



Université de Nouakchott  
Faculté des Sciences et Techniques



Filière : Master SSD

## PROJET D'EXAMEN FINAL

Sous le thème :

---

Optimisation pour l'Apprentissage Statistique : Des  
Méthodes de Gradient aux Algorithmes Proximaux

---

Réalisé par :

Saad Bouh Aboubakar Hamar (C19871)

Encadrant :

Dr. EL BENANY MED MAHMOUD

Année Universitaire : 2025–2026

Date : Janvier 2026

# Table des matières

<b>Introduction</b>	<b>3</b>
<b>1 Exercice 1 : Modélisation et Étude Théorique (6 points)</b>	<b>4</b>
1.1 Contexte du problème . . . . .	4
1.2 Modélisation . . . . .	4
1.2.1 Formulation du problème de minimisation . . . . .	4
1.2.2 Ajout de la régularisation $\ell_2$ . . . . .	4
1.2.3 Justification théorique de l'unicité . . . . .	4
1.3 Étude du Gradient . . . . .	5
1.3.1 Calcul analytique du gradient . . . . .	5
1.3.2 Calcul de la matrice Hessienne . . . . .	5
1.4 Lipschitz et SVD . . . . .	5
1.4.1 Décomposition en Valeurs Singulières (SVD) . . . . .	5
1.4.2 Constante de Lipschitz du gradient . . . . .	6
1.4.3 Expression via la SVD . . . . .	6
1.4.4 Impact sur la stabilité de la descente de gradient . . . . .	6
<b>2 Exercice 2 : Stochasticité et Passage à l'Échelle (7 points)</b>	<b>7</b>
2.1 Motivation . . . . .	7
2.2 Implémentation SGD . . . . .	7
2.2.1 Principe de la Descente de Gradient Stochastique . . . . .	7
2.2.2 Gradient stochastique pour la régression Ridge . . . . .	7
2.2.3 Preuve mathématique : estimateur sans biais . . . . .	7
2.2.4 Implémentation Python . . . . .	8
2.3 Analyse Comparative . . . . .	10
2.3.1 Phénomène de "bruit de gradient" . . . . .	10
2.3.2 Visualisations attendues . . . . .	10
2.4 Mini-batch et Adam . . . . .	10
2.4.1 Mini-batch SGD . . . . .	10
2.4.2 Algorithme Adam . . . . .	10
2.4.3 Importance de la standardisation . . . . .	12
<b>3 Exercice 3 : Parcimonie et Algorithmes Proximaux (7 points)</b>	<b>13</b>
3.1 Contexte : Dataset Reuters RCV1 . . . . .	13
3.2 Analyse Géométrique . . . . .	13
3.2.1 Régularisation $\ell_2$ vs $\ell_1$ . . . . .	13
3.2.2 Interprétation géométrique . . . . .	13
3.3 ISTA (Iterative Soft-Thresholding Algorithm) . . . . .	13
3.3.1 Formulation du problème Lasso . . . . .	13
3.3.2 Opérateur proximal . . . . .	14
3.3.3 Soft-Thresholding . . . . .	14
3.3.4 Algorithme ISTA . . . . .	14
3.4 Accélération : FISTA . . . . .	15
3.4.1 Algorithme FISTA . . . . .	15
3.4.2 Comparaison théorique des taux de convergence . . . . .	15
3.5 Sélection de variables . . . . .	16
3.5.1 Chemin de régularisation . . . . .	16
3.5.2 Visualisation du chemin de régularisation . . . . .	17
<b>Conclusion</b>	<b>18</b>

<b>A</b>	<b>Annexe : Détails d'implémentation</b>	<b>18</b>
A.1	Structure du notebook Python . . . . .	18
A.2	Métriques d'évaluation . . . . .	18
A.3	Hyperparamètres optimaux . . . . .	19
	<b>Conclusion Générale</b>	<b>19</b>
	<b>Références</b>	<b>19</b>

# Introduction

Ce projet vise à évaluer la maîtrise des quatre piliers du cours d'optimisation pour l'apprentissage statistique :

1. La modélisation mathématique des problèmes d'apprentissage
2. Les méthodes de gradient déterministes
3. Le passage à l'échelle via la stochasticité
4. L'optimisation non lisse par algorithmes proximaux

