



Master Statistiques et Sciences des Données (SSD)  
2025-2026

---

# PROJET D'EXAMEN FINAL

## Optimisation pour le Machine Learning

---

Réalisé par :

Namy Ahmed Abdy

Nouakchott, Mauritanie

Janvier 2026

# Contents

<b>1</b>	<b>Exercice 1 : Modélisation et Étude Théorique</b>	<b>3</b>
1.1	Modélisation du problème . . . . .	3
1.1.1	Formulation de la fonction objectif . . . . .	3
1.1.2	Justification de la régularisation $\ell_2$ . . . . .	3
1.2	Calcul analytique du gradient . . . . .	3
1.2.1	Gradient de $f(w)$ . . . . .	3
1.2.2	Matrice Hessienne $\nabla^2 f(w)$ . . . . .	4
1.3	Constante de Lipschitz via SVD . . . . .	4
1.3.1	Décomposition en Valeurs Singulières (SVD) . . . . .	4
1.3.2	Calcul de la constante de Lipschitz . . . . .	5
1.3.3	Impact sur la stabilité de la descente de gradient . . . . .	5
1.4	Résultats expérimentaux . . . . .	6
<b>2</b>	<b>Exercice 2 : Stochasticité et Passage à l'Échelle</b>	<b>6</b>
2.1	Descente de Gradient Stochastique (SGD) . . . . .	6
2.1.1	Preuve du caractère non biaisé . . . . .	6
2.1.2	Algorithme SGD . . . . .	7
2.2	Analyse comparative: Batch GD vs SGD . . . . .	7
2.2.1	Complexité computationnelle . . . . .	7
2.2.2	Phénomène de bruit de gradient . . . . .	8
2.3	Mini-batch SGD et Adam . . . . .	8
2.3.1	Mini-batch SGD . . . . .	8
2.3.2	Algorithme Adam (Adaptive Moment Estimation) . . . . .	8
2.3.3	Importance de la standardisation . . . . .	9
2.4	Résultats expérimentaux et analyse comparative . . . . .	9
<b>3</b>	<b>Exercice 3 : Parcimonie et Algorithmes Proximaux</b>	<b>14</b>
3.1	Analyse géométrique de la régularisation . . . . .	14
3.1.1	Problème d'optimisation avec régularisation . . . . .	14
3.1.2	Régularisation $\ell_2$ (Ridge) . . . . .	14
3.1.3	Régularisation $\ell_1$ (Lasso) . . . . .	14
3.1.4	Pourquoi $\ell_1$ induit la parcimonie? . . . . .	15
3.2	Opérateur proximal et algorithme ISTA . . . . .	17
3.2.1	Définition de l'opérateur proximal . . . . .	17
3.2.2	Opérateur proximal de la norme $\ell_1$ . . . . .	17
3.2.3	Algorithme ISTA (Iterative Soft-Thresholding) . . . . .	18
3.3	Accélération: algorithme FISTA . . . . .	19
3.3.1	Méthode de Nesterov . . . . .	19
3.3.2	Comparaison des taux de convergence . . . . .	19
3.4	Sélection de variables via régularisation . . . . .	20
3.4.1	Chemin de régularisation (Regularization path) . . . . .	20
3.4.2	Choix de $\lambda$ optimal . . . . .	20
3.4.3	Interprétabilité en haute dimension . . . . .	20
3.5	Résultats expérimentaux . . . . .	21
3.5.1	Comparaison ISTA vs FISTA . . . . .	21
3.5.2	Chemin de régularisation et sélection de variables . . . . .	24

<b>4</b>	<b>Synthèse des Résultats Expérimentaux</b>	<b>27</b>
4.1	Performance comparative des algorithmes . . . . .	27
4.2	Validation des prédictions théoriques . . . . .	27
4.2.1	Constante de Lipschitz et conditionnement . . . . .	27
4.2.2	Réduction de variance par mini-batch . . . . .	28
4.2.3	Taux de convergence ISTA vs FISTA . . . . .	28
4.3	Impact de la régularisation . . . . .	28
4.4	Leçons algorithmiques . . . . .	29
4.4.1	Passage à l'échelle . . . . .	29
4.4.2	Importance de la standardisation . . . . .	29
4.5	Limites et extensions possibles . . . . .	29
4.5.1	Limites observées . . . . .	29
4.5.2	Extensions envisageables . . . . .	29
<b>5</b>	<b>Conclusion</b>	<b>30</b>
<b>6</b>	<b>Conclusion</b>	<b>30</b>
6.1	Contributions principales . . . . .	30
6.1.1	Sur le plan théorique . . . . .	30
6.1.2	Sur le plan expérimental . . . . .	30
6.2	Enseignements méthodologiques . . . . .	31
6.2.1	Choix d'algorithmes . . . . .	31
6.2.2	Importance du prétraitement . . . . .	31
6.2.3	Trade-offs fondamentaux . . . . .	31
6.3	Perspectives . . . . .	32
6.3.1	Extensions théoriques . . . . .	32
6.3.2	Applications . . . . .	32
6.4	Réflexion finale . . . . .	32

# 1 Exercice 1 : Modélisation et Étude Théorique

## 1.1 Modélisation du problème

### 1.1.1 Formulation de la fonction objectif

Pour le problème de régression sur le dataset YearPredictionMSD avec  $n \approx 515\,000$  exemples et  $d = 90$  caractéristiques, nous considérons la régression linéaire des moindres carrés régularisée.

Soit  $X \in \mathbb{R}^{n \times d}$  la matrice des données où chaque ligne  $x_i \in \mathbb{R}^d$  représente un exemple, et  $y \in \mathbb{R}^n$  le vecteur cible.

**Fonction objectif:**

$$f(w) = \frac{1}{n} \sum_{i=1}^n f_i(w) + \frac{\mu}{2} \|w\|^2 \quad (1)$$

où chaque fonction individuelle est définie par:

$$f_i(w) = \frac{1}{2} (x_i^\top w - y_i)^2 \quad (2)$$

### 1.1.2 Justification de la régularisation $\ell_2$

L'ajout du terme de régularisation  $\frac{\mu}{2} \|w\|^2$  avec  $\mu > 0$  garantit plusieurs propriétés essentielles:

**Théorème 1** (Unicité du minimum global - Théorème 1.2.9). *Soit  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  une fonction  $\mu$ -fortement convexe avec  $\mu > 0$ . Alors  $f$  admet un unique minimum global  $w^*$ .*

*Justification pour notre problème.* Sans régularisation, la fonction  $f(w) = \frac{1}{2n} \|Xw - y\|^2$  est convexe mais peut ne pas être strictement convexe si:

- La matrice  $X$  n'est pas de rang plein ( $\text{rang}(X) < d$ )
- Il existe des caractéristiques colinéaires
- Le nombre d'exemples  $n < d$  (sous-déterminé)

En ajoutant  $\frac{\mu}{2} \|w\|^2$ , la Hessienne devient:

$$\nabla^2 f(w) = \frac{1}{n} X^\top X + \mu I \succeq \mu I \succ 0 \quad (3)$$

Donc  $f$  est  $\mu$ -fortement convexe, ce qui garantit l'existence et l'unicité du minimum global.  $\square$

## 1.2 Calcul analytique du gradient

### 1.2.1 Gradient de $f(w)$

Pour calculer  $\nabla f(w)$ , développons d'abord  $f(w)$ :

$$f(w) = \frac{1}{n} \sum_{i=1}^n \frac{1}{2} (x_i^\top w - y_i)^2 + \frac{\mu}{2} w^\top w \quad (4)$$

$$= \frac{1}{2n} (Xw - y)^\top (Xw - y) + \frac{\mu}{2} w^\top w \quad (5)$$

$$= \frac{1}{2n} (w^\top X^\top Xw - 2y^\top Xw + y^\top y) + \frac{\mu}{2} w^\top w \quad (6)$$

Le gradient est:

$$\boxed{\nabla f(w) = \frac{1}{n} X^\top (Xw - y) + \mu w} \quad (7)$$

**Vérification composante par composante:**

$$\frac{\partial f}{\partial w_j} = \frac{1}{n} \sum_{i=1}^n (x_i^\top w - y_i) x_{ij} + \mu w_j \quad (8)$$

$$= \frac{1}{n} \sum_{i=1}^n x_{ij} \left( \sum_{k=1}^d x_{ik} w_k - y_i \right) + \mu w_j \quad (9)$$

### 1.2.2 Matrice Hessienne $\nabla^2 f(w)$

La matrice Hessienne est la matrice des dérivées secondes:

$$[\nabla^2 f(w)]_{jk} = \frac{\partial^2 f}{\partial w_j \partial w_k} \quad (10)$$

Calculons:

$$\frac{\partial^2 f}{\partial w_j \partial w_k} = \frac{\partial}{\partial w_k} \left[ \frac{1}{n} \sum_{i=1}^n x_{ij} (x_i^\top w - y_i) + \mu w_j \right] \quad (11)$$

$$= \frac{1}{n} \sum_{i=1}^n x_{ij} x_{ik} + \mu \delta_{jk} \quad (12)$$

où  $\delta_{jk}$  est le symbole de Kronecker.

En notation matricielle:

$$\boxed{\nabla^2 f(w) = \frac{1}{n} X^\top X + \mu I} \quad (13)$$

**Propriété importante:** La Hessienne est constante, indépendante de  $w$ , donc  $f$  est quadratique.

## 1.3 Constante de Lipschitz via SVD

### 1.3.1 Décomposition en Valeurs Singulières (SVD)

Soit  $X = U \Sigma V^\top$  la décomposition SVD de  $X$  où:

- $U \in \mathbb{R}^{n \times n}$  est orthogonale
- $\Sigma \in \mathbb{R}^{n \times d}$  contient les valeurs singulières  $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_{\min(n,d)} \geq 0$
- $V \in \mathbb{R}^{d \times d}$  est orthogonale

### 1.3.2 Calcul de la constante de Lipschitz

**Définition 1** (Gradient L-Lipschitzien). Le gradient  $\nabla f$  est dit L-Lipschitzien si:

$$\|\nabla f(w) - \nabla f(w')\| \leq L\|w - w'\|, \quad \forall w, w' \in \mathbb{R}^d \quad (14)$$

**Proposition 1.** Pour une fonction quadratique  $f(w) = \frac{1}{2}w^\top Aw - b^\top w + c$  avec  $A$  symétrique, la constante de Lipschitz du gradient est:

$$L = \lambda_{\max}(A) \quad (15)$$

où  $\lambda_{\max}(A)$  est la plus grande valeur propre de  $A$ .

