

# Une Analyse des Méthodes de Projections Aléatoires par la Théorie des Matrices Aléatoires

Zhenyu Liao, Romain Couillet

► **To cite this version:**

Zhenyu Liao, Romain Couillet. Une Analyse des Méthodes de Projections Aléatoires par la Théorie des Matrices Aléatoires. Colloque GRETSI'17, 2017, Juan Les Pins, France. hal-01961878

**HAL Id: hal-01961878**

**<https://hal.archives-ouvertes.fr/hal-01961878>**

Submitted on 20 Dec 2018

**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.

# Une Analyse des Méthodes de Projections Aléatoires par la Théorie des Matrices Aléatoires

Zhenyu LIAO, Romain COUILLET \*

Laboratoire des signaux et systèmes, CentraleSupélec, CNRS, Université Paris-Sud  
3, Rue Joliot-Curie, 91192, Gif-sur-Yvette, France

`zhenyu.liao@l2s.centralesupelec.fr, romain.couillet@centralesupelec.fr`

**Résumé** – Dans cet article, nous présentons nos résultats théoriques sur les méthodes de projections aléatoires pour les données de très grandes dimensions, exploitées notamment dans les réseaux de neurones aléatoires. En combinant les avancées récentes de la théorie des matrices aléatoires et l’outil de la concentration des mesures, nous étudions en particulier les performances asymptotiques des dites *extreme learning machines*, permettant ainsi une compréhension plus profonde de ce réseau de neurone aléatoire simple.

**Abstract** – In this paper we present our theoretical findings on *random feature maps* or *random projections* for large dimensional data in large systems, classically found in different types of random neural networks. Combining recent advances in random matrix theory and the concentration of measure phenomenon, we provide a new framework on the theoretical analysis of random feature maps in large dimensional problems. As a concrete example, a study of the asymptotic performance of *extreme learning machines* is provided.

## 1 Introduction

Les réseaux de neurones grands et profonds sont devenus aujourd’hui des incontournables pour l’apprentissage (à partir de grandes bases de données), notamment de classificateurs ou modèles prédictifs non linéaires. Cependant, l’apprentissage de réseaux de neurones profonds reste aujourd’hui coûteux, peu adaptable et présuppose toujours de grandes bases de données d’apprentissage.

Une alternative possible, contre-intuitive au premier abord, consiste à utiliser des méthodes de projections aléatoires afin d’extraire les caractéristiques (*features* en anglais) des données permettant d’effectuer les tâches d’apprentissage. Dans ce nouvel espace projectif, des méthodes d’apprentissage classiques (régression linéaire ou logistique, machine à vecteurs de support, etc.) peuvent alors être utilisées. Ces méthodes apparaissent sous diverses formes dans la littérature : sous le nom des *extreme learning machines* (ELMs) [3] dans le contexte des réseaux de neurones, de *random feature maps* dans le cas d’approximations aléatoires de noyaux [8], de réseaux de neurones récurrents aléatoires dits *echo state networks* [4], etc.

Ce type de réseaux connaît un récent regain d’intérêt ces dernières années, de par ses nombreux avantages : 1) simplicité d’utilisation et adaptabilité à des grands nombres de données, 2) surapprentissage limité sous un nombre faible de données, 3) reparamétrisation facile et adaptabilité aux flux des données (à l’aide, par exemple, de la méthode des moindres carrés récurrents [2]) et 4) une promesse théorique du domaine des pro-

jections aléatoires selon laquelle  $k \log(k)$  projections aléatoires sont suffisantes pour extraire l’information pertinente de données de très grandes dimensions représentables dans une base de dimension  $k$ .

L’étude théorique de ces méthodes de projections aléatoires se limite souvent à l’hypothèse d’un nombre de données  $T$  de dimension  $p$  constants tandis que le nombre de neurones  $n$  tend vers l’infini, comme dans l’exemple populaire des *random fourier features* [8]. Dans cet article, nous nous plaçons dans le régime plus pertinent où  $n, p, T$  sont grands mais commensurables. En se basant sur l’outil de la théorie des grandes matrices aléatoires, nous traitons la non-linéarité délicate intervenant dans ces méthodes au moyen du phénomène de concentration de la mesure [6] ainsi qu’en exploitant les degrés d’indépendance structurels du modèle. La contribution principale du travail porte sur l’analyse spectrale de la matrice corrélation des données dans l’espace des projections aléatoires non-linéaires. Cette étude est ensuite appliquée au problème concret de l’analyse des performances asymptotiques des méthodes de réseaux de neurones de type *extreme learning machine* (ELM) [3].

