



**HAL**  
open science

## Sparsity in neural networks can improve their privacy

Antoine Gonon, Léon Zheng, Clément Lalanne, Quoc-Tung Le, Guillaume Lauga, Can Pouliquen

► **To cite this version:**

Antoine Gonon, Léon Zheng, Clément Lalanne, Quoc-Tung Le, Guillaume Lauga, et al.. Sparsity in neural networks can improve their privacy. 2023. hal-04062317v1

**HAL Id: hal-04062317**

**<https://hal.science/hal-04062317v1>**

Preprint submitted on 7 Apr 2023 (v1), last revised 12 Oct 2023 (v3)

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# La parcimonie des réseaux de neurones peut améliorer leur confidentialité

Antoine GONON\*<sup>1</sup> Léon ZHENG\*<sup>1,2</sup> Clément LALANNE<sup>1</sup> Quoc-Tung LE<sup>1</sup> Guillaume LAUGA†<sup>1</sup> Can POULIQUEN†<sup>1</sup>

<sup>1</sup>Univ Lyon, EnsL, UCBL, CNRS, Inria, LIP, F-69342, LYON Cedex 07, France

<sup>2</sup>valeo.ai, Paris, France

**Résumé** – Cet article mesure l’impact de la parcimonie des réseaux de neurones sur l’efficacité des attaques par inférence d’appartenance. Les résultats obtenus montrent que la parcimonie permet d’améliorer la confidentialité, tout en préservant des performances comparables sur la tâche d’apprentissage considérée. Cette démonstration empirique complète et étend des résultats obtenus précédemment dans la littérature.

**Abstract** – This article measures how sparsity can make neural networks more robust to membership inference attacks. The obtained empirical results show that sparsity improves the privacy of the network, while preserving comparable performances on the task at hand. This empirical study completes and extends existing literature.

## 1 Introduction

Les réseaux de neurones profonds constituent l’état de l’art dans de nombreux problèmes d’apprentissage. En pratique, il est possible d’ajuster les paramètres du réseau considéré afin de parfaitement interpoler les données disponibles [19]. Cette situation de *sur-apprentissage* est intéressante car les modèles ont de bonnes performances dans ce régime [3]. Néanmoins, il présente un risque de confidentialité puisque le modèle mémorise des informations sur les données, au point de les interpoler. Parmi ces informations, certaines sont peut-être confidentielles, et il se pose la question de savoir lesquelles peuvent être retrouvées à partir de la seule connaissance du modèle appris.

Pour détecter une situation de sur-apprentissage, un indicateur est donné par le ratio nombre de paramètres sur nombre de données : plus il y a de paramètres, plus le modèle peut interpoler les données. Afin de contrôler la capacité du modèle à sur-apprendre, et donc à mémoriser des informations confidentielles, cet article étudie le rôle du nombre de paramètres non nuls utilisés. Peut-on trouver un bon compromis entre précision du modèle et confidentialité en ajustant la parcimonie (nombre de paramètres non nuls) des réseaux de neurones ?

Via un type d’attaque nommé « Membership Inference Attack » (MIA), il est possible d’inférer l’appartenance de données au jeu d’entraînement [16]. Cette attaque ne requiert qu’un accès boîte-noire au modèle visé, et peut être problématique selon la confidentialité des données (médicales, etc.). Étant donné un réseau, comment réduire le risque d’une telle attaque, tout en préservant au mieux ses performances ?

De nombreuses procédures ont été proposées pour se défendre contre les MIAs [9]. Ici, l’approche étudiée consiste à diminuer le nombre de paramètres non nuls utilisés par le réseau afin de réduire sa capacité de mémorisation, en préservant autant que possible les performances.

\* , † : Contributions égales. Le travail est en partie soutenu par les projets AllegroAssai ANR-19-CHIA-0009, NuSCAP ANR-20-CE48-0014, SeqALO ANR-20-CHIA-0020-01, MOMIGS du GdR ISIS et la CIFRE N°2020/1643. Les auteurs remercient le Centre Blaise Pascal pour les moyens de calcul. La plateforme exploite SIDUS [14] développée par Emmanuel Quemener.

