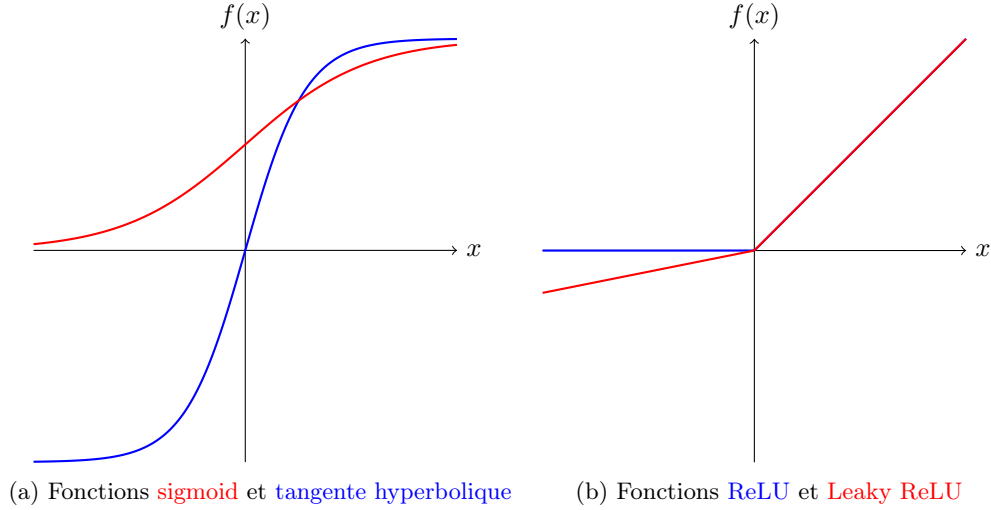


FIGURE 1.2 – Différentes fonctions d'activations non-linéaire

Cependant, la majorité des réseaux moderne utilisent la fonction d'activation ReLU qui n'est pas impaire, comme visible dans la figure (1.2b). L'article *Delving deep into rectifiers : Surpassing human-level performance on ImageNet classification* [He et al., 2015] répond à cette question.

Puisque cette fois f n'est pas impaire mais égale à $f(x) = \max\{0, x\}$, en simplifiant les notations, on a :

$$\begin{aligned}
 \mathbb{E}[x^2] &= \mathbb{E}[f(z)^2] \\
 &= \int_{-\infty}^{+\infty} f(t)^2 p(t) dt \text{ avec } p \text{ la densité de probabilité associé à } z \\
 &= \int_0^{+\infty} t^2 p(t) dt \\
 &= \frac{1}{2} \mathbb{V}[z]
 \end{aligned}$$

Nous avons cette fois :

$$\begin{aligned}
 \forall l \leq L, \forall k \in N_l, \mathbb{V}[z_k^l] &= \frac{n_{\text{input}}}{2} \mathbb{V}[w_{k,j}^l] \mathbb{V}[z_k^{l-1}] \\
 &= \prod_{i=1}^{l-1} \left(\frac{n_{\text{input}}}{2} \mathbb{V}[w_{k,j}^i] \right) \mathbb{V}[z_k^1]
 \end{aligned}$$

A nouveau, une condition suffisante pour que la relation (1.1) soit vérifiée pour la passe forward est :

$$\mathbb{V}[w_{k,j}^l] = \frac{2}{n_{\text{input}}}$$

C'est très proche de la condition quand l'on considère une fonction d'activation impaire, mais cela résulte en un écart significatif lors de l'entraînement. Contrairement à la résolution proposée par Xavier Glorot et Yoshua Bengio, la valeur de la variance pour les poids est obtenu en sélectionnant le *fan in* ou le *fan out*.

Exercice 1.3 (Forme uniforme et normale pour ReLU). *On considère un réseau de neurones avec une fonction d'activation ReLU. Donner l'expression de la variance en respectant la condition précédente, en supposant que les poids suivent une loi normale puis une loi uniforme.*

Solution. Si les poids sont initialisé selon une loi normale $\mathcal{N}(0, \sigma^2)$: $\sigma = \sqrt{\frac{2}{\text{fan mode}}}$ avec *fan mode* valant soit *fan in* si l'on privilégie une conservation maximale des gradients en passe forward, *fan out* pour la passe backward.

Si les poids sont initialisé selon une loi uniforme $\mathcal{U}([-a, a])$: $a = \sqrt{\frac{6}{\text{fan mode}}}$ avec la même définition de *fan mode*. □

Nous savons donc maintenant comment initialiser un réseau de neurones en fonction du choix de la fonction d'activation. Il nous reste à connaître une autre possibilité que l'on rencontre en pratique : le Leaky ReLU. On peut la visualiser avec la figure (1.2b).

Exercice 1.4 (Forme uniforme et normale pour le Leaky ReLU). *On considère un réseau de neurones avec une fonction d'activation Leaky ReLU de paramètre $\alpha \in [0, 1]$. Donner l'expression de la variance en respectant la condition précédente, en supposant que les poids suivent une loi normale puis une loi uniforme.*

Solution. En reprenant la démarche pour la fonction ReLU ainsi que les notations, nous avons :

$$\begin{aligned} \mathbb{E}[x^2] &= \mathbb{E}[f(z)^2] \\ &= \int_{-\infty}^{+\infty} \max\{\alpha t, t\}^2 p(t) dt \\ &= \int_{-\infty}^0 \alpha^2 t^2 p(t) dt + \int_0^{+\infty} t^2 p(t) dt \\ &= \frac{\alpha^2 + 1}{2} \mathbb{V}[z] \end{aligned}$$

Ainsi, si les poids sont initialisé selon une loi normale $\mathcal{N}(0, \sigma^2)$: $\sigma = \sqrt{\frac{2}{(1 + \alpha^2) \times \text{fan mode}}}$ avec *fan mode* valant soit *fan in* si l'on privilégie une conservation maximale des gradients en passe forward, *fan out* pour la passe backward.

Si les poids sont initialisé selon une loi uniforme $\mathcal{U}([-a, a])$: $a = \sqrt{\frac{6}{(1 + \alpha^2) \times \text{fan mode}}}$ avec la même définition de *fan mode*. \square

Notons une certaine similarité entre le cas Leaky ReLU et ReLU : c'est capturé par la documentation PyTorch sous le concepts de *gains*. Pour chaque fonction d'activation est associé une valeur que l'on multiplie à l'expression de la variance pour la stratégie d'initialisation des poids que l'on choisit.

A travers ces problèmes et ses résolutions, nous comprenons l'importance d'une bonne distribution pour l'initialisation des poids et donc d'un bon *départ* dans un réseau de neurone. Essayons de mesurer l'importance du début d'entraînement d'un réseau de neurone.

1.2 Un mauvais départ n'est pas rattrapé

Intuitivement, puisque un réseau peut être amené à s'entraîner longtemps, si l'on lui permet de propager correctement les gradients, il semblerait possible que l'algorithme puisse performer *quoi qu'il arrive*. C'est ce qu'explore l'article *Critical learning periods in deep neural networks* [Achille et al., 2017].

Dans les années 1950, une des premières pistes suivi dans le domaine de l'intelligence artificielle était de s'inspirer du vivant. Plus particulièrement : inspirons-nous du fonctionnement du cerveau humain pour construire une machine capable elle aussi d'apprendre. Pour les raisons techniques et théoriques, comme évoqué plus tôt, les réseaux de neurones sont abandonnés pendant plusieurs décennies avant de devenir incontournable à partir des années 2010.

Le cerveau dans le monde animal exhibe des périodes d'apprentissage critique, notamment à la naissance. Plusieurs études se sont intéressées au sujet dont la mesure l'acuité visuelle des chatons en 1963. Cette étude est adapté au réseau de neurones dans cet article. Que se passerait-il si les premières époques d'apprentissage du réseau de neurones étaient réalisées avec des images floutées, tordues ou incomplète ? Il est montré qu'un exposition trop longue à ces images affecte à long terme les chats, qu'en est-il des réseaux de neurones ?

Nous savons que les enfants apprennent vite, et les adultes moins, grâce au concept de la plasticité cérébrale. Peut-on observer la même chose dans un réseau de neurones ? On peut reprendre notre intuition du début : puisque l'algorithme peut s'entraîner longtemps, il n'y a pas de raison qu'il ne puisse pas exhiber de plasticité, de plus sans être affecté par *l'âge*.

L'ensemble des résultats présentés dans cette section sont issu de l'article [Achille et al., 2017] mais certains ont été reproduit à la main pour ces notes de cours pour des raisons de lisibilité.

La première expérience réalisée consiste à entraîner un réseau de neurones avec des images floues pendant N époques puis l'entraîner avec les images originelle pendant 160 époques supplémentaires comme décrit par la figure (1.3a). Pour comparer, on entraîne un réseau de neurones sans déficit et on montre les performances pour chaque époques jusqu'à 140.

Plus le déficit sera supprimé tard, plus le réseau de neurones se sera entraîné longtemps : on peut imaginer que la performance finale sera donc meilleure. Notamment le réseau qui aura été entraîné avec le déficit pendant 140 époques : il aura 160 époques d'entraînement *en clair* de plus.

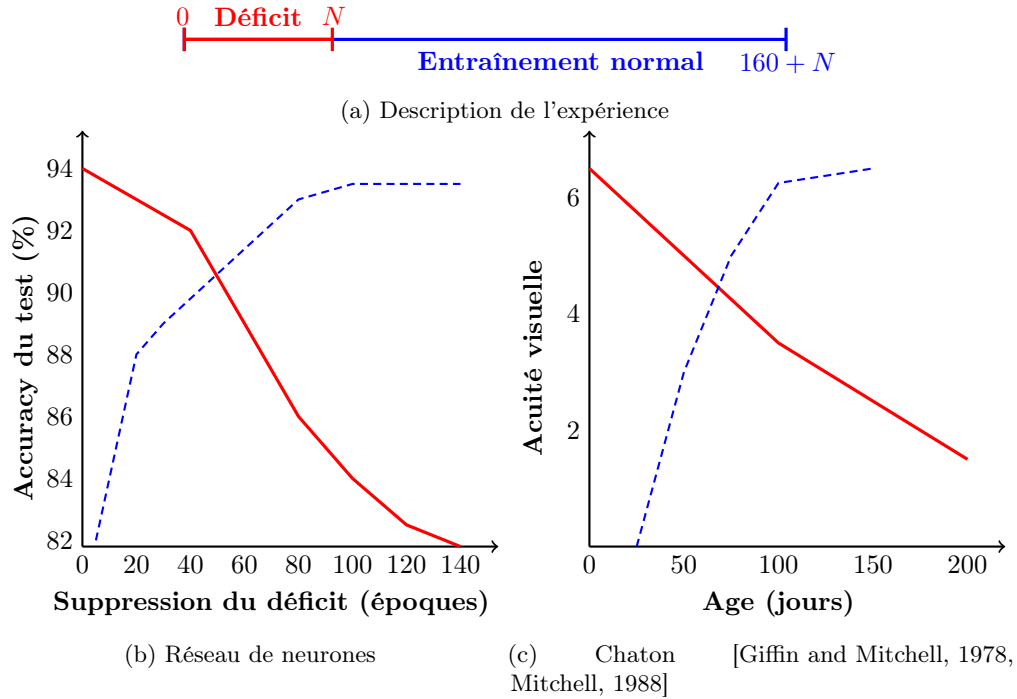


FIGURE 1.3 – Performance/acuité selon la durée du **déficit** par rapport à un développement **normal**

On observe l'inverse ! Si le déficit est supprimé *tôt* alors on peut obtenir une performance proche d'un entraînement sans déficit. Cependant, si le déficit dure, le temps d'entraînement supplémentaire ne semble pas permettre de rattraper ce retard. Ce comportement est identique à celui décrit dans une expérience visuelle réalisée sur des chatons. On comprend donc que la période de départ est importante, portant l'emphase sur les premières époques.

Une deuxième expérience est réalisée en sélectionnant cette fois une fenêtre de taille fixe pour l'entraînement avec du déficit. Le nombre d'époque total est également fixe, c'est le début du déficit qui change cette fois, comme décrit dans (??). Avec l'expérience précédente, on conjecture que lorsque le déficit portera sur les premières époques il y aura un grand écart par rapport à un entraînement sain sur la performance sur le jeu de test. Mais qu'en est-il lorsque le déficit est présent en *fin* d'entraînement ?

Nous retrouvons bien l'importance du début d'entraînement avec une baisse jusqu'à 5% de la performance sur le jeu de test dans le déficit est appliqué sur les premières époques. Notons tout de même que lorsque le déficit est appliqué sur les toutes premières époques, la sensibilité est plus faible que lorsqu'il s'agit des époques 20/40. On observe également que lorsque le déficit est porté sur les dernières époques, l'impact est beaucoup plus faible, mais présent.

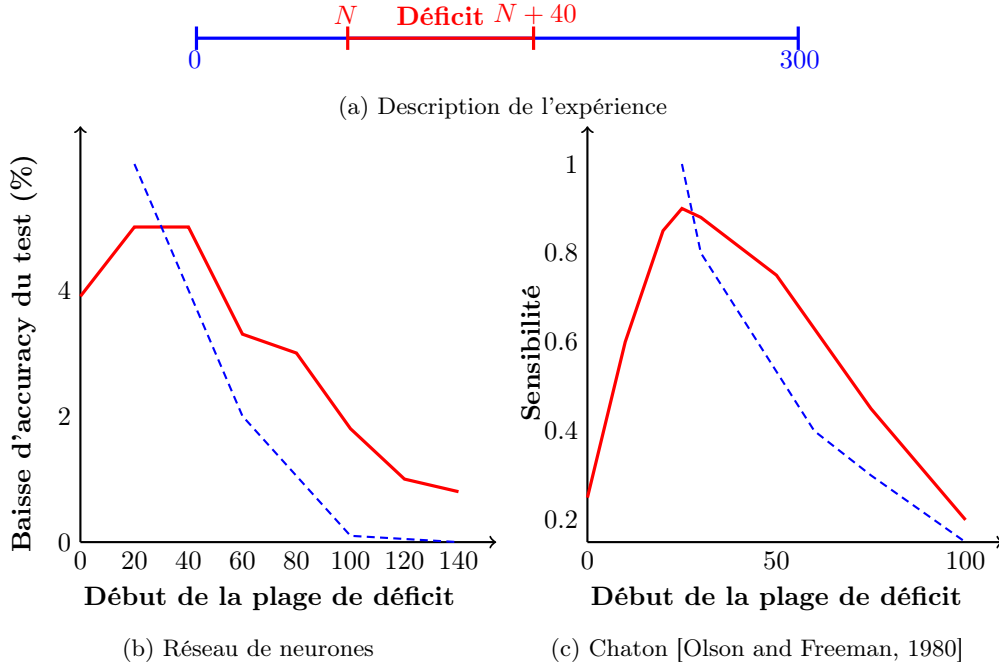


FIGURE 1.4 – Sensibilité selon le début du **déficit** par rapport à un développement **normal**

Les conclusions sont identiques pour les chatons à nouveau. Ce qui montre à nouveau l'importance du début d'entraînement : un mauvais départ n'est pas rattrapable, même avec des époques supplémentaires. Qu'est-ce qui peut expliquer ce phénomène ? Pour répondre à cette question, l'article propose de mesurer l'apprentissage du réseau de neurones en étudiant la force de connexion entre les neurones.

1.2.1 Mesurer la force de connexion

On représente le résultat d'un réseau de neurones comme $\mathbb{P}_w(y | x)$ où w représente l'ensemble des poids qui paramètre le réseau de neurones. Si l'on note $w' = w + \delta w$ où δ représente une petite perturbation, alors on peut mesurer l'impact des poids modifiés par la perturbation dans le réseau de neurones comme *l'écart* entre $\mathbb{P}_w(y | x)$ et $\mathbb{P}_{w'}(y | x)$. Mais comment mesurer *l'écart* entre deux distributions ?

Cette question rejoint la théorie de l'information et plus particulièrement la divergence de Kullback-Leibler :

$$D_{\text{KL}} \left(\overbrace{P}^{\text{Distribution de densité } p} \parallel \underbrace{Q}_{\text{Distribution de densité } q} \right) = \int_{-\infty}^{+\infty} p(x) \ln \left(\frac{p(x)}{q(x)} \right) dx \quad (\text{Kullack-Leibler})$$

Cette divergence mesure ce que l'on peut espérer perdre en remplaçant la **vrai distribution** par une **distribution estimée**. Comme son nom l'indique, la divergence de Kullback-Leibler est une

divergence : donc elle n'est pas symétrique. Puisque c'est une divergence, pour toutes distributions P et Q , la divergence de Kullback-Leibler doit être positive : ce n'est pas évident tel quel. Le prouver est un exercice classique d'analyse présenté dans l'exercice (1.5).

Exercice 1.5 (Inégalité de Gibbs discrète). On considère $P = \{p_1, \dots, p_n\}$ et $Q = \{q_1, \dots, q_n\}$ deux distributions discrètes.

1. Montrer que : $\forall x > 0, \ln x \leq x - 1$ et que l'égalité est vérifiée si, et seulement si, $x=1$.
2. Soit $I = \{i \mid p_i \neq 0\}$. En utilisant le résultat précédent, montrer que :

$$-\sum_{i \in I} p_i \ln \left(\frac{q_i}{p_i} \right) \geq 0$$

3. On rappelle que $\lim_{x \rightarrow 0} x \ln(x) = 0$. Montrer que :

$$-\sum_{i=1}^n p_i \ln(q_i) \geq -\sum_{i=1}^n p_i \ln(p_i)$$

4. Sous quelle condition y a-t-il égalité ?
5. Montrer donc que $D_{KL}(P \parallel Q) \geq 0$

Puisque nous nous intéressons ici à des réseaux de neurones avec énormément de paramètre, l'estimation de la divergence de Kullback-Leibler va être difficile en l'état. Notons que puisque $\lim_{x \rightarrow 0} x \ln(x) = 0$, nous n'avons pas de problème pour des événements impossible pour P mais possible pour Q . Cependant l'inverse n'est pas vrai : nous aurons donc également un problème de définition puisque la divergence vaut, par convention, $+\infty$. En introduisant d'infimes quantités ε à chaque événements (dont ceux impossible pour l'une ou l'autre des distributions) nous sommes capables de calculer la divergence au prix d'un léger écart. Cela ne résout pas notre problème de capacité à estimer la mesure pour un grand réseau de neurones. Il nous faudrait une expression plus simple pour la calculer.

On peut montrer, mais ce n'est pas l'objet de ce cours, que le développement de Taylor à l'ordre deux de la divergence de Kullback-Leibler s'écrit sous la forme suivante :

$$D_{KL}(\mathbb{P}_w(y \mid x) \parallel \mathbb{P}_{w'}(y \mid x)) = \delta w F(\delta w) + o(\delta w^2)$$

Avec F la métrique d'information de Fisher. On peut la comprendre dans ce contexte comme une métrique **locale** de la perturbation de la valeur d'un poids. Concrètement, si un poids a un information de Fisher faible c'est que le poids n'est pas très important pour la sortie du réseau de neurone.

À nouveau la matrice d'information de Fisher est trop grande pour être calculé systématiquement. Ce que l'article propose finalement, est de ne considérer que la force de connexion entre chaque couche : ne considérer que la trace de la matrice d'information de Fisher.

Nous sommes donc maintenant capables d'approcher, localement, la force de la connexion entre les couches d'un réseau de neurone. Notons que puisque la mesure est approchée et largement imparfaite pour des raisons de puissance de calcul, ce sont des réseaux de neurones résiduels qui

ont été utilisés. Ils sont connus pour avoir des propriétés plus lisses que les autres types de réseaux de neurones :

- [Balduzzi et al., 2017] montre que les ResNet ont des gradients plus lisses
- [Li et al., 2018] montre que les ResNet ont une surface de loss plus lisse également
- [Zhang et al., 2019] montre que la particularité des ResNet, les *skip connections*, peuvent être remplacées par une bonne initialisation et tout de même obtenir des résultats compétitifs

Voyons comment cette notion peut nous aider à comprendre l'importance d'un bon départ. Puisque l'on peut interpréter la trace de la matrice de Fisher comme la quantité d'information *apprise* par le modèle, on devrait voir une augmentation continue au fur et à mesure de l'entraînement. Mais peut-être moins quand il y a un déficit.