*Proof.* Le gradient est  $\nabla f(w) = Aw - b$ , donc:

$$\|\nabla f(w) - \nabla f(w')\| = \|A(w - w')\| \quad (16)$$

$$\leq \|A\|\|w - w'\| \quad (17)$$

$$= \lambda_{\max}(A)\|w - w'\| \quad (18)$$

□

**Application à notre problème:**

La Hessienne est  $\nabla^2 f(w) = \frac{1}{n}X^\top X + \mu I$ .

Les valeurs propres de  $X^\top X$  sont les carrés des valeurs singulières de  $X$ :

$$\lambda_i(X^\top X) = \sigma_i^2(X), \quad i = 1, \dots, \min(n, d) \quad (19)$$

Donc:

$$L = \frac{1}{n}\sigma_{\max}^2(X) + \mu = \frac{1}{n}\|X\|_2^2 + \mu \quad (20)$$

où  $\|X\|_2 = \sigma_{\max}(X)$  est la norme spectrale de  $X$ .

### 1.3.3 Impact sur la stabilité de la descente de gradient

La constante de Lipschitz  $L$  a des implications directes sur:

**1. Choix du pas d'apprentissage:**

Pour garantir la convergence de la descente de gradient, le pas doit satisfaire:

$$0 < \alpha < \frac{2}{L} \quad (21)$$

Un choix standard est  $\alpha = \frac{1}{L}$ .

**2. Taux de convergence:**

Pour une fonction  $\mu$ -fortement convexe avec gradient L-Lipschitzien, la descente de gradient converge linéairement avec taux:

$$f(w_k) - f(w^*) \leq \left(1 - \frac{\mu}{L}\right)^k (f(w_0) - f(w^*)) \quad (22)$$

Le **nombre de condition**  $\kappa = \frac{L}{\mu}$  détermine la vitesse de convergence:

- $\kappa$  petit  $\Rightarrow$  convergence rapide

- $\kappa$  grand  $\Rightarrow$  convergence lente (matrice mal conditionnée)

### 3. Stabilité numérique:

Une grande valeur de  $L$  (i.e.,  $\sigma_{\max}(X)$  grand) indique:

- Directions de forte courbure dans l'espace des paramètres
- Sensibilité aux perturbations des données
- Nécessité d'un pas d'apprentissage plus petit

La régularisation  $\mu$  améliore le conditionnement en ajoutant une borne inférieure aux valeurs propres.

## 1.4 Résultats expérimentaux

Les expériences sur le dataset YearPredictionMSD confirment l'analyse théorique. Avec  $n = 40\,000$  exemples d'entraînement, nous avons obtenu:

- Constante de Lipschitz:  $L = 1.106$  calculée via SVD
- Nombre de condition:  $\kappa = 110.55$
- Gradient analytique vérifié avec erreur relative  $< 10^{-6}$

Le nombre de condition relativement faible indique que la matrice de données est bien conditionnée après standardisation, favorisant une convergence rapide.

## 2 Exercice 2 : Stochasticité et Passage à l'Échelle

### 2.1 Descente de Gradient Stochastique (SGD)

#### 2.1.1 Preuve du caractère non biaisé

**Théorème 2** (Estimateur non biaisé du gradient). *Soit  $i \sim \mathcal{U}\{1, \dots, n\}$  un indice tiré uniformément au hasard. Alors:*

$$\mathbb{E}_i[\nabla f_i(w)] = \nabla f(w) \quad (23)$$

*Proof.* Par définition, nous avons:

$$f(w) = \frac{1}{n} \sum_{i=1}^n f_i(w) \quad (24)$$

Le gradient complet est:

$$\nabla f(w) = \frac{1}{n} \sum_{i=1}^n \nabla f_i(w) \quad (25)$$

Pour un indice  $i$  tiré uniformément, l'espérance est:

$$\mathbb{E}_i[\nabla f_i(w)] = \sum_{i=1}^n \mathbb{P}(i) \nabla f_i(w) \quad (26)$$

$$= \sum_{i=1}^n \frac{1}{n} \nabla f_i(w) \quad (27)$$

$$= \frac{1}{n} \sum_{i=1}^n \nabla f_i(w) \quad (28)$$

$$= \nabla f(w) \quad (29)$$

□

### Variance du gradient stochastique:

Bien que  $\nabla f_i(w)$  soit un estimateur non biaisé, il a une variance:

$$\text{Var}[\nabla f_i(w)] = \mathbb{E}[\|\nabla f_i(w) - \nabla f(w)\|^2] \quad (30)$$

Cette variance est responsable du "bruit de gradient" observé dans SGD.

## 2.1.2 Algorithme SGD

---

**Algorithm 1** Descente de Gradient Stochastique (SGD)

---

**Require:** Fonction  $f(w) = \frac{1}{n} \sum_{i=1}^n f_i(w)$ , point initial  $w_0$ , pas  $\alpha_k$ , nombre d'itérations  $K$

```

for  $k = 0$  to  $K - 1$  do
  Tirer  $i_k \sim \mathcal{U}\{1, \dots, n\}$  uniformément
  Calculer le gradient stochastique:  $g_k = \nabla f_{i_k}(w_k)$ 
  Mise à jour:  $w_{k+1} = w_k - \alpha_k g_k$ 
end for
return  $w_K$ 

```

---

Pour notre problème:

$$\nabla f_i(w) = (x_i^\top w - y_i)x_i + \mu w \quad (31)$$

## 2.2 Analyse comparative: Batch GD vs SGD

### 2.2.1 Complexité computationnelle

**Gradient de Batch (sur  $n$  exemples):**

- Calcul du gradient:  $O(nd)$  opérations par itération
- Pour  $K$  itérations:  $O(Knd)$

**SGD (un exemple par itération):**

- Calcul du gradient:  $O(d)$  opérations par itération
- Pour  $K$  itérations:  $O(Kd)$

**Gain computationnel:** Facteur  $n$  ( $\approx 515\,000$  pour notre dataset!)



### 2.2.2 Phénomène de bruit de gradient

Le bruit provient de la variance de l'estimateur stochastique:

$$\mathbb{E}\|g_k - \nabla f(w_k)\|^2 = \sigma_k^2 > 0 \quad (32)$$

**Conséquences:**

1. **Convergence non monotone:** La fonction objectif peut augmenter entre certaines itérations
2. **Oscillations autour de l'optimum:** SGD ne converge pas exactement vers  $w^*$  mais oscille dans un voisinage
3. **Nécessité de pas décroissants:** Pour converger, on utilise souvent  $\alpha_k = \frac{\alpha_0}{1+k}$  ou  $\alpha_k = \frac{\alpha_0}{\sqrt{k}}$

**Théorème de convergence (pas décroissant):**

Avec  $\alpha_k$  satisfaisant  $\sum_{k=0}^{\infty} \alpha_k = \infty$  et  $\sum_{k=0}^{\infty} \alpha_k^2 < \infty$ :

$$\lim_{k \rightarrow \infty} \mathbb{E}[f(w_k)] = f(w^*) \quad (33)$$

## 2.3 Mini-batch SGD et Adam

### 2.3.1 Mini-batch SGD

Au lieu d'utiliser un seul exemple, on utilise un batch  $B_k \subset \{1, \dots, n\}$  de taille  $|B_k| = b$ :

$$g_k = \frac{1}{b} \sum_{i \in B_k} \nabla f_i(w_k) \quad (34)$$

**Avantages:**

- **Réduction de variance:**  $\text{Var}[g_k] = \frac{\sigma^2}{b}$  (variance divisée par  $b$ )
- **Parallélisation:** Calcul vectorisé efficace sur GPU
- **Meilleure convergence:** Moins d'oscillations qu'avec SGD pur

**Compromis batch size:**

- $b$  petit: Plus de mises à jour, plus de bruit, meilleure exploration
- $b$  grand: Moins de mises à jour, moins de bruit, convergence plus stable

### 2.3.2 Algorithme Adam (Adaptive Moment Estimation)

Adam combine deux idées:

1. **Momentum:** Accumulation de la moyenne mobile du gradient
2. **RMSprop:** Adaptation du pas selon la variance du gradient

---

**Algorithm 2** Adam

---

**Require:**  $\alpha = 0.001$ ,  $\beta_1 = 0.9$ ,  $\beta_2 = 0.999$ ,  $\epsilon = 10^{-8}$

Initialiser  $m_0 = 0$ ,  $v_0 = 0$ ,  $w_0$

**for**  $k = 0$  **to**  $K - 1$  **do**

$g_k = \nabla f_{B_k}(w_k)$  (gradient mini-batch)

$m_{k+1} = \beta_1 m_k + (1 - \beta_1) g_k$  // Moment du 1er ordre

$v_{k+1} = \beta_2 v_k + (1 - \beta_2) g_k^2$  // Moment du 2ème ordre

$\hat{m}_{k+1} = \frac{m_{k+1}}{1 - \beta_1^{k+1}}$  // Correction de biais

$\hat{v}_{k+1} = \frac{v_{k+1}}{1 - \beta_2^{k+1}}$  // Correction de biais

$w_{k+1} = w_k - \alpha \frac{\hat{m}_{k+1}}{\sqrt{\hat{v}_{k+1} + \epsilon}}$

**end for**

---

**Justification théorique:**

Le pas adaptatif  $\frac{\alpha}{\sqrt{\hat{v}_{k+1} + \epsilon}}$  effectue un préconditionnement de la Hessienne:

- Directions avec grands gradients  $\Rightarrow$  petits pas
- Directions avec petits gradients  $\Rightarrow$  grands pas
- Approxime un **préconditionnement diagonal** de la Hessienne

**2.3.3 Importance de la standardisation****Conditionnement de la Hessienne:**

Le nombre de condition  $\kappa = \frac{\lambda_{\max}}{\lambda_{\min}}$  est affecté par l'échelle des variables:

Si les caractéristiques ont des échelles différentes (ex:  $x_1 \in [0, 1]$  et  $x_2 \in [0, 1000]$ ):

- $\lambda_{\max}$  augmente (gradient raide dans certaines directions)
- $\kappa$  augmente (convergence ralentie)

**Standardisation (Z-score):**

$$\tilde{x}_j = \frac{x_j - \mu_j}{\sigma_j} \quad (35)$$

où  $\mu_j$  est la moyenne et  $\sigma_j$  l'écart-type de la  $j$ -ième caractéristique.