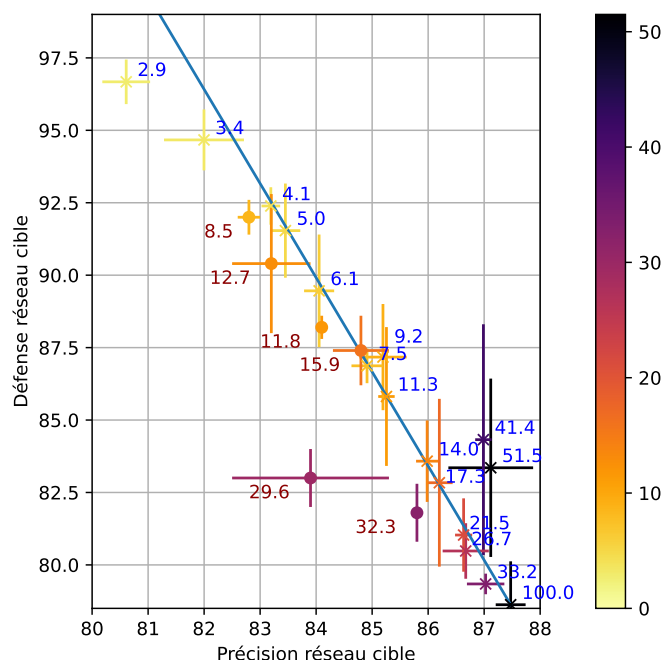


FIGURE 1 : Résultats moyens et écart-types obtenus pour la précision et la défense d’un réseau cible. Le pourcentage de poids non nuls est indiqué en bleu pour IMP (\* p%), en rouge pour Butterfly (\* p%). La couleur des points indique le niveau de parcimonie. La droite est de pente  $-3.25$ .

**Approches similaires.** Les liens entre parcimonie des réseaux de neurones et confidentialité ont déjà été partiellement explorés, mais, à notre connaissance, il n’a pas encore été mis en évidence que la parcimonie permet d’améliorer la confidentialité *sans modification supplémentaire* de l’algorithme d’entraînement. Un positionnement est réalisé en section 4.

**Contributions et résultats.** Les résultats en section 4 corroborent l’hypothèse que la parcimonie permet d’améliorer la défense contre une MIA tout en gardant des performances

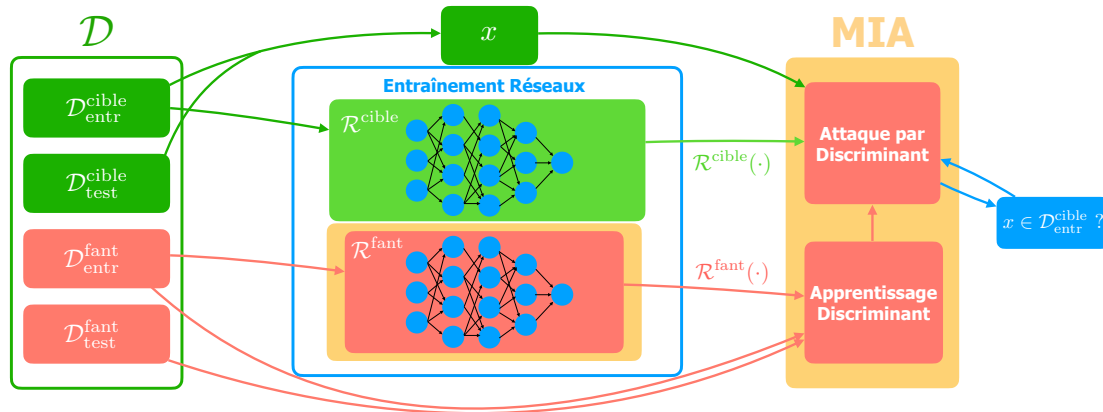


FIGURE 2 : Les expériences réalisées obéissent au même schéma général représenté ici : deux réseaux sont entraînés de la même manière sur  $\mathcal{D}_{\text{entr}}^{\text{cible}}$  et  $\mathcal{D}_{\text{entr}}^{\text{fant}}$  respectivement.  $\mathcal{R}^{\text{fant}}$ ,  $\mathcal{D}_{\text{entr}}^{\text{fant}}$  et  $\mathcal{D}_{\text{test}}^{\text{fant}}$  sont ensuite utilisés pour entraîner un discriminant qui va attaquer  $\mathcal{R}^{\text{cible}}$  en identifiant l'appartenance ou non de  $x$  à  $\mathcal{D}_{\text{entr}}^{\text{cible}}$ .