*Notations* : La norme  $\|\cdot\|$  est la norme euclidienne pour des vecteurs et la norme spectrale pour des matrices, tandis que  $\|\cdot\|_F$  est la norme de Frobenius. La notation  $\xrightarrow{p.s.}$  indique la convergence presque sûre.

## 2 Énoncé du Problème

Les projecteurs aléatoires consistent à construire, à partir des données d’entraînement  $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_T] \in \mathbb{R}^{p \times T}$  une

\* Ce travail a été soutenu par le projet ANR RMT4GRAPH (ANR-14-CE28-0006). Ce travail était présenté partiellement au 42nd International Conference on Acoustics, Speech and Signal Processing (ICASSP’17), New Orleans, USA.

matrice de ‘‘caractéristiques’’ (*features*)  $\Sigma \equiv [\varphi_1, \dots, \varphi_T] \in \mathbb{R}^{n \times T}$  obtenue en multipliant  $\mathbf{X}$  par une matrice de poids aléatoire  $\mathbf{W} \in \mathbb{R}^{n \times p}$ , souvent telle que les  $W_{ij}$  sont i.i.d.  $\mathcal{N}(0, 1)$ , et ensuite en appliquant une fonction d’activation non-linéaire  $\sigma : \mathbb{R} \mapsto \mathbb{R}$  élément par élément à la matrice  $\mathbf{W}\mathbf{X}$ , pour obtenir donc la matrice  $\Sigma = \sigma(\mathbf{W}\mathbf{X}) \in \mathbb{R}^{n \times T}$ .

Dans le cas des *random fourier features* ou bien le cas de l’ELM, les méthodes d’apprentissage sont toujours basées sur le spectre de la matrice de covariance empirique dans l’espace des caractéristiques  $\frac{1}{T}\Sigma^T\Sigma$ , qui devient donc naturellement l’objet central à analyser. La théorie des grandes matrices aléatoires prédit que le spectre de cette matrice est essentiellement relié à la transformée de Stieltjes de la résolvante associée

$$\mathbf{Q}(z) \equiv \left( \frac{1}{T}\Sigma^T\Sigma - z\mathbf{I}_T \right)^{-1}$$

pour  $z \in \mathbb{C} \setminus \mathbb{R}^+$ , qui est par conséquent l’objet clé dans la compréhension de ces méthodes.

Afin d’étudier le comportement asymptotique des projections aléatoires dans le régime où  $n, p, T \rightarrow \infty$ , nous travaillons sous les hypothèses de croissance suivantes.

**Hypothèse 1** (Taux de Croissance). *Lorsque  $n \rightarrow \infty$ ,*

1.  $p/n \rightarrow c_0 \in (0, \infty)$ ,  $T/n \rightarrow c_T \in (0, \infty)$
2.  $\|\mathbf{X}\| = O(1)$
3. *pour évaluer la performance d’ELM, nous avons besoin en plus de la sortie associée  $\mathbf{Y}$  telle que  $\mathbf{Y}_{ij} = O(1)$ .*

Nous avons par ailleurs besoin d’imposer la condition de régularité suivante de la fonction d’activation  $\sigma$ .

**Hypothèse 2** (Fonction d’activation  $\sigma$ ). *La fonction  $\sigma$  est Lipschitzienne avec constante de Lipschitz  $\lambda_\sigma$  indépendante de  $n$ .*

La plus grande difficulté de cette analyse porte sur la non-linéarité de la matrice  $\Sigma$ . Ce problème est résolu en exploitant la caractère Lipschitz de  $\sigma(\cdot)$  qui permet alors d’appliquer des inégalités de concentration de la mesure sur l’application  $\mathbf{W} \mapsto \sigma(\mathbf{W}\mathbf{X})$ . Par ce biais, nous parvenons à démontrer que l’espérance  $\mathbb{E}[\mathbf{Q}]$  de la matrice aléatoire  $\mathbf{Q}$  a un comportement asymptotique dans le régime  $n, p, T \rightarrow \infty$  qui s’exprime en termes d’objets complètement déterministes.