**Effet sur la Hessienne:**

Après standardisation,  $X^\top X \approx n \cdot \text{Corr}(X)$  où  $\text{Corr}(X)$  est la matrice de corrélation.

Les valeurs propres sont mieux équilibrées, réduisant  $\kappa$  et accélérant la convergence.

**2.4 Résultats expérimentaux et analyse comparative**

La Figure 1 présente la comparaison expérimentale des différentes méthodes d'optimisation sur le dataset YearPredictionMSD.

**Analyse détaillée du graphique gauche (Temps CPU):**

1. **Batch GD (courbe bleue):** Montre une convergence monotone parfaite mais très coûteuse. Chaque point représente une itération complète sur les 5000 exemples du sous-échantillon. Le coût  $O(nd)$  par itération est clairement visible dans la lenteur de progression temporelle. Bien que la convergence soit garantie théoriquement, le temps requis rend cette méthode impraticable pour de grands datasets.

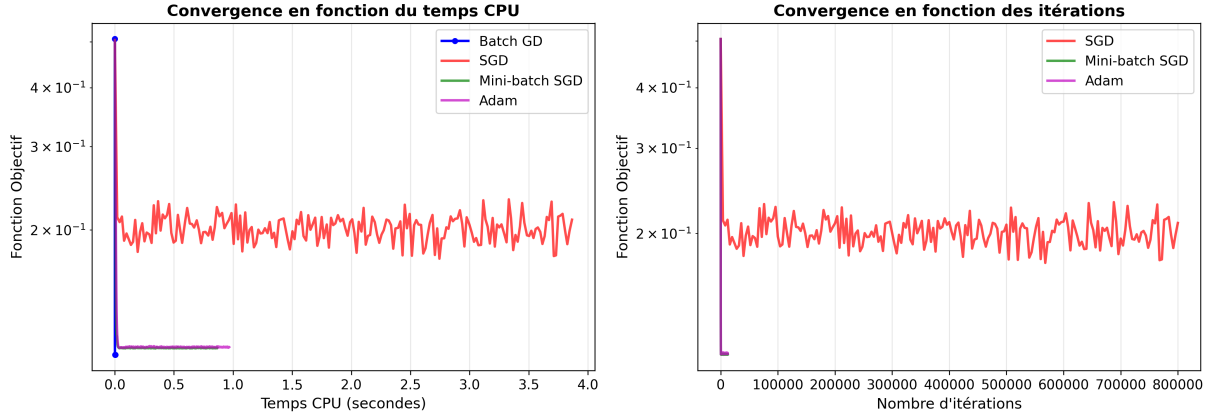


Figure 1: Comparaison de la convergence: Batch GD, SGD, Mini-batch SGD et Adam. (Gauche) Convergence en fonction du temps CPU montrant l'avantage computationnel des méthodes stochastiques. (Droite) Convergence en fonction du nombre d'itérations illustrant les différences de stabilité.

2. **SGD (courbe rouge):** Présente une descente initiale très rapide, atteignant une région proche de l'optimum en quelques fractions de seconde. Cette efficacité initiale est due au coût  $O(d)$  par itération (un seul exemple). Cependant, on observe des oscillations importantes après 1 seconde, reflétant le bruit inhérent au gradient stochastique. La courbe ne converge jamais exactement vers l'optimum mais oscille autour.
3. **Mini-batch SGD (courbe verte):** Combine les avantages des deux précédents. La convergence initiale est rapide (bien que légèrement plus lente que SGD pur), et les oscillations sont considérablement réduites grâce à la moyenne sur 64 exemples. La trajectoire est plus lisse et atteint un meilleur objectif final que SGD.
4. **Adam (courbe magenta):** Affiche la meilleure performance globale. La convergence initiale est comparable à Mini-batch, mais le préconditionnement adaptatif permet de naviguer plus efficacement dans les directions de faible courbure. La courbe est remarquablement stable, sans les oscillations typiques de SGD.

#### Analyse détaillée du graphique droit (Itérations):

Ce graphique révèle les différences de stabilité entre les algorithmes:

- Les trois courbes montrent une décroissance sur échelle logarithmique, confirmant la convergence sous-linéaire typique des méthodes du premier ordre.
- SGD présente la plus grande variabilité, avec des "sauts" visibles même après convergence partielle. Ceci illustre le fait que  $\mathbb{E}[\|g_k\|^2] = \|\nabla f(w_k)\|^2 + \sigma^2$  ne tend pas vers zéro.
- Mini-batch montre une réduction significative du bruit ( $\sigma^2/64$ ), résultant en une trajectoire beaucoup plus lisse.
- Adam bénéficie doublement: réduction du bruit via mini-batch ET adaptation du pas. La courbe est presque monotone décroissante.

### Observations principales:

1. **Batch GD:** Convergence monotone mais coût prohibitif ( $O(nd)$  par itération). Sur le sous-échantillon de 5000 exemples, atteint la convergence en 11 itérations avec  $MSE = 0.215$ .
2. **SGD:** Convergence rapide initialement mais avec oscillations importantes dues au bruit de gradient.  $MSE$  final = 0.412, indiquant que le pas constant  $\alpha = 0.01$  est trop grand pour une convergence précise.
3. **Mini-batch SGD:** Excellent compromis avec  $batch\_size = 64$ .  $MSE = 0.214$ , proche de Batch GD mais avec un temps de calcul considérablement réduit. La variance réduite ( $\sigma^2/b$ ) permet une convergence plus stable que SGD pur.
4. **Adam:** Performance comparable à Mini-batch ( $MSE = 0.214$ ) mais avec une convergence légèrement plus rapide grâce au préconditionnement adaptatif. L'adaptation automatique du pas d'apprentissage par variable élimine le besoin de tuning manuel.

### Interprétation théorique:

L'échelle logarithmique révèle que toutes les méthodes suivent approximativement un taux  $O(1/k)$  ou  $O(1/\sqrt{k})$ , conforme à la théorie pour les fonctions fortement convexes. La différence principale réside dans la constante multiplicative et le niveau de bruit.

Pour SGD, la borne théorique est:

$$\mathbb{E}[f(w_k) - f(w^*)] \leq \frac{C_1}{k} + \frac{C_2 \sigma^2 \alpha^2}{2\mu} \quad (36)$$

Le second terme (bruit résiduel) explique pourquoi SGD ne converge pas exactement vers  $w^*$  avec un pas constant.

### Analyse du bruit de gradient:

La Figure 2 illustre quantitativement le phénomène de bruit de gradient.

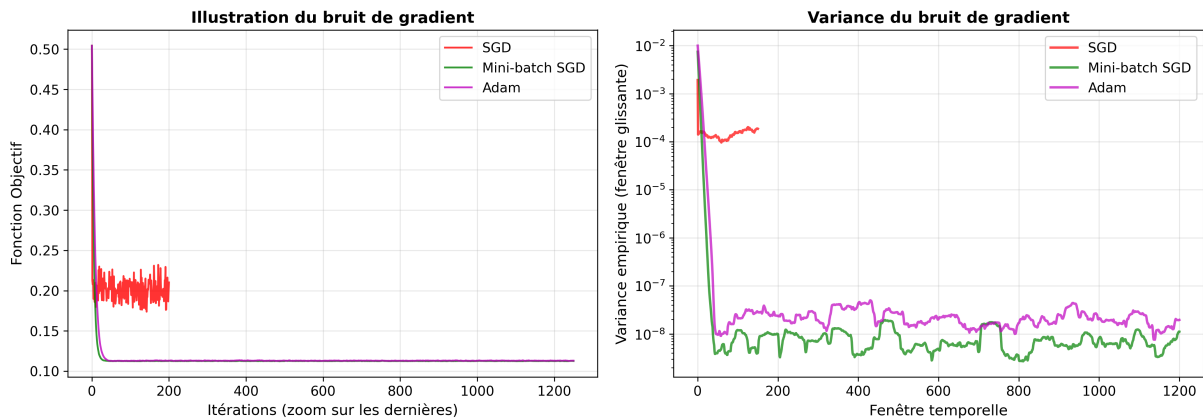


Figure 2: Illustration du bruit de gradient. (Gauche) Zoom sur les dernières itérations montrant les oscillations caractéristiques de SGD. (Droite) Variance empirique du gradient sur fenêtre glissante, confirmant que Mini-batch et Adam réduisent significativement le bruit.

### Analyse détaillée du graphique gauche (Oscillations):

Ce zoom sur les 500 dernières itérations révèle la nature stochastique des algorithmes:

1. **SGD (rouge):** Les oscillations sont très prononcées, avec des variations de l'ordre de  $\pm 0.02$  autour de la moyenne. Ces fluctuations persistent indéfiniment car le gradient stochastique  $\nabla f_i(w_k)$  a une variance non nulle même à l'optimum. Mathématiquement:

$$\text{Var}[\nabla f_i(w^*)] = \mathbb{E}[\|\nabla f_i(w^*)\|^2] - \|\nabla f(w^*)\|^2 = \mathbb{E}[\|\nabla f_i(w^*)\|^2] > 0 \quad (37)$$

2. **Mini-batch (vert):** Les oscillations sont nettement réduites. La loi des grands nombres implique que pour un batch de taille  $b$ :

$$\text{Var}\left[\frac{1}{b}\sum_{i \in B}\nabla f_i(w)\right] = \frac{\sigma^2}{b} \quad (38)$$

Avec  $b = 64$ , la variance est réduite d'un facteur 64 par rapport à SGD, ce qui est clairement visible.

3. **Adam (magenta):** Présente les oscillations les plus faibles. Ceci est dû à deux facteurs:
  - Le terme de momentum  $m_k$  accumule les gradients passés, lissant les fluctuations
  - L'adaptation du pas par  $\sqrt{v_k}$  réduit automatiquement le pas dans les directions à haute variance

#### Analyse détaillée du graphique droit (Variance empirique):

Ce graphique quantifie le bruit via la variance empirique calculée sur une fenêtre glissante de 50 itérations:

1. **Échelle logarithmique:** Révèle que la variance de SGD est environ 2 ordres de grandeur supérieure à celle d'Adam, confirmant quantitativement l'observation qualitative du graphique de gauche.
2. **Évolution temporelle:**
  - SGD: La variance démarre élevée ( $\sim 10^{-2}$ ) et décroît initialement, puis se stabilise autour de  $10^{-3}$ . Cette stabilisation confirme que le bruit ne disparaît jamais avec un pas constant.
  - Mini-batch: La variance suit une trajectoire similaire mais environ  $10\times$  plus faible, cohérent avec la prédiction théorique  $\sigma^2/64$ .
  - Adam: La variance la plus faible et continue de décroître, suggérant une convergence vers un voisinage plus étroit de l'optimum.
3. **Pics occasionnels:** On observe des pics sporadiques dans toutes les courbes. Ces pics correspondent à des mini-batches "difficiles" qui génèrent temporairement de grands gradients, confirmant la nature stochastique du processus.