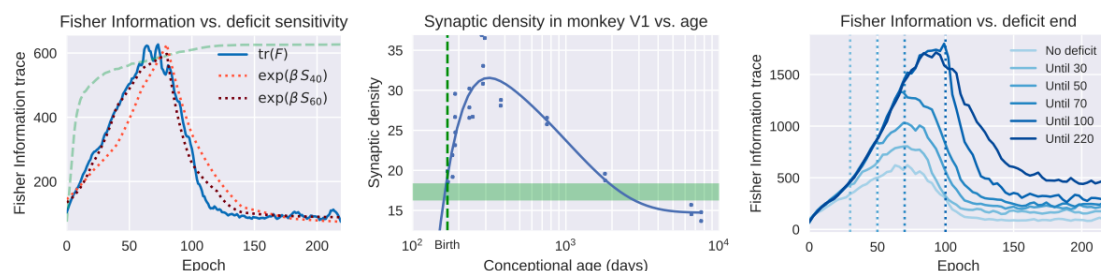


FIGURE 1.5 – Evolutions de la trace de la matrice d'information de Fisher (source : [Achille et al., 2017])

Dans le graphique de **gauche** est tracé en pointillé vert la performance du réseau sur le dataset de test. La courbe bleue correspond à la trace de la matrice d'information de Fisher. On observe deux phases : la trace augmente jusqu'à ce que la performance sur le jeu de test stagne puis s'écroule. On remarque que ce phénomène est commun peu importe qu'il y ait un déficit ou non (graphique *droit*). Cependant, on observe que plus le déficit est maintenu longtemps plus la quantité d'informations stockée est importante et que la phase de *stockage* d'information dure plus longtemps.

On a donc deux phases d'entraînement : une première où le réseau *apprend* et améliore ses performances sur le jeu de test, et une seconde de compression d'informations où les connexions *inutiles* sont supprimées. Remarquons que lorsque les images sont trop difficiles pour être classifiées à cause du déficit introduit, le réseau est forcé de mémoriser des observations.

Voyons comment les différentes couches participent à ce phénomène.

Pour un entraînement normal, le réseau semble *apprendre* avec les couches du milieu plutôt. Ce n'est pas le cas dès que l'on introduit un déficit : les dernières couches sont les plus mises à contributions et stockent énormément d'informations. Ce phénomène semble proportionnel au temps de maintien du déficit.

La puissance de représentation du réseau de neurones n'est pas utilisée puisque les features générées naturellement par les couches intermédiaires n'ont pas lieu : le réseau est forcé de mémoriser les images. Une manière de s'assurer qu'un réseau ne *mémorise* pas les input est de le régulariser. Observons l'impact du *weight decay*¹.

Conformément à notre intuition, la régularisation rend la période critique du début d'entraînement un peu moins sensible au déficit. Cependant, si le déficit est maintenu trop longtemps on

1. La notion est présentée plus en détail dans la section (3.2.1)

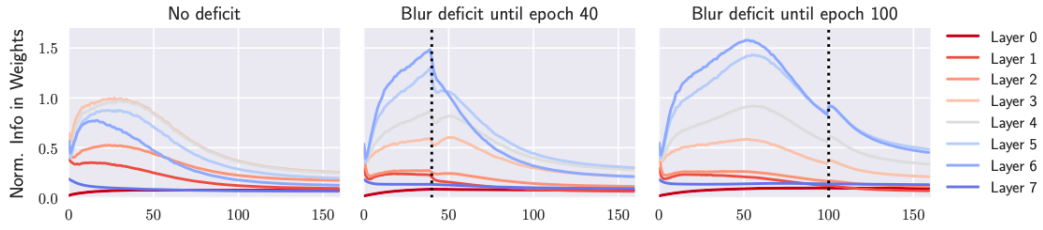


FIGURE 1.6 – Evolutions de la trace de la matrice d’information de Fisher par couche (source : [Achille et al., 2017])

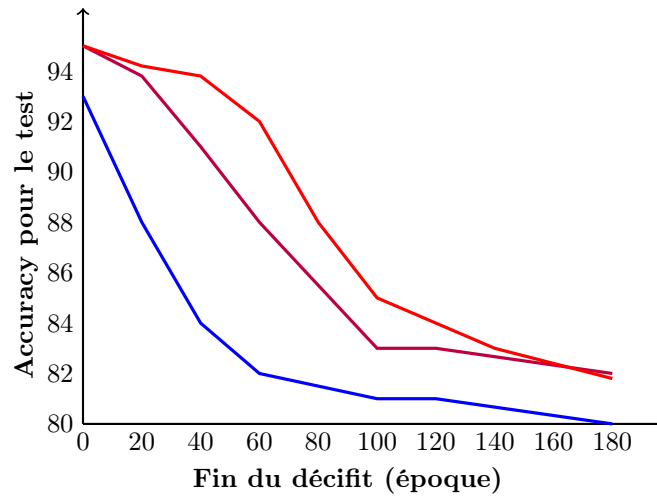


FIGURE 1.7 – Performance sur le dataset de test en fonction du nombre d’époque de maintien du déficit pour un weight decay de $\lambda = 0$, $\lambda = 5 \cdot 10^{-4}$ et $\lambda = 10 \cdot 10^{-4}$

observe tout de même une baisse nette de performance.

Précisons que les graphiques et résultats présentés sont ceux pour des images floutées. Pour des *dégradations* moins forte de l’image comme des rotations ou des symétries, des comportements similaires ne sont pas rapportés par l’article.

Ainsi, de manière surprenante nous comprenons qu’un mauvais départ n’est pas rattrapé et qu’en **contraignant** le réseau de neurones nous pouvons ralentir ce phénomène si le départ n’est pas très bon.

1.3 Bien se lancer pour mieux généraliser

Cette idée rejoint celle de l’article *Dropout Reduces Underfitting* [Liu et al., 2023] en 2023. Le dropout [Hinton et al., 2012b] et [Srivastava et al., 2014] est une méthode de régularisation est rapidement devenu une brique essentielle des réseaux de neurones. Un des facteurs du succès immédiat de la méthode est son utilisation dans le modèle AlexNet [Krizhevsky et al., 2012], vainqueur de la compétition ImageNet en 2012.

Lors de l'entraînement, le dropout va *supprimer* de la chaîne d'entraînement des neurones avec une proportion $1 - p$ à chaque couche. Comme illustré avec la figure 1.8, le dropout s'applique à chaque couches, dont l'input, sauf l'output. On obtient donc à chaque passe forward un sous-ensemble du réseau initial.

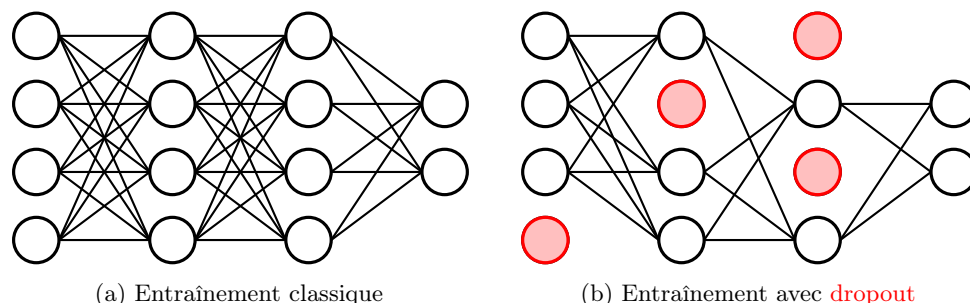


FIGURE 1.8 – Réseau de neurones avec et sans **dropout**

Lors de l'exploitation du réseau de neurones en revanche nous ne pouvons pas avoir du hasard, il faut un comportement déterministe. Remarquons que ce que nous venons de décrire ressemble à du bagging ! Nous avons entraîné tout un ensemble de réseaux de neurones, tous différents (ou presque). Cependant, ils ne sont pas indépendants car ils partagent des poids. Nous devons trouver une manière d'exploiter la force des méthodes ensemblistes.

Classiquement, la prédiction finale d'un ensemble est construite en faisant la moyenne arithmétique des prédictions. Ici, cela coûterait trop en mémoire de conserver l'ensemble des réseaux, donc nous allons approcher ce comportement en faisant la moyenne géométrique. Pour cela, chaque poids est multiplié par la proportion p de dropout.

Le dropout permet d'éviter une sur-spécification de certains neurones, forçant le réseau à être plus redondant et robuste, de la même manière que chaque membre d'une équipe en entreprise doit être capable de se remplacer en cas de congé.

L'amélioration quasi systématique des performances d'un réseau de neurones a rendu cette méthode très populaire mais pas incontournable. Dans l'entraînement des LLM par exemple, aucun dropout n'est ajouté. Cependant, il est fortement conseillé lors d'un fine-tuning. Son utilisation est principalement pour éviter le sur-apprentissage. Il est donc surprenant que le titre du papier [Liu et al., 2023] soit *Dropout Reduces Underfitting*. L'article propose de n'appliquer le dropout que lors du début de l'entraînement et le nomme *early dropout*. En comparant à une baseline, il est noté que la norme des gradients lors de la phase d'early dropout est plus petite que ceux de la baseline. Ainsi, il est probable que la distance parcourue dans l'espace des paramètres soient également plus faible. Cependant, c'est l'inverse qui se produit !

Nous en déduisons que le dropout force le réseau à prendre une direction plus consistante que la baseline originale. C'est confirmé lorsqu'on calcule la variance dans la direction entre le réseau avec early dropout et la baseline. Finalement, contraindre le réseau de neurones permet de le placer dans de meilleures conditions pour la suite de l'entraînement. On observe en effet que la réduction de loss est inférieure tant que le dropout est effectif, et cela s'inverse lorsque l'early dropout s'arrête.

On peut se poser la question de la régularisation en fin d'entraînement. Est donc introduit le *late dropout*, mais il n'est pas très efficace pour régulariser le réseau de neurones. Cependant la méthode *stochastic depth* l'est.

Cette méthode est introduite en 2016 [Huang et al., 2016] et ne s'applique qu'au réseau de neurones résiduels (ResNet) introduite un an plus tôt² [He et al., 2016]. Les réseaux de neurones résiduels sont constitués de nombreux ResBlock dont la constitution est illustrée avec la figure (1.9).

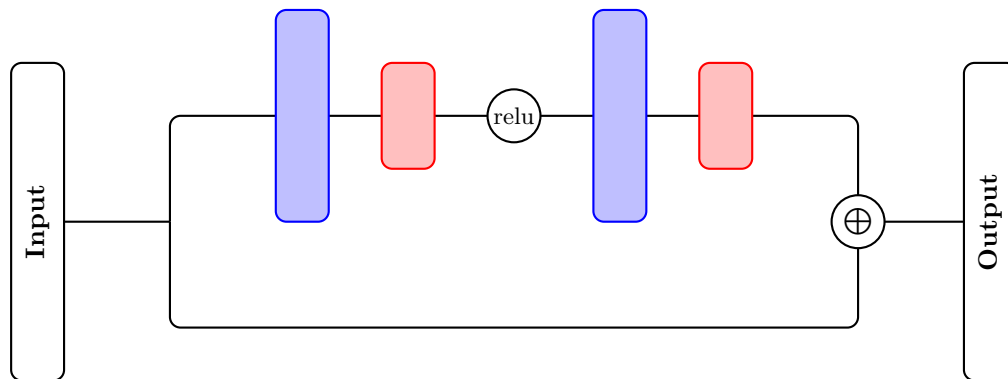


FIGURE 1.9 – ResBlock dans un ResNet avec des couches de **convolutions** et de **Batch Normalisation**

La partie *haute* correspond à une transformation de l'input tandis que la partie *basse* permet de ramener identiquement l'input pour former l'output du ResBlock. Nous étudions la couche de Batch Normalisation plus en détail dans la section 3.3

L'idée de la méthode *stochastic depth* est de ne pas systématiquement réaliser la partie *haute* d'un ResBlock. Il est proposé dans l'article d'associer une probabilité d'évitement du ResBlock de manière linéaire entre le premier ResBlock et le dernier, plutôt que ça soit la même valeur pour tous les blocks, comme dans le dropout.

Cette méthode est unique pour les ResNet, et il n'existe pas d'équivalent³ pour les autres types de réseaux de neurones. On peut noter que GATO de Deepmind [Reed et al., 2022] utilise cet outil.

L'article exploite de nombreuses techniques que l'on définira dans les séances suivantes, et ne traite que des Vision Transformers. Nous verrons en TP si cela reste vrai pour d'autres types de réseaux de neurones⁴. Ce que l'on peut retenir de cet article, s'il est vérifié dans d'autres cas d'usage, est que le dropout devrait plutôt être utilisé comme régularisation du réseau de neurones en début d'entraînement. Si l'on travaille avec un ResNet il peut être intéressant d'ajouter la *stochastic depth* comme mesure de régularisation générale. Dans les deux cas, on peut voir ces méthodes comme des manières d'entraîner une multitude de réseaux de neurones de sorte à créer de la redondance et donc plus de robustesse. Ces deux aspects induisent une régularisation globale du réseau.

2. Cette méthode remporte la compétition ImageNet de 2015.

3. A notre connaissance.

4. Spoiler : Oui, mais il faut bien respecter les hyperparamètres d'architecture de l'article.

Nous avons montré dans cette séance l'importance d'un début d'entraînement. Il conditionne la bonne convergence du réseau à travers une bonne propagation des informations. Nous avons également montré que contraindre le réseau de neurones en début d'entraînement peut le placer dans les meilleures conditions pour atteindre un meilleur optimum. L'importance de cette partie de l'entraînement est telle, que la baisse de performance induite par un mauvais départ ne sera pas rattrapée.

Rien ne sert de courir ; il faut partir à point
— Jean de la Fontaine (1668)

Séance 2

Introduction à la Calibration

Les modèles de Machine Learning sont employés chaque jour un peu plus, en particulier pour des prises de décisions. La performance en termes d'optimisation de la fonction de perte est indispensable. Cependant, nous avons également besoin que le modèle soit capable de renseigner un niveau de confiance **fiable** : prendre une décision avec un niveau de confiance de 99% est plus aisé que de prendre une décision avec un niveau de confiance de 51%. A condition que l'on puisse se fier à ces valeurs !

Cette notion s'appelle la calibration. On dit qu'un prédicteur est calibré quand la probabilité prédite par le modèle est celle que l'on observe réellement. Concrètement, pour modèle parfaitement calibré, si la probabilité qu'il y ait un chat sur une photo est estimé à 0.9, dans 90% des photos où la probabilité était de 0.9 il y avait effectivement un chat. La performance d'optimisation et la calibration sont deux concepts différents, mais pas nécessairement opposés. Il est tout à fait possible d'avoir un modèle très performant mais pas calibré et inversement.

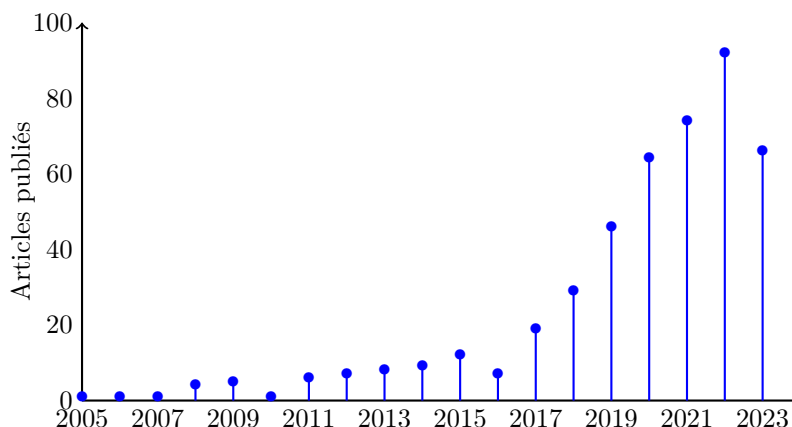


FIGURE 2.1 – Nombres d'articles Arxiv contenant le terme *calibration* dans le titre pour les catégories *cs.AI*, *cs.CL*, *cs.LG*, *math.ST*, *stat.AP*, *stat.CO*, *stat.ML* calculés en août 2023

La demande croissante de garantie sur les modèles déployés a mis petit à petit l'accent sur les travaux autour de la calibration, comme on peut l'observer avec la figure 2.1. Nous avons besoin

de comprendre plus en profondeur ce que cette notion *est* et ce qui l’engendre. La recherche s’articule donc en deux axes majeurs :

1. Comment mesurer la distance à la calibration ?
2. Quel est le lien entre la minimisation d’une fonction de perte et la calibration ?

Nous nous allons nous intéresser à ces deux questions, mais ne traiterons pas l’ensemble des idées proposées. Plus précisément, la séance s’inspire du travail théorique publié en mars et mai 2023 par Jarosław Błasiok, Pariskshi Gopalan, Lunjia Hu et Preetum Nakkiran. Ils proposent un cadre mathématique rigoureux pour répondre à la première question, puis l’exploite pour répondre à la seconde. Nous nous limiterons au cadre d’une classification binaire et ne traiterons pas l’ensemble des métriques proposées pour mesurer la distance à la calibration. Ainsi, nous ne présenterons pas toutes les notions introduites dans les papiers que nous citerons. Nous ne travaillerons que sur celles qui sont utiles pour comprendre les résultats que l’on propose d’étudier. Nous prendrons soin cependant de noter ces écarts à chaque fois que nous le jugerons nécessaire.

2.1 Comment mesurer la distance à la calibration ?


La question de la définition de la distance à la calibration a été le sujet de nombreux articles, qui ont chacun proposé des métriques différentes. Cependant, il est difficile de les relier et donc de les comparer entre elles. De plus, nombre d’entre elles ne parviennent pas à respecter des conditions que l’on souhaiterait intuitivement avoir d’une telle mesure.

L’objectif de l’article *A Unifying Theory of Distance from Calibration* [Błasiok et al., 2023a] publié en mars 2023 par Jarosław Błasiok, Pariskshi Gopalan, Lunjia Hu et Preetum Nakkiran est, comme son titre l’indique, de résoudre ces problèmes et définir un formalisme clair sur ces questions.

Avant de décrire le cheminement de l’article, commençons par formaliser le contexte. Nous avons accès à des observations (x, y) définies comme :

$$(x, y) \sim \mathcal{D} \quad \text{avec } \mathcal{D} \text{ une distribution sur } \mathcal{X} \times \{0, 1\}$$

Domaine discret, espace des inputs



On demande que \mathcal{X} soit discret, éventuellement très grand, pour éviter des difficultés théoriques¹. Cette hypothèse est toujours vérifiée dans la pratique.

On appelle un *prédicteur* une fonction $f : \mathcal{X} \rightarrow [0, 1]$ que l’on interprète comme une estimation de $\mathbb{P}(y = 1|x)$. C’est le résultat de l’entraînement de l’algorithme dont on souhaite mesurer la calibration. On considère la distribution induite par la paire prédiction-label $(f(x), y) \in [0, 1] \times \{0, 1\}$ que l’on notera \mathcal{D}_f . Il ne nous reste plus qu’à définir la parfaite calibration comme [Dawid, 1982].

1. Les auteurs précisent que supposer que f soit mesurable devrait suffire pour traiter le cas où \mathcal{X} est infini.

Définition 2.1 (Parfaite calibration). *On dit qu'une distribution Γ sur $[0, 1] \times \{0, 1\}$ est parfaitement calibrée si, et seulement si, on a :*

$$\mathbb{E}_{(v,y) \sim \Gamma} [y|v] = v$$

On dit qu'un prédicteur f est parfaitement calibré par rapport à \mathcal{D} si \mathcal{D}_f est parfaitement calibrée. On note $\text{cal}(\mathcal{D})$ l'ensemble des prédicteurs parfaitement calibrés par rapport à \mathcal{D} .

Remarquons que l'on peut également définir la parfaite calibration comme l'impossibilité de distinguer $(f(x), y)$ et $(f(x), y')$ où $y' \sim \text{Ber}(f(x))$ quand on ne regarde que les prédictions ! En définissant la parfaite calibration, nous avons également défini l'ensemble des prédicteurs parfaitement calibrés. En sachant que l'on a imposé que \mathcal{X} soit fini, il est légitime de se poser la question pour $\text{cal}(\mathcal{D})$.

Proposition 2.1. *L'ensemble des prédicteurs parfaitement calibrés pour une distribution \mathcal{D} est fini.*

Démonstration. Pour montrer que $\text{cal}(\mathcal{D})$ est fini, montrons qu'il existe une injection entre $\text{cal}(\mathcal{D})$ et l'ensemble des partitions de \mathcal{X} .

Soit $\mathcal{S} = \{S_i \mid \forall i \leq m\}$ une partition de \mathcal{X} . On considère :

$$\forall x \in S_i, g_{\mathcal{S}}(x) = \mathbb{E}[y \mid x \in S_i]$$

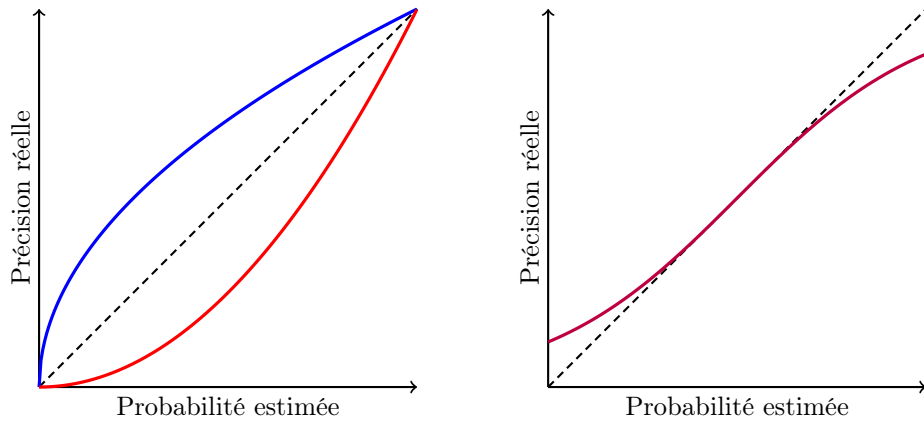
Autrement dit, $g_{\mathcal{S}}$ renvoie la valeur moyenne de y pour chaque partitions $(S_i)_{i \leq m}$ considérées. Par définition (2.1), $g_{\mathcal{S}} \in \text{cal}(\mathcal{D})$.