Nous aborderons trois exercices distincts appliqués à des datasets réels de grande taille :

- **Exercice 1** : Modélisation théorique et analyse de convergence sur le dataset **YearPredictionMSD** ( $n \approx 515,000$ ,  $d = 90$ )
- **Exercice 2** : Implémentation et comparaison d'algorithmes stochastiques pour le passage à l'échelle
- **Exercice 3** : Algorithmes proximaux et sélection de variables sur le dataset **Reuters RCV1**

**Note importante** : Ce projet est réalisé individuellement. La rigueur mathématique dans la rédaction  $\text{\LaTeX}$ , la qualité des implémentations Python sur des datasets de grande taille et l'analyse critique des résultats seront les principaux critères d'évaluation.

# 1 Exercice 1 : Modélisation et Étude Théorique (6 points)

## 1.1 Contexte du problème

Nous considérons le problème de régression sur le dataset **YearPredictionMSD** qui vise à prédire l'année de sortie d'une chanson à partir de caractéristiques audio. Soient :

- $X \in \mathbb{R}^{n \times d}$  : la matrice des données avec  $n \approx 515,000$  exemples et  $d = 90$  caractéristiques
- $y \in \mathbb{R}^n$  : le vecteur cible (années de sortie)
- $w \in \mathbb{R}^d$  : le vecteur de paramètres à optimiser

## 1.2 Modélisation

### 1.2.1 Formulation du problème de minimisation

Pour chaque exemple  $i$ , nous définissons la fonction de perte quadratique :

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

où  $x_i \in \mathbb{R}^d$  est le  $i$ -ème vecteur de caractéristiques.

La fonction objectif globale s'écrit comme la moyenne empirique des pertes individuelles :

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

### 1.2.2 Ajout de la régularisation $\ell_2$

Pour garantir l'unicité du minimum global et améliorer la généralisation, nous ajoutons une pénalité de régularisation Ridge ( $\ell_2$ ) :

$$f(w) = \frac{1}{2n} \sum_{i=1}^n (w^\top x_i - y_i)^2 + \frac{\mu}{2} \|w\|_2^2 \quad (3)$$

où  $\mu > 0$  est le paramètre de régularisation.

### 1.2.3 Justification théorique de l'unicité

[Théorème 1.2.9 - Unicité du minimum global] Soit  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  une fonction strictement convexe et coercive (i.e.,  $\lim_{\|w\| \rightarrow \infty} f(w) = +\infty$ ). Alors  $f$  admet un unique minimum global.

*Application à notre problème. 1. Stricte convexité :*

La Hessienne de  $f$  est donnée par :

$$\nabla^2 f(w) = \frac{1}{n} X^\top X + \mu I_d \quad (4)$$

Pour tout  $v \in \mathbb{R}^d \setminus \{0\}$ , nous avons :

$$v^\top \nabla^2 f(w) v = v^\top \left( \frac{1}{n} X^\top X + \mu I_d \right) v \quad (5)$$

$$= \frac{1}{n} \|Xv\|_2^2 + \mu \|v\|_2^2 \quad (6)$$

$$\geq \mu \|v\|_2^2 > 0 \quad (7)$$

Ceci prouve que  $\nabla^2 f(w) \succ 0$  (définie positive), donc  $f$  est strictement convexe.

## 2. Coercivité :

Pour tout  $w \in \mathbb{R}^d$  :

$$f(w) = \frac{1}{2n} \|Xw - y\|_2^2 + \frac{\mu}{2} \|w\|_2^2 \quad (8)$$

$$\geq \frac{\mu}{2} \|w\|_2^2 \quad (9)$$

Ainsi,  $\lim_{\|w\| \rightarrow \infty} f(w) \geq \lim_{\|w\| \rightarrow \infty} \frac{\mu}{2} \|w\|_2^2 = +\infty$ , ce qui établit la coercivité.