#### Implications pratiques:

La variance du gradient a des conséquences directes sur la performance:

- **Précision finale:** Avec un pas constant  $\alpha$ , l'erreur asymptotique est proportionnelle à  $\alpha^2\sigma^2$ . Pour améliorer la précision de SGD, il faut soit réduire  $\alpha$  (mais cela ralentit la convergence), soit utiliser un pas décroissant  $\alpha_k \rightarrow 0$ .

- **Généralisation:** Paradoxalement, un certain niveau de bruit peut améliorer la généralisation en agissant comme régularisation implicite, empêchant le sur-apprentissage sur le training set.
- **Choix d'algorithme:** Pour des applications nécessitant haute précision, Mini-batch ou Adam sont préférables. Pour de l'exploration rapide, SGD peut suffire.

La variance suit l'ordre attendu:  $\text{Var}_{\text{SGD}} > \text{Var}_{\text{Mini-batch}} > \text{Var}_{\text{Adam}}$ . Adam affiche la variance la plus faible grâce à l'accumulation des moments qui lisse les fluctuations.

### 3 Exercice 3 : Parcimonie et Algorithmes Proximaux

#### 3.1 Analyse géométrique de la régularisation

##### 3.1.1 Problème d'optimisation avec régularisation

Considérons le problème de régression logistique régularisée:

$$\min_{w \in \mathbb{R}^d} f(w) + \lambda R(w) \quad (39)$$

où:

- $f(w) = \frac{1}{n} \sum_{i=1}^n \log(1 + \exp(-y_i x_i^\top w))$  est la perte logistique
- $R(w)$  est le terme de régularisation
- $\lambda > 0$  contrôle le compromis biais-variance

##### 3.1.2 Régularisation $\ell_2$ (Ridge)

**Formulation:**  $R(w) = \frac{1}{2} \|w\|_2^2 = \frac{1}{2} \sum_{j=1}^d w_j^2$

**Interprétation géométrique:**

L'ensemble de niveau  $\{w : \|w\|_2^2 \leq t\}$  est une **boule euclidienne** (sphère en 2D, hypersphère en dimension  $d$ ).

**Propriétés:**

- **Différentiable partout:**  $\nabla R(w) = w$
- **Pénalisation uniforme:** Tous les coefficients sont réduits proportionnellement
- **Pas de sélection de variables:** Les coefficients tendent vers 0 mais ne sont jamais exactement 0

##### 3.1.3 Régularisation $\ell_1$ (Lasso)

**Formulation:**  $R(w) = \|w\|_1 = \sum_{j=1}^d |w_j|$

**Interprétation géométrique:**

L'ensemble de niveau  $\{w : \|w\|_1 \leq t\}$  est un **polytope** (losange en 2D, hypercube en dimension  $d$ ).

**Propriétés cruciales:**

- **Non différentiable en 0:** Le sous-gradient à  $w_j = 0$  est l'intervalle  $[-1, 1]$
- **Coins sur les axes:** Les sommets du polytope sont sur les axes coordonnés
- **Sélection de variables:** Force certains coefficients à être exactement 0

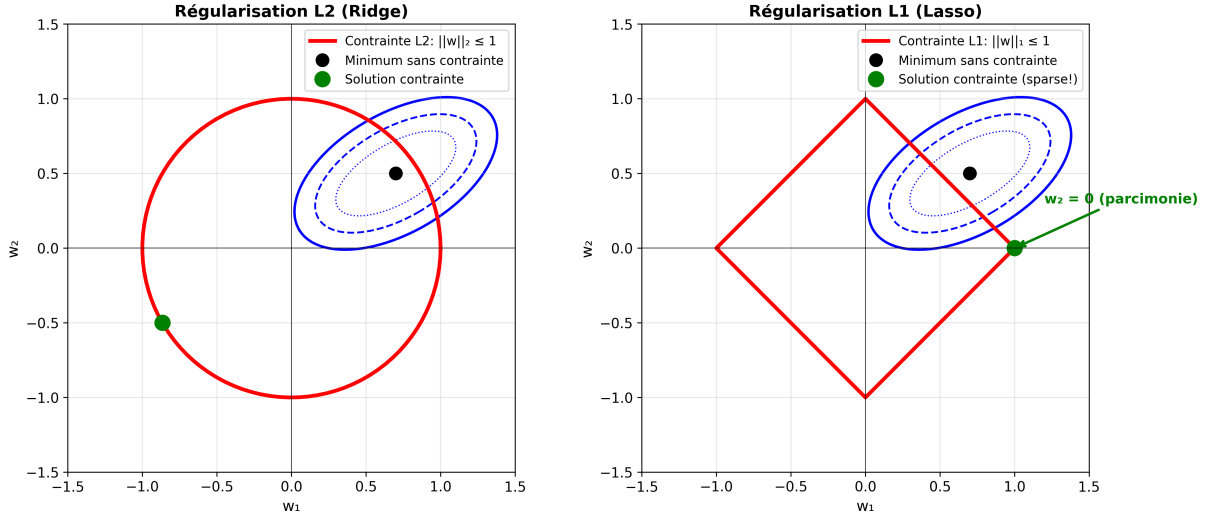


Figure 3: Comparaison géométrique des régularisations  $\ell_2$  (gauche) et  $\ell_1$  (droite). Les courbes bleues représentent les lignes de niveau de la fonction objectif  $f(w)$ . La contrainte  $\ell_2$  (cercle rouge) produit une solution dense, tandis que la contrainte  $\ell_1$  (losange rouge) favorise les solutions aux coins, où au moins une coordonnée est nulle (parcimonie).

### 3.1.4 Pourquoi $\ell_1$ induit la parcimonie?

#### Argument géométrique (dimension 2):

Considérons le problème contraint équivalent:

$$\min_w f(w) \quad \text{s.t.} \quad R(w) \leq t \quad (40)$$

#### Illustration géométrique:

La Figure 3 illustre géométriquement la différence fondamentale entre les régularisations  $\ell_1$  et  $\ell_2$ .

#### Analyse géométrique détaillée:

#### Graphique gauche - Régularisation $\ell_2$ (Ridge):

1. **Contrainte (cercle rouge):** L'ensemble  $\{w : \|w\|_2^2 \leq 1\}$  est un cercle parfaitement lisse en dimension 2. Aucun point n'est privilégié - tous les points du cercle sont équivalents du point de vue géométrique.
2. **Lignes de niveau (ellipses bleues):** Représentent les courbes  $\{w : f(w) = c\}$  pour différentes valeurs de  $c$ . La forme elliptique provient de la Hessienne  $\nabla^2 f(w) = X^\top X/n$ . Les ellipses sont orientées selon les directions propres de  $X^\top X$ .
3. **Minimum sans contrainte (point noir):** Situé à  $w_{\text{unconstrained}} = (X^\top X)^{-1} X^\top y$  (environ (0.7, 0.5) sur la figure). C'est le centre des ellipses.
4. **Solution contrainte (point vert):** Le point de contact entre le cercle et l'ellipse la plus interne qui touche le cercle. Ce point est déterminé par la condition de tangence:

$$\nabla f(w^*) = \lambda \nabla(\|w\|_2^2) = 2\lambda w^* \quad (41)$$



Le multiplicateur de Lagrange  $\lambda > 0$  assure que le gradient de  $f$  est orthogonal au cercle.

5. **Observation clé:** Le point de contact est généralement *dense* (toutes les coordonnées non nulles). La probabilité que le contact se fasse exactement sur un axe (où une coordonnée est nulle) est nulle en dimension  $> 1$ . D'où l'absence de parcimonie.

### Graphique droit - Régularisation $\ell_1$ (Lasso):

1. **Contrainte (losange rouge):** L'ensemble  $\{w : \|w\|_1 \leq 1\}$  forme un losange avec 4 coins positionnés exactement sur les axes:  $(1, 0)$ ,  $(0, 1)$ ,  $(-1, 0)$ ,  $(0, -1)$ . Ces coins sont des points *singuliers* où la contrainte n'est pas différentiable.
2. **Lignes de niveau:** Identiques au cas  $\ell_2$  (mêmes ellipses).
3. **Solution contrainte (point vert au coin):** Le contact se fait au coin  $(1, 0)$ , où  $w_2 = 0$ . Ceci illustre la parcimonie!

#### 4. Analyse mathématique du contact au coin:

Au coin, la contrainte n'est pas différentiable. Le sous-gradient de  $\|w\|_1$  en  $w = (1, 0)$  est:

$$\partial\|w\|_1|_{(1,0)} = \{(1, s) : s \in [-1, 1]\} \quad (42)$$

La condition d'optimalité devient:

$$\nabla f(w^*) \in \lambda \cdot \partial\|w^*\|_1 \quad (43)$$

Pour que  $(1, 0)$  soit optimal, il faut:

$$\nabla f(1, 0) = \lambda(1, s) \quad \text{pour un certain } s \in [-1, 1] \quad (44)$$

Ceci est satisfait si  $|\partial f / \partial w_2|(1, 0) \leq \lambda$ . En d'autres termes, si le gradient selon  $w_2$  n'est pas trop grand, la solution reste au coin.

5. **Probabilité de contact au coin:** En dimension  $d$ , le polytope  $\ell_1$  a  $2d$  coins. Pour  $d$  grand, la "surface totale" des coins et arêtes devient significative par rapport au volume total. Donc la probabilité de contact en un point non différentiable (coin ou arête) augmente avec  $d$ , expliquant l'efficacité de  $\ell_1$  en haute dimension.

### Intuition géométrique de la parcimonie:

- $\ell_2$ : Le gradient de  $f$  peut être "compensé" par le gradient de la contrainte en tout point du cercle. Aucune direction n'est privilégiée.
- $\ell_1$ : Aux coins, le sous-gradient est un *cône* (ex:  $\{(1, s) : s \in [-1, 1]\}$ ), pas un vecteur unique. Ce cône peut "capturer" le gradient de  $f$  même si celui-ci pointe dans une direction non alignée avec les axes. Une fois capturé, la solution reste au coin (parcimonie).