Soit f un prédicteur. Notons $\text{level}(f)$ les niveaux de la fonction f , et partitionnons \mathcal{X} par cet ensemble. Alors,

$$f \in \text{cal}(\mathcal{D}) \iff f = g_{\text{level}(f)}$$

Nous venons de montrer que l'on peut lier un élément de $\text{cal}(\mathcal{D})$ à un seul élément des partitions de \mathcal{X} , d'où l'injectivité. Puisque \mathcal{X} est fini, par injection, $\text{cal}(\mathcal{D})$ est également fini. \square

Il existe une fonction dans scikit-learn [Pedregosa et al., 2011] qui permet de visualiser la calibration d'un algorithme. En traçant la courbe de la probabilité réelle par rapport à celle estimée, avec en ligne de mire la droite $y = x$ qui correspond à un prédicteur parfaitement calibré. Nous avons une manière visuelle de qualifier le comportement de notre algorithme. Voyons schématiquement les différents cas que l'on peut obtenir :



(a) Prédicteur **trop confiant** et **pas assez confiant** (b) Prédicteur **non parfaitement calibré**

FIGURE 2.2 – Comparaison de prédicteurs par rapport à un prédicteur parfaitement calibré (en pointillés)

Il est important lorsqu'on publie un modèle qui sera amené à aider à la décision de connaître le comportement (confiant ou non) pour assurer une prise de décision saine. Il faut noter que ce graphique est présenté ici sans contexte : c'est une courbe générale qui peut se décliner dans autant de sous-domaines du dataset que l'on souhaite. Ainsi, quand on est amené à utiliser un algorithme qui opère sur des domaines variés, il peut être judicieux de présenter la calibration sur chaque domaine. Cela s'appelle la multicalibration, mais nous ne traiterons pas de cette notion dans ce cours.

2.1.1 Vrai distance à la calibration

Notre objectif est d'être capable de mesurer la distance à la calibration sans avoir à le faire visuellement. Nous serons donc amené à comparer les prédicteurs possibles entre eux. On considère la distance ℓ_1 sur cet ensemble :

Ensemble des prédicteurs $f : \mathcal{X} \rightarrow [0, 1]$

$$\forall f, g \in \mathcal{F}_{\mathcal{X}}, \quad \ell_1(f, g) = \mathbb{E}_{\mathcal{D}}[|f(x) - g(x)|]$$

Notons que dans la définition précédente, nous calculons l'espérance par rapport à la distribution \mathcal{D} des observations (x, y) et non les distributions induites par f ou g . Si f et g sont distant alors D_f et D_g le sont aussi ?

Exercice 2.1. Soit $\mathcal{X} = \{0, 1\}$ et on considère une distribution uniforme sur \mathcal{X} .

1. Quel est le meilleur prédicteur que l'on puisse construire ?
2. On considère $f(x) = x$ et $g(x) = 1 - x$ comme prédicteurs. Que dire de D_f et D_g ?
3. Calculer $\ell_1(f, g)$. Commentez.

Solution. On reprend les mêmes notations.

1. Le prédicteur constant valant 0.5 est l'estimateur bayésien.
2. Les prédicteurs f et g étant uniformes sur \mathcal{X} , les labels sont uniformes une fois conditionnés par f et g , donc \mathcal{D}_f et \mathcal{D}_g sont identiques.
3. On a par définition que $\ell_1(f, g) = 1$ donc que f et g sont distants, bien qu'ils engendrent la même distribution paire-label.

□

La distance ainsi définie permet de vraiment cibler la qualité d'un prédicteur, et non la distribution induite, bien que la définition de calibration se fasse par celle-ci. Nous avons donc un moyen bien défini pour comparer deux prédicteurs. Ce qui veut dire que nous avons une manière de comparer tous les prédicteurs entre eux.

Définition 2.2 (Vrai distance à la calibration). *Pour une distribution \mathcal{D} donnée et une famille $\mathcal{F}_{\mathcal{X}}$ donnée, on définit la vraie distance à la calibration comme :*

$$\text{dCE}_{\mathcal{D}}(f) = \inf_{g \in \text{cal}(\mathcal{D})} \ell_1(f, g)$$

Cette distance serait une bonne candidate pour être une mesure de la distance à la calibration. Pourtant nous ne pouvons pas l'utiliser en pratique, voire la définir à cause du résultat suivant :

Proposition 2.2. *Soit $\alpha \in [0, \frac{1}{2}]$. Il existe un domaine \mathcal{X} , un prédicteur $f \in \mathcal{F}_{\mathcal{X}}$ et deux distributions \mathcal{D}_f^1 et \mathcal{D}_f^2 telles que :*

- Les distributions \mathcal{D}_f^1 et \mathcal{D}_f^2 sont identiques
- $\text{dCE}_{\mathcal{D}_f^1}(f) \leq 2\alpha^2$ tandis que $\text{dCE}_{\mathcal{D}_f^2}(f) \geq \alpha$

Nous ne rédigerons pas la preuve de cette proposition dans ce cours, mais elle est présentée dans l'article [Błasiok et al., 2023a]. Si l'on analyse les bornes, visualisées avec la figure (2.3), on s'aperçoit que nous n'avons pas des borne de dCE mais des prédictions contraires.

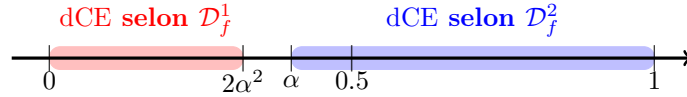


FIGURE 2.3 – Inégalités de la proposition (2.2)

Le coeur du problème soulevé est qu'un même prédicteur, pour des distributions indiscernables, peut avoir des distances à la calibration incohérentes si nous n'avons pas accès à x . Or nous avons fait le choix de travailler avec uniquement la paire prédiction-label $(f(x), y)$.

Comme dCE est intellectuellement simple à comprendre et correspond à ce que l'on cherche à obtenir, nous allons proposer des manières de calculer des approximations de cette distance dans le cadre que l'on s'est défini.

2.1.2 Mesure de calibration consistante

Pour cela, commençons par formaliser ce que l'on attend d'une mesure de calibration : une fonction qui renvoie un nombre entre 0 et 1 à partir d'une distribution \mathcal{D} d'observations et d'un prédicteur f . On note cette valeur $\mu_{\mathcal{D}}(f)$. Au minimum, on souhaite avoir les propriétés suivantes :

$$\begin{aligned} \mu_{\mathcal{D}}(f) &= 0 & \text{si } f \in (\mathcal{D}) & & \text{(Complétude)} \\ \mu_{\mathcal{D}}(f) &> 0 & \text{si } f \notin (\mathcal{D}) & & \text{(Correction)} \end{aligned}$$

Autrement dit, on demande à la mesure de calibration de valoir 0 quand le prédicteur f est parfaitement calibré par rapport à \mathcal{D} , et d'être strictement positif sinon. Ces termes viennent de la logique mathématiques : *Completeness* et *Soundness* en anglais. Concrètement, la complétude garantit que l'on peut prouver ce qui est vrai, et la correction que l'on ne peut pas prouver prouver quelque chose qui est faux.

Cela dit, on préférerait que la valeur soit petite si on est proche de la calibration parfaite, et grande si on en est loin, plutôt qu'une réponse binaire.

Définition 2.3 (Mesure consistante de calibration). *Soit $c \in \mathbb{R}_+$, on dit qu'une mesure de calibration μ satisfait la c -robuste complétude si :*

$$\exists a \in \mathbb{R}_+, \forall \mathcal{D} \sim [0, 1] \times \{0, 1\}, \forall f \in \mathcal{F}_{\mathcal{X}}, \quad \mu_{\mathcal{D}}(f) \leq a (\text{dCE}_{\mathcal{D}}(f))^c$$

Soit $s \in \mathbb{R}_+$, on dit qu'une mesure de calibration μ satisfait la s -robuste correction si :

$$\exists b \in \mathbb{R}_+, \forall \mathcal{D} \sim [0, 1] \times \{0, 1\}, \forall f \in \mathcal{F}_{\mathcal{X}}, \quad \mu_{\mathcal{D}}(f) \geq b (\text{dCE}_{\mathcal{D}}(f))^s$$

On dit que μ est (c, s) -consistante si les deux conditions précédentes sont vérifiées.

Si un prédicteur f est ε proche d'être parfaitement calibré alors par la robuste complétude on a que $\mu_{\mathcal{D}}(f)$ est plus petit que ε^c et tend vers 0 quand ε tend vers 0, de même pour la robuste correction.

Quand la mesure de calibration de f vaut $\eta \gg 0$ alors par la robuste correction est de l'ordre de $\eta^s \gg 0$.

Notons que dans la définition de consistance, nous n'avons pas pris en compte le fait que la mesure de calibration puisse être calculable en ayant accès uniquement aux prédictions et aux labels. Voyons ce que cela impose comme conditions.

Corollaire 2.1. *Soit μ une mesure de calibration (c, s) -consistante que l'on peut calculer à l'aide des prédictions et des labels. Alors $s \geq 2c$*

Démonstration. Il s'agit du corollaire de la proposition (2.2), nous reprenons donc ses notations. Par définition de la complétude robuste pour \mathcal{D}^1 et la correction robuste pour \mathcal{D}^2 , on a :

$$\begin{aligned} \exists a \in \mathbb{R}_+, \quad \mu_{\mathcal{D}^1}(f) &\leq a(2\alpha^2)^c \\ \exists b \in \mathbb{R}_+, \quad \mu_{\mathcal{D}^2}(f) &\geq b\alpha^s \end{aligned}$$

Puisque $\mathcal{D}_f^1 = \mathcal{D}_f^2$ et que μ est calculable à l'aide des prédictions et des labels on a que $\mu_{\mathcal{D}^1}(f) = \mu_{\mathcal{D}^2}(f)$ donc :

$$b\alpha^s \leq a(2\alpha^2)^c \iff \left(\frac{b}{a}\alpha\right)^{\frac{s}{c}} \leq 2\alpha^2$$

On obtient clairement une contradiction quand α tends vers zéro si $\frac{s}{c} < 2$. \square

Nous avons restreint les mesures de calibration éligibles à travers les différents souhaits que nous avons décrit : consistance et capacité d'être calculées dans le cadre prédiction-label. Il nous faut maintenant en exhiber une ! Commençons par la manière la plus intuitive de mesurer la distance à la calibration.

Définition 2.4 (Expected Calibration Error). *Soit \mathcal{D} une distribution sur $[0, 1] \times \{0, 1\}$ et $f : \mathcal{X} \rightarrow [0, 1]$ un prédicteur. On appelle Expected Calibration Error :*

$$ECE_{\mathcal{D}}(f) = \mathbb{E}_{\mathcal{D}} \left[\left| \mathbb{E}_{\mathcal{D}}[y | f(x)] - f(x) \right| \right]$$

On saisit ce que l'ECE veut faire : mesurer la distance moyenne entre la valeur moyenne de y sachant f et la valeur de f . Si l'ECE est égale à zéro, alors le prédicteur est parfaitement calibré. Voyons si l'ECE est consistante.

Lemme 2.1. *Soit $\mathcal{S} = \text{level}(f)$. Alors $ECE(f) = \ell_1(f, g_{\mathcal{S}}) \geq \text{dCE}_{\mathcal{D}}(f)$*

Démonstration. On a que pour tout $x \in f^{-1}(v)$, $f(x) = v$ donc $g_{\mathcal{S}}(x) = \mathbb{E}[y|v]$. Ainsi, on a que :

$$\begin{aligned} \mathbb{E} \left[\left| \mathbb{E}[y|f(x)] - f(x) \right| \right] &= \mathbb{E} [|f(x) - g_{\mathcal{S}}(x)|] \\ ECE(f) &= \ell_1(f, g_{\mathcal{S}}) \end{aligned}$$

Par définition de l'ECE et ℓ_1 . On obtient l'inégalité par la définition de dCE qui minimise la distance ℓ_1 parmi l'ensemble des prédicteurs calibré. \square

Donc l'ECE vérifie la 1-robuste correction. Ce qui veut dire qu'avec l'ECE on ne peut pas prouver quelque chose de faux. Cependant elle ne vérifie pas la robuste complétude.

En effet, on considère une distribution uniforme sur $\mathcal{X} = \{0, 1\}$ où le meilleur prédicteur est $f^*(0) = 0$ et $f^*(1) = 1$. On considère le prédicteur f_{ε} avec $\varepsilon \geq 0$ défini par $f_{\varepsilon}(0) = \frac{1}{2} - \varepsilon$ et $f_{\varepsilon}(1) = \frac{1}{2} + \varepsilon$.

Notons que $f_0(x) = \frac{1}{2}$ qui est parfaitement calibré. Cependant, on a par la définition de l'ECE que $ECE_{\mathcal{D}}(f_{\varepsilon}) = \frac{1}{2} - \varepsilon$ donc $ECE_{\mathcal{D}}(f_0) = \frac{1}{2}$ alors que f_0 est parfaitement calibré !

Nous n'avons donc pas dans ce cas de continuité en 0 en plus d'une non fiabilité dans les valeurs par incomplétude. Autrement dit, il ne nous est pas possible de toujours prouver quelque chose qui est vrai. Nous venons de montrer que l'ECE n'est pas une mesure de calibration consistante bien qu'elle soit largement adopté par la communauté Machine Learning. Nous devons trouver une alternative qui réponde à nos *désiderata*.

2.1.3 Plus petite distance à la calibration

Avant de se lancer dans l'étude d'une nouvelle mesure de calibration, peut-être qu'il serait utile de travailler un peu plus autour de la distance dCE. Nous savons que nous ne pouvons pas l'exploiter directement, mais cette notion est centrale pour la consistance. Nous devons donc obtenir des garanties sur les estimations que nous pourrions faire. Pour cela nous allons définir la distance inférieure à la calibration².

Définition 2.5. Soit Γ une distribution sur $[0, 1] \times \{0, 1\}$. On définit l'ensemble $\text{ext}(\Gamma)$ comme l'ensemble des distributions Π des triplets $(u, v, y) \in [0, 1] \times [0, 1] \times \{0, 1\}$ tel que :

- La distribution marginale de (v, y) est Γ
- La distribution marginale de (u, y) est parfaitement calibré i.e. $\mathbb{E}_{\Pi}[y|u] = u$.

Cet ensemble de distributions n'est pas évident à sentir. Commençons par noter que les deux points ne sont *liés* que par y . Le premier représente une distribution que l'on observe, et le second représente une distribution que l'on souhaiterait obtenir. Finalement, nous aimerions que v et u soient le plus proche possible. Cette intuition est celle de la notion centrale de cette sous-section.

Définition 2.6 (Distance inférieure à la calibration). On définit la distance inférieure à la calibration par :

$$\underline{\text{dCE}}(\Gamma) = \inf_{\Pi \in \text{ext}(\Gamma)} \mathbb{E}_{(u, v, y) \sim \Pi} [|u - v|]$$

Pour une distribution \mathcal{D} et un prédicteur f , on définit $\underline{\text{dCE}}_{\mathcal{D}}(f) = \underline{\text{dCE}}(\mathcal{D}_f)$

Ce qui différencie cette notion de la définition de dCE est que la distributions est *liée* par y comme nous l'avons noté précédemment. Donc nous avons *accès* à toutes les informations possibles dans le cas où l'on observe $f(x)$ et y . Cela règle le problème de dCE. Avec la notation, il semblerait que l'on puisse lier $\underline{\text{dCE}}$ et dCE.

Puisque dCE traite de prédicteurs calibrés, remarquons que tout prédicteur $g \in \text{cal}(\mathcal{D})$ génère une distribution $\Pi \in \text{ext}(\mathcal{D}_f)$ où de $(x, y) \sim \mathcal{D}$ on obtient $(g(x), f(x), y)$. En remarquant que $\mathbb{E}_{(u, v, y) \sim \Pi} [|u - v|] = \ell_1(f, g)$, il suffit de minimiser pour $g \in \text{cal}(\mathcal{D})$ pour obtenir :

$$\underline{\text{dCE}}_{\mathcal{D}}(f) \leq \text{dCE}_{\mathcal{D}}(f)$$

Nous venons de justifier la notation et le nom de $\underline{\text{dCE}}$. Mais nous n'avons pas un lien suffisamment clair encore entre dCE et $\underline{\text{dCE}}$. Dans l'article [Błasiok et al., 2023a], via des raisonnements que nous ne traiterons pas ici, est obtenue l'inégalité suivante :

Proposition 2.3. Pour toute distribution \mathcal{D} et tout prédicteur f , on a :

$$\underline{\text{dCE}}_{\mathcal{D}}(f) \leq \text{dCE}_{\mathcal{D}}(f) \leq 4\sqrt{\underline{\text{dCE}}_{\mathcal{D}}(f)}$$

2. Dans l'article est également défini la distance supérieure à la calibration, mais nous n'en avons pas besoin dans le cadre de ce cours.

Cet encadrement ne permet pas une estimation précise de dCE mais informe que l'écart est au plus quadratique entre dCE et $\underline{\text{dCE}}$. Dans la théorie générale sur la calibration, cela reste vrai dans un grand nombre de cas. De plus, si une mesure de calibration est encadrée par la distance inférieure à la calibration, alors on sera capable de prouver la consistance de cette mesure en exploitant ce résultat. L'intérêt de la proposition (2.3) est plutôt dans la compréhension générale des notions que dans un but concret, pour le moment.

2.1.4 Smooth Calibration

Une alternative proposée à l'ECE est la smooth-calibration. Cette notion est un cas particulier d'une famille plus générale introduite dans l'article *Low-degree multicalibration* [Gopalan et al., 2022].

Définition 2.7 (Calibration pondérée). *Soit W une famille de fonctions $w : [0, 1] \rightarrow \mathbb{R}$. La calibration pondérée d'une distribution Γ sur $[0, 1] \times \{0, 1\}$ est définie comme :*

$$\text{wCE}^W(\Gamma) = \sup_{w \in W} \left| \mathbb{E}_{(v,y) \sim \Gamma} [(y - v)w(v)] \right|$$

Étant donné une distribution \mathcal{D} sur $\mathcal{X} \times \{0, 1\}$ et un prédicteur $f : \mathcal{X} \rightarrow [0, 1]$, on appelle la calibration pondérée de f par rapport à \mathcal{D} :

$$\text{wCE}_{\mathcal{D}}^W(f) = \text{wCE}^W(\mathcal{D}_f) = \sup_{w \in W} \left| \mathbb{E}_{(x,y) \sim \mathcal{D}} [(y - f(x))w(f(x))] \right|$$

La définition précédente est cohérente avec la complétude : si Γ est parfaitement calibrée alors la valeur de wCE^W sera zéro. En effet si l'on conditionne par v :

$$\begin{aligned} \mathbb{E}_{\Gamma} [(y - v)w(v)] &= \mathbb{E}_{\Gamma} \left[\mathbb{E} [(y - v)w(v) \mid v] \right] \\ &= \mathbb{E}_{\Gamma} \left[w(v) \mathbb{E} [y \mid v] - v(w(v)) \right] \\ &= 0 \quad \text{parce que } \Gamma \text{ est parfaitement calibrée, par définition } \mathbb{E} [y \mid v] = v \end{aligned}$$

Une famille importante est celle des fonctions 1-Lipschitzienne bornées. Elle est pour la première fois introduite dans ce contexte dans l'article *Deterministic calibration and Nash equilibrium* [Kakade and Foster, 2004]. Nous suivrons cependant les notations de [Gopalan et al., 2022]. Pour rappel, une fonction $w : [0, 1] \rightarrow [-1, 1]$ est 1-Lipschitzienne si elle vérifie :

$$\forall x, y \in [0, 1], \quad |w(x) - w(y)| \leq |x - y|$$

Cela contraint la vitesse d'évolution de la fonction w selon x . Autrement dit, on demande que w croisse au plus linéairement. Visuellement :

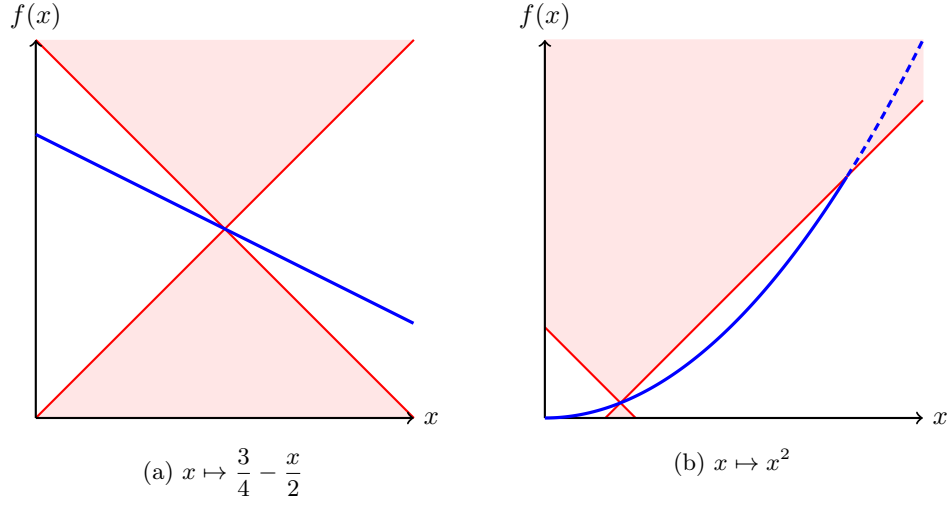


FIGURE 2.4 – Exemple et contre exemple de fonction 1-Lipschitzienne

La première fonction est même $\frac{1}{2}$ -Lipschitzienne, alors que la seconde n'est pas 1-Lipschitzienne. Nous n'avons présenté ici que des fonctions à valeurs dans $[0, 1]$, mais l'intervalle $[-1, 1]$ est autorisé par la définition³. Le rappel étant fait, nous pouvons définir la mesure smooth calibration comme cas particulier de la calibration pondérée.