### 3 Résultats Principaux

Notre analyse commence par un résultat important lié au phénomène de concentration de la mesure [6]. Les vecteurs aléatoires gaussiens satisfont une propriété dite de *concentration normale* qui peut alors se propager à travers toute application Lipschitzienne. Ceci nous permet en particulier d’assurer que, pour  $\mathbf{w} \sim \mathcal{N}(0, \mathbf{I}_p)$  et  $\mathbf{X} \in \mathbb{R}^{p \times T}$  tel que  $\|\mathbf{X}\| = 1$ , la norme de  $\frac{1}{\sqrt{p}}\sigma \equiv \frac{1}{\sqrt{p}}\sigma(\mathbf{X}^T\mathbf{w}) \in \mathbb{R}^T$  est d’ordre  $O(1)$  avec haute probabilité. De plus, conditionnellement au fait que  $\frac{1}{\sqrt{p}}\sigma$

est de norme  $O(1)$ , l’application  $\mathbf{w} \mapsto \frac{1}{p}\sigma^T\mathbf{A}\sigma$  est Lipschitzienne. Avec un contrôle précis des résultats de concentration sous conditionnement, nous obtenons le lemme suivant.

**Lemme 1** (Concentration de la forme quadratique). *Sous l’hypothèse 1–2, pour  $\mathbf{w} \sim \mathcal{N}(0, \mathbf{I}_p)$  et  $\sigma \equiv \sigma(\mathbf{X}^T\mathbf{w}) \in \mathbb{R}^T$ . Nous avons pour toute matrice  $\mathbf{A} \in \mathbb{R}^{T \times T}$  telle que  $\|\mathbf{A}\| \leq 1$  et indépendante de  $\sigma$ ,*

$$P \left( \left| \frac{1}{n}\sigma^T\mathbf{A}\sigma - \frac{1}{n}\text{tr}(\Phi\mathbf{A}) \right| > t \right) \leq Ce^{-cn \min(t, t^2)}$$

pour certains  $c, C > 0$  indépendants de  $n$ , avec  $\Phi \equiv \mathbb{E}[\sigma\sigma^T]$ .

Le lemme 1 nous permet d’estimer les formes quadratiques  $\frac{1}{n}\sigma^T\mathbf{A}\sigma$  par  $\frac{1}{n}\text{tr}(\Phi\mathbf{A})$  avec probabilité exponentiellement proche de 1, qui est essentiellement la version non-linéaire du lemme de la trace (e.g., [1, Lemma B.26]) dans la théorie des grandes matrices aléatoires. Ainsi, en utilisant maintenant ce lemme ainsi que le fait que la matrice

$$\Sigma^T \equiv \sigma(\mathbf{X}^T\mathbf{W}^T) = [\sigma_1, \dots, \sigma_n]$$

avec  $\mathbf{W}^T = [\mathbf{w}_1, \dots, \mathbf{w}_n]$  et  $\sigma_i = \sigma(\mathbf{X}^T\mathbf{w}_i) \in \mathbb{R}^T$ , a des colonnes indépendantes (mais des lignes fortement dépendantes), nous pouvons utiliser des méthodes dorénavant classiques de la théorie des matrices aléatoires afin d’obtenir le résultat suivant.

**Théorème 1** (Équivalent asymptotique de  $\mathbb{E}[\mathbf{Q}]$ ). *Sous l’hypothèse 1–2, pour tout  $\varepsilon > 0$ , il existe  $c > 0$  tel que, pour tout  $n$ ,*

$$\|\mathbb{E}[\mathbf{Q}] - \bar{\mathbf{Q}}\| \leq cn^{\varepsilon - \frac{1}{2}}.$$

où  $\bar{\mathbf{Q}}$  est donnée par

$$\bar{\mathbf{Q}} = \left( \frac{n}{T} \frac{\Phi}{1 + \delta} - z\mathbf{I}_T \right)^{-1}$$

avec  $\delta$  l’unique solution positive de l’équation  $\delta = \frac{1}{T}\text{tr}(\Phi\bar{\mathbf{Q}})$ .