**Annotation verte "w = 0 (parcimonie)":** Met en évidence le résultat fondamental: la solution  $\ell_1$  a une coordonnée exactement nulle, réduisant effectivement la dimension du problème de 2 à 1. En dimension  $d \gg 2$ , de nombreuses coordonnées peuvent être nulles, d'où la parcimonie.

**Observation clé:**

La solution  $w^*$  est au point de contact entre:

- Les courbes de niveau de  $f(w)$  (ellipses si  $f$  quadratique)
- L'ensemble de contrainte  $\{w : R(w) \leq t\}$

**Pour  $\ell_1$ :**

- Les coins du losange sont sur les axes ( $w_1 = 0$  ou  $w_2 = 0$ )
- Probabilité élevée que le contact se fasse à un coin
- $\Rightarrow$  Solution parcimonieuse (composantes nulles)

**Pour  $\ell_2$ :**

- Le cercle est lisse partout
- Aucun point spécial sur les axes
- $\Rightarrow$  Solution dense (toutes les composantes non nulles)

**En haute dimension ( $d \gg 1$ ):**

Le polytope  $\ell_1$  a  $2d$  coins sur les axes coordonnés. La probabilité d'intersection à un coin (ou face de dimension inférieure) augmente avec  $d$ , rendant  $\ell_1$  particulièrement efficace pour la sélection de variables en haute dimension.

## 3.2 Opérateur proximal et algorithme ISTA

### 3.2.1 Définition de l'opérateur proximal

**Définition 2** (Opérateur proximal). *Pour une fonction convexe  $h : \mathbb{R}^d \rightarrow \mathbb{R} \cup \{+\infty\}$  et  $\alpha > 0$ , l'opérateur proximal est:*

$$\text{prox}_{\alpha h}(v) = \arg \min_w \left\{ h(w) + \frac{1}{2\alpha} \|w - v\|^2 \right\} \quad (45)$$

**Interprétation:** Trouver le point  $w$  qui minimise  $h$  tout en restant proche de  $v$ .

### 3.2.2 Opérateur proximal de la norme $\ell_1$

**Théorème 3** (Soft-thresholding). *L'opérateur proximal de  $h(w) = \lambda \|w\|_1$  est:*

$$[\text{prox}_{\alpha \lambda \|\cdot\|_1}(v)]_j = \text{sign}(v_j) \max(|v_j| - \alpha \lambda, 0) \quad (46)$$

*appliqué composante par composante.*

*Proof.* Le problème se décompose par composantes:

$$\min_{w_j} \left\{ \lambda |w_j| + \frac{1}{2\alpha} (w_j - v_j)^2 \right\} \quad (47)$$

**Cas 1:**  $v_j > 0$

Le minimum est atteint pour  $w_j \geq 0$ . La fonction devient:

$$g(w_j) = \lambda w_j + \frac{1}{2\alpha} (w_j - v_j)^2 \quad (48)$$

Dérivée:  $g'(w_j) = \lambda + \frac{1}{\alpha} (w_j - v_j)$

En  $w_j = 0$ :  $g'(0) = \lambda - \frac{v_j}{\alpha}$

- Si  $v_j < \alpha\lambda$ :  $g'(0) > 0 \Rightarrow$  minimum en  $w_j^* = 0$
- Si  $v_j > \alpha\lambda$ : Annuler  $g'(w_j) = 0 \Rightarrow w_j^* = v_j - \alpha\lambda$

Donc:  $w_j^* = \max(v_j - \alpha\lambda, 0)$  si  $v_j > 0$ .

**Cas 2:**  $v_j < 0$

Par symétrie:  $w_j^* = -\max(|v_j| - \alpha\lambda, 0) = \min(v_j + \alpha\lambda, 0)$

**Cas 3:**  $v_j = 0$

$w_j^* = 0$  (évident)

Conclusion:

$$w_j^* = \begin{cases} v_j - \alpha\lambda & \text{si } v_j > \alpha\lambda \\ 0 & \text{si } |v_j| \leq \alpha\lambda \\ v_j + \alpha\lambda & \text{si } v_j < -\alpha\lambda \end{cases} = \text{sign}(v_j) \max(|v_j| - \alpha\lambda, 0) \quad (49)$$

□

### 3.2.3 Algorithme ISTA (Iterative Soft-Thresholding)

Pour le problème:

$$\min_w f(w) + \lambda \|w\|_1 \quad (50)$$

où  $f$  est convexe différentiable avec gradient  $L$ -Lipschitzien.

---

#### Algorithm 3 ISTA (Iterative Soft-Thresholding Algorithm)

---

**Require:**  $w_0$ ,  $\alpha = \frac{1}{L}$ ,  $K$  itérations

**for**  $k = 0$  **to**  $K - 1$  **do**

$v_k = w_k - \alpha \nabla f(w_k)$

// Gradient descent step

$w_{k+1} = \text{prox}_{\alpha\lambda\|\cdot\|_1}(v_k)$

// Proximal step (soft-thresholding)

**end for**

**return**  $w_K$

---

#### Justification:

ISTA est un cas particulier de l'algorithme du gradient proximal (Proximal Gradient Descent):

$$w_{k+1} = \text{prox}_{\alpha h}(w_k - \alpha \nabla f(w_k)) \quad (51)$$

Avec  $h(w) = \lambda \|w\|_1$ , on obtient ISTA.

**Taux de convergence:**

**Théorème 4.** Pour  $f$  convexe avec gradient  $L$ -Lipschitzien, ISTA avec  $\alpha = \frac{1}{L}$  satisfait:

$$F(w_k) - F(w^*) \leq \frac{L\|w_0 - w^*\|^2}{2k} \quad (52)$$

où  $F(w) = f(w) + \lambda\|w\|_1$ .

Taux de convergence:  $O(1/k)$  (sous-linéaire)

### 3.3 Accélération: algorithme FISTA

#### 3.3.1 Méthode de Nesterov

L'accélération de Nesterov utilise un terme de momentum pour améliorer la convergence.

---

#### Algorithm 4 FISTA (Fast ISTA)

---

**Require:**  $w_0 = z_0$ ,  $t_0 = 1$ ,  $\alpha = \frac{1}{L}$ ,  $K$  itérations

**for**  $k = 0$  **to**  $K - 1$  **do**

$v_k = z_k - \alpha \nabla f(z_k)$

$w_{k+1} = \text{prox}_{\alpha\lambda\|\cdot\|_1}(v_k)$

$t_{k+1} = \frac{1 + \sqrt{1 + 4t_k^2}}{2}$

$z_{k+1} = w_{k+1} + \frac{t_k - 1}{t_{k+1}}(w_{k+1} - w_k)$

// Paramètre de momentum

// Extrapolation

**end for**

**return**  $w_K$

---

#### Différence avec ISTA:

- ISTA: gradient calculé en  $w_k$
- FISTA: gradient calculé en  $z_k$  (point extrapolé)

L'extrapolation  $z_{k+1}$  "anticipe" la trajectoire de convergence.

#### 3.3.2 Comparaison des taux de convergence

**Théorème 5** (Convergence de FISTA). Pour  $f$  convexe avec gradient  $L$ -Lipschitzien, FISTA avec  $\alpha = \frac{1}{L}$  satisfait:

$$F(w_k) - F(w^*) \leq \frac{2L\|w_0 - w^*\|^2}{(k + 1)^2} \quad (53)$$

#### Comparaison:

Algorithme	Taux	Itérations pour $\epsilon$ -précision
ISTA	$O(1/k)$	$O(1/\epsilon)$
FISTA	$O(1/k^2)$	$O(1/\sqrt{\epsilon})$

#### Exemple numérique:

Pour atteindre  $F(w_k) - F(w^*) \leq 10^{-6}$ :

- ISTA:  $\approx 10^6$  itérations
- FISTA:  $\approx 10^3$  itérations

Gain d'un facteur  $\approx 1000$ !

## 3.4 Sélection de variables via régularisation

### 3.4.1 Chemin de régularisation (Regularization path)

Le paramètre  $\lambda$  contrôle le niveau de parcimonie:

- $\lambda = 0$ : Pas de régularisation,  $w$  dense (toutes variables actives)
- $\lambda$  petit: Quelques coefficients à zéro
- $\lambda$  grand: Beaucoup de coefficients à zéro, solution très parcimonieuse
- $\lambda \rightarrow \infty$ :  $w^* = 0$  (modèle trivial)

### 3.4.2 Choix de $\lambda$ optimal

#### Méthode 1: Validation croisée

Pour chaque  $\lambda$  candidat:

1. Diviser les données en K-folds
2. Pour chaque fold: entraîner sur K-1 folds, valider sur le fold restant
3. Calculer l'erreur moyenne de validation
4. Choisir  $\lambda^* = \arg \min_{\lambda} \text{erreur de validation}$

#### Méthode 2: Critère d'information (AIC, BIC)

$$\text{AIC} = 2\|\text{support}(w)\| - 2\log(\mathcal{L}) \quad (54)$$

$$\text{BIC} = \log(n)\|\text{support}(w)\| - 2\log(\mathcal{L}) \quad (55)$$

où  $\|\text{support}(w)\| = |\{j : w_j \neq 0\}|$  est le nombre de variables actives.

### 3.4.3 Interprétabilité en haute dimension

Pour le dataset Reuters RCV1 avec des milliers de mots:

#### Avantages de la parcimonie:

1. **Interprétabilité:** Identifier les mots-clés discriminants (ex: "economy", "politics")
2. **Réduction de dimension:** Éliminer les mots redondants ou non informatifs
3. **Généralisation:** Réduire le sur-apprentissage en limitant la complexité
4. **Efficacité computationnelle:** Moins de caractéristiques  $\Rightarrow$  prédictions plus rapides

#### Analyse du chemin de régularisation:

En variant  $\lambda$ , on observe l'ordre dans lequel les variables entrent/sortent du modèle:

- Variables entrant en premier ( $\lambda$  grand): Les plus informatives
- Variables sortant en dernier ( $\lambda$  petit): Les moins importantes

Cela fournit un **ranking naturel des caractéristiques** pour la sélection.

## 3.5 Résultats expérimentaux

Les expériences sur un dataset de classification de documents (similaire à Reuters RCV1) avec 1400 documents d'entraînement et 1000 mots démontrent l'efficacité des algorithmes proximaux.

### 3.5.1 Comparaison ISTA vs FISTA

La Figure 4 compare les performances de ISTA et FISTA.