Définition 2.8 (Smooth calibration). *Soit L la famille des fonctions 1-Lipschitzienne $w : [0, 1] \rightarrow [-1, 1]$. L'erreur de smooth calibration d'une distribution Γ sur $[0, 1] \times \{0, 1\}$ est définie comme :*

$$\text{smCE}(\Gamma) = \text{wCE}^L(\Gamma)$$

De manière similaire, pour une distribution \mathcal{D} et un prédicteur f on définit :

$$\text{smCE}_{\mathcal{D}}(f) = \text{smCE}(\mathcal{D}_f) = \text{wCE}^L(\mathcal{D}_f)$$

Il reste à se poser toujours la même question : smCE est-elle consistante ? Nous avons déjà montré que la calibration pondérée en général vérifie la complétude. Il reste à identifier avec quel coefficient et savoir si elle vérifie la correction robuste. Pour cela, lions la smooth-calibration à la plus petite distance à la calibration.

Théorème 2.1. *Soit une distribution Γ sur $[0, 1] \times \{0, 1\}$. Alors,*

$$\frac{1}{2} \underline{\text{dCE}}(\Gamma) \leq \text{smCE}(\Gamma) \leq 2 \underline{\text{dCE}}(\Gamma)$$

Nous avons écrit plus tôt que si nous étions en capacité de borner une mesure de calibration avec dCE alors nous serons capable de prouver la consistance. Nous avons un corollaire immédiat de cet encadrement :

Corollaire 2.2. *La mesure de calibration smCE est (1, 2)-consistante.*

3. Cette distinction fera encore plus sens dans la seconde partie du chapitre.

Démonstration. Soit Γ une distribution sur $[0, 1] \times \{0, 1\}$. Avec la partie droite de l'inégalité du théorème 2.1 :

$$\text{smCE}(\Gamma) \leq 2\text{dCE}(\Gamma) \leq 2 \times \text{dCE}(\Gamma)$$

D'où la 1-robuste complétude. Avec la partie gauche de l'inégalité du théorème 2.1 :

$$\text{smCE}(\Gamma) \geq \frac{1}{2}\text{dCE}(\Gamma)$$

En exploitant la partie droite de l'inégalité de la proposition 2.3 :

$$\text{dCE}_{\mathcal{D}}(f) \leq 4\sqrt{\text{dCE}_{\mathcal{D}}(f)}$$

En combinant les deux précédentes inégalités, on a :

$$\begin{aligned} \text{smCE}(\Gamma) \geq \frac{1}{2}\text{dCE}(\Gamma) &\implies \sqrt{\text{smCE}(\Gamma)} \geq \sqrt{\frac{1}{2}}\sqrt{\text{dCE}(\Gamma)} \\ &\implies 4\sqrt{\text{smCE}(\Gamma)} \geq \sqrt{\frac{1}{2}}\text{dCE}(\Gamma) \\ &\implies \text{smCE}(\Gamma) \geq \frac{1}{32}\text{dCE}(\Gamma)^2 \end{aligned}$$

D'où la 2-robuste correction et donc la $(1, 2)$ -consistance. \square

Remarquons que la condition imposée par le corollaire 2.1 est bien vérifiée et même optimale puisque $\frac{s}{c} = 2$. Il nous reste à prouver le théorème 2.1. Nous ne prouverons que la borne supérieure ici, la borne inférieure nécessitant beaucoup plus de matériel technique pour être traitée dans ce cours.

Démonstration du théorème 2.1, borne supérieure. Soit Π une distribution sur $[0, 1] \times [0, 1] \times \{0, 1\}$ et soit $w : [0, 1] \rightarrow [-1, 1]$ une fonction 1-Lipschitzienne. On a :

$$\begin{aligned} \left| \mathbb{E}_{(u,v,y) \sim \Pi} [(y-u)w(u)] - \mathbb{E}_{(u,v,y) \sim \Pi} [(y-v)w(v)] \right| &= \left| \mathbb{E}_{(u,v,y) \sim \Pi} [(y-u)w(u) - (y-v)w(v)] \right| \\ &\leq \mathbb{E}_{(u,v,y) \sim \Pi} [| (y-u)(w(u) - w(v)) |] + \mathbb{E}_{(u,v,y) \sim \Pi} [| (u-v)w(v) |] \\ &\leq 2 \mathbb{E}_{(u,v,y) \sim \Pi} [| (u-v) |] \end{aligned}$$

Parce que $|y - u| \leq 1$ et par définition de w . Ainsi si $\Pi \in \text{ext}(\Gamma)$:

$$\begin{aligned} \left| \mathbb{E}_{(u,v,y) \sim \Pi} [(y-v)w(v)] \right| &\leq 2 \mathbb{E}_{(u,v,y) \sim \Pi} [|u-v|] + \left| \mathbb{E}_{(u,v,y) \sim \Pi} [(y-u)w(u)] \right| \\ &\leq 2 \mathbb{E}_{(u,v,y) \sim \Pi} [|u-v|] \end{aligned}$$

Puisque par définition de $\text{ext}(\Gamma)$ la distribution de (u, y) est parfaitement calibrée. En prenant l'infimum sur $\Pi \in \text{ext}(\Gamma)$ on obtient l'inégalité souhaitée. \square

Au terme de cette section, nous savons définir une mesure de calibration consistante et en avons exhibé une que nous avons étudiée plus en détail. Il est important de noter que plusieurs autres mesures de calibration existent et que les principales ont été étudiées dans l'article référence de ce chapitre [Błasiok et al., 2023a]. Si l'on poursuit l'étude de cet article, on comprendra que le cadre théorique introduit permet d'étudier en profondeur les mesures proposées dans la littérature, et les améliorer quand elles présentent des faiblesses.

Cependant, après ce travail pour qualifier la calibration, nous ne savons pas comment *l'obtenir*. Est-on capable d'avoir des garanties sur la calibration d'un algorithme ? Et si oui, quels sont les critères qui permettent de s'en assurer ? La question est vaste, et nous nous concentrerons dans la deuxième partie de ce chapitre au critère le plus intuitif qui pourrait être lié à la calibration.

2.2 Comment induire la calibration ?

Les algorithmes de classifications sont souvent interprétés comme des estimateurs de $\mathbb{P}(y|x)$. Ainsi, si l'optimisation de la fonction de perte est optimale, on devrait être proche de la calibration. Pour la régression logistique, qui a une interprétation naturelle avec les probabilités, ce lien semble évident comme on peut le lire dans le guide utilisateur de scikit-learn :

Properly regularized logistic regression is well calibrated by default thanks to the use of the log-loss

— Guide utilisateur scikit-learn (2007)

Cette affirmation n'est pas prouvée dans le guide, mais semble plausible. Si l'on généralise, l'utilisation d'une fonction de perte, peu importe l'algorithme, doit induire la même qualité de calibration à condition que l'optimal soit atteint. Cependant, en pratique, on observe qu'un réseau de neurones est souvent mieux calibré qu'une régression logistique sur les mêmes tâches alors qu'ils s'entraînent avec la même fonction de perte.

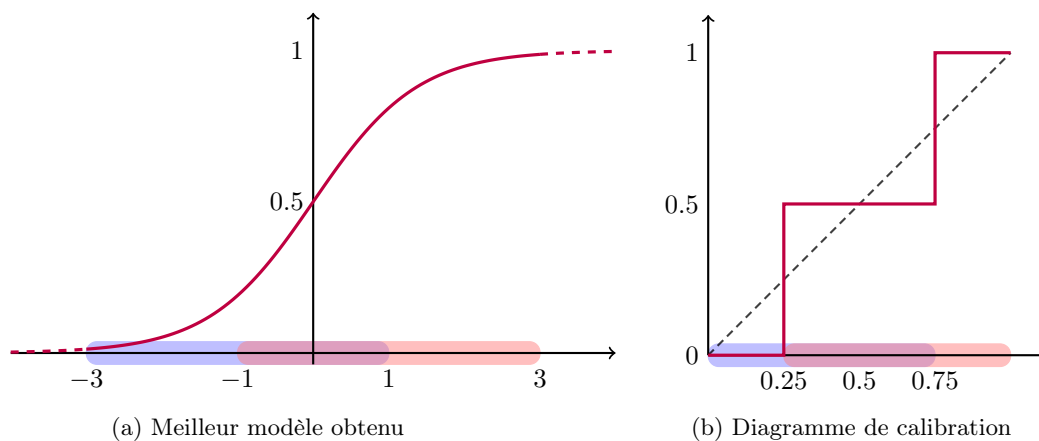


FIGURE 2.5 – Régression logistique pour la distribution $(x, y) \sim \mathcal{D}$ avec $y = 0$ pour $x \sim \mathcal{U}([-3, 1])$ et $y = 1$ pour $x \sim \mathcal{U}([-1, 3])$

Avec l'exemple (2.5) on voit que l'on a réussi à obtenir le meilleur prédicteur en terme de minimisation de loss et pourtant, l'estimateur n'est pas calibré. Notons que le prédicteur constant de valeur 0.5 est parfaitement calibré. Notre intuition de départ serait plus correcte si nous

optimisations parmi l'ensemble des prédicteurs possibles. Hors ce n'est pas ce que nous faisons : nous avons développé des algorithmes qui sont une partie des prédicteurs possibles. Dans l'exemple (2.5) nous avons choisi la régression logistique, mais il se pourrait qu'un autre algorithme réussisse à mieux combiner calibration et optimisation de la fonction de perte.

Le lien entre performance d'optimisation de la fonction de perte et calibration n'est finalement pas si évident et contre-intuitif. C'est ce travail qui a été réalisé avec l'article *When Does Optimizing a Proper Loss Yield Calibration?* [Błasiok et al., 2023b] publié en mai 2023, qui suit l'article *A Unifying Theory of Distance from Calibration* publié en mars 2023.

Comme dans la première section, nous ne présenterons pas l'ensemble de la théorie développée dans l'article, mais nous nous concentrerons sur les concepts centraux qui permettent d'avoir une piste solide pour apprécier le lien entre performance en fonction de perte et la calibration.

Commençons par éclaircir ce qu'on appelle *proper loss*. Dans l'ensemble de la section on considère une distribution \mathcal{D} sur $\mathcal{X} \rightarrow \{0, 1\}$ que nous omettrons de préciser pour simplifier la lecture, même dans la notation de smCE par exemple. On considère un prédicteur $f : \mathcal{X} \rightarrow [0, 1]$, sauf indication contraire.

Définition 2.9 (Fonction de perte bien définie). *Soit $V \subseteq [0, 1]$ un interval non vide. On dit qu'une fonction de perte $\mathcal{L} : \{0, 1\} \times [0, 1]$ est bien définie si pour tout $v \in V$, on a que :*

$$v \in \arg \min_{v' \in V} \mathbb{E}_{y \sim \text{Ber}(v)} [\mathcal{L}(y, v')]$$

Autrement dit, on demande qu'une fonction de perte permette effectivement d'atteindre le minimum que l'on cherche. On peut montrer que la *Squared loss* \mathcal{L}_{sq} est bien définie.

$$\forall y \in \{0, 1\}, \forall v \in V \subseteq [0, 1], \quad \mathcal{L}_{\text{sq}}(y, v) = (y - v)^2$$

Comme souvent en recherche, continuons d'étudier ce cas particulier pour voir ce que l'on pourrait obtenir quand on généralisera.

2.2.1 Cas particulier de la *Squared loss*

La notion centrale qu'introduit l'article est le *post-processing gap* :

Définition 2.10 (Post-processing gap). *Soit K une famille de fonctions de post-processing $\kappa : [0, 1] \rightarrow [0, 1]$ telle que la fonction de mise à jour $\eta(v) = \kappa(v) - v$ soit 1-Lipschitzienne. On appelle post-processing gap de f par rapport à \mathcal{D} :*

$$\text{pGap}_{\mathcal{D}}(f) = \mathbb{E}[\mathcal{L}_{\text{sq}}(y, f(x))] - \inf_{\kappa \in K} \mathbb{E}[\mathcal{L}_{\text{sq}}(y, \kappa \circ f(x))]$$

Le post-processing gap mesure le gain que l'on peut obtenir en performance sur la loss si l'on modifiait les prédictions par : $f(x) \mapsto f(x) + \eta(f(x))$ avec η une fonctions 1-Lipschitziennes. Le choix de la famille de fonction 1-Lipschitzienne, nous fait penser qu'on pourrait relier le post-processing gap avec la smooth calibration définie plus haut : c'est l'objet du théorème 2.2.

Théorème 2.2. Pour la fonction de perte \mathcal{L}_{sq} :

$$\text{smCE}(f)^2 \leq \text{pGap}_{\mathcal{D}}(f) \leq 2 \text{smCE}(f)$$

Ce résultat permet de lier la performance en terme de calibration à la performance en terme de fonction de perte. Plus précisément, à quel point on peut gagner si l'on modifie un peu les prédictions que l'on fait. Si l'on modifie très peu la performance alors pGap est faible et donc smCE l'est aussi donc nous avons une bonne calibration ! Inversement, si pGap est grand, c'est qu'en modifiant les prédictions on peut obtenir un vrai gain de performance. En utilisant la borne droite de l'inégalité de (2.2) on obtient une calibration de faible qualité.

Intuitivement, si nous donnions plus de *puissance* à notre algorithme, il pourrait être en mesure *d'apprendre* lui-même cette fonction de post-processing et ainsi optimiser encore plus ses performances. Cela nous amène à nous poser la question du choix de l'algorithme, paramètre que nous n'avons pas pris en compte pour le moment.

Avant d'intégrer cette réflexion dans l'analyse, le résultat précédent doit être plus général. En effet, cela ne fonctionne que pour la fonction de perte précisée. Prenons une autre fonction de perte, l'entropie croisée :

Observation positive

$$\mathcal{L}_{\text{xent}}(y, v) = - \left[y \ln(v) + (1 - y) \ln(1 - v) \right]$$

Observation négative

On considère $\mathcal{X} = \{x_0, x_1\}$ et \mathcal{D} une distribution sur $\mathcal{X} \rightarrow \{0, 1\}$ telle que :

- $\mathbb{E}_{(x,y) \in \mathcal{D}} [y \mid x = x_0] = 0.1$
- $\mathbb{E}_{(x,y) \in \mathcal{D}} [y \mid x = x_1] = 0.9$

Soit $f_\varepsilon : \mathcal{X} \rightarrow [0, 1]$ le prédicteur défini par $f_\varepsilon(x_0) = \varepsilon$ et $f_\varepsilon(x_1) = 1 - \varepsilon$ pour $\varepsilon \in]0, 0.1[$.

Par la définition, on a que $\lim_{\varepsilon \rightarrow 0} \text{smCE}(f_\varepsilon) = 0.05$. Pourtant, puisque $\lim_{\varepsilon \rightarrow 0} \mathcal{L}_{\text{xent}}(1, \varepsilon) = \lim_{\varepsilon \rightarrow 0} \mathcal{L}_{\text{xent}}(0, 1 - \varepsilon) = +\infty$ invalidant le résultat obtenu précédemment pour \mathcal{L}_{sq} .

Nous ne pouvons donc pas *simplement* remplacer la fonction de perte. Il faut trouver une approche plus générale, en regardant la notion de fonction de perte d'une manière *différente*.

2.2.2 Généralisation du lien entre pGap et smCE : par la dualité

Les travaux de Leonard Savage en 1971 [Savage, 1971] et ceux de Tilman Gneiting et Adrian Raftery en 2007 [Gneiting and Raftery, 2007] ont construit un pont entre les fonctions de pertes bien définies et les fonctions convexes. Ce genre de pont est un exemple d'un grand concept : la dualité. On peut penser à la dualité onde-particule en physique quantique par exemple. Ici, plus modestement, nous allons traduire l'ensemble des notions précédentes dans le monde dual adapté. Le pont entre les deux mondes que nous considérerons se réalise via le résultat suivant :

Proposition 2.4. Soit $V \subseteq [0, 1]$ un intervalle non vide. Soit $\mathcal{L} : \{0, 1\} \times V \rightarrow \mathbb{R}$ une fonction de perte bien définie. On a :

Fonction convexe $\psi : \mathbb{R} \rightarrow \mathbb{R}$

$$\forall y \in \{0, 1\}, \forall v \in V, \quad \mathcal{L}(y, v) = \psi \left(\text{dual}(v) \right) - y \text{dual}(v)$$

$\forall v \in V, \text{dual}(v) = \mathcal{L}(0, v) - \mathcal{L}(1, v)$

De plus, si ψ est différentiable, alors $\nabla \psi(t) \in [0, 1]$

Cette proposition peut faire penser au théorème de représentation de Riesz : on explicite comment lier l'espace primal au dual. Dans notre cas, ce résultat indique que d'une fonction de perte bien définie s'induit une fonction convexe ainsi qu'une fonction qui relie une prédiction $v \in V$ à sa prédiction duale $\text{dual}(v) \in \mathbb{R}$. Nous verrons plus concrètement avec l'exercice 2.2 comment cela se décline, mais avant traduisons une fonction de perte dans le monde dual.

Définition 2.11 (Fonction de perte duale). Soit une fonction $\psi : \mathbb{R} \rightarrow \mathbb{R}$. On définit une fonction de perte $\mathcal{L}^{(\psi)} : \{0, 1\} \times \mathbb{R} \rightarrow \mathbb{R}$ comme :

$$\forall y \in \{0, 1\}, \forall t \in \mathbb{R}, \quad \mathcal{L}^{(\psi)}(y, t) = \psi(t) - yt$$

Ainsi, si une fonction de perte $\mathcal{L} : \{0, 1\} \times V \rightarrow \mathbb{R}$ satisfait les conditions de la proposition 2.4 pour $V \subseteq [0, 1]$ un intervalle non vide et $\text{dual} : V \rightarrow \mathbb{R}$, alors :

$$\forall y \in \{0, 1\}, \forall v \in V, \quad \mathcal{L}(y, v) = \mathcal{L}^{(\psi)}(y, \text{dual}(v))$$

La définition de la fonction de perte duale met en lumière le rôle de la fonction dual introduite dans la proposition 2.4. Cette fonction permet de faire la traduction entre la prédiction v dans le monde réel et la prédiction t dans le monde dual. Notons que si $v \in [0, 1]$, t peut lui valoir l'ensemble des valeurs réelles. Manipulons ces notions sur l'entropie croisée.

Exercice 2.2 (Entropie croisée). On considère la fonction de perte entropie croisée définie comme :

$$\mathcal{L}_{\text{xent}}(y, v) = -y \ln(v) - (1 - y) \ln(1 - v)$$

On souhaite obtenir l'ensemble des informations nécessaires pour expliciter le mapping vers le monde dual. Identifier la fonction de perte duale, ψ et dual.

Solution. Par la définition de dual dans la proposition 2.4, on obtient que $v = \sigma(t) = \frac{1}{1 + e^{-t}}$. La fonction de perte duale de l'entropie croisée est la fonction de perte logistique :

$$\mathcal{L}^{(\psi)}(y, t) = \ln(1 + e^t) - yt$$

On a par définition de la fonction de perte duale que $\psi(t) = \ln(1 + e^t)$. On peut visualiser ces changements avec la figure (2.6).