comparables sur la tâche d'apprentissage. Néanmoins les écarts-types observés nuancent les résultats et suggèrent qu'il est nécessaire de réaliser des expériences à plus grande échelle avant de pouvoir confirmer cette tendance. La Figure 1 montre que le compromis entre robustesse aux attaques et précision du réseau est similaire entre, d'une part, des parcimonies non structurées obtenues par une procédure itérative d'élagage des poids par magnitude et de ré-ajustement (« Iterative Magnitude Pruning », IMP) [5], et, d'autre part, des parcimonies structurées dites « butterfly », où les matrices de poids sont contraintes pour admettre une certaine factorisation creuse structurée [13, 4]. À notre connaissance, la structure « butterfly » n'a pas été étudiée dans la littérature dans ce contexte. Cette structure atteint des compromis comparables à IMP, avec l'avantage de permettre une implémentation efficace de la multiplication matrice-vecteur. Ceci est remarquable, car la structure est fixée indépendamment des données.

Les expériences réalisées sur CIFAR-10 montrent que lorsque le pourcentage des poids non nuls de ResNet-20 est entre 3.4% et 17.3%, une perte relative de  $p\%$  par rapport au réseau dense entraîné<sup>1</sup> en précision mène à un gain relatif de  $3.6 \times p\%$  en défense contre la MIA, voir la Figure 1.

La section 2 introduit le modèle d'attaque par réseau fantôme utilisé pour les expériences. La section 3 décrit les types de parcimonie utilisés pour se défendre contre les MIAs. Les expériences sont présentées en section 4.

## 2 Attaque par réseau fantôme

Soient  $\mathcal{D}$  un jeu de données et un sous-ensemble  $\mathcal{D}_{\text{entr}} \subset \mathcal{D}$ . La fonction d'appartenance  $a_{\mathcal{D}_{\text{entr}}, \mathcal{D}}$  associée est définie par :

$$a_{\mathcal{D}_{\text{entr}}, \mathcal{D}} : x \in \mathcal{D} \mapsto \begin{cases} 1 & \text{si } x \in \mathcal{D}_{\text{entr}}, \\ 0 & \text{sinon.} \end{cases}$$

Étant donné un jeu de donnée  $\mathcal{D}^{\text{cible}}$ , un réseau cible  $\mathcal{R}^{\text{cible}}$  entraîné sur  $\mathcal{D}_{\text{entr}}^{\text{cible}} \subset \mathcal{D}^{\text{cible}}$ , une MIA consiste à retrouver la fonction d'appartenance  $a_{\text{cible}} := a_{\mathcal{D}_{\text{entr}}^{\text{cible}}, \mathcal{D}^{\text{cible}}}$  associée en ayant seulement un accès *boîte noire* à la fonction  $x \mapsto \mathcal{R}^{\text{cible}}(x)$ . La plupart des attaques connues se basent sur une observation de la sortie du modèle  $\mathcal{R}^{\text{cible}}$ , localement autour de  $x$  [9]. En

général, ces attaques cherchent à mesurer la confiance du modèle en ses prédictions réalisées localement autour de  $x$ . Si la mesure de confiance est suffisamment élevée, alors l'attaquant répond oui à la question d'appartenance.

En pratique, l'attaque la plus performante [9] consiste à entraîner un modèle *discriminant* qui prend une décision en fonction d'informations locales sur  $\mathcal{R}^{\text{cible}}$  autour de  $x$ . Ce discriminant est entraîné à partir d'un réseau *fantôme* [9], comme expliqué ci-dessous (voir aussi Figure 2).