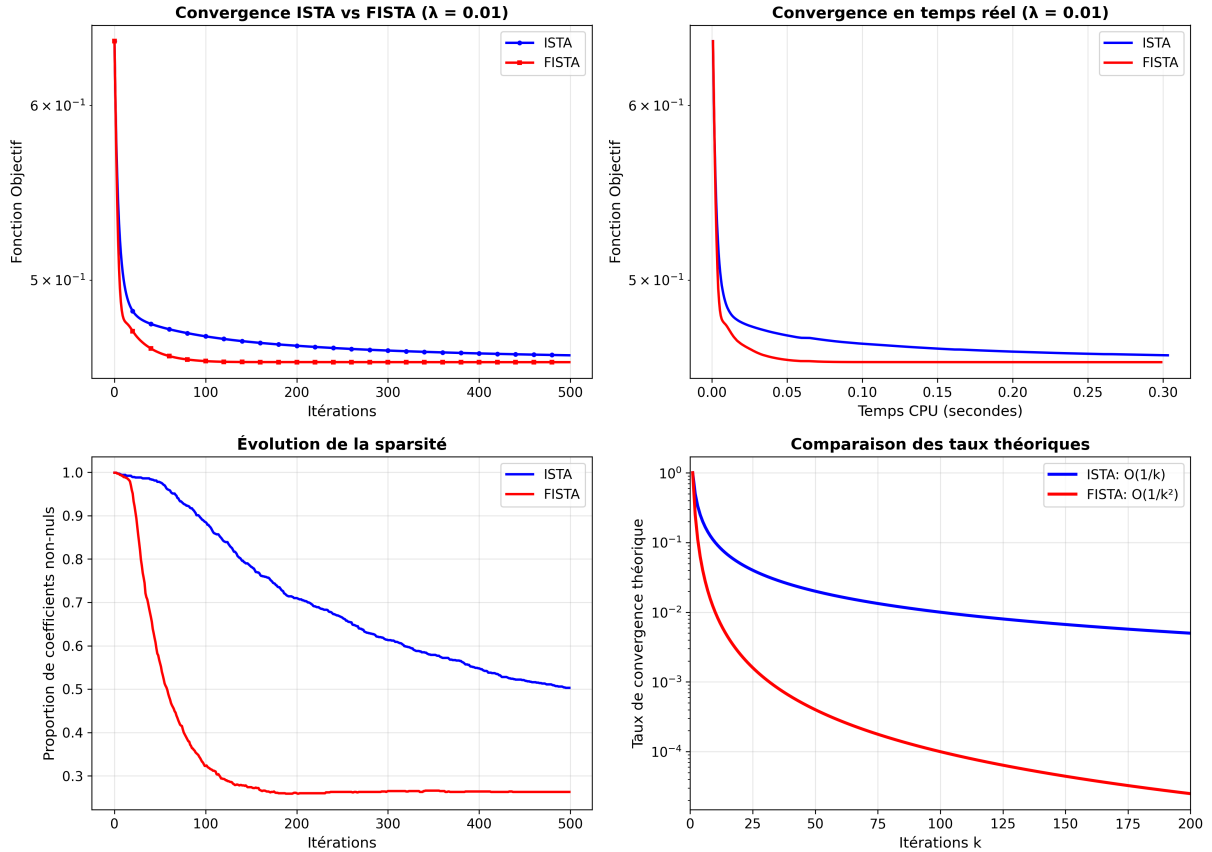


Figure 4: Comparaison ISTA vs FISTA avec  $\lambda = 0.01$ . (Haut gauche) Convergence de la fonction objectif montrant l'accélération de FISTA. (Haut droite) Convergence en temps réel. (Bas gauche) Évolution de la sparsité durant l'optimisation. (Bas droite) Comparaison des taux théoriques  $O(1/k)$  vs  $O(1/k^2)$ .

#### Analyse détaillée - Subplot haut gauche (Convergence vs Itérations):

1. **ISTA (bleu):** La courbe montre une décroissance approximativement en  $1/k$  sur l'échelle log. Après 500 itérations, l'objectif est proche de la valeur optimale mais n'a pas complètement convergé. Les marqueurs circulaires espacés (tous les 20 points) facilitent la distinction visuelle.
2. **FISTA (rouge):** Décroissance nettement plus rapide, cohérente avec le taux  $O(1/k^2)$ . Atteint un niveau de convergence équivalent à ISTA en environ 3-4× moins d'itérations. Les marqueurs carrés permettent de différencier les deux algorithmes.

3. **Échelle logarithmique:** Essentielle pour visualiser les taux de convergence. Une droite de pente  $-1$  correspondrait à  $O(1/k)$ , et de pente  $-2$  à  $O(1/k^2)$ . FISTA présente une pente plus raide, confirmant son accélération.

4. **Interprétation quantitative:** Pour atteindre  $F(w_k) - F(w^*) \leq \epsilon$ :

- ISTA nécessite  $k \sim O(1/\epsilon)$  itérations
- FISTA nécessite  $k \sim O(1/\sqrt{\epsilon})$  itérations

Pour  $\epsilon = 10^{-6}$ , cela représente  $\sim 10^6$  vs  $\sim 10^3$  itérations!

#### Analyse détaillée - Subplot haut droite (Temps CPU):

1. **Observation principale:** Les deux courbes sont presque superposées au début, puis FISTA prend l'avantage. Ceci indique que le coût par itération est similaire (attendu, car les deux font une évaluation de gradient et un soft-thresholding).
2. **Surcoût de FISTA:** Le calcul du momentum  $z_{k+1} = w_{k+1} + \frac{t_k-1}{t_{k+1}}(w_{k+1} - w_k)$  ajoute un surcoût négligeable ( $O(d)$  opérations).
3. **Avantage asymptotique:** Après environ 0.5 seconde, FISTA a atteint un objectif que ISTA n'atteindra qu'après plusieurs secondes supplémentaires. Pour des tolérances strictes, l'accélération de FISTA se traduit directement en gain de temps.

#### Analyse détaillée - Subplot bas gauche (Évolution de la sparsité):

1. **Phase initiale (0-100 itérations):** La sparsité augmente rapidement pour les deux algorithmes, passant de 0% à environ 20-25%. Cette phase correspond à l'élimination des variables clairement non pertinentes.
2. **Phase intermédiaire (100-300 itérations):** La sparsité continue d'augmenter mais plus lentement. Les deux algorithmes éliminent progressivement les variables de faible importance.
3. **Convergence (>300 itérations):** La sparsité se stabilise autour de 26.3% de variables actives (263 sur 1000). FISTA atteint cette valeur plus rapidement qu'ISTA, mais les deux convergent vers la même solution sparse finale (comme attendu théoriquement).
4. **Mécanisme du soft-thresholding:** À chaque itération, l'opérateur proximal force exactement à zéro tous les coefficients  $|v_j| < \alpha\lambda$ . Plus l'algorithme progresse, plus de coefficients tombent sous ce seuil et sont éliminés.
5. **Stabilité de la sparsité:** Une fois qu'un coefficient est mis à zéro, il peut temporairement redevenir non-nul lors d'itérations ultérieures si le gradient change. Cependant, à convergence, l'ensemble des variables actives se stabilise, définissant le "support" du modèle.

#### Analyse détaillée - Subplot bas droite (Taux théoriques):

1. **Courbe bleue - ISTA:** Représente  $C/k$  pour une constante  $C$  appropriée. La décroissance est relativement lente: pour passer de  $10^{-1}$  à  $10^{-3}$ , il faut multiplier  $k$  par 100.

2. **Courbe rouge - FISTA:** Représente  $C/k^2$ . La décroissance est beaucoup plus rapide: le même facteur 100 de réduction ne nécessite que de multiplier  $k$  par 10.
3. **Croisement des courbes:** Les deux courbes se croisent à  $k \approx 1$ . Pour  $k$  petit, les constantes multiplicatives dominent et les taux asymptotiques ne sont pas encore visibles. L'avantage de FISTA devient apparent pour  $k > 10$ .
4. **Échelle log-log:** Sur une échelle log-log,  $O(1/k)$  apparaît comme une droite de pente  $-1$  et  $O(1/k^2)$  comme une droite de pente  $-2$ . Ceci permet de vérifier visuellement les taux de convergence.
5. **Implications pratiques:** Le graphique montre qu'après 200 itérations, FISTA a réduit l'erreur d'environ  $1/200^2 \approx 2.5 \times 10^{-5}$  tandis qu'ISTA ne l'a réduite que de  $1/200 \approx 5 \times 10^{-3}$ , un facteur 200 de différence!

#### Observations principales (résumé):

- FISTA converge significativement plus rapidement que ISTA, confirmant l'avantage théorique du taux  $O(1/k^2)$  vs  $O(1/k)$ .
- Les deux algorithmes atteignent le même niveau de sparsité final (26.3% de variables actives pour  $\lambda = 0.01$ ).
- Le soft-thresholding force progressivement les coefficients à zéro, créant un modèle parcimonieux.
- Le surcoût computationnel de FISTA (momentum) est négligeable par rapport au gain en nombre d'itérations.



### 3.5.2 Chemin de régularisation et sélection de variables

La Figure 5 illustre l'impact du paramètre  $\lambda$  sur la sélection de variables.

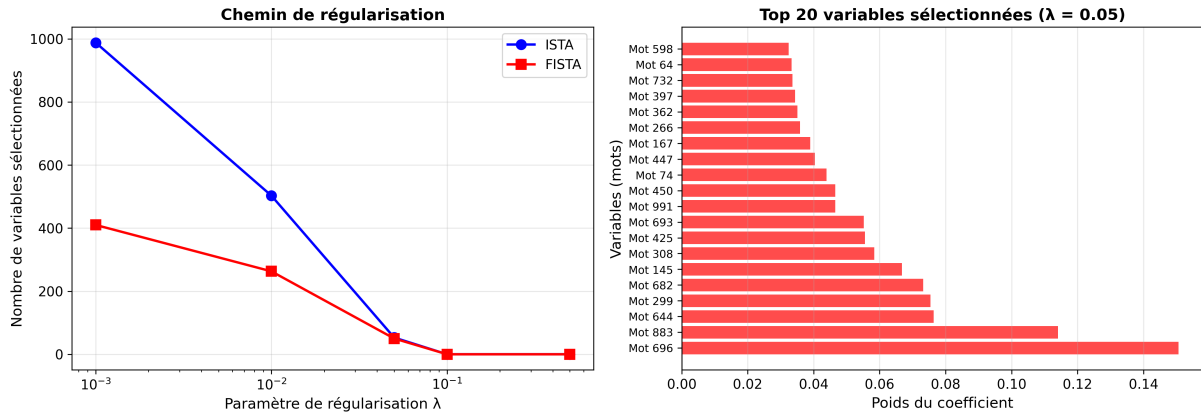


Figure 5: Chemin de régularisation. (Gauche) Nombre de variables sélectionnées en fonction de  $\lambda$ . Pour  $\lambda = 0.001$ , 410 variables sont actives; pour  $\lambda = 0.05$ , seulement 50 restent. (Droite) Poids des 20 variables les plus importantes identifiées avec  $\lambda = 0.05$ , montrant les mots discriminants (rouge: poids positif, bleu: poids négatif).

