

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

Colloque GRETSI 2017

Zhenyu Liao, Romain Couillet

CentraleSupélec
Université Paris-Saclay

Juan-Les-Pins, 05 Septembre, 2017



CentraleSupélec

1 Introduction

2 Résultats Principaux

3 Conclusion

1 Introduction

2 Résultats Principaux

3 Conclusion

Motivation: Extraire des features en utilisant des projections aléatoires

Apprentissage = Représentation + Évaluation + Optimisation.¹

¹Domingos, Pedro. "A few useful things to know about machine learning." Communications of the ACM 55.10 (2012): 78-87.

Motivation: Extraire des features en utilisant des projections aléatoires

Apprentissage = Représentation + Évaluation + Optimisation.¹

Comment trouver les “bonnes” features?

- ▶ apprendre les features à partir des données:

¹Domingos, Pedro. "A few useful things to know about machine learning." Communications of the ACM 55.10 (2012): 78-87.

Motivation: Extraire des features en utilisant des projections aléatoires

Apprentissage = Représentation + Évaluation + Optimisation.¹

Comment trouver les “bonnes” features?

- ▶ apprendre les features à partir des données:
 - ▶ réseaux de neurones: backprop coûteux

¹Domingos, Pedro. "A few useful things to know about machine learning." Communications of the ACM 55.10 (2012): 78-87.

Motivation: Extraire des features en utilisant des projections aléatoires

Apprentissage = Représentation + Évaluation + Optimisation.¹

Comment trouver les “bonnes” features?

- ▶ apprendre les features à partir des données:
 - ▶ réseaux de neurones: backprop coûteux
 - ▶ problème de sur-apprentissage: jeux de données très grands

¹Domingos, Pedro. "A few useful things to know about machine learning." Communications of the ACM 55.10 (2012): 78-87.

Motivation: Extraire des features en utilisant des projections aléatoires

Apprentissage = Représentation + Évaluation + Optimisation.¹

Comment trouver les “bonnes” features?

- ▶ apprendre les features à partir des données:
 - ▶ réseaux de neurones: backprop coûteux
 - ▶ problème de sur-apprentissage: jeux de données très grands
 - ▶ absence de compréhension théorique

¹Domingos, Pedro. "A few useful things to know about machine learning." Communications of the ACM 55.10 (2012): 78-87.

Motivation: Extraire des features en utilisant des projections aléatoires

Apprentissage = Représentation + Évaluation + Optimisation.¹

Comment trouver les “bonnes” features?

- ▶ apprendre les features à partir des données:
 - ▶ réseaux de neurones: backprop coûteux
 - ▶ problème de sur-apprentissage: jeux de données très grands
 - ▶ absence de compréhension théorique
- ▶ méthodes aléatoires: projections aléatoires, compressive sensing

¹Domingos, Pedro. "A few useful things to know about machine learning." Communications of the ACM 55.10 (2012): 78-87.

Motivation: Extraire des features en utilisant des projections aléatoires

Apprentissage = Représentation + Évaluation + Optimisation.¹

Comment trouver les “bonnes” features?

- ▶ apprendre les features à partir des données:
 - ▶ réseaux de neurones: backprop coûteux
 - ▶ problème de sur-apprentissage: jeux de données très grands
 - ▶ absence de compréhension théorique
- ▶ méthodes aléatoires: projections aléatoires, compressive sensing
 - ▶ simples et rapides

¹Domingos, Pedro. "A few useful things to know about machine learning." Communications of the ACM 55.10 (2012): 78-87.

Motivation: Extraire des features en utilisant des projections aléatoires

Apprentissage = Représentation + Évaluation + Optimisation.¹

Comment trouver les “bonnes” features?

- ▶ apprendre les features à partir des données:
 - ▶ réseaux de neurones: backprop coûteux
 - ▶ problème de sur-apprentissage: jeux de données très grands
 - ▶ absence de compréhension théorique
- ▶ méthodes aléatoires: projections aléatoires, compressive sensing
 - ▶ simples et rapides
 - ▶ adaptées aux applications en ligne (e.g., moindres carrés récursifs)

¹Domingos, Pedro. "A few useful things to know about machine learning." Communications of the ACM 55.10 (2012): 78-87.

Motivation: Extraire des features en utilisant des projections aléatoires

Apprentissage = Représentation + Évaluation + Optimisation.¹

Comment trouver les “bonnes” features?

- ▶ apprendre les features à partir des données:
 - ▶ réseaux de neurones: backprop coûteux
 - ▶ problème de sur-apprentissage: jeux de données très grands
 - ▶ absence de compréhension théorique
- ▶ méthodes aléatoires: projections aléatoires, compressive sensing
 - ▶ simples et rapides
 - ▶ adaptées aux applications en ligne (e.g., moindres carrés récursifs)
 - ▶ analyse théorique possible

¹Domingos, Pedro. "A few useful things to know about machine learning." Communications of the ACM 55.10 (2012): 78-87.

Motivation: Extraire des features en utilisant des projections aléatoires

Apprentissage = Représentation + Évaluation + Optimisation.¹

Comment trouver les “bonnes” features?

- ▶ apprendre les features à partir des données:
 - ▶ réseaux de neurones: backprop coûteux
 - ▶ problème de sur-apprentissage: jeux de données très grands
 - ▶ absence de compréhension théorique
- ▶ méthodes aléatoires: projections aléatoires, compressive sensing
 - ▶ simples et rapides
 - ▶ adaptées aux applications en ligne (e.g., moindres carrés récursifs)
 - ▶ analyse théorique possible
 - ▶ mais difficultés théoriques dans le cas non linéaire

¹Domingos, Pedro. "A few useful things to know about machine learning." Communications of the ACM 55.10 (2012): 78-87.

Motivation: Extraire des features en utilisant des projections aléatoires

Apprentissage = Représentation + Évaluation + Optimisation.¹

Comment trouver les “bonnes” features?

- ▶ apprendre les features à partir des données:
 - ▶ réseaux de neurones: backprop coûteux
 - ▶ problème de sur-apprentissage: jeux de données très grands
 - ▶ absence de compréhension théorique
- ▶ méthodes aléatoires: projections aléatoires, compressive sensing
 - ▶ simples et rapides
 - ▶ adaptées aux applications en ligne (e.g., moindres carrés récursifs)
 - ▶ analyse théorique possible
 - ▶ mais difficultés théoriques dans le cas non linéaire

Exemples dans la littérature:

- ▶ *random features maps*: projections non linéaires aléatoires

¹Domingos, Pedro. "A few useful things to know about machine learning." Communications of the ACM 55.10 (2012): 78-87.

Motivation: Extraire des features en utilisant des projections aléatoires

Apprentissage = Représentation + Évaluation + Optimisation.¹

Comment trouver les “bonnes” features?

- ▶ apprendre les features à partir des données:
 - ▶ réseaux de neurones: backprop coûteux
 - ▶ problème de sur-apprentissage: jeux de données très grands
 - ▶ absence de compréhension théorique
- ▶ méthodes aléatoires: projections aléatoires, compressive sensing
 - ▶ simples et rapides
 - ▶ adaptées aux applications en ligne (e.g., moindres carrés récurrents)
 - ▶ analyse théorique possible
 - ▶ mais difficultés théoriques dans le cas non linéaire

Exemples dans la littérature:

- ▶ *random features maps*: projections non linéaires aléatoires
- ▶ *extreme learning machine (ELM)*: réseaux de neurones aléatoires simples

¹Domingos, Pedro. "A few useful things to know about machine learning." Communications of the ACM 55.10 (2012): 78-87.

Motivation: Extraire des features en utilisant des projections aléatoires

Apprentissage = Représentation + Évaluation + Optimisation.¹

Comment trouver les “bonnes” features?