**Réseau Fantôme.** Supposons que l'attaquant ait accès à un jeu de données  $\mathcal{D}^{\text{fant}}$  issu de la même distribution que  $\mathcal{D}^{\text{cible}}$ . Il entraîne alors son propre réseau fantôme  $\mathcal{R}^{\text{fant}}$  sur un sous-ensemble  $\mathcal{D}_{\text{entr}}^{\text{fant}} \subset \mathcal{D}^{\text{fant}}$  des données qu'il possède. Idéalement,  $\mathcal{R}^{\text{fant}}$  est entraîné dans les mêmes conditions que  $\mathcal{R}^{\text{cible}}$  (même structure et même algorithme d'optimisation). L'attaquant a alors un triplet  $(\mathcal{R}^{\text{fant}}, \mathcal{D}_{\text{entr}}^{\text{fant}}, \mathcal{D}_{\text{test}}^{\text{fant}})$  ayant des similarités avec le triplet  $(\mathcal{R}^{\text{cible}}, \mathcal{D}_{\text{entr}}^{\text{cible}}, \mathcal{D}_{\text{test}}^{\text{cible}})$ , et la connaissance de la fonction d'appartenance  $a_{\text{fant}} := a_{\mathcal{D}_{\text{entr}}^{\text{fant}}, \mathcal{D}^{\text{fant}}}$ .

**Discriminant.** L'attaquant peut ensuite entraîner un discriminant afin d'approcher  $a_{\text{fant}}$  à partir du seul accès boîte noire de  $\mathcal{R}^{\text{fant}}$ . Ce discriminant peut alors servir à approcher  $a_{\text{cible}}$  à partir du seul accès boîte noire de  $\mathcal{R}^{\text{cible}}$ . Le modèle pour le discriminant peut être n'importe quel classificateur classique (régression logistique, réseau de neurones, etc.) [9].

## 3 Défense et parcimonie des réseaux

L'entraînement de réseaux de neurones parcimonieux est d'abord motivé par des besoins de frugalité en ressources [8] (mémoire, temps d'inférence, temps d'entraînement, etc.).

Ici, l'hypothèse suivante est étudiée : la parcimonie peut limiter le sur-apprentissage, et ainsi limiter la capacité du modèle à mémoriser des informations confidentielles sur les données qu'il a vues. Un réseau parfaitement confidentiel n'a rien appris de ses données et n'a donc pas d'intérêt en pratique. Un compromis entre confidentialité et précision du réseau est à réaliser en fonction de la tâche considérée.

### 3.1 Parcimonie non structurée via IMP

Dans le premier cas, aucune structure spécifique n'est imposée sur l'ensemble des poids non nuls. Les poids nuls sont

<sup>1</sup>Le réseau dense est le réseau d'origine, avec 100% des poids non nuls.

sélectionnés par un processus itératif d'élagage par amplitude (« Iterative Magnitude Pruning », IMP) [5] qui consiste à : (i) entraîner un réseau de manière usuelle, (ii) élaguer (mettre à zéro)  $p\%$  des poids ayant l'amplitude la plus faible, (iii) ajuster les poids restants en ré-entraînant le réseau (les poids ayant été élagués sont masqués et ne sont plus mis à jour), puis revenir à l'étape (ii) tant que la parcimonie souhaitée n'est pas atteinte. Cette procédure permet de trouver des sous-réseaux ayant empiriquement de bonnes propriétés statistiques [5, 6].