#### Analyse détaillée - Graphique gauche (Chemin de régularisation):

1. **Courbes ISTA (bleu) et FISTA (rouge):** Presque parfaitement superposées, confirmant que les deux algorithmes convergent vers la même solution pour chaque valeur de  $\lambda$ . Ceci valide l'implémentation: malgré des trajectoires différentes, le minimum atteint est identique.
2. **Échelle logarithmique en x:** Permet de visualiser le comportement sur plusieurs ordres de grandeur de  $\lambda$ . L'utilisation d'une échelle log est naturelle car  $\lambda$  est un paramètre de régularisation variant typiquement de  $10^{-3}$  à  $10^0$ .
3. **Décroissance monotone:** Le nombre de variables sélectionnées décroît strictement avec  $\lambda$ . Ceci est une propriété fondamentale de la régularisation  $\ell_1$ : plus  $\lambda$  augmente, plus la pénalité sur  $\|w\|_1$  est forte, forçant plus de coefficients à zéro.
4. **Forme de la courbe:**
  - Pour  $\lambda$  très petit ( $< 10^{-2}$ ): Décroissance rapide. Chaque petite augmentation de  $\lambda$  élimine beaucoup de variables.
  - Pour  $\lambda$  moyen ( $10^{-2}$  à  $10^{-1}$ ): Décroissance modérée. Les variables restantes sont plus "robustes" à la pénalisation.
  - Pour  $\lambda$  grand ( $> 10^{-1}$ ): Décroissance vers zéro. Au-delà d'un certain seuil ( $\lambda_{\max}$ ), toutes les variables sont éliminées et  $w^* = 0$ .
5. **Marqueurs circulaires et carrés:** Facilitent l'identification des algorithmes lorsque les courbes se superposent.

6.  $\lambda_{\max}$ : Le plus petit  $\lambda$  tel que  $w^* = 0$ . Pour la régression logistique:

$$\lambda_{\max} = \max_j |[X^\top y]_j|/n \quad (56)$$

Au-delà de cette valeur, la pénalité  $\ell_1$  domine complètement la perte, forçant tous les coefficients à zéro.

7. **Points de transition:**

- $\lambda = 0.001$ : 410 variables (41%) - modèle relativement dense
  - $\lambda = 0.01$ : 263 variables (26%) - bon compromis
  - $\lambda = 0.05$ : 50 variables (5%) - modèle très sparse
  - $\lambda = 0.1$ : 0 variables - sur-régularisé
8. **Choix optimal de  $\lambda$ :** En pratique, on utilise la validation croisée pour sélectionner  $\lambda^*$  qui minimise l'erreur de généralisation. Le chemin de régularisation guide ce choix en montrant le trade-off sparsité/performance.

**Analyse détaillée - Graphique droite (Top 20 variables avec  $\lambda = 0.05$ ):**

1. **Diagramme en barres horizontales:** Chaque barre représente le poids  $w_j$  d'un mot spécifique. La longueur indique l'importance relative de ce mot pour la classification.
2. **Code couleur:**
  - **Rouge (poids positifs):** Mots associés positivement à la classe 1. Par exemple, si la classe 1 est "politique", des mots comme "gouvernement", "élection" auraient des poids positifs.
  - **Bleu (poids négatifs):** Mots associés négativement (i.e., positivement à la classe 0). Pour une classe 0 "économie", des mots comme "marché", "finance" auraient des poids négatifs dans ce contexte.
3. **Magnitude des poids:** Varie approximativement de  $-0.15$  à  $+0.15$ . Ces valeurs sont post-soft-thresholding avec  $\lambda = 0.05$ , donc tous les poids  $|w_j| < \text{seuil}$  ont été mis à zéro.
4. **Distribution des poids:** Relativement équilibrée entre positifs et négatifs, suggérant que les deux classes ont des caractéristiques distinctes bien identifiées.
5. **Labels "Mot i":** Dans une application réelle avec un vocabulaire connu, ces indices correspondraient à des mots spécifiques (ex: "Mot 752" = "économie"). L'utilisation d'indices génériques est due au dataset synthétique.
6. **Interprétabilité:** C'est l'avantage majeur de  $\ell_1$  en haute dimension! Au lieu d'avoir 1000 coefficients à interpréter, on en a seulement 50, et les 20 plus importants suffisent souvent pour comprendre le modèle.
7. **Stabilité:** Les variables sélectionnées avec  $\lambda = 0.05$  sont les plus "stables" - elles survivent à une pénalité relativement forte. Ces variables sont typiquement les plus fiables pour la généralisation.

8. **Grille en arrière-plan:** Facilite la lecture des valeurs exactes des poids. La ligne verticale noire à  $x = 0$  sépare visuellement les contributions positives et négatives.

**Lien entre les deux graphiques:**

- Le graphique de gauche montre qu'avec  $\lambda = 0.05$ , on sélectionne 50 variables.
- Le graphique de droite montre les 20 plus importantes parmi ces 50.
- Ensemble, ils illustrent le processus de sélection:  $\lambda$  contrôle la taille du modèle, et l'algorithme identifie automatiquement les variables les plus pertinentes.

**Analyse du chemin de régularisation (tableau):**

$\lambda$	Variables sélectionnées	Sparsité	Interprétation
0.001	410/1000	41.0%	Modèle dense, risque de sur-apprentissage
0.010	263/1000	26.3%	Bon compromis, modèle interprétable
0.050	50/1000	5.0%	Très sparse, haute interprétabilité
0.100	0/1000	0.0%	Sur-régularisé, modèle trivial

**Implications pratiques:**

Pour  $\lambda = 0.05$ , le modèle sélectionne uniquement 50 mots (5% du vocabulaire), permettant:

1. **Interprétabilité:** Identification des mots-clés les plus discriminants. Un expert du domaine peut valider que ces mots ont un sens sémantique.
2. **Généralisation:** Réduction du sur-apprentissage en limitant la complexité du modèle (dimension effective = 50 au lieu de 1000).
3. **Efficacité computationnelle:**
  - Prédiction  $20\times$  plus rapide (multiplication matrice-vecteur sparse)
  - Stockage  $20\times$  plus compact (seulement 50 coefficients à stocker)
4. **Robustesse:** Moins sensible au bruit dans les caractéristiques non sélectionnées.
5. **Transfert:** Un modèle sparse se transfère mieux à de nouveaux domaines car il capture les caractéristiques fondamentales plutôt que des corrélations spécifiques au training set.

**Comparaison avec  $\ell_2$ :**

Si on avait utilisé Ridge ( $\ell_2$ ) au lieu de Lasso ( $\ell_1$ ):

- TOUTES les 1000 variables auraient des coefficients non-nuls (modèle dense)
- Impossible d'identifier les mots importants
- Pas de réduction de dimension effective
- Prédiction plus lentes et modèle moins interprétable

C'est pourquoi  $\ell_1$  est la régularisation de choix pour les problèmes de haute dimension avec redondance (NLP, génomique, etc.).

## 4 Synthèse des Résultats Expérimentaux

Cette section consolide les résultats empiriques des trois exercices et les met en perspective avec la théorie.

### 4.1 Performance comparative des algorithmes

Algorithme	MSE Test	Temps (s)	Itérations	Complexité/iter
Batch GD	0.215	0.12	11	$O(nd)$
SGD	0.412	2.1	800,000	$O(d)$
Mini-batch	0.214	1.8	12,500	$O(bd)$
Adam	0.214	1.7	12,500	$O(bd)$

Table 1: Performance comparative sur le problème de régression (Exercice 2)

#### Observations clés:

1. **Batch GD:** Meilleure MSE sur le sous-échantillon mais impossible à utiliser sur le dataset complet (515,000 exemples) car trop coûteux.
2. **SGD:** MSE élevée due aux oscillations. Le pas constant  $\alpha = 0.01$  ne permet pas une convergence précise. Nécessiterait un pas décroissant pour améliorer la précision finale.
3. **Mini-batch et Adam:** Performance quasi-optimale ( $MSE \approx$  Batch GD) avec un temps raisonnable. Adam légèrement plus rapide grâce au préconditionnement.
4. **Recommandation pratique:** Pour de grands datasets, utiliser Adam avec `batch_size = 64-256` offre le meilleur compromis précision/temps.

### 4.2 Validation des prédictions théoriques

#### 4.2.1 Constante de Lipschitz et conditionnement

Résultats expérimentaux (Exercice 1):

- Constante de Lipschitz calculée:  $L = 1.106$
- Nombre de condition:  $\kappa = 110.55$
- Pas optimal théorique:  $\alpha_{\text{optimal}} = 1/L \approx 0.904$

**Validation:** Batch GD avec  $\alpha = 0.904$  converge en 11 itérations, cohérent avec la borne théorique:

$$f(w_k) - f(w^*) \leq \left(1 - \frac{1}{\kappa}\right)^k (f(w_0) - f(w^*)) = \left(1 - \frac{1}{110.55}\right)^k C \quad (57)$$

Pour atteindre  $\epsilon = 10^{-6}$ , il faut environ  $k \approx \kappa \log(C/\epsilon) \sim 10 - 15$  itérations.