*Trame de la Preuve.* Nous introduisons  $\alpha \equiv \frac{1}{T}\text{tr}(\Phi\mathbb{E}[\mathbf{Q}])$  et  $\tilde{\mathbf{Q}} \equiv \left( \frac{n}{T} \frac{\Phi}{1 + \alpha} - z\mathbf{I}_T \right)^{-1}$ . Nous déduisons de l’identité de la résolvante  $\mathbf{A}^{-1} - \mathbf{B}^{-1} = \mathbf{A}^{-1}(\mathbf{B} - \mathbf{A})\mathbf{B}^{-1}$  que

$$\begin{aligned} \mathbb{E}[\mathbf{Q}] - \tilde{\mathbf{Q}} &= \mathbb{E}[\mathbf{Q} - \tilde{\mathbf{Q}}] = \mathbb{E}[\mathbf{Q}] \frac{n}{T} \frac{\Phi}{1 + \alpha} \tilde{\mathbf{Q}} - \mathbb{E} \left[ \mathbf{Q} \frac{1}{T} \Sigma^T \Sigma \right] \tilde{\mathbf{Q}} \\ &= \mathbb{E}[\mathbf{Q}] \frac{n}{T} \frac{\Phi}{1 + \alpha} \tilde{\mathbf{Q}} - \frac{1}{T} \sum_{i=1}^n \mathbb{E} [\mathbf{Q} \sigma_i \sigma_i^T] \tilde{\mathbf{Q}}. \end{aligned}$$

Avec la formule de Sherman-Morrison, pour

$$\mathbf{Q}_{-i} \equiv \left( \frac{1}{T} \Sigma^T \Sigma - \frac{1}{T} \sigma_i \sigma_i^T - z\mathbf{I}_T \right)^{-1}$$

nous obtenons

$$\begin{aligned} \mathbb{E}[\mathbf{Q}] - \tilde{\mathbf{Q}} &= \mathbb{E}[\mathbf{Q}] \frac{n}{T} \frac{\Phi}{1 + \alpha} \tilde{\mathbf{Q}} - \frac{1}{T} \sum_{i=1}^n \mathbb{E} \left[ \frac{\mathbf{Q}_{-i} \sigma_i \sigma_i^T}{1 + \frac{1}{T} \sigma_i^T \mathbf{Q}_{-i} \sigma_i} \right] \tilde{\mathbf{Q}} \\ &= \mathbb{E}[\mathbf{Q}] \frac{n}{T} \frac{\Phi}{1 + \alpha} \tilde{\mathbf{Q}} - \frac{1}{T} \frac{1}{1 + \alpha} \sum_{i=1}^n \mathbb{E} [\mathbf{Q}_{-i} \sigma_i \sigma_i^T] \tilde{\mathbf{Q}} \\ &\quad + \frac{1}{T} \frac{1}{1 + \alpha} \sum_{i=1}^n \mathbb{E} \left[ \frac{\mathbf{Q}_{-i} \sigma_i \sigma_i^T \left( \frac{1}{T} \sigma_i^T \mathbf{Q}_{-i} \sigma_i - \alpha \right)}{1 + \frac{1}{T} \sigma_i^T \mathbf{Q}_{-i} \sigma_i} \right] \tilde{\mathbf{Q}}. \end{aligned}$$