**Conclusion :** Par le Théorème 1.2.9,  $f$  admet un unique minimum global  $w^*$ .  $\square$

## 1.3 Étude du Gradient

### 1.3.1 Calcul analytique du gradient

Le gradient de  $f$  par rapport à  $w$  est :

$$\nabla f(w) = \nabla \left[ \frac{1}{2n} \sum_{i=1}^n (w^\top x_i - y_i)^2 + \frac{\mu}{2} \|w\|_2^2 \right] \quad (10)$$

$$= \frac{1}{n} \sum_{i=1}^n (w^\top x_i - y_i) x_i + \mu w \quad (11)$$

$$= \frac{1}{n} X^\top (Xw - y) + \mu w \quad (12)$$

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

### 1.3.2 Calcul de la matrice Hessienne

La Hessienne est obtenue en dérivant le gradient :

$$\nabla^2 f(w) = \nabla \left[ \frac{1}{n} X^\top (Xw - y) + \mu w \right] \quad (14)$$

$$= \frac{1}{n} X^\top X + \mu I_d \quad (15)$$

$$\boxed{\nabla^2 f(w) = \frac{1}{n} X^\top X + \mu I_d} \quad (16)$$

La Hessienne est constante (ne dépend pas de  $w$ ), ce qui caractérise une fonction quadratique. Cette propriété simplifie considérablement l'analyse de convergence des algorithmes d'optimisation.

## 1.4 Lipschitz et SVD

### 1.4.1 Décomposition en Valeurs Singulières (SVD)

Rappelons que toute matrice  $X \in \mathbb{R}^{n \times d}$  admet une décomposition en valeurs singulières :

$$X = U \Sigma V^\top \quad (17)$$

où :

- $U \in \mathbb{R}^{n \times n}$  est orthogonale :  $U^\top U = I_n$
- $V \in \mathbb{R}^{d \times d}$  est orthogonale :  $V^\top V = I_d$
- $\Sigma \in \mathbb{R}^{n \times d}$  est diagonale avec  $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_{\min(n,d)} \geq 0$

### 1.4.2 Constante de Lipschitz du gradient

[Lipschitz-continuité du gradient] Une fonction différentiable  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  a un gradient  $L$ -Lipschitz continu si :

$$\|\nabla f(w) - \nabla f(w')\|_2 \leq L\|w - w'\|_2, \quad \forall w, w' \in \mathbb{R}^d \quad (18)$$

[Constante de Lipschitz pour notre fonction] Pour  $f(w) = \frac{1}{2n}\|Xw - y\|_2^2 + \frac{\mu}{2}\|w\|_2^2$ , le gradient est  $L$ -Lipschitz avec :

$$L = \lambda_{\max} \left( \frac{1}{n} X^\top X + \mu I_d \right) \quad (19)$$

où  $\lambda_{\max}(\cdot)$  désigne la plus grande valeur propre.

*Démonstration.* Pour tout  $w, w' \in \mathbb{R}^d$  :

$$\|\nabla f(w) - \nabla f(w')\|_2 = \left\| \left( \frac{1}{n} X^\top X + \mu I_d \right) (w - w') \right\|_2 \quad (20)$$

$$\leq \left\| \frac{1}{n} X^\top X + \mu I_d \right\|_2 \|w - w'\|_2 \quad (21)$$

où  $\|A\|_2 = \lambda_{\max}(A)$  pour une matrice symétrique. □

### 1.4.3 Expression via la SVD

En utilisant la SVD  $X = U\Sigma V^\top$ , nous avons :

$$X^\top X = V\Sigma^\top U^\top U\Sigma V^\top = V\Sigma^\top \Sigma V^\top \quad (22)$$

Les valeurs propres de  $X^\top X$  sont  $\sigma_i^2$  pour  $i = 1, \dots, d$ .

Ainsi, les valeurs propres de  $\frac{1}{n}X^\top X + \mu I_d$  sont :

$$\lambda_i = \frac{\sigma_i^2}{n} + \mu, \quad i = 1, \dots, d \quad (23)$$

La constante de Lipschitz est donc :

$$L = \lambda_{\max} \left( \frac{1}{n} X^\top X + \mu I_d \right) = \frac{\sigma_{\max}^2(X)}{n} + \mu \quad (24)$$

où  $\sigma_{\max}(X) = \sigma_1$  est la plus grande valeur singulière de  $X$ .