### 3.2 Parcimonie structurée butterfly

Dans le second cas, la parcimonie est structurée : les matrices de poids des couches du réseau de neurones ont une factorisation « butterfly », qui permet d'implémenter efficacement la multiplication matrice-vecteur [4]. Selon [12], une matrice  $\mathbf{W}$  carrée de taille  $N := 2^L$  admet une factorisation butterfly lorsqu'il s'écrit comme le produit exact  $\mathbf{W} = \mathbf{X}^{(1)} \dots \mathbf{X}^{(L)}$  de  $L$  facteurs carrés de taille  $N$ , où chaque facteur satisfait la contrainte de support<sup>2</sup>  $\text{supp}(\mathbf{X}^{(\ell)}) \subseteq \text{supp}(\mathbf{S}_{\text{bf}}^{(\ell)})$ , avec  $\mathbf{S}_{\text{bf}}^{(\ell)} := \mathbf{I}_{2^{\ell-1}} \otimes \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix} \otimes \mathbf{I}_{N/2^{\ell}}$ . Une illustration des supports est donnée en Figure 3. Les facteurs ont au plus deux coefficients non nuls par ligne et par colonne. La multiplication matrice-vecteur a pour complexité  $\mathcal{O}(N \log N)$  en utilisant la forme factorisée, contre  $\mathcal{O}(N^2)$  en général.

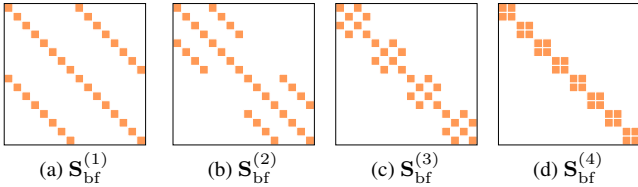


FIGURE 3 : Supports des facteurs butterfly de taille  $N = 16$ .

Pour imposer la structure butterfly dans un réseau de neurones, les matrices de poids  $\mathbf{W}$  sont paramétrisées sous la forme  $\mathbf{W} = \mathbf{X}^{(1)} \dots \mathbf{X}^{(L)}$ , et seuls les coefficients non nuls des facteurs  $\mathbf{X}^{(1)}, \dots, \mathbf{X}^{(L)}$  sont ajustés pour minimiser la fonction de coût au cours de l'entraînement. L'initialisation de ces coefficients au début de l'entraînement est aléatoire.

Dans le cas d'une couche de convolution, la matrice  $\mathbf{W}$  pour laquelle on impose une factorisation butterfly correspond à la concaténation des noyaux de convolution [13]. Dans nos expériences, pour une taille de  $\mathbf{W}$  et un nombre de facteurs  $L$  fixés, la factorisation butterfly généralisée au cas rectangulaire est paramétrée selon une chaîne *monotone* en suivant [13]. Parmi tous les choix de chaînes possibles, celle avec le nombre minimal de paramètres est sélectionnée pour imposer la parcimonie butterfly.

Les réseaux butterfly obtenus atteignent des performances empiriques comparables à un réseau dense sur des tâches de classification d'images [4, 13].

## 4 Résultats expérimentaux

Tous les hyperparamètres (y compris l'architecture du discriminant) ont été déterminés suite à une recherche sur grille où

<sup>2</sup> $\text{supp}(\cdot)$  est le support d'une matrice, c'est-à-dire l'ensemble des indices de la matrice correspondants à des coefficients non nuls,  $\mathbf{I}_N$  est la matrice identité de taille  $N$ , et  $\otimes$  est le produit de Kronecker.

TABLE 1 : Hyperparamètres pour l'entraînement des réseaux cibles et fantômes.

| Réseau                       | % de paramètres                 | Pas initial | Weight decay |
|------------------------------|---------------------------------|-------------|--------------|
| ResNet-20 dense              | 100 %                           | 0.03        | 0.005        |
| Butterfly ( $S = 1, L = 2$ ) | 32.3 %                          | 0.3         | 0.0005       |
| Butterfly ( $S = 1, L = 3$ ) | 29.6 %                          | 0.3         | 0.0001       |
| Butterfly ( $S = 2, L = 2$ ) | 15.9 %                          | 0.3         | 0.0005       |
| Butterfly ( $S = 2, L = 3$ ) | 12.9 %                          | 0.1         | 0.001        |
| Butterfly ( $S = 3, L = 2$ ) | 11.8 %                          | 0.3         | 0.0005       |
| Butterfly ( $S = 3, L = 3$ ) | 8.5 %                           | 0.1         | 0.001        |
| IMP avec $k$ élagages        | $\approx 100 \times (0.8)^k \%$ | 0.03        | 0.005        |

l'aléa a été moyenné sur trois expériences.