□

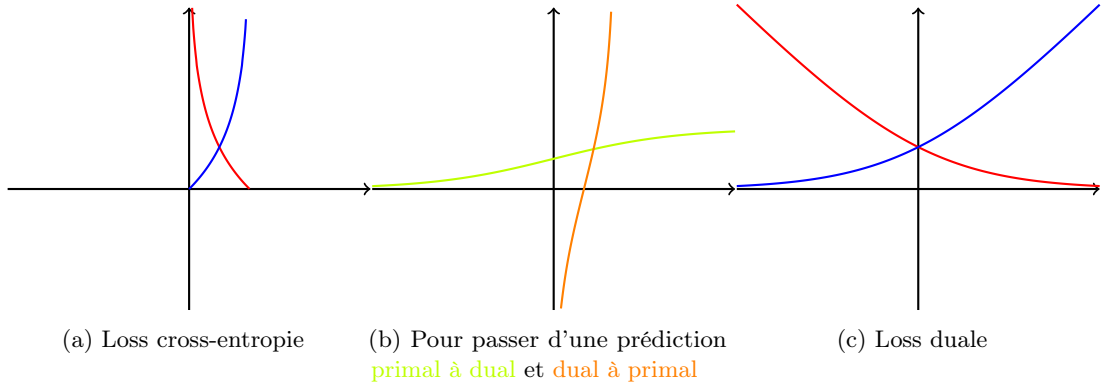


FIGURE 2.6 – Courbes des fonctions de pertes pour $y = 0$ et $y = 1$ ainsi que les fonctions de transferts

Rappelons que notre objectif est de généraliser les conclusions obtenue sur le cas particulier de la square loss. Puisque ce résultat impliquait la smooth calibration et le post processing gap du monde primal, nous devons définir ces notions dans le monde dual. Cette fois, les notions feront intervenir les prédictions duales. Pour un prédicteur $f : \mathcal{X} \rightarrow V$, ses prédictions duales sont : $g(x) = \text{dual}(f(x))$ par rapport à la fonction de perte $\mathcal{L}^{(\psi)}$.

Définition 2.12 (Gap de post-processing dual). On note K_λ pour $\lambda > 0$ l'ensemble des fonctions de post-processing $\kappa : \mathbb{R} \rightarrow \mathbb{R}$ telle qu'il existe une fonction $\eta : \mathbb{R} \rightarrow [-\frac{1}{\lambda}, \frac{1}{\lambda}]$ qui est 1-Lipschitzienne et qui vérifie pour tout $t \in \mathbb{R}$, $\kappa(t) = t + \eta(t)$.

On définit le gap de post-processing de $g : \mathcal{X} \rightarrow \mathbb{R}$:

$$\text{pGap}^{(\psi, \lambda)}(g) = \mathbb{E} \left[\underbrace{\mathcal{L}^{(\psi)}(y, g(x))}_{\text{Fonction de perte duale}} \right] - \inf_{\kappa \in K_\lambda} \mathbb{E} \left[\mathcal{L}^{(\psi)}(y, \kappa \circ g(x)) \right]$$

Notons que pour définir le gap de post-processing dual, nous n'avons eu qu'à décliner et revoir les ensembles de définitions des fonctions de la définition du gap de post-processing primal. L'adaptation de chaque ensemble de définition peut se deviner avec la figure (2.6) précédente.

Définition 2.13 (Smooth calibration dual). On note H_λ pour $\lambda > 0$ l'ensemble des fonctions $\eta : \mathbb{R} \rightarrow [-1, 1]$ qui sont λ -Lipschitzienne. Soit $\psi : \mathbb{R} \rightarrow \mathbb{R}$ une fonction différentiable avec pour tout $t \in \mathbb{R}$, $\nabla \psi(t) \in [0, 1]$.

Pour une fonction $g : \mathcal{X} \rightarrow \mathbb{R}$, on définit le prédicteur $f : \mathcal{X} \rightarrow [0, 1]$ tel que pour tout $x \in \mathcal{X}$, $f(x) = \nabla \psi(g(x))$.

On appelle smooth calibration duale de g :

$$\text{smCE}^{(\psi, \lambda)}(g) = \sup_{\eta \in H_\lambda} \left| \mathbb{E} [(y - f(x)) \eta \circ g(x)] \right|$$

Nous pouvons faire le même commentaire que pour la définition du gap de post-processing dual pour cette définition. Cette fois cependant nous apprenons comment d'un prédicteur dual

on construit un prédicteur primal. Préciser que $\nabla\psi(t) \in [0, 1]$ systématiquement permet à ce sous-gradient de faire le lien entre le prédicteur f du monde *réel* à g du monde dual. Avant de continuer les développements, nous devons nous assurer que nous sommes capables de relier la définition classique de la smooth calibration à celle du monde dual.

Proposition 2.5. *Soit $\psi : \mathbb{R} \rightarrow \mathbb{R}$ une fonction convexe différentiable qui vérifie pour tout $t \in \mathbb{R}$, $\nabla\psi(t) \in [0, 1]$. Soit $\lambda > 0$, on suppose que ψ est λ -smooth :*

$$\forall x, y \in \mathbb{R}, \quad |\nabla\psi(x) - \nabla\psi(y)| \leq \lambda|x, y|$$

Pour $g : \mathcal{X} \rightarrow \mathbb{R}$ on définit $f : \mathcal{X} \rightarrow [0, 1]$ telle que $f(x) = \nabla\psi(g(x))$ pour $x \in \mathcal{X}$. On a :

$$\text{smCE}(f) \leq \text{smCE}^{(\psi, \lambda)}(g)$$

Ainsi, avoir une petite valeur de smooth calibration dans le monde dual garantit une plus petite valeur (ou égale) dans le monde réel. Il ne nous reste plus qu'à obtenir une inégalité dans le monde dual que l'on pourra exploiter dans le monde réel.

Théorème 2.3. *Soit $\psi : \mathbb{R} \rightarrow \mathbb{R}$ une fonction convexe différentiable qui vérifie pour tout $t \in \mathbb{R}$, $\nabla\psi(t) \in [0, 1]$. Soit $\lambda > 0$, on suppose que ψ est λ -smooth. Alors pour toute fonction $g : \mathcal{X} \rightarrow \mathbb{R}$ et toute distribution \mathcal{D} sur $\mathcal{X} \times \{0, 1\}$:*

$$\frac{1}{2} \text{smCE}^{(\psi, \lambda)}(g)^2 \leq \lambda \text{pGap}^{(\psi, \lambda)}(g) \leq \text{smCE}^{(\psi, \lambda)}(g)$$

L'écart entre les deux notions que l'on observe est à nouveau au plus quadratique, comme pour dCE et dCE. Il est précisé que les constantes sont optimales, nous ne présenterons pas les exemples qui le montrent.

Maintenant que nous avons une borne dans le monde dual, développons-là dans le monde réel :

Qualité de la calibration

$$\text{smCE}(f)^2 \leq 2\lambda \text{pGap}^{(\psi, \lambda)}(g)$$

Gain d'un post-processing

Nous avons lié la distance à la calibration d'un prédicteur à un gain potentiel de fonction de perte à la suite d'un post-processing par une fonction 1-Lipschitzienne. Ainsi, nous comprenons pourquoi il est difficile pour une régression logistique ou un arbre d'être parfaitement calibré : en plus d'apprendre la structure des données pour bien prédire, elle doit être capable de simuler une fonction 1-Lipschitzienne ! Cela est en revanche plus simple pour un réseau de neurone puisqu'une couche supplémentaire permettrait de le faire. Nous revenons à notre discussion sur le choix de l'algorithme dans la performance de calibration.

Cette remarque met encore plus l'emphasis sur l'importance de la profondeur dans un réseau de neurones ainsi que du rôle de la composition. Mais si l'on pousse la réflexion un peu plus loin, on peut se questionner sur l'impact de la méthode de descente de gradient. S'il suffit d'une nouvelle couche, peut-être suffit-il également d'une mise à jour supplémentaire des gradients par SGD ? A ce stade nous n'avons pas la réponse, mais cela reste une piste de recherche.

Également, nous savons que l'implémentation de la régression logistique dans scikit-learn est accompagnée par défaut d'une régularisation \mathcal{L}_2 . Il est montré dans l'article pour la fonction de perte *loss square* que l'on peut conserver les conclusions précédentes à partir du moment où la famille de prédicteur est fermée pour la composition par les fonctions 1-Lipschitzienne et que la mesure de complexité n'évolue pas trop fortement avec la composition. Ce n'est pas le cas des arbres, SVM ou régressions logistiques par exemple alors que ça l'est pour un réseau de neurones. Ce qui expliquerait pourquoi, même avec de nombreuses régularisations structurelles, ils peuvent être calibrés.

Nous avons présenté dans ce chapitre deux nouveaux exemples où l'intuition faillit à comprendre en profondeur une notion. A travers ce travail théorique autour de la définition de la distance à la calibration, nous avons compris en quoi l'ECE, bien qu'elle semble être la solution idoine, est incomplète. Partant de constat, nous avons proposé une version plus robuste. Nous avons également montré que la relation entre performance et calibration est plus complexe qu'elle n'y paraît. Cependant, comme discuté rapidement, la performance de calibration ne semble pas se limiter à la seule performance d'optimisation.

Cette question reste en suspens, comme les autres questions soulevées sans réponses dans ce cours. Elles, et probablement d'autres encore, nous guideront vers une compréhension plus globale et fine de ce domaine. Cela passera également par la notion de multi-calibration que nous avons rapidement évoquée : être calibré selon plusieurs sous-ensemble du domaine. Cette notion n'est pas traitée dans ce cours, déjà dense et pourtant incomplet, signe d'un domaine de recherche riche en enseignements.

Annexes

Annexe B

Convexité : rappel et dualité avec les fonctions de perte bien définies

B.1 Définition et propriétés

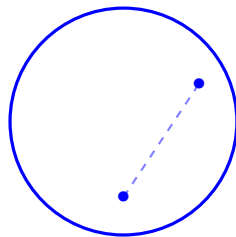
B.1.1 Ensemble et fonction convexe

Dans l'ensemble de la section on travaillera toujours en dimension $d \in \mathbb{N}^*$. Commençons par définir ce que l'on appelle un ensemble convexe :

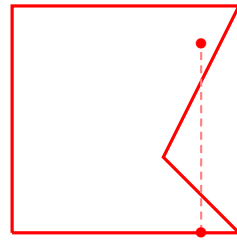
Définition B.1. Soit A un sous-ensemble de \mathbb{R}^d . On dit que A est un ensemble convexe si :

$$\forall a, b \in A, \forall t \in [0, 1], ta + (1 - t)b \in A$$

On appelle donc un ensemble convexe un ensemble dans lequel on peut relier deux points linéairement avec uniquement des points de l'ensemble. On peut illustrer cette définition avec la figure (B.1).



(a) Ensemble convexe



(b) Ensemble non convexe

FIGURE B.1 – Exemple et contre exemple d'ensemble convexe

L'ensemble **bleu** répond parfaitement à l'exigence de la définition : chaque point de l'ensemble est joignable linéairement avec uniquement des points de l'ensemble. L'ensemble **rouge** n'est pas convexe parce qu'entre a et b le *chemin* entre les deux points n'est pas entièrement contenu dans

l'ensemble.

Avec cette notion, nous sommes capable de définir une fonction convexe :

Définition B.2. Soit A un sous-ensemble convexe de \mathbb{R}^d . Une fonction $f : A \mapsto \mathbb{R}$ est convexe si et seulement si :

$$\forall a, b \in A, \forall t \in [0, 1], f(ta + (1 - t)b) \leq tf(a) + (1 - t)f(b)$$

On dira que f est strictement convexe si l'inégalité est stricte.

La définition de fonction convexe ressemble à la définition d'un ensemble convexe. Mais la différence réside dans l'utilisation d'une inégalité ici, et que l'on traite avec les images des points de l'ensemble convexe A . A nouveau, on peut visualiser cette définition avec la figure (B.2).

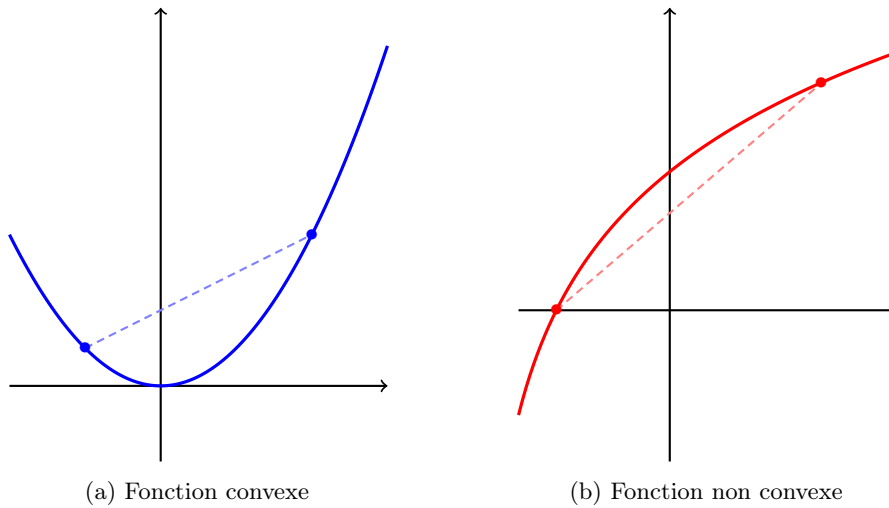


FIGURE B.2 – Exemple et contre exemple de fonction convexe

Exercice B.1. Donner des exemples de fonctions qui répondent aux critères suivants.

1. Une fonction strictement convexe.
2. Une fonction convexe qui n'est pas strictement convexe.
3. Une fonction convexe qui n'est pas strictement convexe ni affine.

Solution. On propose ici une possibilité, mais il y en a bien sur beaucoup plus.

1. La fonction $x \mapsto x^4$ est strictement convexe.
2. N'importe quelle fonction affine est convexe mais pas strictement convexe.
3. La fonction $x \mapsto \max\{0, x\}$ est convexe mais pas strictement convexe ni affine.

□

B.1.2 Caractérisation du premier et deuxième ordre

Il peut parfois être difficile de prouver qu'une fonction est convexe avec la définition que l'on vient de donner. Il nous faudrait une caractérisation plus simple d'utilisation :

Théorème B.1 (Caractérisation des fonctions convexes). *Soit A un sous-ensemble convexe de \mathbb{R}^d et soit $f : A \mapsto \mathbb{R}$ deux fois différentiable. Alors les propriétés suivantes sont équivalentes :*

1. f est convexe
2. $\forall x, y \in A, f(y) \geq f(x) + \nabla f(x) \cdot (y - x)$
3. $\forall x \in A, \nabla^2 f(x) \succeq 0$

On sait déjà ce que la première propriété veut dire, il nous reste à comprendre les deux suivantes.

Dans la deuxième condition, on reconnaît un développement de Taylor à l'ordre 1, ce qui nous dit que chaque tangente de la fonction f est un sous-estimateur global. On peut le visualiser avec deux exemples dans la figure (B.3).

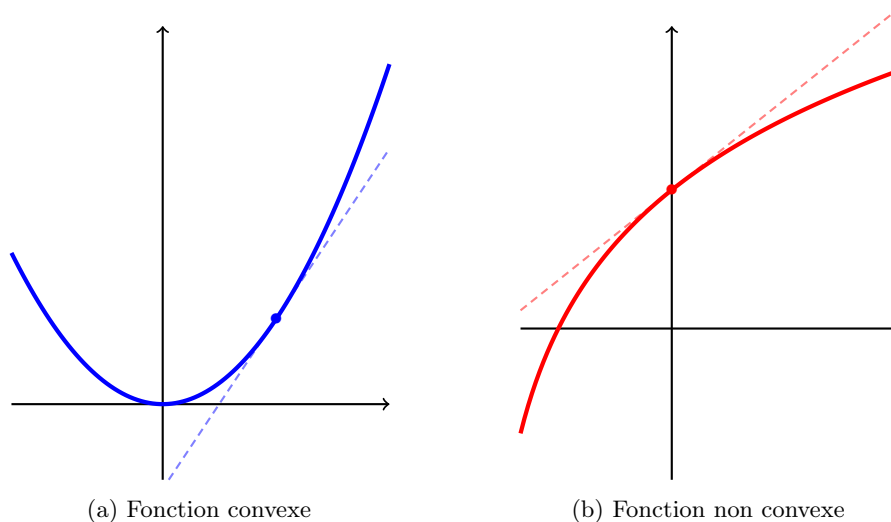


FIGURE B.3 – Illustration de la propriété de sous-estimateur pour une fonction **convexe** et contre exemple pour une fonction **non convexe**

La troisième propriété veut dire qu'il n'y a pas de courbure négative dans la courbe de la fonction f . Autrement dit, que la dérivée de la fonction f est croissante. En dimension une, cela veut dire que la dérivée seconde est toujours positive ou nulle.

Il s'agit maintenant de prouver le théorème (B.1)

Démonstration. Commençons par montrer que si f est convexe, alors $\forall x, y \in A, f(y) \geq f(x) + \nabla f(x) \cdot (y - x)$.

Soit $x, y \in A$, par définition on a que :

$$\begin{aligned}\forall t \in [0, 1], f(tx + (1-t)y) &\leq tf(x) + (1-t)f(y) \\ \forall t \in [0, 1], f(x + t(y-x)) &\leq f(x) + t(f(y) - f(x)) \\ \forall t \in [0, 1], f(y) - f(x) &\geq \frac{f(x + t(y-x)) - f(x)}{t}\end{aligned}$$

Ainsi, en prenant la limite pour $t \downarrow 0$, on a que :

$$\forall x, y \in A, f(y) \geq f(x) + \nabla f(x)(y-x)$$

D'où le résultat souhaité. Montrons maintenant que si on a le résultat précédent, alors f est convexe. Soit $x, y \in A$ et on définit $z = tx + (1-t)y$. Par la propriété que l'on suppose, on a :

$$\begin{aligned}f(x) &\geq f(z) + \nabla f(z)(x-z) \\ f(y) &\geq f(z) + \nabla f(z)(y-z)\end{aligned}$$

En multipliant la première équation par $t \in [0, 1]$ et la seconde équation par $(1-t)$, on obtient :

$$\begin{aligned}tf(x) + (1-t)f(y) &\geq f(z) + \nabla f(z)(tx + (1-t)y - z) \\ &= f(z) \\ f(tx + (1-t)y) &\leq tf(x) + (1-t)f(y)\end{aligned}$$

On a donc montré que les deux premières propositions sont équivalentes. Montrons maintenant que les deux dernières propriétés sont équivalentes en dimension 1 pour simplifier les calculs.

Supposons que $\forall x \in A, f''(x) \geq 0$, alors par le théorème de la valeur moyenne de Taylor on a que :

$$\begin{aligned}\exists z \in [x, y], f(y) &= f(x) + f'(x)(y-x) + \frac{1}{2}f''(z)(y-x)^2 \\ \Rightarrow f(y) &\geq f(x) + f'(x)(y-x)\end{aligned}$$

Finalement, supposons que $\forall x, y \in A, f(y) \geq f(x) + \nabla f(x)(y-x)$. Soit $x, y \in A$ tels que $y > x$. On a que :

$$\begin{aligned}f(y) &\geq f(x) + f'(x)(y-x) \\ f(x) &\geq f(y) + f'(y)(x-y)\end{aligned}$$

On en déduit donc que :

$$f'(x)(y-x) \leq f(y) - f(x) \leq f'(y)(y-x)$$

Donc en divisant par $(y-x)^2$ puis en prenant la limite pour $y \rightarrow x$, on obtient bien que la propriété souhaité. \square

Ce résultat conclut notre section de présentation des notions de convexité. Intéressons-nous maintenant à l'utilisation de cette notion pour l'optimisation.

B.2 Résultats d'optimisations

On considère un problème d'optimisation sans contraintes avec une fonction $f : \mathbb{R}^d \mapsto \mathbb{R}$ différentiable :

$$x^* = \arg \min_{x \in \mathbb{R}^d} f(x) \quad (\text{B.1})$$

Notons qu'ici nous n'avons pas encore spécifié que f est convexe. Dans le cadre général, on sait qu'une condition nécessaire pour que x soit une solution de ce problème est que $\nabla f(x) = 0$. Mais ce n'est pas une condition suffisante ! De plus, si cette condition nécessaire est vraie, et qu'il s'agit d'un minimum, alors on ne peut pas dire plus que " x est un minimum local" avec ces informations.

Exercice B.2. Donner un exemple de fonction qui répond à chaque critère :

1. Une fonction où il existe $x \in \mathbb{R}$ tel que $f'(x) = 0$ et que x est un minimum local.
2. Une fonction où il existe $x \in \mathbb{R}$ tel que $f'(x) = 0$ mais que x n'est ni un minimum local ni un maximum local.
3. Une fonction où il existe une infinité de $x \in \mathbb{R}$ tel que $f'(x) = 0$ qui sont tous des minimum locaux.

Solution. On propose ici une possibilité, mais il y en a bien sûr beaucoup plus.

1. La fonction $x \mapsto x^3 - x$ possède un minimum local (et un maximum local).
2. La fonction $x \mapsto x^3$ en $x = 0$.
3. La fonction $x \mapsto \cos(x)$ contient une infinité de point x tel que $f'(x) = 0$ (tous espacés de π) mais seulement *la moitié* sont des minimums.