- ▶ apprendre les features à partir des données:
 - ▶ réseaux de neurones: backprop coûteux
 - ▶ problème de sur-apprentissage: jeux de données très grands
 - ▶ absence de compréhension théorique
- ▶ méthodes aléatoires: projections aléatoires, compressive sensing
 - ▶ simples et rapides
 - ▶ adaptées aux applications en ligne (e.g., moindres carrés récurrents)
 - ▶ analyse théorique possible
 - ▶ mais difficultés théoriques dans le cas non linéaire

Exemples dans la littérature:

- ▶ *random features maps*: projections non linéaires aléatoires
- ▶ *extreme learning machine (ELM)*: réseaux de neurones aléatoires simples
- ▶ *echo state nets*: réseaux de neurones récurrents simples

¹Domingos, Pedro. "A few useful things to know about machine learning." Communications of the ACM 55.10 (2012): 78-87.

Motivation: Extraire des features en utilisant des projections aléatoires

Apprentissage = Représentation + Évaluation + Optimisation.¹

Comment trouver les “bonnes” features?

- ▶ apprendre les features à partir des données:
 - ▶ réseaux de neurones: backprop coûteux
 - ▶ problème de sur-apprentissage: jeux de données très grands
 - ▶ absence de compréhension théorique
- ▶ méthodes aléatoires: projections aléatoires, compressive sensing
 - ▶ simples et rapides
 - ▶ adaptées aux applications en ligne (e.g., moindres carrés récurrents)
 - ▶ analyse théorique possible
 - ▶ mais difficultés théoriques dans le cas non linéaire

Exemples dans la littérature:

- ▶ *random features maps*: projections non linéaires aléatoires
- ▶ *extreme learning machine (ELM)*: réseaux de neurones aléatoires simples
- ▶ *echo state nets*: réseaux de neurones récurrents simples

Dans ce travail:

- ▶ analyse spectrale des matrices de projections aléatoires non linéaire en grande dimension

¹Domingos, Pedro. "A few useful things to know about machine learning." Communications of the ACM 55.10 (2012): 78-87.

Motivation: Extraire des features en utilisant des projections aléatoires

Apprentissage = Représentation + Évaluation + Optimisation.¹

Comment trouver les “bonnes” features?

- ▶ apprendre les features à partir des données:
 - ▶ réseaux de neurones: backprop coûteux
 - ▶ problème de sur-apprentissage: jeux de données très grands
 - ▶ absence de compréhension théorique
- ▶ méthodes aléatoires: projections aléatoires, compressive sensing
 - ▶ simples et rapides
 - ▶ adaptées aux applications en ligne (e.g., moindres carrés récurrents)
 - ▶ analyse théorique possible
 - ▶ mais difficultés théoriques dans le cas non linéaire

Exemples dans la littérature:

- ▶ *random features maps*: projections non linéaires aléatoires
- ▶ *extreme learning machine (ELM)*: réseaux de neurones aléatoires simples
- ▶ *echo state nets*: réseaux de neurones récurrents simples

Dans ce travail:

- ▶ analyse spectrale des matrices de projections aléatoires non linéaire en grande dimension
- ▶ application à l'analyse des performances asymptotiques de l'ELM

¹Domingos, Pedro. "A few useful things to know about machine learning." Communications of the ACM 55.10 (2012): 78-87.

Projections aléatoires: un exemple

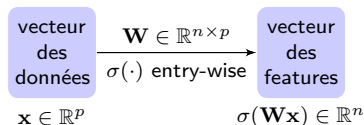


Figure: Illustration de la méthode de projections aléatoires

avec la matrice de poids $\mathbf{W} \in \mathbb{R}^{n \times p}$ telle que les \mathbf{W}_{ij} 's sont des variables aléatoires i.i.d. et la fonction d'activation $\sigma(\cdot)$ appliquée élément par élément.

Projections aléatoires: un exemple

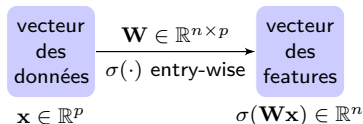


Figure: Illustration de la méthode de projections aléatoires

avec la matrice de poids $\mathbf{W} \in \mathbb{R}^{n \times p}$ telle que les \mathbf{W}_{ij} 's sont des variables aléatoires i.i.d. et la fonction d'activation $\sigma(\cdot)$ appliquée élément par élément.

Exemple: Classification des spams

$$\underbrace{\begin{bmatrix} 0 & 0 & 0 & 1 & 1 \\ 0 & 1 & 0 & 0 & 0 \\ 1 & 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 1 & 0 \end{bmatrix}}_{\mathbf{W}} \underbrace{\begin{bmatrix} \text{Cadeaux} \\ \text{gratuits} \\ \text{pour} \\ \text{votre} \\ \text{famille} \end{bmatrix}}_{\mathbf{x}} \xrightarrow{\sigma(t)=t} \underbrace{\begin{bmatrix} \text{votre+famille} \\ \text{Cadeaux} \\ \text{Cadeaux+gratuits} \\ \text{pour+famille} \end{bmatrix}}_{\sigma(\mathbf{W}\mathbf{x})}$$

Quand $n, p \rightarrow \infty$, on obtient de “bonnes” features, ensuite un prédicteur peut être appris de manière **supervisée**, e.g., en appliquant une régression normalisée.

Comprendre la méthode de projections aléatoires

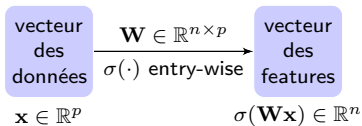


Figure: Illustration de la méthode de projections aléatoires

Comprendre la méthode de projections aléatoires

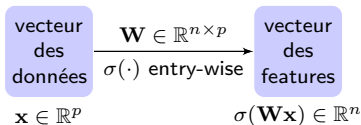


Figure: Illustration de la méthode de projections aléatoires

- Pour matrice de données $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_T] \in \mathbb{R}^{p \times T}$, on note la **matrice de feature** $\mathbf{\Sigma} \equiv \sigma(\mathbf{W}\mathbf{X}) = [\sigma(\mathbf{W}\mathbf{x}_1), \dots, \sigma(\mathbf{W}\mathbf{x}_T)] \in \mathbb{R}^{n \times T}$

Comprendre la méthode de projections aléatoires

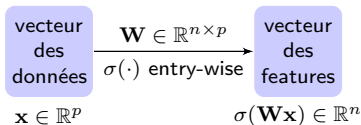


Figure: Illustration de la méthode de projections aléatoires

- Pour matrice de données $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_T] \in \mathbb{R}^{p \times T}$, on note la **matrice de feature** $\mathbf{\Sigma} \equiv \sigma(\mathbf{W}\mathbf{X}) = [\sigma(\mathbf{W}\mathbf{x}_1), \dots, \sigma(\mathbf{W}\mathbf{x}_T)] \in \mathbb{R}^{n \times T}$
- Objet clé: la **matrice de Gram** $\mathbf{G} = \frac{1}{T} \mathbf{\Sigma}^T \mathbf{\Sigma}$: matrice de **corrélation** dans l'espace "features"

Comprendre la méthode de projections aléatoires

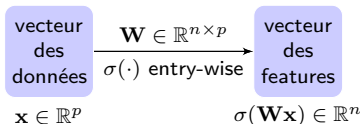


Figure: Illustration de la méthode de projections aléatoires

- ▶ Pour matrice de données $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_T] \in \mathbb{R}^{p \times T}$, on note la **matrice de feature** $\mathbf{\Sigma} \equiv \sigma(\mathbf{W}\mathbf{X}) = [\sigma(\mathbf{W}\mathbf{x}_1), \dots, \sigma(\mathbf{W}\mathbf{x}_T)] \in \mathbb{R}^{n \times T}$
- ▶ Objet clé: la **matrice de Gram** $\mathbf{G} = \frac{1}{T} \mathbf{\Sigma}^T \mathbf{\Sigma}$: matrice de **corrélation dans l'espace "features"**
- ▶ RMT \Rightarrow analyse de la **résolvante** associée

$$\mathbf{Q}(z) = (\mathbf{G} - z\mathbf{I}_T)^{-1} = \left(\frac{1}{T} \mathbf{\Sigma}^T \mathbf{\Sigma} - z\mathbf{I}_T \right)^{-1}, \quad z \in \mathbb{C} \setminus \mathbb{R}^+$$

donne accès à:

Comprendre la méthode de projections aléatoires

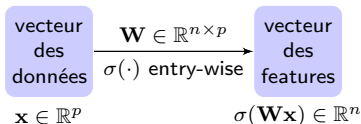


Figure: Illustration de la méthode de projections aléatoires

- Pour matrice de données $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_T] \in \mathbb{R}^{p \times T}$, on note la **matrice de feature** $\Sigma \equiv \sigma(\mathbf{W}\mathbf{X}) = [\sigma(\mathbf{W}\mathbf{x}_1), \dots, \sigma(\mathbf{W}\mathbf{x}_T)] \in \mathbb{R}^{n \times T}$
- Objet clé: la **matrice de Gram** $\mathbf{G} = \frac{1}{T} \Sigma^T \Sigma$: matrice de **corrélation dans l'espace "features"**
- RMT \Rightarrow analyse de la **résolvante** associée

$$\mathbf{Q}(z) = (\mathbf{G} - z\mathbf{I}_T)^{-1} = \left(\frac{1}{T} \Sigma^T \Sigma - z\mathbf{I}_T \right)^{-1}, \quad z \in \mathbb{C} \setminus \mathbb{R}^+$$

donne accès à:

- valeurs propres de \mathbf{G} (via transformée de Stieltjes $\frac{1}{T} \text{tr } \mathbf{Q}$)

Comprendre la méthode de projections aléatoires

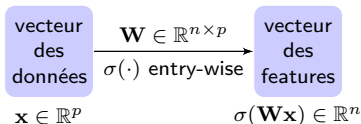


Figure: Illustration de la méthode de projections aléatoires

- ▶ Pour matrice de données $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_T] \in \mathbb{R}^{p \times T}$, on note la **matrice de feature** $\Sigma \equiv \sigma(\mathbf{W}\mathbf{X}) = [\sigma(\mathbf{W}\mathbf{x}_1), \dots, \sigma(\mathbf{W}\mathbf{x}_T)] \in \mathbb{R}^{n \times T}$
- ▶ Objet clé: la **matrice de Gram** $\mathbf{G} = \frac{1}{T} \Sigma^T \Sigma$: matrice de **corrélation dans l'espace "features"**
- ▶ RMT \Rightarrow analyse de la **résolvante** associée

$$\mathbf{Q}(z) = (\mathbf{G} - z\mathbf{I}_T)^{-1} = \left(\frac{1}{T} \Sigma^T \Sigma - z\mathbf{I}_T \right)^{-1}, \quad z \in \mathbb{C} \setminus \mathbb{R}^+$$

donne accès à:

- ▶ valeurs propres de \mathbf{G} (via transformée de Stieltjes $\frac{1}{T} \text{tr } \mathbf{Q}$)
- ▶ vecteurs propres isolés de \mathbf{G} (via intégrale complexe sur \mathbf{Q})

Comprendre la méthode de projections aléatoires

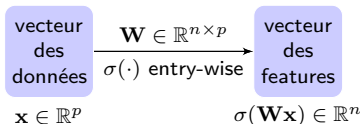


Figure: Illustration de la méthode de projections aléatoires

- Pour matrice de données $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_T] \in \mathbb{R}^{p \times T}$, on note la **matrice de feature** $\Sigma \equiv \sigma(\mathbf{W}\mathbf{X}) = [\sigma(\mathbf{W}\mathbf{x}_1), \dots, \sigma(\mathbf{W}\mathbf{x}_T)] \in \mathbb{R}^{n \times T}$
- Objet clé: la **matrice de Gram** $\mathbf{G} = \frac{1}{T} \Sigma^T \Sigma$: matrice de **corrélation** dans l'espace "features"
- RMT \Rightarrow analyse de la **résolvante** associée

$$\mathbf{Q}(z) = (\mathbf{G} - z\mathbf{I}_T)^{-1} = \left(\frac{1}{T} \Sigma^T \Sigma - z\mathbf{I}_T \right)^{-1}, \quad z \in \mathbb{C} \setminus \mathbb{R}^+$$

donne accès à:

- valeurs propres de \mathbf{G} (via transformée de Stieltjes $\frac{1}{T} \text{tr } \mathbf{Q}$)
- vecteurs propres isolés de \mathbf{G} (via intégrale complexe sur \mathbf{Q})
- dans le régime où $n, p, T \rightarrow \infty$, on cherche un **équivalent déterministe asymptotique** de \mathbf{Q}

1 Introduction

2 Résultats Principaux

3 Conclusion

Hypothèse 1: Taux de Croissance

Lorsque $n \rightarrow \infty$, on a

- ▶ $\frac{p}{n} \rightarrow c_0 \in (0, \infty)$

Hypothèse 1: Taux de Croissance

Lorsque $n \rightarrow \infty$, on a

- ▶ $\frac{p}{n} \rightarrow c_0 \in (0, \infty)$
- ▶ $\frac{T}{n} \rightarrow c_T \in (0, \infty)$

Hypothèse 1: Taux de Croissance

Lorsque $n \rightarrow \infty$, on a

- ▶ $\frac{p}{n} \rightarrow c_0 \in (0, \infty)$
- ▶ $\frac{T}{n} \rightarrow c_T \in (0, \infty)$
- ▶ $\|\mathbf{X}\| = O(1)$ [norme opérateur]

Hypothèse 1: Taux de Croissance

Lorsque $n \rightarrow \infty$, on a

- ▶ $\frac{p}{n} \rightarrow c_0 \in (0, \infty)$
- ▶ $\frac{T}{n} \rightarrow c_T \in (0, \infty)$
- ▶ $\|\mathbf{X}\| = O(1)$ [norme opérateur]

Hypothèse 2: Fonction d'activation

La fonction σ est Lipschitzienne avec constante de Lipschitz λ_σ indépendante de n .

Hypothèse 1: Taux de Croissance

Lorsque $n \rightarrow \infty$, on a

- ▶ $\frac{p}{n} \rightarrow c_0 \in (0, \infty)$
- ▶ $\frac{T}{n} \rightarrow c_T \in (0, \infty)$
- ▶ $\|\mathbf{X}\| = O(1)$ [norme opérateur]

Hypothèse 2: Fonction d'activation

La fonction σ est Lipschitzienne avec constante de Lipschitz λ_σ indépendante de n .

Hypothèse 3: \mathbf{W} Gaussienne

La matrice \mathbf{W} a ses entrées i.i.d. $\mathcal{N}(0, 1)$.

Équivalent déterministe de \mathbf{Q} (1)

Theorem (Équivalent déterministe de \mathbf{Q})

Sous hypothèses 1 - 3, lorsque $n \rightarrow \infty$, on a

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

où $\bar{\mathbf{Q}}$ est donnée par

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

avec δ l'unique solution positive $\delta = \frac{1}{T} \text{tr}(\Phi \bar{\mathbf{Q}})$ et $\Phi \equiv \mathbb{E}[\sigma(\mathbf{X}^\top \mathbf{w}) \sigma(\mathbf{w}^\top \mathbf{X})]$ pour $\mathbf{w} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_p)$.

Équivalent déterministe de \mathbf{Q} (1)

Theorem (Équivalent déterministe de \mathbf{Q})

Sous hypothèses 1 - 3, lorsque $n \rightarrow \infty$, on a

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

où $\bar{\mathbf{Q}}$ est donnée par

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

avec δ l'unique solution positive $\delta = \frac{1}{T} \text{tr}(\Phi \bar{\mathbf{Q}})$ et $\Phi \equiv \mathbb{E}[\sigma(\mathbf{X}^\top \mathbf{w}) \sigma(\mathbf{w}^\top \mathbf{X})]$ pour $\mathbf{w} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_p)$.