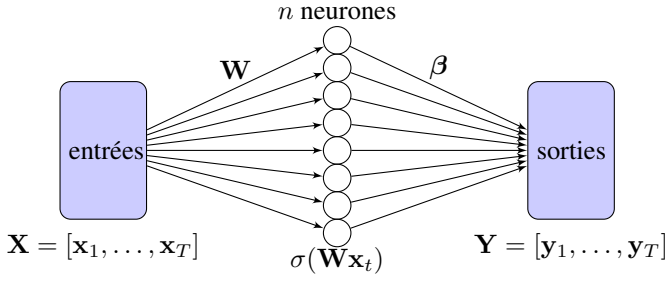


FIGURE 1 – Illustration d l'ELM

Puisque  $\mathbf{Q}_{-i}$  est indépendant de  $\sigma_i$ , nous obtenons alors

$$\begin{aligned} \mathbb{E}[\mathbf{Q}] - \tilde{\mathbf{Q}} &= \mathbb{E}[\mathbf{Q}] \frac{n}{T} \frac{\Phi}{1+\alpha} \tilde{\mathbf{Q}} - \frac{1}{T} \frac{1}{1+\alpha} \sum_{i=1}^n \mathbb{E}[\mathbf{Q}_{-i}] \Phi \tilde{\mathbf{Q}} \\ &+ \frac{1}{T} \frac{1}{1+\alpha} \mathbb{E}[\mathbf{Q} \Sigma^T \mathbf{D} \Sigma] \tilde{\mathbf{Q}} \end{aligned}$$

avec  $\mathbf{D} = \text{diag}(\{\frac{1}{T} \sigma_i^T \mathbf{Q}_{-i} \sigma_i - \alpha\}_{i=1}^n)$  qui est telle que  $\|\mathbf{D}\| \rightarrow 0$  quand  $n \rightarrow \infty$  d'après le lemme 1. De plus, nous montrons que  $\|\mathbb{E}[\mathbf{Q} - \mathbf{Q}_{-i}]\| = O(n^{-1})$ , ce qui nous dit que

$$\|\mathbb{E}[\mathbf{Q}] - \tilde{\mathbf{Q}}\| \rightarrow 0,$$

ensuite en observant le fait que

$$\alpha = \frac{1}{T} \text{tr}(\Phi \mathbb{E}[\mathbf{Q}]) = \frac{1}{T} \text{tr}\left(\Phi \left(\frac{n}{T} \frac{\Phi}{1+\alpha} - z \mathbf{I}_T\right)^{-1}\right) + o(1)$$

avec  $\delta \equiv \frac{1}{T} \text{tr}(\Phi \tilde{\mathbf{Q}})$  nous concluons la preuve.  $\square$

## 4 Exemple Concret de l'ELM

Dans cette section, nous nous concentrons sur l'exemple concret de l'ELM, qui est une simple régression linéaire normalisée de  $\Sigma$  par rapport à une sortie désirée  $\mathbf{Y} = [y_1, \dots, y_T] \in \mathbb{R}^{d \times T}$  associée à  $\mathbf{X}$  donnée par le régresseur

$$\beta = \frac{1}{T} \Sigma \left( \frac{1}{T} \Sigma^T \Sigma + \gamma \mathbf{I}_T \right)^{-1} \mathbf{Y}^T = \frac{1}{T} \Sigma \mathbf{Q} (-\gamma) \mathbf{Y}^T \quad (1)$$

avec  $\gamma$  le facteur de régularisation strictement positif, comme illustré dans la Figure 1.

L'erreur d'entraînement de l'ELM pour les données d'entraînement  $\mathbf{X}$  est ainsi donnée par

$$E_{\text{train}} = \frac{1}{T} \left\| \mathbf{Y}^T - \Sigma^T \beta \right\|_F^2 = \frac{\gamma^2}{T} \text{tr}(\mathbf{Y}^T \mathbf{Y} \mathbf{Q}^2). \quad (2)$$

Pendant la *phase de test*, nous considérons l'ensemble de données de test :  $\hat{\mathbf{X}} \in \mathbb{R}^{p \times \hat{T}}$  de taille  $\hat{T}$ , dont les sorties correspondantes sont notées  $\hat{\mathbf{Y}} \in \mathbb{R}^{d \times \hat{T}}$ . En appliquant le régresseur  $\beta$  défini par (1) (qui ne dépend que de  $\Sigma$  et  $\mathbf{Y}$ ), nous obtenons l'erreur de test

$$E_{\text{test}} = \frac{1}{T} \left\| \hat{\mathbf{Y}}^T - \hat{\Sigma}^T \beta \right\|_F^2 \quad (3)$$

avec  $\hat{\Sigma} \equiv \sigma(\mathbf{W} \hat{\mathbf{X}})$ . Pour la convenance de l'étude de  $E_{\text{test}}$  dans la suite, nous étendons la définition de  $\Phi$  dans Lemme 1 pour toute paire de matrices  $(\mathbf{A}, \mathbf{B})$  à

$$\Phi_{\mathbf{AB}} \equiv \mathbb{E}[\sigma(\mathbf{w}^T \mathbf{A})^T \sigma(\mathbf{w}^T \mathbf{B})]$$

avec  $\mathbf{w} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_p)$  et notons

$$\Psi_{\mathbf{AB}} \equiv \frac{n}{T} \frac{\Phi_{\mathbf{AB}}}{1+\delta}.$$

Pour simplifier les notations, nous notons  $\Phi = \Phi_{\mathbf{XX}}$  et  $\Psi = \Psi_{\mathbf{XX}}$ .