□

On voit donc que nous n'avons pas de critères simples et clairs dans le cas général sur l'existence et l'unicité d'un minimum global. Voyons ce qu'il en est quand la fonction f est convexe.

Proposition B.1. Soit un problème d'optimisation sans contraintes comme présenté dans (B.1) avec f une fonction convexe et différentiable. Alors, chaque point x qui vérifie $\nabla f(x) = 0$ est un minimum global.

Pour une fonction convexe différentiable, la condition $\nabla f(x) = 0$ est une condition nécessaire et suffisante pour caractériser un minimum global.

Exercice B.3. Prouver la proposition (B.1) à l'aide du théorème (B.1).

Solution. Comme $f : A \mapsto \mathbb{R}$ est convexe et différentiable, d'après le théorème (B.1) on a :

$$\forall x, y \in A, f(y) \geq f(x) + \nabla f(x)(y - x)$$

Donc en particulier pour \bar{x} défini comme $\nabla f(\bar{x}) = 0$:

$$\begin{aligned}\forall y \in A, f(y) &\geq f(\bar{x}) + \nabla f(\bar{x})(y - \bar{x}) \\ \forall y \in A, f(y) &\geq f(\bar{x})\end{aligned}$$

Qui est bien la définition d'un minimum global d'une fonction. \square

Mais nous n'avons toujours pas l'unicité d'un minimum à ce stade. Pour l'obtenir, nous avons besoin d'avoir la stricte convexité.

Proposition B.2. *Soit un problème d'optimisation sans contraintes comme présenté dans (B.1) avec f une fonction strictement convexe et différentiable. Si $\nabla f(\bar{x}) = 0$, alors \bar{x} est l'unique minimum global de f .*

Nous obtenons cette fois l'unicité à l'aide de la stricte convexité. Voyons comment.

Exercice B.4. *Prouver la proposition (B.2) en raisonnant par l'absurde. On suppose donc qu'il existe deux minimaux globaux et on aboutit à une absurdité en exploitant la stricte convexité.*

Solution. On suppose qu'il existe $a, b \in A$ qui minimisent f tel quel que $a \neq b$. Prenons $z = \frac{a+b}{2} \in A$ comme combinaison convexe de deux points du domaine (avec $t = \frac{1}{2}$). Alors :

$$f(z) < \frac{1}{2}f(a) + \frac{1}{2}f(b) = f(a) = f(b)$$

On vient de trouver un nouveau nombre z qui minimise encore mieux la fonction f que les deux meilleurs minimiseurs : absurde, d'où l'unicité. \square

Avec ces deux derniers résultats, nous comprenons pourquoi il est important de travailler avec des fonctions de perte convexes : elles nous garantissent qu'en suivant une descente de gradient, nous atteindrons bien un minimum global. Voyons à présent à quelle vitesse.

B.3 Vitesse de convergence pour la descente de gradient

Dans cette section nous allons donner des preuves de vitesse de convergence pour deux hypothèses sur la fonction f .

B.3.1 Pour une fonction Lipschitzienne

La plus petite supposition qui nous permet de garantir la convergence vers le minimum global est que la fonction soit Lipschitzienne.

Définition B.3 (Fonction L -Lipschitzienne). *Soit $f : \mathcal{D} \rightarrow \mathbb{R}$ une fonction convexe. On dit que f est L -Lipschitzienne si il existe un nombre L tel que :*

$$|f(y) - f(x)| \leq L\|y - x\|$$

Pour comprendre visuellement cette définition, on peut s'appuyer sur la figure (B.4).

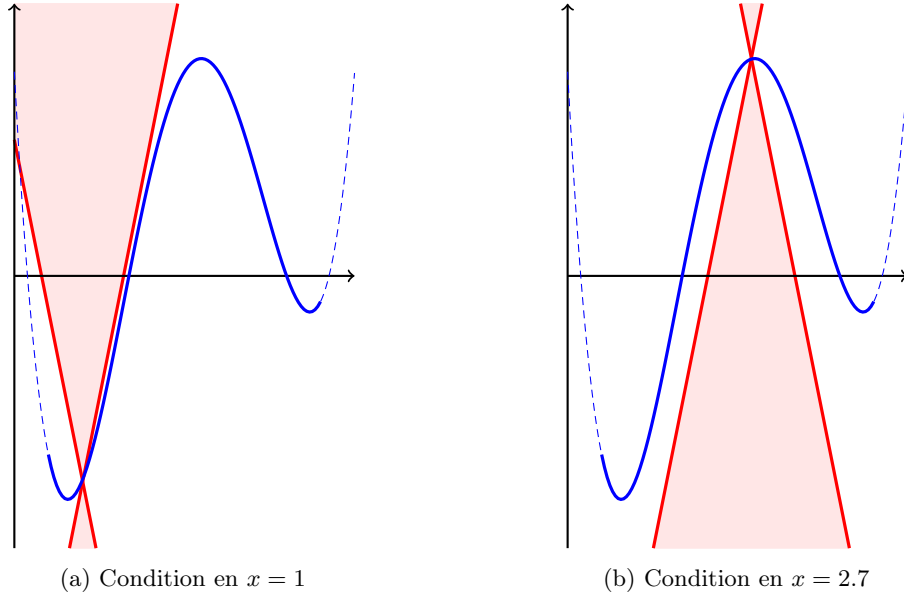


FIGURE B.4 – Illustration de la condition lipschitzienne pour une fonction

Dire qu'une fonction est L -lipschitzienne revient à dire que ses variations seront toujours hors des cônes rouge, autrement dit on contraint la progression de la fonction. Ici, on remarque que la fonction n'est pas lipschitzienne pour ce L -là sur $[0, 5]$ parce que la fonction passe dans les cônes rouge. En revanche, elle l'est pour l'intervalle $[0.5, 4.5]$. Remarquons également que nous avons une fonction qui peut être lipschitzienne sans être convexe. Supposer les deux, permet d'avoir le résultat suivant.

Proposition B.3. Soit $x^* \in \mathbb{R}^d$ le minimum de la fonction $f : \mathbb{R}^d \mapsto \mathbb{R}$ convexe et L -Lipschitzienne. Soit $\varepsilon > 0$ l'erreur que l'on accorde pour arrêter la descente de gradient. On choisit $\eta = \frac{\varepsilon}{L^2}$. Alors en $T = \frac{L^2 \|x^* - x_0\|^2}{\varepsilon^2}$ itérations, l'algorithme converge vers x^* avec une erreur de ε .

Nous avons donc la garantie que nous sommes capables de converger vers le minimum de la fonction f avec un nombre d'itérations que l'on maîtrise. Si l'on note $\tilde{x} \in \mathbb{R}^d$ le point que l'on obtient à la suite de la descente de gradient, on a $f(\tilde{x}) \leq f(x^*) + \varepsilon$.

Pour faire la preuve de ce résultat, nous avons besoin d'un résultat intermédiaire.

Lemme B.1. Pour $(x_t)_{t \in \mathbb{N}}$ la suite définie pour une descente de gradient qui cherche à minimiser une fonction f convexe et L -Lipschitzienne, on a :

$$\|x_{t+1} - x^*\|^2 \leq \|x_t - x^*\|^2 - 2\eta(f(x_t) - f(x^*)) + \eta^2 L^2$$

Exercice B.5. Prouver le lemme (B.1).

Solution.

$$\begin{aligned}
\|x_{t+1} - x^*\|^2 &= \|x_t - x^* - \eta \nabla f(x)\|^2 \\
&= \|x_t - x^*\|^2 - 2\eta (\nabla f(x))^t (x_t - x^*) + \eta^2 \|\nabla f(x_t)\|^2 \\
&\leq \|x_t - x^*\|^2 - 2\eta (f(x_t) - f(x^*)) + \eta^2 L^2
\end{aligned}$$

En exploitant le fait que f soit convexe et L -Lipschizienne pour la dernière inégalité. \square

Nous avons donc maintenant tous les outils pour démontrer le résultat annoncé dans la proposition (B.3).

Démonstration. Soit $\Phi(t) = \|x_t - x^*\|^2$. Alors avec $\eta = \frac{\varepsilon}{L^2}$ et le précédent lemme on a :

$$\Phi(t) - \Phi(t+1) > 2\eta\varepsilon - \eta^2 L^2 = \frac{\varepsilon}{L^2}$$

Puisque par définition, $\Phi(0) = \|x^* - x_0\|^2$ et $\Phi(t) \geq 0$, la précédente équation ne peut pas être vérifiée pour tous $0 \leq t \leq \frac{L^2 \|x^* - x_0\|^2}{\varepsilon^2}$ d'où le résultat. \square

B.3.2 Fonction β -smooth

Une manière d'accélérer cette convergence est de faire plus d'hypothèse sur la fonction en demandant qu'elle soit β -smooth.

Définition B.4 (Fonction β -smooth). Soit $f : \mathcal{D} \mapsto \mathbb{R}$ une fonction convexe différentiable. On dit que f est β -smooth si :

$$\forall x, y \in \mathcal{D}, f(y) \leq f(x) + \langle \nabla f(x), y - x \rangle + \frac{\beta}{2} \|y - x\|^2$$

Avec l'exercice suivant, on peut découvrir une interprétation de cette nouvelle notion.

Exercice B.6 (Caractérisation d'une fonction β -smooth). Montrer que f est une fonction β -smooth si, et seulement si, le gradient de f est β -Lipschitz.

Commençons par nous assurer que l'on aura bien toujours une *descente* :

$$\begin{aligned}
f(x_{t+1}) &\leq f(x_t) - \langle \nabla f(x_t), x_{t+1} - x_t \rangle + \frac{\beta}{2} \|x_{t+1} - x_t\|^2 \\
&\leq f(x_t) - \eta_t \langle \nabla f(x_t), \nabla f(x_t) \rangle + \frac{\beta}{2} \eta_t^2 \|\nabla f(x_t)\|^2 \\
&\leq f(x_t) - \left(\eta_t - \frac{\beta \eta_t^2}{2} \right) \|\nabla f(x_t)\|^2
\end{aligned}$$

Ainsi, si $\eta_t \leq \frac{1}{\beta}$ on a bien une descente ! Voyons comment, avec cette condition supplémentaire, nous sommes en capacité d'avoir une convergence plus rapide.

Proposition B.4. Soit $x^* \in \mathbb{R}^d$ le minimum de la fonction $f : \mathbb{R}^d \mapsto \mathbb{R}$ une fonction convexe et β -smooth. On choisit $\eta = \frac{1}{\beta}$.

Alors pour toutes itérations $t \leq T$ avec $T \geq 1$, on a :

$$f(x_t) - f(x^*) \leq \mathcal{O}\left(\frac{1}{T}\right)$$

Nous ne prouverons pas ici le résultat et renvoyons vers le très complet livre de Sébastien Bubeck *Convex optimization : Algorithms and complexity* publié en 2015, pour avoir le détail de la preuve.

Dans ce même article, nous pouvons trouver d'autres hypothèses, comme par exemple que f soit fortement convexe, pour accélérer encore la convergence.

Connaître ces propriétés permet parfois de choisir une fonction plutôt qu'une autre quand elles remplissent le même rôle dans l'entraînement d'un modèle. C'est particulièrement vrai dans l'entraînement des réseaux de neurones, que nous ne traitons pas ici, où un très grand nombre de paramètres est à apprendre. Avoir une vitesse de convergence plus rapide parce que nous avons choisi une fonction à optimiser la plus régulière possible permet de faire gagner parfois des jours voire semaines de calculs. C'est ce qui explique le très grand nombre de différentes méthodes de descente de gradient qui ont été développées.

B.4 La fonction de perte des réseaux de neurones n'est pas en général convexe

L'ensemble de cette annexe est à propos des fonctions convexe, et nous avons en particulier traité de la vitesse de convergence pour des objectifs avec des fonctions convexe. Dans les réseaux de neurones, la descente de gradient¹ est clé pour l'algorithme de back-propagation. Ainsi, nous pourrions tirer des sections de cette annexes des garanties de vitesses de convergence pour des réseaux de neurones.

Pour cela, nous devons nous assurer que la fonction de perte d'un réseau de neurone est convexe. Cela semble plausible puisque chaque neurones est une combinaison affine de ses inputs, et que la plupart des fonctions d'activation (par exemple ReLU) sont convexes.

Exercice B.7 (Combinaison de fonctions convexes). Soient $U, V \subseteq \mathbb{R}$ et $f : U \rightarrow V$ et $g : V \rightarrow \mathbb{R}$ deux fonctions convexes.

1. Montrer que si g est croissante, alors $g \circ f$ est une fonction convexe.
2. Montrer à l'aide d'un contre-exemple que si g est strictement décroissante, alors $g \circ f$ n'est pas une fonction convexe.

Solution. On reprend les notations de l'exercice.

1. Plus précisément ses variantes en mini-batch, avec momentum ou pas adaptatif.

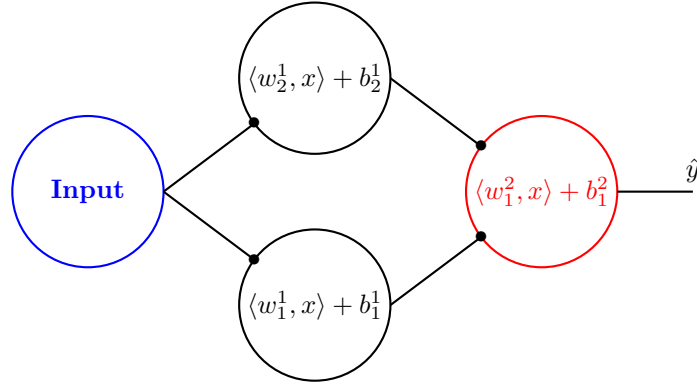


FIGURE B.5 – Réseau de neurone avec une couche caché de deux neurones, fonction d'activation ReLU

1. Soit $x, y \in U$ et $\lambda \in [0, 1]$. Puisque f est convexe, on a :

$$\begin{aligned}
 f(\lambda x + (1 - \lambda)y) &\leq \lambda f(x) + (1 - \lambda)f(y) \quad \text{par définition} \\
 (g \circ f)(\lambda x + (1 - \lambda)y) &\leq g(\lambda f(x) + (1 - \lambda)f(y)) \quad \text{car } g \text{ est croissante} \\
 &\leq \lambda(g \circ f)(x) + (1 - \lambda)(g \circ f)(y) \quad \text{car } g \text{ est convexe}
 \end{aligned}$$

D'où $g \circ f$ est une fonction convexe.

2. Prenons $f(x) = x^2$ et $g(x) = e^{-x}$. Il s'agit bien de deux fonctions convexe puisque pour tout $x \in \mathbb{R}$, on a $f''(x) \geq 0$ et $g''(x) \geq 0$ avec le théorème B.1 de caractérisation des fonctions convexes. g est également strictement décroissante. Ici, $g \circ f(x) = e^{-x^2}$ et cette fonction n'est clairement pas convexe.

□

Avec l'exercice précédent, tout porte à croire que l'on a des chances d'avoir un réseau de neurones qui ait une fonction de perte convexe.

Exercice B.8. On considère le réseau de neurones défini à la figure B.5. Pour simplifier les calculs, on fixe les valeurs de l'ensemble des biais $b_1^1 = b_2^1 = b_1^2 = 0$.

On considère le dataset $\mathcal{D} = \{(-1, -1), (1, 1)\}$ et la fonction de perte $\mathcal{L}(\theta) = \frac{1}{2} \sum_{i=1}^2 (\hat{y}_i - y_i)^2$.

1. Montrer que pour $\theta_1 : w_1^1 = 1, w_2^1 = -1, w_1^2 = (1, -1)$ on a $\mathcal{L}(\theta_1) = 0$
2. Montrer que pour $\theta_2 : w_1^1 = -1, w_2^1 = 1, w_1^2 = (-1, 1)$ on a $\mathcal{L}(\theta_2) = 0$
3. Conclure sur la nature convexe ou non de la fonction de perte.

Solution. Les deux premières questions sont immédiates. Si la loss est convexe, alors n'importe quelle combinaison convexe de ces paramètres donnera une loss inférieure ou égale à 0 (puisque $\mathcal{L}(\theta_1) = \mathcal{L}(\theta_2) = 0$).

En particulier, la moyenne arithmétique de θ_1 et θ_2 donne : $w_1^1 = 0, w_2^1 = 0, w_1^2 = (0, 0)$. On obtient alors $\mathcal{L}\left(\frac{\theta_1 + \theta_2}{2}\right) = 1 > 0$. Donc la fonction de perte n'est pas convexe. \square

Nous avons donc qu'avec un réseau de neurones aussi simple que le réseau de la figure B.5 la fonction de perte n'est pas convexe. Il est ainsi probable que pour des architectures de réseaux de neurones plus conséquent en taille, la fonction de perte soit encore plus non-convexe. Nous ne pouvons donc plus nous appuyer sur les résultats théoriques obtenu pour la descente de gradient dans cette annexe.

Annexe C

Dualité entre les fonctions convexes et les fonctions de pertes bien définies

La notion de dérivée est largement utilisé dans la partie précédente, pour caractériser une fonction convexe par exemple. Ce concept est largement utilisé en mathématiques plus généralement pour être capable de quantifier le comportement locale d'une fonction. Cependant, une fonction convexe n'est pas forcément différentiable partout. Il nous faudrait un outil plus général que la dérivation mais qui permette d'atteindre un but similaire. C'est la notion de sous-différentiel.

Définition C.1 (Sous-différentiel). Soit $\varphi : V \rightarrow \mathbb{R}$ une fonction définie pour $V \subseteq \mathbb{R}$. On dit que $t \in \mathbb{R}$ est un sous-gradient de φ en $v \in V$ si :

$$\forall v' \in V, \quad \varphi(v') \geq \varphi(v) + (v' - v)t$$

Ou de manière équivalente :

$$vt - \phi(v) = \max_{v' \in V} v't - \varphi(v')$$

On appelle l'ensemble des sous-gradients le sous-différentiel.

Pour comprendre la définition, considérons la fonction valeur absolue. Elle est convexe, mais pas différentiable en 0.

Avec la figure C.1 on a l'exemple de deux sous-gradient qui appartiennent au sous-différentiel. On a aussi un contre-exemple : la droite tracé n'est pas toujours une minorante comme demandé dans la première définition. Ainsi, le sous-différentiel de la valeur absolue en 0 est l'intervalle $[-1, 1]$.

On comprend que si le sous-différentiel est réduit à un singleton, alors c'est que la fonction considéré est différentiable en ce point et que la valeur du singleton est la dérivée en ce point. De même, si zéro appartient au différentiel d'un point, alors c'est que la fonction atteint son minimum en ce point. Nous avons donc bien une généralisation de la notion de différentiabilité !

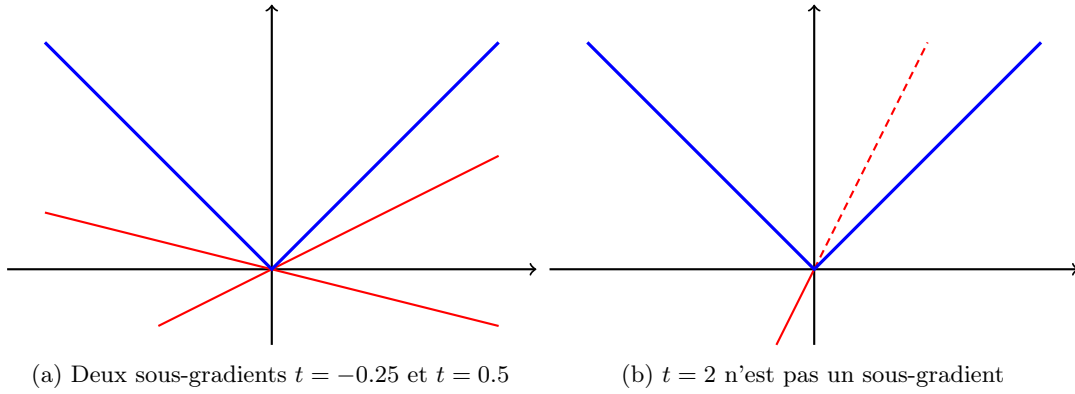


FIGURE C.1 – Fonction $x \mapsto |x|$ avec exemple et contre-exemple de sous-gradient

Définition C.2 (Divergence de Fenchel-Young). Soit $V, T \subseteq \mathbb{R}$. On définit la divergence de Fenchel-Young $D_{\varphi, \psi} : V \times T \rightarrow \mathbb{R}$ pour une paire de fonction (φ, ψ) par :

$$\forall v \in V, \forall t \in T, \quad D_{\varphi, \psi}(v, t) = \overbrace{\varphi}^{\varphi : V \rightarrow \mathbb{R}}(v) + \underbrace{\psi(t)}_{\psi : T \rightarrow \mathbb{R}} - vt$$

Considérons une fonction $\varphi : V \rightarrow \mathbb{R}$ et $\psi : T \rightarrow \mathbb{R}$ définie avec $V, T \subseteq \mathbb{R}$. Essayons de lier les deux précédentes notions en commençant par supposer que $t \in T$ est un sous-gradient de φ en $v \in V$:

$$\begin{aligned} vt - \varphi(v) = \max_{v' \in V} (v't - \varphi(v')) &\iff vt - \varphi(v) - \psi(t) = \max_{v' \in V} (v't - \varphi(v') - \psi(t)) \\ &\iff \varphi(v) + \psi(t) - vt = \min_{v' \in V} (v't - \varphi(v') - \psi(t)) \\ &\iff D_{\varphi, \psi}(v, t) = \min_{v' \in V} D_{\varphi, \psi}(v', t) \end{aligned}$$

Ainsi, $t \in T$ est un sous-gradient de φ en $v \in V$ si, et seulement si, $D_{\varphi, \psi}(v, t) = \min_{v' \in V} D_{\varphi, \psi}(v', t)$.