Intuition:

- pour $\sigma(t) = t$, [Silverstein et Bai'95]: $\|\mathbb{E}[\mathbf{Q}] - \bar{\mathbf{Q}}\| \rightarrow 0$, $\bar{\mathbf{Q}} \equiv \left(\frac{n}{T} \frac{\mathbf{X}^\top \mathbf{X}}{1 + \delta} - z \mathbf{I}_T \right)^{-1}$

Équivalent déterministe de \mathbf{Q} (1)

Theorem (Équivalent déterministe de \mathbf{Q})

Sous hypothèses 1 - 3, lorsque $n \rightarrow \infty$, on a

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

où $\bar{\mathbf{Q}}$ est donnée par

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

avec δ l'unique solution positive $\delta = \frac{1}{T} \text{tr}(\Phi \bar{\mathbf{Q}})$ et $\Phi \equiv \mathbb{E}[\sigma(\mathbf{X}^\top \mathbf{w}) \sigma(\mathbf{w}^\top \mathbf{X})]$ pour $\mathbf{w} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_p)$.

Intuition:

- ▶ pour $\sigma(t) = t$, [Silverstein et Bai'95]: $\|\mathbb{E}[\mathbf{Q}] - \bar{\mathbf{Q}}\| \rightarrow 0$, $\bar{\mathbf{Q}} \equiv \left(\frac{n}{T} \frac{\mathbf{X}^\top \mathbf{X}}{1 + \delta} - z \mathbf{I}_T \right)^{-1}$
- ▶ repose essentiellement sur **lemme de la trace**: pour \mathbf{A} de norme borné et indépendante de \mathbf{W} ,

$$\left| \frac{1}{n} \mathbf{w}_i^\top \mathbf{A} \mathbf{w}_i - \frac{1}{n} \text{tr} \mathbf{A} \right| \rightarrow 0$$

Équivalent déterministe de \mathbf{Q} (1)

Theorem (Équivalent déterministe de \mathbf{Q})

Sous hypothèses 1 - 3, lorsque $n \rightarrow \infty$, on a

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

où $\bar{\mathbf{Q}}$ est donnée par

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

avec δ l'unique solution positive $\delta = \frac{1}{T} \text{tr}(\Phi \bar{\mathbf{Q}})$ et $\Phi \equiv \mathbb{E}[\sigma(\mathbf{X}^\top \mathbf{w}) \sigma(\mathbf{w}^\top \mathbf{X})]$ pour $\mathbf{w} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_p)$.

Intuition:

- ▶ pour $\sigma(t) = t$, [Silverstein et Bai'95]: $\|\mathbb{E}[\mathbf{Q}] - \bar{\mathbf{Q}}\| \rightarrow 0$, $\bar{\mathbf{Q}} \equiv \left(\frac{n}{T} \frac{\mathbf{X}^\top \mathbf{X}}{1 + \delta} - z \mathbf{I}_T \right)^{-1}$
- ▶ repose essentiellement sur **lemme de la trace**: pour \mathbf{A} de norme borné et indépendante de \mathbf{W} ,

$$\left| \frac{1}{n} \mathbf{w}_i^\top \mathbf{A} \mathbf{w}_i - \frac{1}{n} \text{tr} \mathbf{A} \right| \rightarrow 0$$

- ✓ \mathbf{w}_i^\top (colonne de \mathbf{W}) indépendants

Équivalent déterministe de \mathbf{Q} (1)

Theorem (Équivalent déterministe de \mathbf{Q})

Sous hypothèses 1 - 3, lorsque $n \rightarrow \infty$, on a

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

où $\bar{\mathbf{Q}}$ est donnée par

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

avec δ l'unique solution positive $\delta = \frac{1}{T} \text{tr}(\Phi \bar{\mathbf{Q}})$ et $\Phi \equiv \mathbb{E}[\sigma(\mathbf{X}^\top \mathbf{w}) \sigma(\mathbf{w}^\top \mathbf{X})]$ pour $\mathbf{w} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_p)$.

Intuition:

- ▶ pour $\sigma(t) = t$, [Silverstein et Bai'95]: $\|\mathbb{E}[\mathbf{Q}] - \bar{\mathbf{Q}}\| \rightarrow 0$, $\bar{\mathbf{Q}} \equiv \left(\frac{n}{T} \frac{\mathbf{X}^\top \mathbf{X}}{1 + \delta} - z \mathbf{I}_T \right)^{-1}$
- ▶ repose essentiellement sur **lemme de la trace**: pour \mathbf{A} de norme borné et indépendante de \mathbf{W} ,

$$\left| \frac{1}{n} \mathbf{w}_i^\top \mathbf{A} \mathbf{w}_i - \frac{1}{n} \text{tr} \mathbf{A} \right| \rightarrow 0$$

- ✓ \mathbf{w}_i^\top (colonne de \mathbf{W}) indépendants
- ✓ \mathbf{w}_i à entrées i.i.d. $\mathcal{N}(0, 1)$,

Équivalent déterministe de \mathbf{Q} (1)

Theorem (Équivalent déterministe de \mathbf{Q})

Sous hypothèses 1 - 3, lorsque $n \rightarrow \infty$, on a

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

où $\bar{\mathbf{Q}}$ est donnée par

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

avec δ l'unique solution positive $\delta = \frac{1}{T} \text{tr}(\Phi \bar{\mathbf{Q}})$ et $\Phi \equiv \mathbb{E}[\sigma(\mathbf{X}^\top \mathbf{w}) \sigma(\mathbf{w}^\top \mathbf{X})]$ pour $\mathbf{w} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_p)$.

Intuition:

- ▶ pour $\sigma(t) = t$, [Silverstein et Bai'95]: $\|\mathbb{E}[\mathbf{Q}] - \bar{\mathbf{Q}}\| \rightarrow 0$, $\bar{\mathbf{Q}} \equiv \left(\frac{n}{T} \frac{\mathbf{X}^\top \mathbf{X}}{1 + \delta} - z \mathbf{I}_T \right)^{-1}$
- ▶ repose essentiellement sur **lemme de la trace**: pour \mathbf{A} de norme borné et indépendante de \mathbf{W} ,

$$\left| \frac{1}{n} \mathbf{w}_i^\top \mathbf{A} \mathbf{w}_i - \frac{1}{n} \text{tr} \mathbf{A} \right| \rightarrow 0$$

- ✓ \mathbf{w}_i^\top (colonne de \mathbf{W}) indépendants
- ✓ \mathbf{w}_i à entrées i.i.d. $\mathcal{N}(0, 1)$,
- ▶ Ici, pour $\frac{1}{n} \sigma(\mathbf{w}_i^\top \mathbf{X}) \mathbf{A} \sigma(\mathbf{X}^\top \mathbf{w}_i)$,

Équivalent déterministe de \mathbf{Q} (1)

Theorem (Équivalent déterministe de \mathbf{Q})

Sous hypothèses 1 - 3, lorsque $n \rightarrow \infty$, on a

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

où $\bar{\mathbf{Q}}$ est donnée par

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

avec δ l'unique solution positive $\delta = \frac{1}{T} \text{tr}(\Phi \bar{\mathbf{Q}})$ et $\Phi \equiv \mathbb{E}[\sigma(\mathbf{X}^\top \mathbf{w}) \sigma(\mathbf{w}^\top \mathbf{X})]$ pour $\mathbf{w} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_p)$.