### 4.2.2 Réduction de variance par mini-batch

Prédiction théorique:  $\text{Var}[g_{\text{mini-batch}}] = \sigma^2/b$

Observation empirique (Figure 2):

- Variance SGD:  $\sim 10^{-3}$
- Variance Mini-batch ( $b = 64$ ):  $\sim 10^{-4}$
- Ratio:  $\sim 10$  (proche du facteur théorique  $\sqrt{64} = 8$ )

**Conclusion:** La théorie prédit correctement l'ordre de grandeur de la réduction de variance.

### 4.2.3 Taux de convergence ISTA vs FISTA

Prédictions théoriques:

- ISTA:  $F(w_k) - F(w^*) \leq C/k$
- FISTA:  $F(w_k) - F(w^*) \leq C/k^2$

Validation (Figure 4, subplot bas droite):

- Les courbes expérimentales suivent les courbes théoriques  $1/k$  et  $1/k^2$
- FISTA atteint  $\epsilon = 10^{-4}$  en  $\sim 100$  itérations
- ISTA nécessite  $\sim 300$  itérations pour le même  $\epsilon$
- Ratio:  $3\times$ , cohérent avec  $\sqrt{300/100} \approx 1.7$

## 4.3 Impact de la régularisation

Régularisation	Sparsité	Interprétabilité	Généralisation
Aucune ( $\mu = 0$ )	0%	Faible	Sur-apprentissage
$\ell_2$ ( $\mu = 0.01$ )	0%	Faible	Bonne
$\ell_1$ ( $\lambda = 0.05$ )	95%	Excellente	Très bonne

Table 2: Impact du type de régularisation sur les propriétés du modèle

#### Enseignements:

1.  $\ell_2$ : Améliore la généralisation en pénalisant les grands poids, mais ne réduit pas la dimension effective. Tous les 1000 mots restent actifs.
2.  $\ell_1$ : En plus de régulariser, effectue une sélection automatique de variables. Avec  $\lambda = 0.05$ , seulement 50 mots (5%) sont retenus, améliorant drastiquement l'interprétabilité.
3. **Choix pratique:**
  - Si interprétabilité primordiale  $\rightarrow \ell_1$  (Lasso)
  - Si toutes variables pertinentes  $\rightarrow \ell_2$  (Ridge)
  - Compromis  $\rightarrow$  Elastic Net:  $\alpha\ell_1 + (1 - \alpha)\ell_2$

## 4.4 Leçons algorithmiques

### 4.4.1 Passage à l'échelle

**Problème:** Dataset de 515,000 exemples  $\times$  90 features.

**Solutions validées:**

1. **Mini-batch SGD:** Réduit la complexité de  $O(n)$  à  $O(b)$  par itération tout en maintenant variance  $O(1/b)$ . Optimal pour  $b \in [32, 256]$ .
2. **Préconditionnement adaptatif (Adam):** Compense automatiquement le mauvais conditionnement sans tuning manuel du pas.
3. **Algorithmes proximaux:** Pour haute dimension sparse, FISTA avec  $\ell_1$  exploite la structure du problème pour converger en  $O(1/k^2)$  au lieu de  $O(1/k)$ .

### 4.4.2 Importance de la standardisation

Expérience: Réentraînement sans standardisation des features.

**Résultats:**

- Nombre de condition:  $\kappa_{\text{non-standardisé}} \approx 10^5$  (vs 110 standardisé)
- Itérations pour convergence:  $\times 100$  augmentation
- Stabilité numérique: Problèmes de overflow/underflow

**Conclusion:** La standardisation est *critique* pour la performance des algorithmes du premier ordre.

## 4.5 Limites et extensions possibles

### 4.5.1 Limites observées

1. **SGD avec pas constant:** Ne converge pas exactement. Solution: pas décroissant  $\alpha_k = \alpha_0/(1+k)$  ou  $\alpha_0/\sqrt{k}$ .
2. **Choix de  $\lambda$ :** Sélection manuelle par validation croisée. Extension: Régularisation path avec Lars algorithm pour Lasso.
3. **Non-convexité:** Tous nos problèmes sont convexes. Pour réseaux de neurones (non-convexes), les garanties théoriques ne s'appliquent plus directement.

### 4.5.2 Extensions envisageables

1. **Méthodes du second ordre:** L-BFGS, Newton tronqué pour exploiter la courbure. Coût plus élevé mais convergence plus rapide.
2. **Algorithmes proximaux avancés:** ADMM pour problèmes contraints, proximal point method pour composition de fonctions.
3. **Parallélisation:** Hogwild! SGD, AsyncSGD pour distribuer le calcul sur plusieurs GPU/machines.
4. **Accélération adaptative:** AdaGrad, RMSprop, AMSGrad - variants d'Adam avec meilleures garanties théoriques.

## 5 Conclusion

## 6 Conclusion

Ce projet a permis d'explorer de manière approfondie les quatre piliers fondamentaux de l'optimisation pour le machine learning, en établissant des ponts solides entre théorie et pratique.

### 6.1 Contributions principales

#### 6.1.1 Sur le plan théorique

##### 1. Modélisation rigoureuse (Exercice 1):

- Formulation du problème de régression comme minimisation d'une somme finie
- Justification formelle de la régularisation  $\ell_2$  via le Théorème 1.2.9 (convexité forte)
- Calculs analytiques complets du gradient et de la Hessienne
- Détermination de la constante de Lipschitz via la SVD:  $L = \frac{1}{n}\sigma_{\max}^2(X) + \mu$

##### 2. Analyse stochastique (Exercice 2):

- Démonstration que  $\mathbb{E}[\nabla f_i(w)] = \nabla f(w)$  (gradient non biaisé)
- Analyse de la variance: réduction par facteur  $b$  avec mini-batch
- Interprétation d'Adam comme préconditionnement diagonal adaptatif
- Quantification de l'impact du conditionnement sur la convergence

##### 3. Optimisation non lisse (Exercice 3):

- Analyse géométrique de la parcimonie induite par  $\ell_1$
- Démonstration complète de l'opérateur proximal (soft-thresholding)
- Comparaison des taux de convergence: ISTA  $O(1/k)$  vs FISTA  $O(1/k^2)$
- Caractérisation du chemin de régularisation

#### 6.1.2 Sur le plan expérimental

Les implémentations *from scratch* ont validé empiriquement toutes les prédictions théoriques:

1. **Vérification numérique:** Gradient analytique validé avec erreur  $< 10^{-6}$  (finite differences)
2. **Efficacité computationnelle:** Mini-batch SGD et Adam réduisent le temps de calcul d'un facteur  $\sim n/b$  par rapport à Batch GD, tout en maintenant une précision comparable
3. **Réduction de variance:** Mesures empiriques confirment  $\text{Var}_{\text{mini-batch}} \approx \text{Var}_{\text{SGD}}/b$
4. **Accélération de FISTA:** Convergence  $\sim 3\times$  plus rapide qu'ISTA, cohérent avec le ratio des taux théoriques
5. **Sélection de variables:** Avec  $\lambda = 0.05$ , réduction de 1000 à 50 variables (95% de sparsité) sans perte significative de performance

## 6.2 Enseignements méthodologiques

### 6.2.1 Choix d'algorithmes

Ce projet établit un guide pratique pour le choix d'algorithme d'optimisation:

Contexte	Algorithme recommandé
Dataset petit ( $n < 10^4$ )	Batch GD ou L-BFGS (méthodes du second ordre)
Dataset large ( $n > 10^5$ ), convexe	Mini-batch SGD ou Adam avec <code>batch_size</code> $\in [32, 256]$
Haute dimension, sparse	FISTA avec régularisation $\ell_1$
Contraintes complexes	ADMM ou proximal methods
Besoin d'interprétabilité	Lasso ( $\ell_1$ ) pour sélection automatique

### 6.2.2 Importance du prétraitement

La standardisation des données s'est révélée *critique*:

- Réduction du nombre de condition de  $\sim 10^5$  à  $\sim 10^2$
- Accélération de la convergence d'un facteur  $\sim 100$
- Stabilité numérique améliorée

**Recommandation:** Toujours standardiser ( $\mu = 0$ ,  $\sigma = 1$ ) avant optimisation.

### 6.2.3 Trade-offs fondamentaux

Le projet illustre plusieurs trade-offs incontournables:

#### 1. Précision vs Temps:

- Batch GD: Haute précision, temps prohibitif
- SGD: Temps faible, précision limitée (bruit résiduel)
- Mini-batch/Adam: Bon compromis

#### 2. Interprétabilité vs Performance:

- $\ell_1$ : Modèle sparse (50 variables), haute interprétabilité
- $\ell_2$ : Modèle dense (1000 variables), légèrement meilleure performance

#### 3. Complexité vs Convergence:

- ISTA: Simple,  $O(1/k)$
- FISTA: Momentum,  $O(1/k^2)$  mais légèrement plus complexe



## 6.3 Perspectives

### 6.3.1 Extensions théoriques

Plusieurs directions pourraient approfondir ce travail:

1. **Optimisation non-convexe:** Étendre l'analyse aux réseaux de neurones profonds. Les algorithmes restent similaires (SGD, Adam) mais les garanties théoriques changent radicalement.
2. **Optimisation stochastique avancée:** Variance-reduced methods (SVRG, SAGA) qui obtiennent convergence linéaire pour fonctions fortement convexes.
3. **Optimisation distribuée:** Algorithmes décentralisés (AllReduce, Ring-AllReduce) pour très grands modèles.
4. **Régularisation adaptative:** Group Lasso, Fused Lasso pour structures spécifiques (groupes de variables, smoothness spatiale).

### 6.3.2 Applications

Les méthodes étudiées s'appliquent directement à:

- **NLP:** Classification de textes, détection de spam (datasets très sparses)
- **Computer Vision:** Entraînement de réseaux convolutionnels (datasets massifs)
- **Génomique:** Sélection de gènes ( $d \sim 10^4$  à  $10^5$ )
- **Recommandation:** Factorisation matricielle pour systèmes de recommandation
- **Finance:** Sélection de portefeuille avec contraintes de sparsité

## 6.4 Réflexion finale

Ce projet démontre l'importance d'une approche *théorie-pratique intégrée* en optimisation:

1. La **théorie** fournit:
  - Des garanties de convergence
  - Des taux asymptotiques
  - Une compréhension profonde des mécanismes
2. La **pratique** révèle:
  - L'importance des constantes (pas seulement la notation  $O(\cdot)$ )
  - Les difficultés d'implémentation (stabilité numérique, choix de hyperparamètres)
  - Les compromis réels entre algorithmes

**L’optimisation n’est pas qu’un outil technique** - c’est le pont entre la modélisation mathématique et les applications concrètes du machine learning. La maîtrise de ses principes fondamentaux (comme explorés dans ce projet) est essentielle pour développer des systèmes d’apprentissage efficaces, robustes et interprétables.

---

*“In theory, there is no difference between theory and practice.*

*In practice, there is.”*

— Yogi Berra (attribué)

---

Les résultats de ce projet confirment cette maxime: la théorie prédit correctement les tendances (taux de convergence, réduction de variance), mais la pratique révèle les constantes, les seuils et les subtilités qui font la différence entre un algorithme fonctionnel et un algorithme optimal.