Exercice C.1. En s'inspirant de la démonstration précédente, montrer que v est un sous-gradient de ψ en $t \in T$ si, et seulement si, $D_{\varphi, \psi}(v, t) = \min_{t' \in T} D_{\varphi, \psi}(v, t')$

Avec le résultat précédent et celui de l'exercice, on comprend que si l'on suppose que $D_{\varphi, \psi}(v', t') \geq 0$ pour tout $v' \in V$ et $t' \in T$ mais qu'il existe $v \in V$ et $t \in T$ tels que $D_{\varphi, \psi}(v, t) = 0$ alors t est un sous-gradient de φ en v et v est un sous-gradient de ψ en t .

On saisit l'intérêt de l'introduction de la divergence de Fenchel-Young : elle permet de lier la notion de sous-gradient d'une paire de fonction. C'est ce lien que l'on va exploiter pour montrer comment on peut induire une fonction convexe depuis une fonction de perte bien définie.

Proposition C.1. Soit $V \subseteq [0, 1]$ un intervalle non vide. Soit $\mathcal{L} : \{0, 1\} \times V \mapsto \mathbb{R}$ une fonction de perte bien définie. Pour tout $v \in V$ on définit $\text{dual}(v) = \mathcal{L}(0, v) - \mathcal{L}(1, v)$. Alors, il existe une fonction convexe $\varphi : V \rightarrow \mathbb{R}$ et une fonction convexe $\psi : \mathbb{R} \rightarrow \mathbb{R}$ telles que :

$$\begin{aligned} \forall y \in \{0, 1\}, \forall v \in V, & \quad \mathcal{L}(y, v) = \psi(\text{dual}(v)) - y \text{dual}(v) \\ \forall v \in V, \forall t \in \mathbb{R}, & \quad D_{\varphi, \psi}(v, t) \geq 0 \\ \forall v \in V, & \quad D_{\varphi, \psi}(v, \text{dual}(v)) = 0 \\ \forall (t_1, t_2) \in \mathbb{R}, t_1 \neq t_2, & \quad 0 \leq \frac{\psi(t_2) - \psi(t_1)}{t_2 - t_1} \leq 1 \end{aligned}$$

La démonstration de ce résultat montre comment sont construite les fonctions φ et ψ . Elle s'appuie également sur une réécriture de la fonction de perte à l'aide de la fonction duale. Nous proposons de nous intéresser à une autre manière décrire la fonction de perte \mathcal{L} avec les deux fonctions induites φ et ψ .

Exercice C.2 (Réécriture d'une fonction de perte bien définie). Montrer, avec les mêmes hypothèses et notation que la proposition (C.1), que l'on peut écrire la fonction de perte \mathcal{L} bien définie comme :

$$\forall y \in \{0, 1\}, \forall v \in V, \quad \mathcal{L}(y, v) = -\varphi(v) + (v - y) \text{dual}(v)$$

Solution. Soit $v \in V$. Avec la proposition (C.1), et les mêmes notations, on a :

$$\begin{aligned} D_{\varphi, \psi}(v, \text{dual}(v)) &= 0 \\ \varphi(v) + \psi(\text{dual}(v)) - v \text{dual}(v) &= 0 \end{aligned}$$

Mais on sait aussi que pour tout $y \in \{0, 1\}$:

$$\mathcal{L}(y, v) = \psi(\text{dual}(v)) - y \text{dual}(v) \iff \psi(\text{dual}(v)) = \mathcal{L}(y, v) + y \text{dual}(v)$$

Finalement, on a :

$$\begin{aligned} \varphi(v) + (\mathcal{L}(y, v) + y \text{dual}(v)) - v \text{dual}(v) &= 0 \iff \varphi(v) + \mathcal{L}(y, v) + (y - v) \text{dual}(v) = 0 \\ &\iff \mathcal{L}(y, v) = -\varphi(v) + (v - y) \text{dual}(v) \end{aligned}$$

□

Il nous reste à voir comment à partir de fonctions convexes on peut induire une fonction de perte bien définie.

Proposition C.2. Soit $V \subseteq [0, 1]$ un interval non vide. Soit $\varphi : V \rightarrow \mathbb{R}$, $\psi : \mathbb{R} \rightarrow \mathbb{R}$ et $\text{dual} : V \rightarrow \mathbb{R}$ des fonctions telles que :

$$\begin{aligned}\forall v \in V, \forall t \in \mathbb{R}, D_{\varphi, \psi}(v, t) &\geq 0 \\ \forall v \in V, D_{\varphi, \psi}(v, \text{dual}(v)) &= 0\end{aligned}$$

On définit la fonction de perte $\mathcal{L} : \{0, 1\} \times V \rightarrow \mathbb{R}$ qui vérifie

$$\forall y \in \{0, 1\}, \forall v \in V, \quad \mathcal{L}(y, v) = \psi(\text{dual}(v)) - y \text{dual}(v)$$

Alors, on a :

1. \mathcal{L} est une fonction de perte bien définie
2. φ est convexe
3. Pour tout $v \in V$, $\text{dual}(v)$ est un sous-gradient de φ en v , et v est un sous-gradient de ψ en $\text{dual}(v)$
4. $\varphi(v) = v \text{dual}(v) - \psi(\text{dual}(v)) = - \mathbb{E}_{y \sim \text{Ber}(v)} [\mathcal{L}(y, v)]$
5. Si on définit $\psi'(t) = \sup_{v \in V} vt - \varphi(v)$ alors ψ' est convexe et $\psi'(\text{dual}(v)) = \psi(\text{dual}(v))$ pour tout $v \in V$. De plus, pour t_1, t_2 deux réels distincts on a :

$$0 \leq \frac{\psi(t_2) - \psi(t_1)}{t_2 - t_1} \leq 1$$

De ces deux propositions nous voyons la correspondance entre une fonction de perte bien définie $\mathcal{L} : \{0, 1\} \times V \rightarrow \mathbb{R}$ et $(\varphi, \psi, \text{dual})$ qui vérifient $D_{\varphi, \psi}(v, t) \geq 0$ et $D_{\varphi, \psi}(v, \text{dual}(v)) = 0$. D'après la proposition (C.2) en supposant que ψ est différentiable, on a que $v = \nabla \psi(\text{dual}(v))$.

Ceci conclut le lien qui a permis de mieux comprendre le lien entre performance d'optimisation et qualité de calibration d'un algorithme de Machine Learning.

Annexe D

Fléau de la dimension

Georg Cantor est un mathématicien allemand du 19e siècle qui est particulièrement connu pour son travail sur la théorie des ensembles et plus spécifiquement sur ses résultats concernant l'infini. L'ensemble des nombres entiers est infini par construction, mais l'ensemble des nombres entiers relatifs également. Et intuitivement, nous nous disons que ces infinis ne sont pas vraiment les mêmes, puisqu'il semblerait que \mathbb{Z} soit plus grand que \mathbb{N} ! Georg Cantor montre que ces deux infinis sont en fait les mêmes : il y a autant de nombres dans \mathbb{Z} que dans \mathbb{N} . Plus fort encore, il démontre la puissance de l'infini qui nous donne un résultat encore plus contre intuitif : il y a autant de nombres dans l'intervalle $[0, 1]$ que dans \mathbb{R} tout entier ! Alors qu'il venait de le démontrer, il a envoyé une lettre à son ami mathématicien Dedekind :

Tant que vous ne m'aurez pas approuvé, je ne puis que dire : je le vois mais je ne le crois pas.

— Georg Cantor (1877)

Il a prouvé quelque chose que la communauté pensait intuitivement fausse, et lui-même n'y croyait pas. Nous sommes toujours mis en difficulté quand il s'agit de traiter avec l'infini, ou des grandes quantités. Cette remarque nous amène donc à nous questionner sur l'impact d'un grand nombre d'informations quand nous entraînons un modèle de Machine Learning.

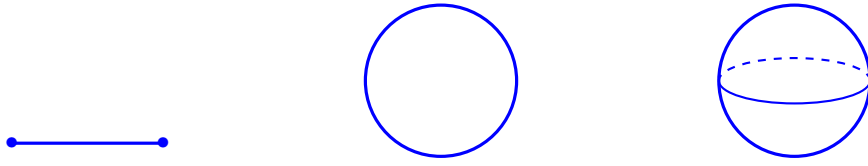
Le **fléau de la dimension** est une notion connue en statistiques et en Machine Learning. Ce terme rassemble tout un ensemble de phénomènes qui se produit en très grande dimension, mais pas dans une dimension plus petite. Nous proposons dans cette annexe d'illustrer quelques-uns des phénomènes étranges de la grande dimension et ses impacts en Machine Learning.

D.1 Volume d'une hypersphère

Pour essayer de sentir les problèmes de la très grande dimension, on s'intéresse au volume d'une hypersphère.

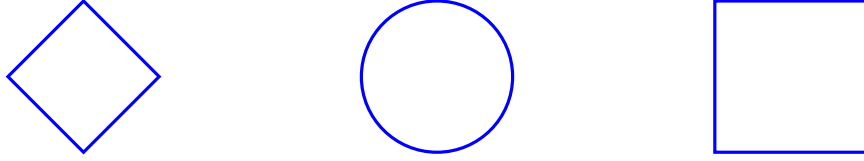
Avec la figure (D.1) nous avons l'intuition que le volume augmente avec la dimension. Donc pour une hypersphère de très grande dimension, on devrait avoir un très grand volume. Nous avons tracé et mesuré le volume pour la distance euclidienne classique, mais nous pouvons aussi utiliser d'autres distances (autre norme) comme montré dans la figure (D.2).

Formalisons le problème et généralisons-le pour calculer le volume d'une hypersphère en n'importe quelle dimension et pour n'importe quelle norme. Soit $n \in \mathbb{N}^*$ la dimension de l'espace, on appelle *boule* ou hypersphère l'objet défini par :



(a) En dimension 1, volume = 2 (b) En dimension 2, volume = π (c) En dimension 3, volume = $\frac{4}{3}\pi$

FIGURE D.1 – Représentation et volume d’une hypersphère de rayon 1 dans 3 espaces de dimensions différentes



(a) Avec la norme 1

(b) Avec la norme 2

(c) Avec la norme infinie

FIGURE D.2 – Représentation d’une hypersphère de rayon 1 en dimension 2 pour 3 normes différentes

$$\begin{aligned} B_n^p(R) &= \{(x_1, \dots, x_n) \in \mathbb{R}^n, \sum_{i=1}^n x_i^p \leq R^p\} \\ &= \{u \in \mathbb{R}^n, \|u\|_p^p \leq R^p\} \end{aligned}$$

Avec $\|u\|_p$ la norme p définie comme $\|u\|_p^p = \sum_{i=1}^n x_i^p$. $B_n^p(R)$ est la boule de dimension n avec une p -norme de rayon R . On définit $V_n^p(R)$ le volume de la boule $B_n^p(R)$ *i.e.* la mesure de $B_n^p(R)$ pour la mesure de Lebesgue dans \mathbb{R}^n . Formellement :

$$V_n^p(R) = \int_{B_n^p(R)} \bigotimes_{i=1}^n dx_i$$

Proposition D.1 (Volume d'une hypersphère). *Avec les notations précédentes, on a :*

$$\forall R > 0, \forall n \geq 2, \forall p \geq 1, \quad V_n^p(R) = \frac{\left(2R\Gamma\left(\frac{1}{p} + 1\right)\right)^n}{\Gamma\left(\frac{n}{p} + 1\right)}$$

Et son équivalent quand n tend vers l'infini :

$$V_n^p(R) \sim \sqrt{\frac{p}{2\pi n}} \left[2R\Gamma\left(\frac{1}{p} + 1\right) \left(\frac{pe}{n}\right)^{\frac{1}{p}}\right]^n$$

Avec la fonction Γ définie comme :

$$\Gamma(x) = \int_0^{+\infty} e^{-t} t^{x-1} dt$$

Démonstration. Soit la fonction ϕ définie comme :

$$\begin{aligned} \phi : \quad B_n^p(1) &\rightarrow B_n^p(R) \\ (x_1, \dots, x_n) &\mapsto (Rx_1, \dots, Rx_n) \end{aligned}$$

Par définition, si $u \in B_n^p(1)$, alors $\|u\|_p^p \leq 1 \iff \|Ru\|_p^p \leq R^p$ donc $\phi(u) \in B_n^p(R)$, en supposant que $R > 0$. Autrement dit, on ne considère pas le cas dégénéré. Par équivalence, on a que $\phi(B_n^p(1)) = \phi(B_n^p(R))$. ϕ est donc une bijection, ainsi par changement de variable : $V_n^p(R) = R^n V_n^p(1)$.

Par le théorème de Fubini :

$$\begin{aligned} V_n^p(1) &= \int_{-1}^1 V_{n-1}^p\left(\sqrt[p]{1-|x|^p}\right) dx \\ &= V_{n-1}^p(1) \int_{-1}^1 (1-|x|^p)^{\frac{n-1}{p}} dx \\ &= 2V_{n-1}^p(1) \int_0^1 (1-x^p)^{\frac{n-1}{p}} dx \\ &= \frac{2}{p} V_{n-1}^p(1) \int_0^1 y^{\frac{1}{p}-1} (1-y)^{\frac{n-1}{p}} dy \end{aligned}$$

Rappelons la définition de la fonction Beta, définie pour $x, y \in \mathbb{R}_+^*$:

$$B(x, y) = \int_0^1 t^{x-1} (1-t)^{y-1} dt$$

La fonction Beta est liée à la fonction Gamma par l'identité suivante :

$$B(x, y) = \frac{\Gamma(x)\Gamma(y)}{\Gamma(x+y)}$$

On peut alors écrire :

$$\begin{aligned}
V_n^p(1) &= \frac{2}{p} V_{n-1}^p(1) \frac{\Gamma\left(\frac{1}{p}\right) \Gamma\left(\frac{n-1}{p} + 1\right)}{\Gamma\left(\frac{n}{p} + 1\right)} \\
&= V_{n-1}^p(1) \left[\frac{2}{p} \frac{\Gamma\left(\frac{1}{p}\right) \Gamma\left(\frac{n-1}{p} + 1\right)}{\Gamma\left(\frac{n}{p} + 1\right)} \right] \\
&= V_1^p(1) 2^{n-1} \Gamma\left(\frac{1}{p} + 1\right)^{n-1} \frac{\Gamma\left(\frac{1}{p} + 1\right)}{\Gamma\left(\frac{n}{p} + 1\right)} \\
&= V_1^p(1) 2^{n-1} \frac{\Gamma\left(\frac{1}{p} + 1\right)^n}{\Gamma\left(\frac{n}{p} + 1\right)}
\end{aligned}$$

Puisque $V_1^p(1) = 2$, on obtient finalement :

$$\forall n \geq 2, \forall p \geq 1, \quad V_n^p(1) = \frac{\left(2\Gamma\left(\frac{1}{p} + 1\right)\right)^n}{\Gamma\left(\frac{n}{p} + 1\right)}$$

Avec un rayon $R > 0$, plus généralement :

$$\forall R > 0, \forall n \geq 2, \forall p \geq 1, \quad V_n^p(R) = \frac{\left(2R\Gamma\left(\frac{1}{p} + 1\right)\right)^n}{\Gamma\left(\frac{n}{p} + 1\right)} \quad (\text{D.1})$$

On déduit l'équivalent directement. □

Ce que ce résultat exhibe, c'est que le volume d'une hypersphère en grande dimension tend exponentiellement vite vers 0, c'est complètement contre intuitif! Visualisons les courbes de cette fonction avec la figure (??).

On retrouve bien le comportement en hausse que nous avons observé, mais on comprend que le comportement ultime est que le volume tende vers 0 très rapidement. Avant de discuter de ce que ce résultat implique, regardons un autre résultat contre intuitif.

Exercice D.1 (Concentration dans l'hypersphère). *Soit $\varepsilon > 0$. On considère une hypersphère de rayon R . Montrer que :*

$$\frac{V_n^p(R - \varepsilon)}{V_n^p(R)} = \left(1 - \frac{\varepsilon}{R}\right)^n$$

C'est encore plus étrange : les points semblent se concentrer proche des frontières de l'hypersphère, donc en ayant un *centre* vide. Cela veut dire que plus la dimension augmente, plus le volume tend vers 0 et que dans le même temps les données se rapprochent des frontières.

Donc si l'on distribue des points uniformément dans une sphère, la distribution des distances entre les points ne sera pas informative du tout. Ces intuitions sont confirmées par la figure (D.4).

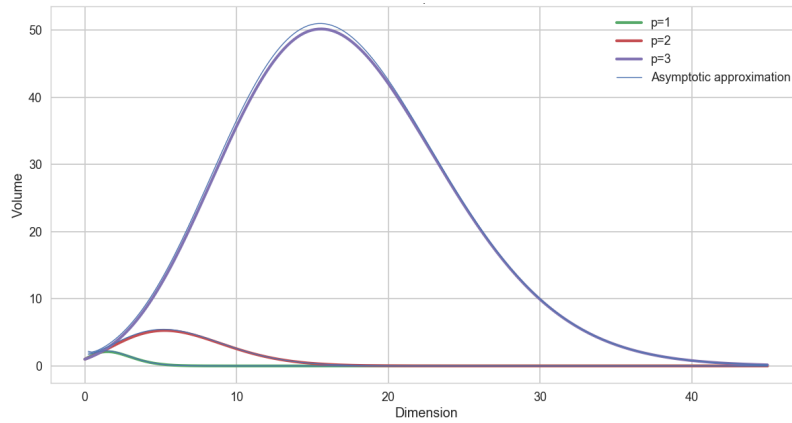


FIGURE D.3 – Volume de la boule unité en fonction de la dimension de son espace pour trois p -normes

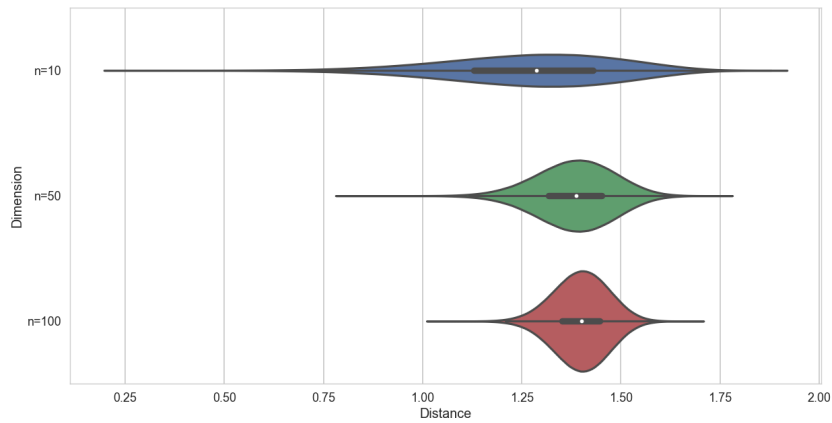


FIGURE D.4 – Distribution des distances entre chaque point en fonction de la dimension de l'espace

Ainsi, plus la dimension est grande, moins la notion de distance a du sens. Au-delà de l'aspect combinatoire et de stockage des données, avoir un modèle qui a moins d'indicateurs pour s'entraîner aura plus de chance d'être performant et *utile*. Par exemple, l'ensemble des méthodes de clustering par exemple seront largement impactées par une très grande dimension. Finalement, on peut remettre en cause l'exactitude d'une idée répandue : "*Avec plus d'informations les modèles sont meilleurs*". Ce qui est plus exact est qu'avec les informations utiles, les modèles sont meilleurs. C'est tout l'enjeu de la phase exploratoire et d'augmentation des données pour répondre à un problème de Machine Learning.

D.2 Orthogonalité à la surface d'une hypersphère

Nous avons donc montré que n'importe quelle distance issue d'une norme \mathcal{L}_p était soumise au fléau de la dimension. En NLP *classique*, il est fréquent d'utiliser la distance cosinus, où le cadre est très souvent en très grande dimension. Les résultats semblent montrer que le fléau de la dimension n'affecte pas la distance cosinus. Vérifions.

On peut définir le produit scalaire entre $x, y \in \mathbb{R}^n$ comme :

$$\langle x, y \rangle = \|x\|_2 \|y\|_2 \cos(\theta)$$

Avec θ l'angle entre le représentant de x et le représentant de y à l'origine comme défini dans la figure (D.5).