Intuition:

- ▶ pour $\sigma(t) = t$, [Silverstein et Bai'95]: $\|\mathbb{E}[\mathbf{Q}] - \bar{\mathbf{Q}}\| \rightarrow 0$, $\bar{\mathbf{Q}} \equiv \left(\frac{n}{T} \frac{\mathbf{X}^\top \mathbf{X}}{1 + \delta} - z \mathbf{I}_T \right)^{-1}$
- ▶ repose essentiellement sur **lemme de la trace**: pour \mathbf{A} de norme borné et indépendante de \mathbf{W} ,

$$\left| \frac{1}{n} \mathbf{w}_i^\top \mathbf{A} \mathbf{w}_i - \frac{1}{n} \text{tr} \mathbf{A} \right| \rightarrow 0$$

- ✓ \mathbf{w}_i^\top (colonne de \mathbf{W}) indépendants
- ✓ \mathbf{w}_i à entrées i.i.d. $\mathcal{N}(0, 1)$,
- ▶ Ici, pour $\frac{1}{n} \sigma(\mathbf{w}_i^\top \mathbf{X}) \mathbf{A} \sigma(\mathbf{X}^\top \mathbf{w}_i)$,
 - ✓ les colonnes de $\Sigma \equiv \sigma(\mathbf{W} \mathbf{X})$: $\sigma(\mathbf{X}^\top \mathbf{w}_i)$ i.i.d.

Équivalent déterministe de \mathbf{Q} (1)

Theorem (Équivalent déterministe de \mathbf{Q})

Sous hypothèses 1 - 3, lorsque $n \rightarrow \infty$, on a

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

où $\bar{\mathbf{Q}}$ est donnée par

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

avec δ l'unique solution positive $\delta = \frac{1}{T} \text{tr}(\Phi \bar{\mathbf{Q}})$ et $\Phi \equiv \mathbb{E}[\sigma(\mathbf{X}^\top \mathbf{w}) \sigma(\mathbf{w}^\top \mathbf{X})]$ pour $\mathbf{w} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_p)$.

Intuition:

- ▶ pour $\sigma(t) = t$, [Silverstein et Bai'95]: $\|\mathbb{E}[\mathbf{Q}] - \bar{\mathbf{Q}}\| \rightarrow 0$, $\bar{\mathbf{Q}} \equiv \left(\frac{n}{T} \frac{\mathbf{X}^\top \mathbf{X}}{1 + \delta} - z \mathbf{I}_T \right)^{-1}$
- ▶ repose essentiellement sur **lemme de la trace**: pour \mathbf{A} de norme borné et indépendante de \mathbf{W} ,

$$\left| \frac{1}{n} \mathbf{w}_i^\top \mathbf{A} \mathbf{w}_i - \frac{1}{n} \text{tr} \mathbf{A} \right| \rightarrow 0$$

- ✓ \mathbf{w}_i^\top (colonne de \mathbf{W}) indépendants
- ✓ \mathbf{w}_i à entrées i.i.d. $\mathcal{N}(0, 1)$,
- ▶ Ici, pour $\frac{1}{n} \sigma(\mathbf{w}_i^\top \mathbf{X}) \mathbf{A} \sigma(\mathbf{X}^\top \mathbf{w}_i)$,
 - ✓ les colonnes de $\Sigma \equiv \sigma(\mathbf{W} \mathbf{X})$: $\sigma(\mathbf{X}^\top \mathbf{w}_i)$ i.i.d.
 - ✗ mais $\sigma(\mathbf{X}^\top \mathbf{w}_i)$ à entrées **non-indépendants**

Équivalent déterministe de \mathbf{Q} (2)

Theorem (Équivalent déterministe de \mathbf{Q})

Sous hypothèses 1 - 3, lorsque $n \rightarrow \infty$, on a

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

où $\bar{\mathbf{Q}}$ est donnée par

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

avec δ l'unique solution positive $\delta = \frac{1}{T} \text{tr}(\Phi \bar{\mathbf{Q}})$ et $\Phi \equiv \mathbb{E}[\sigma(\mathbf{X}^\top \mathbf{w}) \sigma(\mathbf{w}^\top \mathbf{X})]$ pour $\mathbf{w} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_p)$.

Équivalent déterministe de \mathbf{Q} (2)

Theorem (Équivalent déterministe de \mathbf{Q})

Sous hypothèses 1 - 3, lorsque $n \rightarrow \infty$, on a

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

où $\bar{\mathbf{Q}}$ est donnée par

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

avec δ l'unique solution positive $\delta = \frac{1}{T} \text{tr}(\Phi \bar{\mathbf{Q}})$ et $\Phi \equiv \mathbb{E}[\sigma(\mathbf{X}^\top \mathbf{w}) \sigma(\mathbf{w}^\top \mathbf{X})]$ pour $\mathbf{w} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_p)$.

Intuition:

- ▶ On remplace le lemme de la trace par **concentration de la mesure**:

$$P \left(\left| \frac{1}{n} \sigma(\mathbf{w}_i^\top \mathbf{X}) \mathbf{A} \sigma(\mathbf{X}^\top \mathbf{w}_i) - \frac{1}{n} \text{tr}(\Phi \mathbf{A}) \right| > t \right) \leq C e^{-cn \max(t, t^2)}$$

avec $\Phi \equiv \mathbb{E}_{\mathbf{w}}[\sigma(\mathbf{X}^\top \mathbf{w}) \sigma(\mathbf{w}^\top \mathbf{X})]$ pour $\mathbf{w} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_p)$.

Application: Analyse de la performance de l'ELM

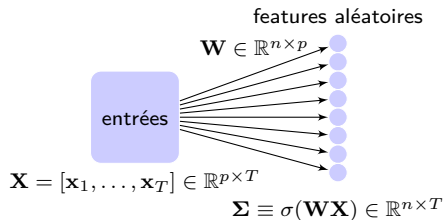


Figure: Illustration de l'ELM

Phase d'apprentissage:

Application: Analyse de la performance de l'ELM

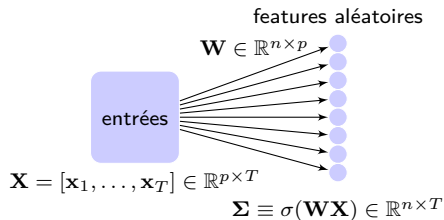


Figure: Illustration de l'ELM

Phase d'apprentissage:

- génération de features aléatoires pour les **données d'entraînement** \mathbf{X} : $\Sigma \equiv \sigma(\mathbf{W}\mathbf{X})$

Application: Analyse de la performance de l'ELM

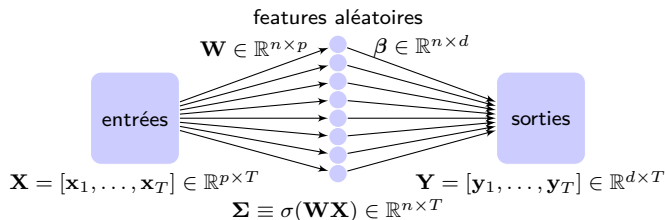


Figure: Illustration de l'ELM

Phase d'apprentissage:

- ▶ génération de features aléatoires pour les **données d'entraînement** \mathbf{X} : $\Sigma \equiv \sigma(\mathbf{W}\mathbf{X})$
- ▶ régression normalisée de Σ aux sorties associées \mathbf{Y} :
$$\beta = \frac{1}{T} \Sigma \left(\frac{1}{T} \Sigma^T \Sigma + \gamma \mathbf{I}_T \right)^{-1} \mathbf{Y}^T = \frac{1}{T} \Sigma \mathbf{Q}(-\gamma) \mathbf{Y}^T$$

Application: Analyse de la performance de l'ELM

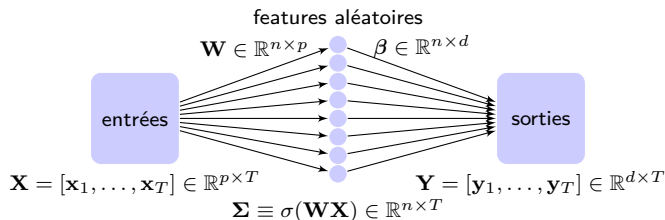


Figure: Illustration de l'ELM

Phase d'apprentissage:

- ▶ génération de features aléatoires pour les données d'entraînement \mathbf{X} : $\Sigma \equiv \sigma(\mathbf{W}\mathbf{X})$
- ▶ régression normalisée de Σ aux sorties associées \mathbf{Y} :
$$\beta = \frac{1}{T} \Sigma \left(\frac{1}{T} \Sigma^T \Sigma + \gamma \mathbf{I}_T \right)^{-1} \mathbf{Y}^T = \frac{1}{T} \Sigma \mathbf{Q}(-\gamma) \mathbf{Y}^T$$

Phase de test:

Application: Analyse de la performance de l'ELM

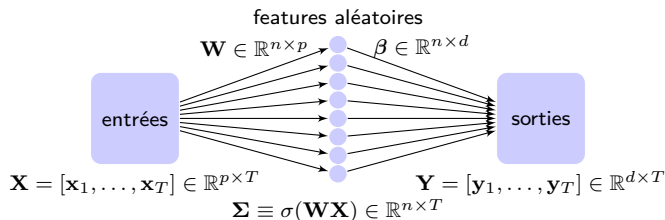


Figure: Illustration de l'ELM

Phase d'apprentissage:

- ▶ génération de features aléatoires pour les **données d'entraînement** \mathbf{X} : $\Sigma \equiv \sigma(\mathbf{W}\mathbf{X})$
- ▶ régression normalisée de Σ aux sorties associées \mathbf{Y} :
$$\beta = \frac{1}{T} \Sigma \left(\frac{1}{T} \Sigma^T \Sigma + \gamma \mathbf{I}_T \right)^{-1} \mathbf{Y}^T = \frac{1}{T} \Sigma \mathbf{Q}(-\gamma) \mathbf{Y}^T$$

Phase de test:

- ▶ appliquer le régresseur aux **nouvelles données** $\hat{\mathbf{X}}$, sorties $\hat{\mathbf{Y}}$ inconnues

Application: Analyse de la performance de l'ELM

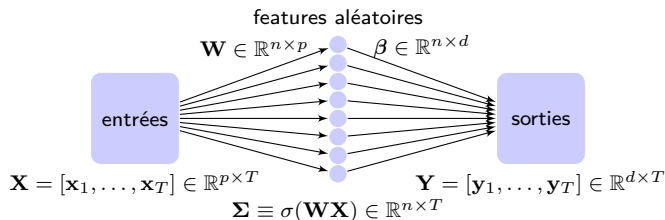


Figure: Illustration de l'ELM

Phase d'apprentissage:

- ▶ génération de features aléatoires pour les **données d'entraînement** \mathbf{X} : $\mathbf{\Sigma} \equiv \sigma(\mathbf{W}\mathbf{X})$
- ▶ régression normalisée de $\mathbf{\Sigma}$ aux sorties associées \mathbf{Y} :
$$\beta = \frac{1}{T} \mathbf{\Sigma} \left(\frac{1}{T} \mathbf{\Sigma}^T \mathbf{\Sigma} + \gamma \mathbf{I}_T \right)^{-1} \mathbf{Y}^T = \frac{1}{T} \mathbf{\Sigma} \mathbf{Q}(-\gamma) \mathbf{Y}^T$$

Phase de test:

- ▶ appliquer le régresseur aux **nouvelles données** $\hat{\mathbf{X}}$, sorties $\hat{\mathbf{Y}}$ inconnues
- ▶ sortie de l'ELM donnée par $\mathbf{Z} = \beta^T \hat{\mathbf{\Sigma}}$, avec $\hat{\mathbf{\Sigma}} \equiv \sigma(\mathbf{W}\hat{\mathbf{X}})$

Compréhension la performance de l'ELM

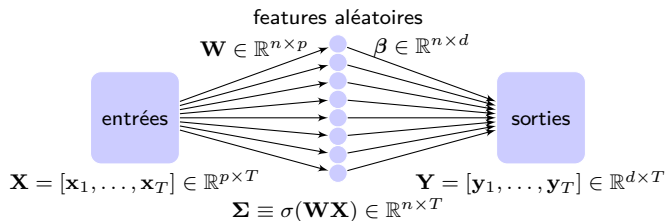


Figure: Illustration of ELM

avec un *ridge-regressor* $\beta = \frac{1}{T} \Sigma \left(\frac{1}{T} \Sigma^T \Sigma + \gamma \mathbf{I}_T \right)^{-1} \mathbf{Y}^T = \frac{1}{T} \Sigma \mathbf{Q} \mathbf{Y}^T$.

Compréhension la performance de l'ELM

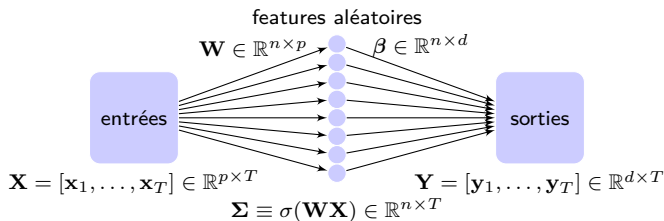


Figure: Illustration of ELM

avec un *ridge-regressor* $\beta = \frac{1}{T} \Sigma \left(\frac{1}{T} \Sigma^T \Sigma + \gamma \mathbf{I}_T \right)^{-1} \mathbf{Y}^T = \frac{1}{T} \Sigma \mathbf{Q} \mathbf{Y}^T$.

Pour comprendre le MSE (mean-square error) d'apprentissage:

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

Compréhension la performance de l'ELM

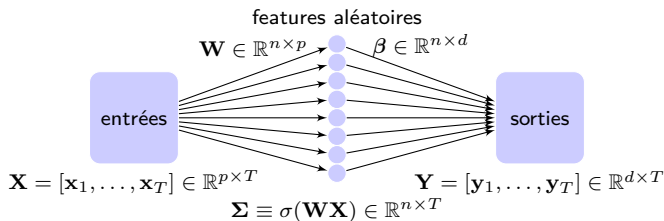


Figure: Illustration of ELM

avec un *ridge-regressor* $\beta = \frac{1}{T} \Sigma \left(\frac{1}{T} \Sigma^T \Sigma + \gamma \mathbf{I}_T \right)^{-1} \mathbf{Y}^T = \frac{1}{T} \Sigma \mathbf{Q} \mathbf{Y}^T$.

Pour comprendre le MSE (mean-square error) d'apprentissage:

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

aussi la performance de **test** pour les nouvelles données $\hat{\mathbf{X}}$ et les sorties associées $\hat{\mathbf{Y}}$:

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

Theorem (Équivalent déterministe de E_{train})

Sous hypothèses 1-3, lorsque $n \rightarrow \infty$, on a

$$|E_{\text{train}} - \bar{E}_{\text{train}}| \rightarrow 0$$

avec

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

Theorem (Équivalent déterministe de E_{train})

Sous hypothèses 1-3, lorsque $n \rightarrow \infty$, on a

$$|E_{\text{train}} - \bar{E}_{\text{train}}| \rightarrow 0$$

avec

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

Dernière mission: évaluer $\Phi_{\mathbf{a}\mathbf{b}} \equiv \mathbb{E}_{\mathbf{w}}[\sigma(\mathbf{w}^\top \mathbf{a})\sigma(\mathbf{w}^\top \mathbf{b})]$ pour des $\sigma(\cdot)$ différents:

Équivalent déterministe de MSE

Theorem (Équivalent déterministe de E_{train})

Sous hypothèses 1-3, lorsque $n \rightarrow \infty$, on a

$$|E_{\text{train}} - \bar{E}_{\text{train}}| \rightarrow 0$$

avec

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

Dernière mission: évaluer $\Phi_{\mathbf{a}\mathbf{b}} \equiv \mathbb{E}_{\mathbf{w}}[\sigma(\mathbf{w}^\top \mathbf{a})\sigma(\mathbf{w}^\top \mathbf{b})]$ pour des $\sigma(\cdot)$ différents:

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

Table: Valeurs de $\Phi_{\mathbf{a}\mathbf{b}}$ pour $w \sim \mathcal{N}(0, \mathbf{I}_p)$, avec $\angle(\mathbf{a}, \mathbf{b}) \equiv \frac{\mathbf{a}^\top \mathbf{b}}{\|\mathbf{a}\| \|\mathbf{b}\|}$.