**Données.** Les expériences sont réalisées sur le problème de classification CIFAR-10 (60000 images  $32 \times 32 \times 3$ , 10 classes). Les données sont aléatoirement (uniformément) partitionnées en 4 sous-ensembles  $\mathcal{D}_{\text{entr}}^{\text{cible}}, \mathcal{D}_{\text{test}}^{\text{cible}}, \mathcal{D}_{\text{entr}}^{\text{fant}}, \mathcal{D}_{\text{test}}^{\text{fant}}$  de 15000 images, respectivement utilisés pour entraîner et tester les réseaux cibles et fantômes. Les questions d'appartenance se posent pour  $\mathcal{D}^{\text{cible}} := \mathcal{D}_{\text{entr}}^{\text{cible}} \cup \mathcal{D}_{\text{test}}^{\text{cible}}$  et  $\mathcal{D}^{\text{fant}} := \mathcal{D}_{\text{entr}}^{\text{fant}} \cup \mathcal{D}_{\text{test}}^{\text{fant}}$ . Pour le réseau cible et fantôme, parmi leurs 15000 données d'entraînement, 1000 sont choisies aléatoirement et fixées pour toutes nos expériences comme ensemble de validation (utilisé pour choisir les hyperparamètres, et pour le critère d'arrêt).

**Entraînement des réseaux cibles et fantômes.** Les réseaux cibles et fantômes sont des ResNet-20 [7] (272474 paramètres) entraînés pour minimiser l'entropie croisée par descente de gradient stochastique (avec momentum 0.9 et sans accélération Nesterov) sur leur jeu d'entraînement respectif pendant 300 époques, avec des batchs de taille 256. Les données sont augmentées avec retournement horizontal aléatoire et recadrage aléatoire. Le pas d'apprentissage initial est divisé par 10 au bout de 150 époques, puis à nouveau par 10 au bout de 225 époques. Les poids des réseaux de neurones sont initialisés avec la méthode standard par défaut sur Pytorch selon une loi uniforme sur  $(-1/\sqrt{n}, 1/\sqrt{n})$  où  $n =$  dimension entrée pour une couche linéaire, et  $n =$  dimension entrée  $\times$  largeur noyau  $\times$  hauteur noyau pour une convolution. Les valeurs du pas d'apprentissage initial et du « weight decay » sont reportées dans le tableau 1. Elles permettent de reproduire les résultats de [7] lorsque les réseaux sont entraînés sur l'ensemble du jeu d'entraînement de CIFAR-10 avec ResNet-20.

Pour IMP, 24 élagages et ré-ajustements des poids sont réalisés. Chaque ré-ajustement consiste en un entraînement comme ci-dessus (300 époques, etc.). Avant chaque élagage, les poids sont rembobinés à leurs valeurs associées à l'époque ayant la précision maximale sur le jeu de validation lors des dernières 300 époques.

Pour l'entraînement de ResNet-20 avec structure butterfly, les matrices de poids originelles de certaines couches de convolution sont substituées par des matrices ayant une factorisation, avec un nombre  $L = 2, 3$  de facteurs, suivant une chaîne monotone minimisant le nombre de paramètres dans la factorisation, comme décrit en section 3.2. Les couches substituées sont celle des  $S = 1, 2, 3$  derniers segments de ResNet-20.

**Entraînement du discriminant.** Un discriminant prend en entrée la classe  $i$  de  $x$ , la prédiction  $\mathcal{R}(x)$  réalisée par un réseau  $\mathcal{R}$  (cible ou fantôme), ainsi que  $\frac{1}{\epsilon} \mathbb{E}(|\mathcal{R}(x) - \mathcal{R}(x + \epsilon \mathcal{N})|)$

( $\epsilon = 0,001$  et  $\mathcal{N}$  un vecteur gaussien indépendant centré réduit) encodant des informations locales du premier ordre sur  $\mathcal{R}$  autour de  $x$ . L'espérance est estimée par moyenne de Monte-Carlo (5 échantillons). Pour chaque couple de réseaux ( $\mathcal{R}^{\text{cible}}, \mathcal{R}^{\text{fant}}$ ), sont entraînés trois discriminants (perceptrons) à respectivement 1, 2, 3 couche(s) cachée(s) avec respectivement 30, 30, 100 neurones sur chaque couche cachée. L'entropie croisée est minimisée avec Adam, sans weight decay avec des pas dans  $\{0.01, 0.001, 0.0001\}$  sur 80 époques.