En étendant le lemme 1, nous pouvons démontrer que les erreurs  $E_{\text{train}}$  et  $E_{\text{test}}$  convergent (donc vers leur espérance). Concernant  $\mathbb{E}[E_{\text{train}}]$ , nous avons notamment

$$\mathbb{E}[E_{\text{train}}] = \frac{\gamma^2}{T} \text{tr}(\mathbf{Y}^T \mathbf{Y} \mathbb{E}[\mathbf{Q}^2])$$

qui nous conduit à l'analyse asymptotique de  $\mathbb{E}[\mathbf{Q}^2]$  ou plus généralement de  $\mathbb{E}[\mathbf{QAQ}]$  pour toute matrice  $\mathbf{A}$  déterministe, qui nécessite plus d'efforts et est annoncé sans preuve ici. Nous renvoyons les lecteurs à [7] pour la démonstration complète.

**Proposition 1** (Équivalent asymptotique de  $\mathbb{E}[\mathbf{QAQ}]$ ). *Soit  $\mathbf{A}$  la matrice  $\Phi$  ou une matrice symétrique de norme bornée. Alors, sous l'hypothèse 1-2, avec la matrice  $\tilde{\mathbf{Q}}$  donnée dans le théorème 1, pour tout  $\varepsilon > 0$ , il existe une constante  $c$  telle que*

$$\left\| \mathbb{E}[\mathbf{QAQ}] - \tilde{\mathbf{Q}} \mathbf{A} \tilde{\mathbf{Q}} - \tilde{\mathbf{Q}} \Psi \tilde{\mathbf{Q}} \frac{\frac{1}{n} \text{tr}(\Psi \tilde{\mathbf{Q}} \mathbf{A} \tilde{\mathbf{Q}})}{1 - \frac{1}{n} \text{tr}(\Psi^2 \tilde{\mathbf{Q}}^2)} \right\| \leq c n^{\varepsilon - \frac{1}{2}}.$$

En prenant  $\mathbf{A}$  la matrice d'identité, nous accédons à l'équivalent déterministe de l'erreur d'entraînement  $E_{\text{train}}$  dans le régime où  $n, p, T \rightarrow \infty$ . La dérivation de l'erreur de test nécessite plus d'efforts et sera annoncée sans démonstration dans le théorème suivant.

**Théorème 2** (Performance de l'ELM). *Sous l'hypothèse 1-2, pour tout  $\varepsilon > 0$ ,*

$$\begin{aligned} n^{\frac{1}{2}-\varepsilon} (E_{\text{train}} - \bar{E}_{\text{train}}) &\xrightarrow{p.s.} 0 \\ n^{\frac{1}{2}-\varepsilon} (E_{\text{test}} - \bar{E}_{\text{test}}) &\xrightarrow{p.s.} 0 \end{aligned}$$

où

$$\bar{E}_{\text{train}} = \frac{\gamma^2}{T} \text{tr} \left( \mathbf{Y}^T \mathbf{Y} \tilde{\mathbf{Q}} \left[ \frac{\frac{1}{n} \text{tr}(\tilde{\mathbf{Q}} \Psi \tilde{\mathbf{Q}})}{1 - \frac{1}{n} \text{tr}(\Psi^2 \tilde{\mathbf{Q}}^2)} \Psi + \mathbf{I}_T \right] \tilde{\mathbf{Q}} \right) \quad (4)$$

et  $\bar{E}_{\text{test}}$  est donnée par l'équation (6).