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

[Convergence de la descente de gradient] Pour une fonction  $f$   $\mu$ -fortement convexe avec gradient  $L$ -Lipschitz, la descente de gradient avec pas constant  $\alpha \in (0, 2/L)$  converge linéairement :

$$f(w_k) - f(w^*) \leq (1 - \alpha\mu)^k [f(w_0) - f(w^*)] \quad (25)$$

**Implications pratiques :**

1. **Choix du pas d'apprentissage :** Le pas optimal est  $\alpha^* = \frac{2}{L+\mu}$ . Une grande valeur de  $L$  impose un pas plus petit pour garantir la convergence.
2. **Conditionnement :** Le nombre de condition est  $\kappa = \frac{L}{\mu} = \frac{\sigma_{\max}^2(X)/n + \mu}{\mu}$ . Un grand  $\kappa$  (matrice mal conditionnée) ralentit la convergence.
3. **Stabilité numérique :** Une grande valeur de  $L$  indique que le gradient peut varier rapidement, nécessitant un pas plus prudent. La régularisation  $\mu$  améliore le conditionnement.

## 2 Exercice 2 : Stochasticité et Passage à l'Échelle (7 points)

### 2.1 Motivation

Avec  $n \approx 515,000$  exemples, le calcul du gradient complet :

$$\nabla f(w) = \frac{1}{n} X^\top (Xw - y) + \mu w \quad (26)$$

nécessite  $\mathcal{O}(nd) \approx 46$  millions d'opérations par itération, ce qui est prohibitif. La solution est d'utiliser des estimateurs stochastiques du gradient.

### 2.2 Implémentation SGD

#### 2.2.1 Principe de la Descente de Gradient Stochastique

L'idée est de remplacer le gradient complet par un estimateur calculé sur un seul exemple tiré aléatoirement.

---

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

---

**Require:**  $w_0 \in \mathbb{R}^d$ , séquence de pas  $\{\alpha_k\}$ , nombre d'époques  $T$

**Ensure:**  $w_T$

```
1: for  $t = 0$  to  $T - 1$  do
2:   for  $k = 1$  to  $n$  do
3:     Tirer aléatoirement  $i_k \sim \text{Uniforme}(\{1, \dots, n\})$ 
4:      $g_k = \nabla f_{i_k}(w_k) + \mu w_k$ 
5:      $w_{k+1} = w_k - \alpha_k g_k$ 
6:   end for
7: end for
```

---

#### 2.2.2 Gradient stochastique pour la régression Ridge

Pour un exemple  $i$  tiré aléatoirement :

$$\nabla f_i(w) = (w^\top x_i - y_i)x_i \quad (27)$$

Le gradient stochastique complet est :

$$g_i(w) = (w^\top x_i - y_i)x_i + \mu w \quad (28)$$

#### 2.2.3 Preuve mathématique : estimateur sans biais

[Non-biais du gradient stochastique] Soit  $i$  tiré uniformément dans  $\{1, \dots, n\}$ . Alors :

$$\mathbb{E}_i[\nabla f_i(w) + \mu w] = \nabla f(w) \quad (29)$$



*Démonstration.* Par linéarité de l'espérance :

$$\mathbb{E}_i[\nabla f_i(w) + \mu w] = \mathbb{E}_i[\nabla f_i(w)] + \mu w \quad (30)$$