**Précision et défense** La *précision* d'un réseau est le pourcentage de données dont la classe est celle prédite avec le plus de probabilité par le réseau. La *défense*  $D$  d'un réseau vis-à-vis d'un discriminant ayant une précision  $P\%$  est définie comme  $D = 200 - 2P$ . Par exemple, si un discriminant a une précision d'attaque  $P = 50 + x$ , alors la défense est de  $D = 100 - 2x$ . Dans notre cas, il y a autant de données vues que non vues par le réseau (cible ou fantôme) durant l'entraînement. Idéalement, le discriminant ne peut pas faire mieux que deviner au hasard, ayant alors une précision de 50%.

**Résultats** Les réseaux cibles et fantômes denses atteignent en moyenne 87.5% de précision sur l'ensemble test. Cette précision diminue avec la parcimonie, voir la Figure 1. Un gain (ou perte) en défense est significatif si l'intervalle donné par la moyenne plus ou moins l'écart-type est disjoint de l'intervalle correspondant au réseau dense entraîné. De manière significative, une amélioration de la défense est observée pour une proportion de poids entre 0% et 17.3%, pour 41.4% et 51.5%. Entre 0% et 17.3%, une perte relative de  $p\%$  en précision, par rapport au réseau dense entraîné, mène à un gain relatif de  $3.6 \times p\%$  :  $3.6 \simeq \frac{|\text{défense} - \text{défense dense}|}{\text{défense dense}} \frac{\text{précision dense}}{|\text{précision} - \text{précision dense}|}$ .

**Positionnement** Les résultats expérimentaux de [18] suggèrent à l'inverse que l'entraînement d'un réseau avec régularisation parcimonieuse par IMP *dégrade* la confidentialité. Mais ces résultats n'ont pas été moyennés sur plusieurs expériences pour diminuer la variabilité de l'aléa. Les expériences de [18] sont réalisées sur un modèle ayant 40 fois plus de poids que ResNet-20, et pour des proportions de poids non nuls  $> 50\%$ . Étant donné les écart-types observés en Figure 1 pour des niveaux de parcimonie  $> 20\%$  sur ResNet-20, il convient de rester prudent sur l'interprétabilité des résultats de [18].

L'article [17] fixe un niveau de parcimonie, et cherche les paramètres qui minimisent la fonction de perte du problème d'apprentissage, pénalisé par la plus grande précision d'attaque MIA atteignable contre ces paramètres. Néanmoins, ce terme de pénalisation n'est en général pas calculable explicitement, et difficile à minimiser. Sans comparaison avec le cas non pénalisé [17], on ne peut conclure sur la nécessité de cette pénalisation. Ici, la confidentialité est améliorée sans cette pénalisation. De plus, [17] ne présente pas quelle confidentialité est atteinte pour chaque niveau de parcimonie, mais seulement au niveau de parcimonie ayant la plus petite fonction de perte pénalisée. La Figure 1 montre quel est l'effet de la parcimonie sur la confidentialité.

Enfin, il a été observé qu'imposer la parcimonie durant l'entraînement des réseaux de neurones avec DP-SGD (« Differentially Private Stochastic Gradient Descent ») [1, 2] améliore leurs performances, à garanties égales de « Differential Pri-

vacy » (donnant des garanties fortes de confidentialité) [10, 2]. L'utilisation de la DP-SGD a cependant un coût en performances et en ressources [15, 11] prohibitif pour des expériences à grandes échelles. Ici, l'amélioration de la confidentialité se fait à un coût moindre (en performance, en ressources car SGD est utilisé) mais n'apporte pas de garantie théorique.