Notons que  $\tilde{\mathbf{Q}}$  qui dépend implicitement de la matrice  $\Phi$  est au cœur du Théorème 2. L'évaluation de la matrice  $\Phi_{\mathbf{AB}}$  pour deux matrices arbitraires  $\mathbf{A}, \mathbf{B}$  requiert naturellement l'évaluation de ses entrées et donc revient à calculer, pour deux vecteurs arbitraires  $\mathbf{a}, \mathbf{b} \in \mathbb{R}^p$ , la valeur de  $\Phi(\mathbf{a}, \mathbf{b}) \equiv \mathbb{E}[\sigma(\mathbf{w}^T \mathbf{a}) \sigma(\mathbf{w}^T \mathbf{b})]$ , avec  $\mathbf{w} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_p)$ . Les résultats pour quelques fonctions d'activation  $\sigma$  couramment utilisées sont fournis dans Table 1. Encore une fois, nous renvoyons les lecteurs à [7] dans lequel tous les détails de calculs sont donnés.

La matrice  $\Phi$  joue ici le rôle de la matrices à noyau qui serait obtenue dans la limite  $n \rightarrow \infty$  seul comme prédit dans [8]. Il est facile de voir en effet ici que pour  $T/n, p/n \rightarrow 0$ , l'ELM devient équivalent à une régression linéaire par noyau. Ces liens entre matrice à noyau et projections aléatoires sont aussi mis en évidence dans le contexte d'autres méthodes d'apprentissage, ce qui confère à cette étude une dimension d'ouverture plus large que seulement restreintes aux ELM.

$$\bar{E}_{\text{test}} = \frac{1}{\hat{T}} \left\| \hat{\mathbf{Y}}^\top - \Psi_{\mathbf{X}\mathbf{X}}^\top \bar{\mathbf{Q}} \mathbf{Y}^\top \right\|_F^2 + \frac{\frac{1}{n} \text{tr}(\mathbf{Y}^\top \mathbf{Y} \bar{\mathbf{Q}} \Psi \bar{\mathbf{Q}})}{1 - \frac{1}{n} \text{tr}(\Psi^2 \bar{\mathbf{Q}}^2)} \left[ \frac{1}{\hat{T}} \text{tr} \Psi_{\hat{\mathbf{X}}\hat{\mathbf{X}}} - \frac{1}{\hat{T}} \text{tr}(\mathbf{I}_T + \gamma \bar{\mathbf{Q}})(\Psi_{\mathbf{X}\mathbf{X}} \Psi_{\hat{\mathbf{X}}\hat{\mathbf{X}}} \bar{\mathbf{Q}}) \right]. \quad (6)$$

$\sigma(t)$	$\Phi_{\mathbf{a}\mathbf{b}}$
$t$	$\mathbf{a}^\top \mathbf{b}$
$\max(t, 0)$	$\frac{1}{2\pi} \ \mathbf{a}\  \ \mathbf{b}\  \left( \angle(\mathbf{a}, \mathbf{b}) \arccos(-\angle(\mathbf{a}, \mathbf{b})) + \sqrt{1 - \angle(\mathbf{a}, \mathbf{b})^2} \right)$
$\text{erf}(t)$	$\frac{2}{\pi} \arcsin \left( \frac{2\mathbf{a}^\top \mathbf{b}}{\sqrt{(1+2\ \mathbf{a}\ ^2)(1+2\ \mathbf{b}\ ^2)}} \right)$
$\cos(t)$	$\exp\left(-\frac{1}{2}(\ \mathbf{a}\ ^2 + \ \mathbf{b}\ ^2)\right) \cosh(\mathbf{a}^\top \mathbf{b})$
$\sin(t)$	$\exp\left(-\frac{1}{2}(\ \mathbf{a}\ ^2 + \ \mathbf{b}\ ^2)\right) \sinh(\mathbf{a}^\top \mathbf{b})$

TABLE 1 – Valeur de  $\Phi_{\mathbf{a}\mathbf{b}}$  pour  $w \sim \mathcal{N}(0, \mathbf{I}_p)$ ,  $\angle(\mathbf{a}, \mathbf{b}) \equiv \frac{\mathbf{a}^\top \mathbf{b}}{\|\mathbf{a}\| \|\mathbf{b}\|}$ .