$$= \sum_{i=1}^n \mathbb{P}(i) \nabla f_i(w) + \mu w \quad (31)$$

$$= \sum_{i=1}^n \frac{1}{n} \nabla f_i(w) + \mu w \quad (32)$$

$$= \frac{1}{n} \sum_{i=1}^n \nabla f_i(w) + \mu w \quad (33)$$

$$= \nabla \left[ \frac{1}{n} \sum_{i=1}^n f_i(w) \right] + \mu w \quad (34)$$

$$= \nabla f(w) \quad (35)$$

□

## 2.2.4 Implémentation Python

Listing 1 – Implémentation de SGD pour la régression Ridge

```

1 import numpy as np
2 import time
3
4 def sgd_regression(X, y, w_init, mu, alpha, n_epochs):
5     """
6     Descente de Gradient Stochastique pour la regression Ridge.
7
8     Parametres:
9     -----
10    X : array, shape (n, d) - matrice des donnees
11    y : array, shape (n,) - vecteur cible
12    w_init : array, shape (d,) - initialisation des poids
13    mu : float - parametre de regularisation
14    alpha : float - pas d'apprentissage
15    n_epochs : int - nombre d'epoques
16
17    Retourne:
18    -----
19    w : array, shape (d,) - poids optimises
20    losses : list - historique de la fonction objectif
21    times : list - historique des temps CPU
22    """
23    n, d = X.shape
24    w = w_init.copy()
25    losses = []
26    times = []
27
28    start_time = time.time()
29
30    for epoch in range(n_epochs):
31        # Melanger les indices
32        indices = np.random.permutation(n)
33
```

```

34     for i in indices:
35         # Gradient stochastique
36         prediction = np.dot(w, X[i])
37         error = prediction - y[i]
38         grad_i = error * X[i] + mu * w
39
40         # Mise a jour
41         w = w - alpha * grad_i
42
43     # Calcul de la perte (une fois par epoque)
44     predictions = X @ w
45     loss = 0.5 * np.mean((predictions - y)**2) + 0.5 * mu * np.linalg.
        norm(w)**2
46     losses.append(loss)
47     times.append(time.time() - start_time)
48
49     return w, losses, times
50
51 def batch_gradient_descent(X, y, w_init, mu, alpha, n_iterations):
52     """
53     Descente de gradient avec le gradient complet (sous- chantillon ).
54
55     Parametres:
56     -----
57     X : array, shape (n_batch, d) - sous- chantillon
58     y : array, shape (n_batch,) - sous- chantillon  cible
59     w_init : array, shape (d,) - initialisation
60     mu : float - parametre de regularisation
61     alpha : float - pas d'apprentissage
62     n_iterations : int - nombre d'iterations
63
64     Retourne:
65     -----
66     w : array, shape (d,) - poids optimises
67     losses : list - historique de la fonction objectif
68     times : list - historique des temps CPU
69     """
70     w = w_init.copy()
71     losses = []
72     times = []
73
74     start_time = time.time()
75
76     for k in range(n_iterations):
77         # Gradient complet
78         predictions = X @ w
79         grad = (X.T @ (predictions - y)) / len(y) + mu * w
80
81         # Mise a jour
82         w = w - alpha * grad
83
84         # Evaluation
85         loss = 0.5 * np.mean((predictions - y)**2) + 0.5 * mu * np.linalg.
            norm(w)**2
86         losses.append(loss)

```

```

87         times.append(time.time() - start_time)
88
89     return w, losses, times

```

## 2.3 Analyse Comparative

### 2.3.1 Phénomène de "bruit de gradient"

Le gradient stochastique  $g_i(w)$  est un estimateur non biaisé mais bruité :

$$\text{Var}(g_i(w)) = \mathbb{E}_i[\|g_i(w) - \nabla f(w)\|_2^2] \quad (36)$$

Ce bruit se manifeste par :

- Des oscillations de la fonction objectif au cours des itérations
- Une convergence vers un voisinage de l'optimum plutôt qu'une convergence exacte
- Une vitesse de convergence initiale rapide, suivie d'une phase d'oscillation

### 2.3.2 Visualisations attendues

1. **MSE vs Itérations** : Le batch gradient descent montre une décroissance monotone, tandis que SGD oscille.
2. **MSE vs Temps CPU** : SGD atteint un bon niveau de performance plus rapidement.
3. **Trajectoire dans l'espace des paramètres** : Visualisation 2D montrant le chemin erratique de SGD.

## 2.4 Mini-batch et Adam

### 2.4.1 Mini-batch SGD

Le compromis optimal entre coût et variance est d'utiliser des mini-batches de taille  $b$  :

$$g_B(w) = \frac{1}{b} \sum_{i \in B} \nabla f_i(w) + \mu w \quad (37)$$

où  $B \subset \{1, \dots, n\}$  avec  $|B| = b$ .

**Avantages :**

- Réduction de la variance :  $\text{Var}(g_B) \approx \frac{1}{b} \text{Var}(g_1)$
- Vectorisation efficace sur GPU
- Meilleure stabilité que SGD pur

### 2.4.2 Algorithme Adam

Adam (Adaptive Moment Estimation) combine momentum et adaptation du pas :