Validation numérique sur la base de données MNIST²

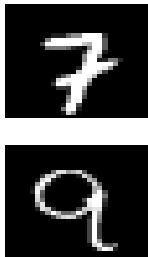


Figure: Exemples de la base de données MNIST

²Les images des chiffres manuscrits de taille 28×28 .

Validation numérique sur la base de données MNIST²

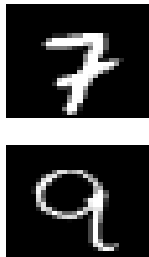
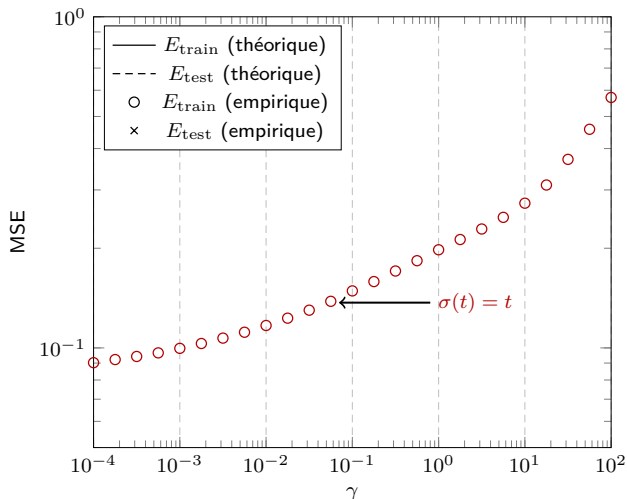


Figure: Exemples de la base de données MNIST



²Les images des chiffres manuscrits de taille 28×28 .

Validation numérique sur la base de données MNIST²

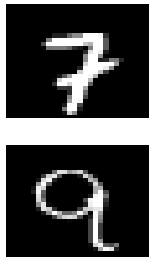
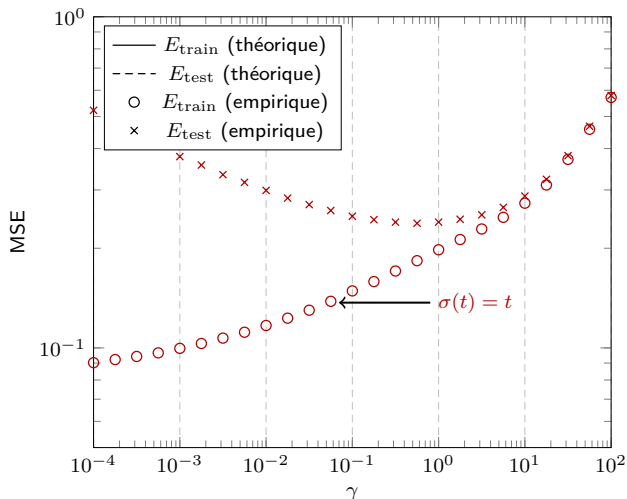


Figure: Exemples de la base de données MNIST



²Les images des chiffres manuscrits de taille 28×28 .

Validation numérique sur la base de données MNIST²



Figure: Exemples de la base de données MNIST

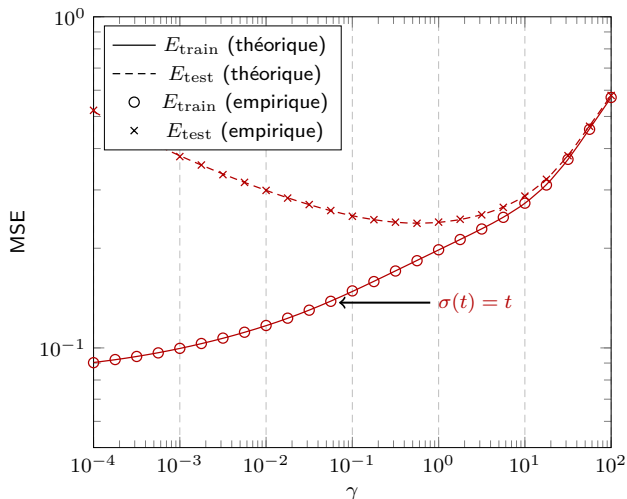


Figure: Performance de l'ELM pour σ Lipschitzienne, en fonction de γ , pour 2 classes de données de MNIST (sept et neuf), $n = 512$, $T = \hat{T} = 1024$, $p = 784$.

²Les images des chiffres manuscrits de taille 28×28 .

Validation numérique sur la base de données MNIST²

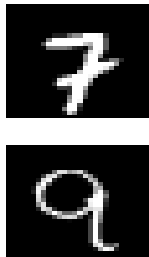
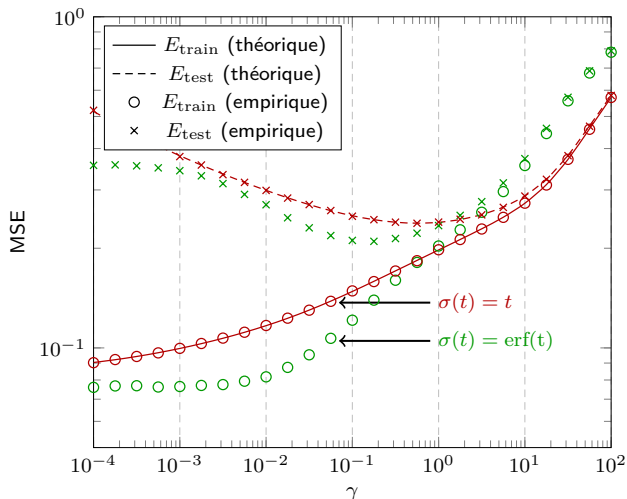


Figure: Exemples de la base de données MNIST



²Les images des chiffres manuscrits de taille 28×28 .

Validation numérique sur la base de données MNIST²

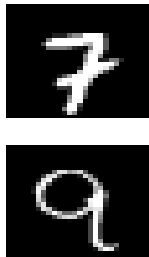


Figure: Exemples de la base de données MNIST

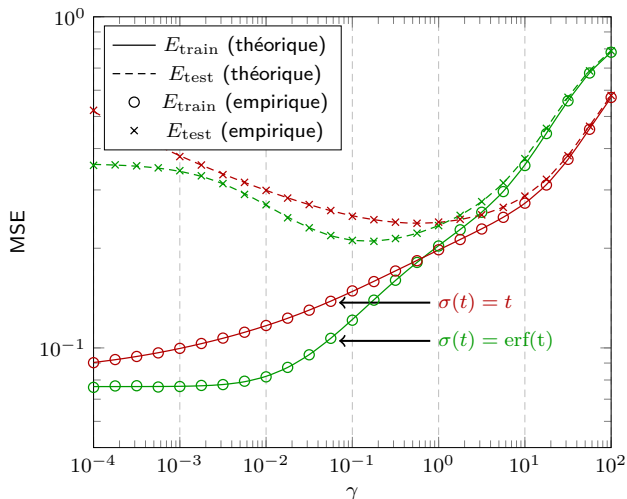


Figure: Performance de l'ELM pour σ Lipschitzienne, en fonction de γ , pour 2 classes de données de MNIST (sept et neuf), $n = 512$, $T = \hat{T} = 1024$, $p = 784$.

²Les images des chiffres manuscrits de taille 28×28 .