## 5 Conclusion

Les résultats obtenus font pencher la balance pour la thèse suivante : face aux attaques par inférence d'appartenance, la parcimonie agit en tout point comme une défense ; elle diminue l'efficacité des attaques avec un coût relativement faible sur les performances des réseaux, y compris pour de la parcimonie structurée butterfly, à notre connaissance jamais explorée dans ce contexte dans la littérature.

Une extension à un ensemble plus riche de modèles et de données permettrait de valider l'intérêt de la parcimonie et de définir un point de départ fiable pour de futurs entraînements visant à diminuer ces risques de perte de confidentialité.

## Références

- [1] M. ABADI, A. CHU, I. GOODFELLOW, H. B. MCMAHAN, I. MIRONOV, K. TALWAR et Li. ZHANG : Deep learning with differential privacy. *In SIGSAC*, 2016.
- [2] K. ADAMCZEWSKI et M. PARK : Differential privacy meets neural network pruning. *Preprint*, 2023.
- [3] M. BELKIN, D. HSU, S. MA et S. MANDAL : Reconciling modern machine-learning practice and the classical bias-variance trade-off. *National Academy of Sciences USA*, 2019.
- [4] T. DAO, B. CHEN, N. S. SOHONI, A. DESAI, M. POLI, J. GROGAN, A. LIU, A. RAO, A. RUDRA et C. RÉ : Monarch : Expressive structured matrices for efficient and accurate training. *In ICML*, 2022.
- [5] J. FRANKLE et M. CARBIN : The lottery ticket hypothesis : Finding sparse, trainable neural networks. *In ICLR*, 2019.
- [6] J. FRANKLE, G. K. DZIUGAITE, D. ROY et M. CARBIN : Pruning neural networks at initialization : Why are we missing the mark? *In ICLR*, 2021.
- [7] K. HE, X. ZHANG, Sh. REN et J. SUN : Deep residual learning for image recognition. *In CVPR*. IEEE, 2016.
- [8] T. HOEFLE, D. ALISTARH, T. BEN-NUN, N. DRYDEN et A. PESTE : Sparsity in deep learning : Pruning and growth for efficient inference and training in neural networks. *Journal of Machine Learning Research*, 2021.
- [9] H. HU, Z. SALCIC, L. SUN, G. DOBBIE, P. S. YU et X. ZHANG : Membership inference attacks on machine learning : A survey. *ACM Computing Surveys*, 2022.
- [10] Y. HUANG, Y. SU, S. RAVI, Z. SONG, S. ARORA et K. LI : Privacy-preserving learning via deep net pruning. *Preprint*, 2020.
- [11] C. LALANNE, A. GARIVIER et R. GRIBONVAL : On the statistical complexity of estimation and testing under privacy constraints. *Preprint*, 2022.
- [12] Q.-T. LE, L. ZHENG, E. RICCIETTI et R. GRIBONVAL : Fast learning of fast transforms, with guarantees. *In ICASSP*. IEEE, 2022.
- [13] R. LIN, J. RAN, K. H. CHIU, G. CHESI et N. WONG : Deformable butterfly : A highly structured and sparse linear transform. *In NeurIPS*, 2021.
- [14] E. QUEMENER et M. CORVELLEC : SIDUS—the Solution for Extreme Deduplication of an Operating System. *Linux Journal*, 2013.
- [15] T. SANDER, P. STOCK et A. SABLAYROLLES : Tan without a burn : Scaling laws of dp-sgd. *Preprint*, 2022.
- [16] R. SHOKRI, M. STRONATI, C. SONG et V. SHMATIKOV : Membership inference attacks against machine learning models. *In SP*. IEEE, 2017.
- [17] Y. WANG, C. WANG, Z. WANG, S. ZHOU, H. LIU, J. BI, C. DING et S. RAJASEKARAN : Against membership inference attack : Pruning is all you need. *IJCAI*, 2021.
- [18] X. YUAN et L. ZHANG : Membership inference attacks and defenses in neural network pruning. *In USENIX*. IEEE, 2022.
- [19] C. ZHANG, S. BENGIO, M. HARDT, B. RECHT et O. VINYALS : Understanding deep learning (still) requires rethinking generalization. *ACM*, 2021.