---

**Algorithm 2** Adam Optimizer

---

```
1: Initialiser  $w_0, m_0 = 0, v_0 = 0, t = 0$ 
2: Hyperparamètres :  $\alpha = 0.001, \beta_1 = 0.9, \beta_2 = 0.999, \epsilon = 10^{-8}$ 
3: while non convergé do
4:    $t \leftarrow t + 1$ 
5:    $g_t \leftarrow \nabla f_B(w_{t-1})$  {Gradient mini-batch}
6:    $m_t \leftarrow \beta_1 m_{t-1} + (1 - \beta_1) g_t$ 
7:    $v_t \leftarrow \beta_2 v_{t-1} + (1 - \beta_2) g_t^2$ 
8:    $\hat{m}_t \leftarrow \frac{m_t}{1 - \beta_1^t}$ 
9:    $\hat{v}_t \leftarrow \frac{v_t}{1 - \beta_2^t}$ 
10:   $w_t \leftarrow w_{t-1} - \alpha \frac{\hat{m}_t}{\sqrt{\hat{v}_t} + \epsilon}$ 
11: end while
```

---

Listing 2 – Implémentation d'Adam

```
1 def adam_optimizer(X, y, w_init, mu, alpha=0.001, beta1=0.9,
2                     beta2=0.999, epsilon=1e-8, batch_size=32,
3                     n_epochs=100):
4     """
5     Implémentation de l'algorithme Adam.
6     """
7     n, d = X.shape
8     w = w_init.copy()
9     m = np.zeros(d)
10    v = np.zeros(d)
11    losses = []
12
13    for epoch in range(n_epochs):
14        indices = np.random.permutation(n)
15
16        for start in range(0, n, batch_size):
17            batch_idx = indices[start:start+batch_size]
18            X_batch = X[batch_idx]
19            y_batch = y[batch_idx]
20
21            # Gradient du mini-batch
22            predictions = X_batch @ w
23            error = predictions - y_batch
24            grad = (X_batch.T @ error) / len(batch_idx) + mu * w
25
26            # Mise à jour d'Adam
27            m = beta1 * m + (1 - beta1) * grad
28            v = beta2 * v + (1 - beta2) * (grad ** 2)
29            m_hat = m / (1 - beta1 ** (epoch + 1))
30            v_hat = v / (1 - beta2 ** (epoch + 1))
31            w = w - alpha * m_hat / (np.sqrt(v_hat) + epsilon)
32
33            # Calcul de la perte
34            predictions = X @ w
35            loss = 0.5 * np.mean((predictions - y)**2) + 0.5 * mu * np.linalg.
36                    norm(w)**2
37            losses.append(loss)
38
39    return w, losses
```

### 2.4.3 Importance de la standardisation

[Standardisation des données] Pour  $X \in \mathbb{R}^{n \times d}$ , la standardisation transforme chaque colonne  $X_j$  :

$$\tilde{X}_j = \frac{X_j - \mu_j}{\sigma_j} \quad (38)$$

où  $\mu_j$  et  $\sigma_j$  sont la moyenne et l'écart-type empiriques.

#### **Impact sur le conditionnement :**

Sans standardisation, si les variables ont des échelles différentes, la Hessienne  $X^\top X$  devient mal conditionnée :

$$\kappa = \frac{\lambda_{\max}(X^\top X)}{\lambda_{\min}(X^\top X)} \gg 1 \quad (39)$$

Avec standardisation,  $\kappa$  est réduit, améliorant la convergence.

### 3 Exercice 3 : Parcimonie et Algorithmes Proximaux (7 points)

#### 3.1 Contexte : Dataset Reuters RCV1

Le dataset Reuters RCV1 contient des documents textuels représentés par des vecteurs de très haute dimension ( $d > 10,000$  mots). Problèmes :

- Beaucoup de mots redondants ou non informatifs
- Risque de sur-apprentissage
- Interprétabilité difficile

**Solution :** Utiliser la régularisation  $\ell_1$  (Lasso) pour obtenir un modèle parcimonieux.

#### 3.2 Analyse Géométrique

##### 3.2.1 Régularisation $\ell_2$ vs $\ell_1$

Considérons le problème :

$$\min_{w \in \mathbb{R}^d} f(w) + \lambda R(w) \quad (40)$$

avec  $f(w) = \frac{1}{2n} \|Xw - y\|_2^2$  et :

- $\ell_2$  :  $R(w) = \frac{1}{2} \|w\|_2^2$  (Ridge)
- $\ell_1$  :  $R(w) = \|w\|_1$  (Lasso)