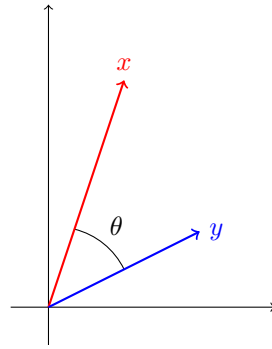


FIGURE D.5 – Définition de la métrique cosinus

Nous pouvons donc naturellement définir la métrique cosinus comme :

$$\text{cosine}(x, y) = \frac{\langle x, y \rangle}{\|x\|_2 \|y\|_2}$$

On comprend que la distance cosinus sera bornée dans $[-1, 1]$ contrairement aux restes des distances qui ne sont pas bornées. Par définition, la métrique cosinus est liée à la distance euclidienne, et on remarque que si l'on prend x et y deux vecteurs unitaires, on obtient :

$$\begin{aligned} \|x - y\|_2^2 &= \|x\|_2^2 + \|y\|_2^2 - 2 \langle x, y \rangle \\ &= \|x\|_2^2 + \|y\|_2^2 - 2 \text{cosine}(x, y) \|x\|_2 \|y\|_2 \\ &= 2 [1 - \text{cosine}(x, y)] \end{aligned}$$

Il serait donc très surprenant qu'une distance avec un lien aussi fort avec la distance euclidienne, ne souffre pas du fléau de la dimension. Nous avons un résultat intéressant :

Lemme D.1. Soit x, y deux vecteurs choisis indépendamment à la surface d'une hypersphère. Alors, avec une probabilité supérieure à $1 - \frac{1}{n}$:

$$|\text{cosine}(x, y)| = O\left(\sqrt{\frac{\log n}{n}}\right)$$

Autrement dit, en prenant deux vecteurs aléatoirement à la surface d'une boule en dimension n , on a avec une très grande probabilité que ces deux vecteurs sont orthogonaux. Cela rend inutilisable la métrique cosinus en très grande dimension.

Démonstration. Soit $a \in \mathbb{R}^n$ tel que $\|a\|_2 = 1$. Soit $x \in S_n^2$ avec $S_n^p = \{x \in \mathbb{R}^n \mid \|x\|_p = 1\}$ la surface de l'hyperboule unitaire. Soit $x \in \mathbb{R}^n$ un vecteur construit avec chacune de ses coordonnées sélectionnées aléatoirement entre -1 et 1 . Puis on normalise le vecteur x . Soit $X = \langle a, x \rangle$, alors il est simple de montrer que $\mathbb{E}[X] = 0$ et $\mathbb{E}[X^2] = \frac{1}{n}$. Ainsi, en exploitant l'inégalité de Chernoff on a :

$$\mathbb{P}(|X| \geq t) \leq 2e^{-\frac{nt^2}{4}}$$

Donc pour $\varepsilon = 2e^{-\frac{nt^2}{4}} \iff t = \sqrt{\frac{-4 \log(\frac{\varepsilon}{2})}{n}}$ et si l'on choisit $\varepsilon = \frac{1}{n}$, on obtient :

$$\mathbb{P}\left(|\cosine(x, y)| \geq \sqrt{\frac{4 \log(2n)}{n}}\right) \leq \frac{1}{n}$$

□

Cette preuve complète la présentation du second comportement étrange que l'on mentionne concernant les hyperboules et hypersphères.

Un deuxième problème majeur en grande dimension est que le nombre de données à obtenir pour être capable d'avoir des garanties statistiques sur la qualité de l'apprentissage est colossal, c'est exponentiel. Ainsi, on peut se poser la question de la capacité des algorithmes à *apprendre* en grande dimension.

D.3 Interpolation et extrapolation

L'ensemble du Machine Learning tel qu'on l'a présenté correspond à de l'interpolation et à essayer de faire en sorte que cette interpolation puisse être capable d'extrapoler correctement. Randall Balestrieri, Jerome Pesenti et Yann Le Cun ont publié en 2021 l'article *Learning in High Dimension always amount to extrapolation* [Balestrieri et al., 2021] dont voici le résumé :

The notion of interpolation and extrapolation is fundamental in various fields from deep learning to function approximation. Interpolation occurs for a sample x whenever this sample falls inside or on the boundary of the given dataset's convex hull. Extrapolation occurs when x falls outside of that convex hull. One fundamental (mis)conception is that state-of-the-art algorithms work so well because of their ability to correctly interpolate training data. A second (mis)conception is that interpolation happens throughout tasks and datasets, in fact, many intuitions and theories rely on that assumption. We empirically and theoretically argue against those two points and demonstrate that on any high-dimensional (>100) dataset, interpolation almost surely never happens. Those results challenge the validity of our current interpolation/extrapolation definition as an indicator of generalization performances.

— Randall Balestrieri, Jerome Pesenti et Yann Le Cun (2021)

L'objet de l'article est de montrer que l'on comprend et définit mal les notions d'interpolation et d'extrapolation en Machine Learning. Cela a des impacts théoriques et donc pratiques sur notre conception et les garanties mathématiques que l'on peut avoir sur les comportements des algorithmes présentés en très grande dimension. Nous invitons à lire en détail cet article pour en apprendre plus sur le sujet en lui-même, mais également pour voir qu'un domaine qui semble plutôt bien établi et en constante expansion se pose encore des questions sur ses fondements.

En résumé, le fléau de la dimension met en lumière les limites de notre intuition humaine et nous amène à nous questionner encore aujourd'hui sur les fondements communément acceptés. Être capable de répondre à ces questions nous permettrait d'être plus précis et plus complet sur notre approche de l'apprentissage en grande dimension.

De manière plus pragmatique, un data scientist doit être au courant que ces questions existent et que le fléau de la dimension va impacter son travail. D'où les techniques de réduction de dimension qui aident à résoudre le problème, mais ne le résolvent clairement pas par construction.

Bibliographie

- [Achille et al., 2017] Achille, A., Rovere, M., and Soatto, S. (2017). Critical learning periods in deep neural networks. *arXiv preprint arXiv :1711.08856*.
- [Anil et al., 2023] Anil, R., Dai, A. M., Firat, O., Johnson, M., Lepikhin, D., Passos, A., Shakeri, S., Taropa, E., Bailey, P., Chen, Z., et al. (2023). Palm 2 technical report. *arXiv preprint arXiv :2305.10403*.
- [Ba et al., 2016] Ba, J. L., Kiros, J. R., and Hinton, G. E. (2016). Layer normalization. *arXiv preprint arXiv :1607.06450*.
- [Balduzzi et al., 2017] Balduzzi, D., Frean, M., Leary, L., Lewis, J., Ma, K. W.-D., and McWilliams, B. (2017). The shattered gradients problem : If resnets are the answer, then what is the question? In *International Conference on Machine Learning*. PMLR.
- [Balestrierio et al., 2021] Balestrierio, R., Pesenti, J., and LeCun, Y. (2021). Learning in high dimension always amounts to extrapolation. *arXiv preprint arXiv :2110.09485*.
- [Barron, 2017] Barron, J. T. (2017). Continuously differentiable exponential linear units. *arXiv preprint arXiv :1704.07483*.
- [Barron, 2021] Barron, J. T. (2021). Squareplus : A softplus-like algebraic rectifier. *arXiv preprint arXiv :2112.11687*.
- [Belkin et al., 2019] Belkin, M., Hsu, D., Ma, S., and Mandal, S. (2019). Reconciling modern machine-learning practice and the classical bias-variance trade-off. *Proceedings of the National Academy of Sciences*.
- [Błasiok et al., 2023a] Błasiok, J., Gopalan, P., Hu, L., and Nakkiran, P. (2023a). A unifying theory of distance from calibration. In *Proceedings of the 55th Annual ACM Symposium on Theory of Computing*, pages 1727–1740.
- [Błasiok et al., 2023b] Błasiok, J., Gopalan, P., Hu, L., and Nakkiran, P. (2023b). When does optimizing a proper loss yield calibration? *arXiv preprint arXiv :2305.18764*.
- [Brown et al., 2020] Brown, T., Mann, B., Ryder, N., Subbiah, M., Kaplan, J. D., Dhariwal, P., Neelakantan, A., Shyam, P., Sastry, G., Askell, A., et al. (2020). Language models are few-shot learners. *Advances in Neural Information Processing Systems*.
- [Chen et al., 2023] Chen, X., Liang, C., Huang, D., Real, E., Wang, K., Liu, Y., Pham, H., Dong, X., Luong, T., Hsieh, C.-J., et al. (2023). Symbolic discovery of optimization algorithms. *arXiv preprint arXiv :2302.06675*.

- [Chowdhery et al., 2022] Chowdhery, A., Narang, S., Devlin, J., Bosma, M., Mishra, G., Roberts, A., Barham, P., Chung, H. W., Sutton, C., Gehrmann, S., et al. (2022). Palm : Scaling language modeling with pathways. *arXiv preprint arXiv :2204.02311*.
- [Dauphin et al., 2016] Dauphin, Y. N., Fan, A., Auli, M., and Grangier, D. (2016). Language modeling with gated convolutional networks. In *International conference on machine learning*, pages 933–941. PMLR.
- [Dawid, 1982] Dawid, A. P. (1982). The well-calibrated bayesian. *Journal of the American Statistical Association*.
- [Duchi et al., 2011] Duchi, J., Hazan, E., and Singer, Y. (2011). Adaptive subgradient methods for online learning and stochastic optimization. *Journal of machine learning research*.
- [Fan et al., 2018] Fan, A., Lewis, M., and Dauphin, Y. (2018). Hierarchical neural story generation. *arXiv preprint arXiv :1805.04833*.
- [Giffin and Mitchell, 1978] Giffin, F. and Mitchell, D. E. (1978). The rate of recovery of vision after early monocular deprivation in kittens. *The Journal of Physiology*.
- [Glorot and Bengio, 2010] Glorot, X. and Bengio, Y. (2010). Understanding the difficulty of training deep feedforward neural networks. In *Proceedings of the thirteenth international conference on artificial intelligence and statistics*, pages 249–256. JMLR Workshop and Conference Proceedings.
- [Gneiting and Raftery, 2007] Gneiting, T. and Raftery, A. E. (2007). Strictly proper scoring rules, prediction, and estimation. *Journal of the American statistical Association*, 102(477) :359–378.
- [Gopalan et al., 2022] Gopalan, P., Kim, M. P., Singhal, M. A., and Zhao, S. (2022). Low-degree multicalibration. In *Conference on Learning Theory*, pages 3193–3234. PMLR.
- [He et al., 2015] He, K., Zhang, X., Ren, S., and Sun, J. (2015). Delving deep into rectifiers : Surpassing human-level performance on imagenet classification. In *Proceedings of the IEEE international conference on computer vision*, pages 1026–1034.
- [He et al., 2016] He, K., Zhang, X., Ren, S., and Sun, J. (2016). Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*.
- [Hendrycks and Gimpel, 2016] Hendrycks, D. and Gimpel, K. (2016). Gaussian error linear units (gelus). *arXiv preprint arXiv :1606.08415*.
- [Hinton et al., 2012a] Hinton, G., Srivastava, N., and Swersky, K. (2012a). Neural networks for machine learning lecture 6a overview of mini-batch gradient descent. *Cited on*.
- [Hinton et al., 2012b] Hinton, G. E., Srivastava, N., Krizhevsky, A., Sutskever, I., and Salakhutdinov, R. R. (2012b). Improving neural networks by preventing co-adaptation of feature detectors. *arXiv preprint arXiv :1207.0580*.
- [Hoffmann et al., 2022] Hoffmann, J., Borgeaud, S., Mensch, A., Buchatskaya, E., Cai, T., Rutherford, E., Casas, D. d. L., Hendricks, L. A., Welbl, J., Clark, A., et al. (2022). Training compute-optimal large language models. *arXiv preprint arXiv :2203.15556*.
- [Holtzman et al., 2019] Holtzman, A., Buys, J., Du, L., Forbes, M., and Choi, Y. (2019). The curious case of neural text degeneration. *arXiv preprint arXiv :1904.09751*.

- [Huang et al., 2016] Huang, G., Sun, Y., Liu, Z., Sedra, D., and Weinberger, K. Q. (2016). Deep networks with stochastic depth. In *Computer Vision–ECCV 2016 : 14th European Conference, Amsterdam, The Netherlands, October 11–14, 2016, Proceedings, Part IV* 14.
- [Ioffe and Szegedy, 2015] Ioffe, S. and Szegedy, C. (2015). Batch normalization : Accelerating deep network training by reducing internal covariate shift. *International conference on machine learning*.
- [Ishida et al., 2020] Ishida, T., Yamane, I., Sakai, T., Niu, G., and Sugiyama, M. (2020). Do we need zero training loss after achieving zero training error ? *arXiv preprint arXiv :2002.08709*.
- [Kakade and Foster, 2004] Kakade, S. M. and Foster, D. P. (2004). Deterministic calibration and nash equilibrium. In *International Conference on Computational Learning Theory*, pages 33–48. Springer.
- [Kaplan et al., 2020] Kaplan, J., McCandlish, S., Henighan, T., Brown, T. B., Chess, B., Child, R., Gray, S., Radford, A., Wu, J., and Amodei, D. (2020). Scaling laws for neural language models. *arXiv preprint arXiv :2001.08361*.
- [Kingma and Ba, 2015] Kingma, D. P. and Ba, J. (2015). Adam : A method for stochastic optimization. In *International Conference on Learning Representations (ICLR)*.
- [Klambauer et al., 2017] Klambauer, G., Unterthiner, T., Mayr, A., and Hochreiter, S. (2017). Self-normalizing neural networks. *Advances in neural information processing systems*.
- [Krizhevsky et al., 2012] Krizhevsky, A., Sutskever, I., and Hinton, G. E. (2012). Imagenet classification with deep convolutional neural networks. *Advances in neural information processing systems*.
- [Krogh and Hertz, 1991] Krogh, A. and Hertz, J. (1991). A simple weight decay can improve generalization. *Advances in neural information processing systems*.
- [LeCun et al., 1998] LeCun, Y. A., Bottou, L., Orr, G. B., and Müller, K.-R. (1998). Efficient backprop. *Tricks of the Trade*.
- [Lee and Cherkassky, 2022] Lee, E. H. and Cherkassky, V. (2022). Vc theoretical explanation of double descent. *arXiv preprint arXiv :2205.15549*.
- [Li et al., 2018] Li, H., Xu, Z., Taylor, G., Studer, C., and Goldstein, T. (2018). Visualizing the loss landscape of neural nets. *Advances in neural information processing systems*, 31.
- [Liu et al., 2018] Liu, W., Lin, R., Liu, Z., Liu, L., Yu, Z., Dai, B., and Song, L. (2018). Learning towards minimum hyperspherical energy. *Advances in neural information processing systems*, 31.
- [Liu et al., 2023] Liu, Z., Xu, Z., Jin, J., Shen, Z., and Darrell, T. (2023). Dropout reduces underfitting. *arXiv preprint arXiv :2303.01500*.
- [Loshchilov and Hutter, 2016] Loshchilov, I. and Hutter, F. (2016). Sgdr : Stochastic gradient descent with warm restarts. *arXiv preprint arXiv :1608.03983*.
- [Loshchilov and Hutter, 2019] Loshchilov, I. and Hutter, F. (2019). Decoupled weight decay regularization. *International Conference on Learning Representations (ICLR)*.

- [Maas et al., 2013] Maas, A. L., Hannun, A. Y., Ng, A. Y., et al. (2013). Rectifier nonlinearities improve neural network acoustic models. In *Proc. icml*.
- [Misra, 2019] Misra, D. (2019). Mish : A self regularized non-monotonic activation function. *arXiv preprint arXiv :1908.08681*.
- [Mitchell, 1988] Mitchell, D. E. (1988). The extent of visual recovery from early monocular or binocular visual deprivation in kittens. *The Journal of physiology*.
- [Nakkiran, 2021] Nakkiran, P. (2021). *Towards an Empirical Theory of Deep Learning*. PhD thesis, Harvard University.
- [Olson and Freeman, 1980] Olson, C. R. and Freeman, R. (1980). Profile of the sensitive period for monocular deprivation in kittens. *Experimental Brain Research*.
- [Pascanu et al., 2013] Pascanu, R., Mikolov, T., and Bengio, Y. (2013). On the difficulty of training recurrent neural networks. In *International conference on machine learning*, pages 1310–1318. PMLR.
- [Pedregosa et al., 2011] Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., Blondel, M., Prettenhofer, P., Weiss, R., Dubourg, V., et al. (2011). Scikit-learn : Machine learning in python. *The journal of machine learning research*.
- [Penedo et al., 2023] Penedo, G., Malartic, Q., Hesslow, D., Cojocaru, R., Cappelli, A., Alobeidli, H., Pannier, B., Almazrouei, E., and Launay, J. (2023). The refinedweb dataset for falcon llm : outperforming curated corpora with web data, and web data only. *arXiv preprint arXiv :2306.01116*.
- [Rae et al., 2021] Rae, J. W., Borgeaud, S., Cai, T., Millican, K., Hoffmann, J., Song, F., Aslanides, J., Henderson, S., Ring, R., Young, S., et al. (2021). Scaling language models : Methods, analysis & insights from training gopher. *arXiv preprint arXiv :2112.11446*.
- [Ramachandran et al., 2017] Ramachandran, P., Zoph, B., and Le, Q. V. (2017). Searching for activation functions. *arXiv preprint arXiv :1710.05941*.
- [Reed et al., 2022] Reed, S., Zolna, K., Parisotto, E., Colmenarejo, S. G., Novikov, A., Barth-Maron, G., Gimenez, M., Sulsky, Y., Kay, J., Springenberg, J. T., et al. (2022). A generalist agent. *arXiv preprint arXiv :2205.06175*.
- [Savage, 1971] Savage, L. J. (1971). Elicitation of personal probabilities and expectations. *Journal of the American Statistical Association*, 66(336) :783–801.
- [Schmidt et al., 2021] Schmidt, R. M., Schneider, F., and Hennig, P. (2021). Descending through a crowded valley-benchmarking deep learning optimizers. In *International Conference on Machine Learning*. PMLR.
- [Shazeer, 2020] Shazeer, N. (2020). Glu variants improve transformer. *arXiv preprint arXiv :2002.05202*.
- [Smith, 2017] Smith, L. N. (2017). Cyclical learning rates for training neural networks. In *2017 IEEE winter conference on applications of computer vision (WACV)*, pages 464–472. IEEE.
- [Srivastava et al., 2014] Srivastava, N., Hinton, G., Krizhevsky, A., Sutskever, I., and Salakhutdinov, R. (2014). Dropout : a simple way to prevent neural networks from overfitting. *The journal of machine learning research*.

- [Tan et al., 2022] Tan, C., Gao, Z., Wu, L., Li, S., and Li, S. Z. (2022). Hyperspherical consistency regularization. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7244–7255.
- [Thoppilan et al., 2022] Thoppilan, R., De Freitas, D., Hall, J., Shazeer, N., Kulshreshtha, A., Cheng, H.-T., Jin, A., Bos, T., Baker, L., Du, Y., et al. (2022). Lamda : Language models for dialog applications. *arXiv preprint arXiv :2201.08239*.
- [Touvron et al., 2023a] Touvron, H., Lavril, T., Izacard, G., Martinet, X., Lachaux, M.-A., Lacroix, T., Rozière, B., Goyal, N., Hambro, E., Azhar, F., et al. (2023a). Llama : Open and efficient foundation language models. *arXiv preprint arXiv :2302.13971*.
- [Touvron et al., 2023b] Touvron, H., Martin, L., Stone, K., Albert, P., Almahairi, A., Babaei, Y., Bashlykov, N., Batra, S., Bhargava, P., Bhosale, S., et al. (2023b). Llama 2 : Open foundation and fine-tuned chat models. *arXiv preprint arXiv :2307.09288*.
- [Vaswani et al., 2017] Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, Ł., and Polosukhin, I. (2017). Attention is all you need. *Advances in neural information processing systems*, 30.
- [Xiong et al., 2020] Xiong, R., Yang, Y., He, D., Zheng, K., Zheng, S., Xing, C., Zhang, H., Lan, Y., Wang, L., and Liu, T. (2020). On layer normalization in the transformer architecture. In *International Conference on Machine Learning*. PMLR.
- [Zhang and Sennrich, 2019] Zhang, B. and Sennrich, R. (2019). Root mean square layer normalization. *Advances in Neural Information Processing Systems*, 32.
- [Zhang et al., 2019] Zhang, H., Dauphin, Y. N., and Ma, T. (2019). Fixup initialization : Residual learning without normalization. *ICLR*.
- [Zhou et al., 2023] Zhou, C., Liu, P., Xu, P., Iyer, S., Sun, J., Mao, Y., Ma, X., Efrat, A., Yu, P., Yu, L., et al. (2023). Lima : Less is more for alignment. *arXiv preprint arXiv :2305.11206*.