Validation numérique sur la base de données MNIST²



Figure: Exemples de la base de données MNIST

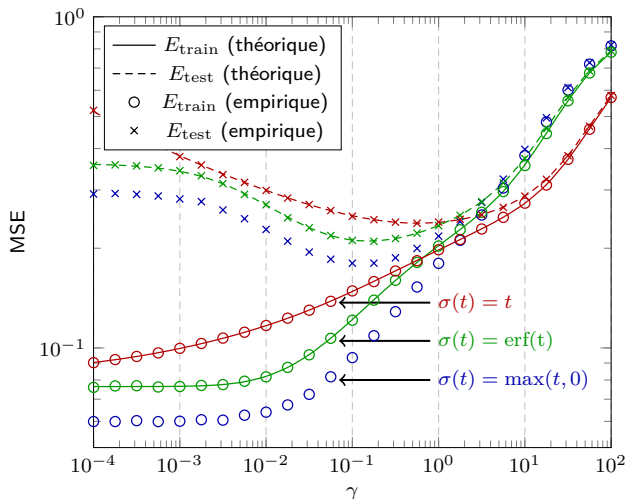


Figure: Performance de l'ELM pour σ Lipschitzienne, en fonction de γ , pour 2 classes de données de MNIST (sept et neuf), $n = 512$, $T = \hat{T} = 1024$, $p = 784$.

²Les images des chiffres manuscrits de taille 28×28 .

Validation numérique sur la base de données MNIST²



Figure: Exemples de la base de données MNIST

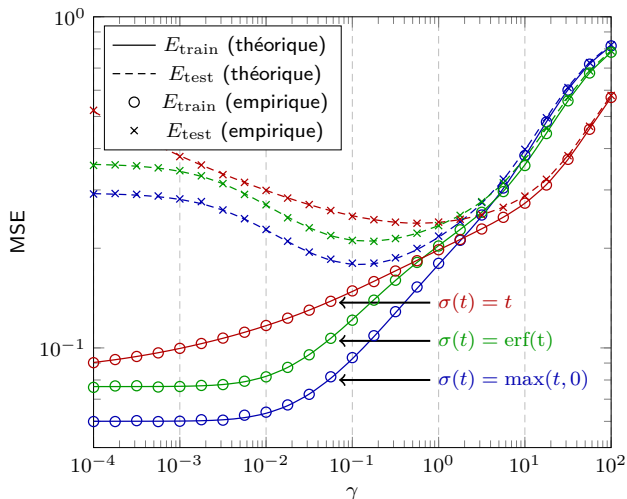


Figure: Performance de l'ELM pour σ Lipschitzienne, en fonction de γ , pour 2 classes de données de MNIST (sept et neuf), $n = 512$, $T = \hat{T} = 1024$, $p = 784$.

²Les images des chiffres manuscrits de taille 28×28 .

1 Introduction

2 Résultats Principaux

3 Conclusion

Conclusion

Messages:

- ▶ étude de la matrice de Gram **non-linéaire**: $\mathbf{G} \equiv \frac{1}{T} \mathbf{\Sigma}^T \mathbf{\Sigma}$

Conclusion

Messages:

- ▶ étude de la matrice de Gram **non-linéaire**: $\mathbf{G} \equiv \frac{1}{T} \mathbf{\Sigma}^T \mathbf{\Sigma}$
- ▶ **objet clé**: la résolvante \mathbf{Q} et son **équivalent déterministe asymptotique** $\bar{\mathbf{Q}}$

Conclusion

Messages:

- ▶ étude de la matrice de Gram **non-linéaire**: $\mathbf{G} \equiv \frac{1}{T} \mathbf{\Sigma}^T \mathbf{\Sigma}$
- ▶ **objet clé**: la résolvante \mathbf{Q} et son **équivalent déterministe asymptotique** $\bar{\mathbf{Q}}$
- ▶ permet la compréhension de la performance de l'ELM et projections aléatoires

Conclusion

Messages:

- ▶ étude de la matrice de Gram **non-linéaire**: $\mathbf{G} \equiv \frac{1}{T} \mathbf{\Sigma}^T \mathbf{\Sigma}$
- ▶ **objet clé**: la résolvante \mathbf{Q} et son **équivalent déterministe asymptotique** $\bar{\mathbf{Q}}$
- ▶ permet la compréhension de la performance de l'ELM et projections aléatoires
- ▶ optimisation théorique des hyper-paramètres (ici γ)

Conclusion

Messages:

- ▶ étude de la matrice de Gram **non-linéaire**: $\mathbf{G} \equiv \frac{1}{T} \mathbf{\Sigma}^T \mathbf{\Sigma}$
- ▶ **objet clé**: la résolvante \mathbf{Q} et son **équivalent déterministe asymptotique** $\bar{\mathbf{Q}}$
- ▶ permet la compréhension de la performance de l'ELM et projections aléatoires
- ▶ optimisation théorique des hyper-paramètres (ici γ)

Travaux à venir:

- ▶ analyse des méthodes plus sophistiquées basées sur les projections aléatoires, e.g., projections aléatoires + régression logistique, projections aléatoires + SVM, etc.

Conclusion

Messages:

- ▶ étude de la matrice de Gram **non-linéaire**: $\mathbf{G} \equiv \frac{1}{T} \mathbf{\Sigma}^T \mathbf{\Sigma}$
- ▶ **objet clé**: la résolvante \mathbf{Q} et son **équivalent déterministe asymptotique** $\bar{\mathbf{Q}}$
- ▶ permet la compréhension de la performance de l'ELM et projections aléatoires
- ▶ optimisation théorique des hyper-paramètres (ici γ)

Travaux à venir:

- ▶ analyse des méthodes plus sophistiquées basées sur les projections aléatoires, e.g., projections aléatoires + régression logistique, projections aléatoires + SVM, etc.
- ▶ analyse de \mathbf{Q} pour un modèle statistique des données: choisir $\sigma(\cdot)$ en fonction des données

Conclusion

Messages:

- ▶ étude de la matrice de Gram **non-linéaire**: $\mathbf{G} \equiv \frac{1}{T} \mathbf{\Sigma}^T \mathbf{\Sigma}$
- ▶ **objet clé**: la résolvante \mathbf{Q} et son **équivalent déterministe asymptotique** $\bar{\mathbf{Q}}$
- ▶ permet la compréhension de la performance de l'ELM et projections aléatoires
- ▶ optimisation théorique des hyper-paramètres (ici γ)

Travaux à venir:

- ▶ analyse des méthodes plus sophistiquées basées sur les projections aléatoires, e.g., projections aléatoires + régression logistique, projections aléatoires + SVM, etc.
- ▶ analyse de \mathbf{Q} pour un modèle statistique des données: choisir $\sigma(\cdot)$ en fonction des données
- ▶ combinaison de plusieurs types de non-linéarités, e.g., random Fourier features ($\sin + \cos \approx$ noyau Gaussien)

Conclusion

Messages:

- ▶ étude de la matrice de Gram **non-linéaire**: $\mathbf{G} \equiv \frac{1}{T} \mathbf{\Sigma}^T \mathbf{\Sigma}$
- ▶ **objet clé**: la résolvante \mathbf{Q} et son **équivalent déterministe asymptotique** $\bar{\mathbf{Q}}$
- ▶ permet la compréhension de la performance de l'ELM et projections aléatoires
- ▶ optimisation théorique des hyper-paramètres (ici γ)

Travaux à venir:

- ▶ analyse des méthodes plus sophistiquées basées sur les projections aléatoires, e.g., projections aléatoires + régression logistique, projections aléatoires + SVM, etc.
- ▶ analyse de \mathbf{Q} pour un modèle statistique des données: choisir $\sigma(\cdot)$ en fonction des données
- ▶ combinaison de plusieurs types de non-linéarités, e.g., random Fourier features ($\sin + \cos \approx$ noyau Gaussien)

References:

C. Louart, Z. Liao, R. Couillet, "A Random Matrix Approach to Neural Networks", (à paraître) Annals of Applied Probability, 2017.

Merci!