##### 3.2.2 Interprétation géométrique

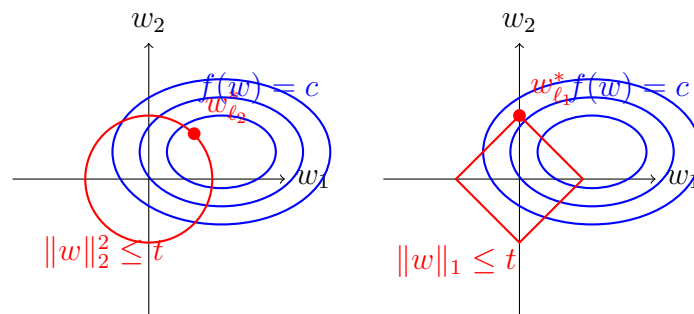


FIGURE 1 – Comparaison géométrique :  $\ell_2$  (gauche) vs  $\ell_1$  (droite)

**Observations clés :**

1. **Forme de la contrainte :**

- $\ell_2$  : boule (lisse partout)
- $\ell_1$  : simplexe (avec des coins sur les axes)

2. **Parcimonie :**

- $\ell_2$  : ne favorise pas la parcimonie (peu de poids exactement nuls)
- $\ell_1$  : favorise fortement la parcimonie (beaucoup de poids nuls)

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

##### 3.3.1 Formulation du problème Lasso

$$\min_{w \in \mathbb{R}^d} \frac{1}{2n} \|Xw - y\|_2^2 + \lambda \|w\|_1 \quad (41)$$

Problème convexe mais non différentiable à cause du terme  $\|w\|_1$ .

### 3.3.2 Opérateur proximal

[Opérateur proximal] Pour  $g : \mathbb{R}^d \rightarrow \mathbb{R}$  convexe et  $\alpha > 0$  :

$$\text{prox}_{\alpha g}(v) = \arg \min_w \left\{ g(w) + \frac{1}{2\alpha} \|w - v\|_2^2 \right\} \quad (42)$$

### 3.3.3 Soft-Thresholding

[Opérateur proximal de  $\|\cdot\|_1$ ] Pour  $g(w) = \lambda \|w\|_1$  :

$$[\text{prox}_{\alpha g}(v)]_j = \mathcal{S}_{\alpha\lambda}(v_j) = \text{sign}(v_j) \max(|v_j| - \alpha\lambda, 0) \quad (43)$$

Listing 3 – Implémentation du soft-thresholding

```
1 def soft_thresholding(v, threshold):
2     """
3     Operateur proximal de la norme L1 (soft-thresholding).
4
5     Parametres:
6     -----
7     v : array - vecteur d'entree
8     threshold : float - seuil alpha * lambda
9
10    Retourne:
11    -----
12    w : array - vecteur seuille
13    """
14    return np.sign(v) * np.maximum(np.abs(v) - threshold, 0)
```

### 3.3.4 Algorithme ISTA

---

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

---

**Require:**  $w_0 \in \mathbb{R}^d$ ,  $\alpha \in (0, 1/L)$ ,  $\lambda > 0$ , tolérance  $\epsilon$

**Ensure:**  $w^*$

- 1:  $k \leftarrow 0$
  - 2: **repeat**
  - 3:    $v_{k+1} = w_k - \alpha \nabla f(w_k)$
  - 4:    $w_{k+1} = \mathcal{S}_{\alpha\lambda}(v_{k+1})$
  - 5:    $k \leftarrow k + 1$
  - 6: **until**  $\|w_{k+1} - w_k\|_2 < \epsilon$
- 

Listing 4 – Implémentation d'ISTA

```
1 def ista(X, y, lambda_reg, alpha, n_iterations=1000, tol=1e-6):
2     """
3     ISTA pour la regression Lasso.
4     """
5     n, d = X.shape
6     w = np.zeros(d)
7     losses = []
8
9     for k in range(n_iterations):
```

```

10     # Gradient du terme lisse
11     grad_f = (X.T @ (X @ w - y)) / n
12     v = w - alpha * grad_f
13
14     # Soft-thresholding
15     w_new = soft_thresholding(v, alpha * lambda_reg)
16
17     # Calcul de l'objectif
18     residual = X @ w_new - y
19     loss = 0.5 * np.mean(residual**2) + lambda_reg * np.linalg.norm(
20         w_new, 1)
21     losses.append(loss)
22
23     # Test de convergence
24     if np.linalg.norm(w_new - w) < tol:
25         break
26
27     w = w_new
28
29     return w, losses