## 5 Illustrations Numériques

Nous testons maintenant nos résultats théoriques sur la base de données MNIST [5]. Ici, nous considérons une tâche de classification binaire des images des chiffres manuscrits de taille  $28 \times 28$ , chaque image étant représentée par un vecteur de taille  $p = 784$ , avec une ELM de  $n = 512$  unités de neurones et  $\mathbf{W}$  de type Gaussien standard. La base de donnée est composée de  $T = \hat{T} = 1024$  échantillons : 1024 images du chiffre sept et 1024 du chiffre neuf sont extraits aléatoirement et repartis en 512 échantillons d’entraînement et 512 échantillons de test.

Dans la figure 2 nous évaluons la performance de l’ELM pour certaines fonctions Lipschitziennes  $\sigma$  (linéaire,  $\text{erf}(t)$ ) ainsi que  $\text{ReLU}(t) \equiv \max(t, 0)$  en fonction de l’hyper-paramètre  $\gamma$ . Nous comparons l’erreur quadratique moyenne empirique à l’approximant asymptotique donné par le théorème 2. Nous pouvons constater une précision impressionnante de nos résultats théoriques dans ce contexte, quand bien même limité à assez peu de données.

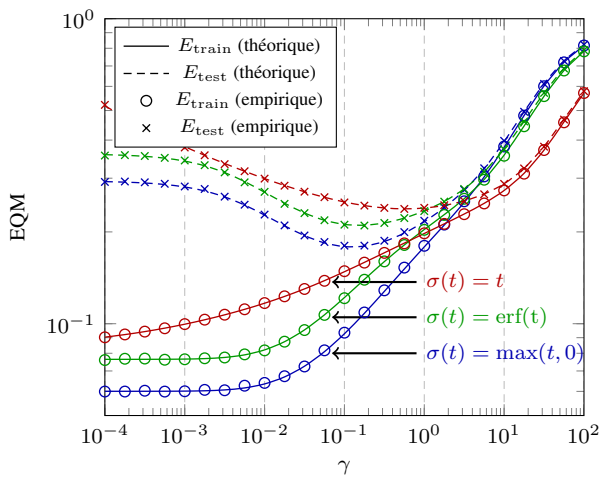


FIGURE 2 – Performance de l’ELM pour  $\sigma$  Lipschitzienne, en fonction de  $\gamma$ , pour 2-classe données de MNIST (sept et neuf),  $n = 512$ ,  $T = \hat{T} = 1024$ ,  $p = 784$ .

## 6 Conclusion

Ce travail apporte une nouvelle compréhension des méthodes de projections aléatoires à l’aide de l’outil de la théorie des grandes matrices aléatoires. Il ouvre également la voie à l’analyse et l’amélioration d’algorithmes standards dans le contexte plus large de l’apprentissage automatisé des grands systèmes, dont les réseaux de neurones, qui sont au centre des questions modernes du traitement du signal et des données.

## Références

- [1] Zhidong Bai and Jack W Silverstein. *Spectral analysis of large dimensional random matrices*, volume 20. Springer, 2010.
- [2] Monson H. Hayes. *Statistical Digital Signal Processing and Modeling*. John Wiley & Sons, Inc., New York, NY, USA, 1st edition, 1996.
- [3] Guang-Bin Huang, Hongming Zhou, Xiaojian Ding, and Rui Zhang. Extreme learning machine for regression and multiclass classification. *IEEE Transactions on Systems, Man, and Cybernetics, Part B (Cybernetics)*, 42(2) :513–529, 2012.
- [4] Herbert Jaeger. The “echo state” approach to analysing and training recurrent neural networks—with an erratum note. *Bonn, Germany : German National Research Center for Information Technology GMD Technical Report*, 148(34) :13, 2001.
- [5] Yann LeCun, Corinna Cortes, and Christopher JC Burges. The MNIST database of handwritten digits, 1998.
- [6] Michel Ledoux. *The concentration of measure phenomenon*. Number 89. American Mathematical Soc., 2005.
- [7] Cosme Louart, Zhenyu Liao, and Romain Couillet. A random matrix approach to neural networks. *arXiv preprint arXiv :1702.05419*, 2017.
- [8] Ali Rahimi, Benjamin Recht, et al. Random Features for Large-Scale Kernel Machines. In *NIPS*, volume 3, page 5, 2007.