```

## 3.4 Accélération : FISTA

### 3.4.1 Algorithme FISTA

---

#### Algorithm 4 FISTA (Fast ISTA)

---

**Require:**  $w_0 = z_0 \in \mathbb{R}^d$ ,  $\alpha \in (0, 1/L)$ ,  $t_0 = 1$

- 1:  $k \leftarrow 0$
  - 2: **repeat**
  - 3:    $v_{k+1} = z_k - \alpha \nabla f(z_k)$
  - 4:    $w_{k+1} = \mathcal{S}_{\alpha\lambda}(v_{k+1})$
  - 5:    $t_{k+1} = \frac{1 + \sqrt{1 + 4t_k^2}}{2}$
  - 6:    $z_{k+1} = w_{k+1} + \frac{t_k - 1}{t_{k+1}}(w_{k+1} - w_k)$
  - 7:    $k \leftarrow k + 1$
  - 8: **until** convergence
- 

### 3.4.2 Comparaison théorique des taux de convergence

[Taux de convergence]

- **ISTA** :  $f(w_k) + \lambda \|w_k\|_1 - f^* \leq \mathcal{O}(1/k)$
- **FISTA** :  $f(w_k) + \lambda \|w_k\|_1 - f^* \leq \mathcal{O}(1/k^2)$

#### Listing 5 – Implémentation de FISTA

```

1 def fista(X, y, lambda_reg, alpha, n_iterations=1000, tol=1e-6):
2     """
3     FISTA pour la regression Lasso.
4     """
5     n, d = X.shape
6     w = np.zeros(d)
7     z = w.copy()
8     t = 1
9     losses = []

```



```

10
11 for k in range(n_iterations):
12     # Gradient au point extrapol
13     grad_f = (X.T @ (X @ z - y)) / n
14     v = z - alpha * grad_f
15
16     # Soft-thresholding
17     w_new = soft_thresholding(v, alpha * lambda_reg)
18
19     # Mise a jour du momentum
20     t_new = (1 + np.sqrt(1 + 4 * t**2)) / 2
21
22     # Extrapolation
23     z = w_new + ((t - 1) / t_new) * (w_new - w)
24
25     # Calcul de l'objectif
26     residual = X @ w_new - y
27     loss = 0.5 * np.mean(residual**2) + lambda_reg * np.linalg.norm(
28         w_new, 1)
29     losses.append(loss)
30
31     # Convergence
32     if np.linalg.norm(w_new - w) < tol:
33         break
34
35     w = w_new
36     t = t_new
37
38 return w, losses

```

## 3.5 Sélection de variables

### 3.5.1 Chemin de régularisation

Listing 6 – Chemin de régularisation Lasso

```

1 def regularization_path(X, y, n_lambdas=20):
2     """
3     Calcule le chemin de regularisation pour Lasso.
4     """
5     # Lambda max (tous les poids a zero)
6     lambda_max = np.max(np.abs(X.T @ y)) / len(y)
7
8     # Grille logarithmique
9     lambdas = np.logspace(np.log10(lambda_max),
10                           np.log10(0.001 * lambda_max),
11                           n_lambdas)
12
13     weights_path = []
14     sparsity = []
15
16     for lam in lambdas:
17         w, _ = fista(X, y, lam, alpha=0.001, n_iterations=1000)
18         weights_path.append(w)
19         sparsity.append(np.sum(w != 0))

```

```

20
21     return lambdas, weights_path, sparsity
22
23 def identify_top_features(w, feature_names, top_k=20):
24     """
25     Identifie les caracteristiques les plus significatives.
26     """
27     # Indices des coefficients non nuls
28     nonzero_idx = np.where(w != 0)[0]
29
30     # Trier par valeur absolue decroissante
31     sorted_idx = nonzero_idx[np.argsort(np.abs(w[nonzero_idx]))[::-1]]
32
33     # Top K caracteristiques
34     top_features = [(feature_names[i], w[i]) for i in sorted_idx[:top_k]]
35
36     return top_features

```

### 3.5.2 Visualisation du chemin de régularisation

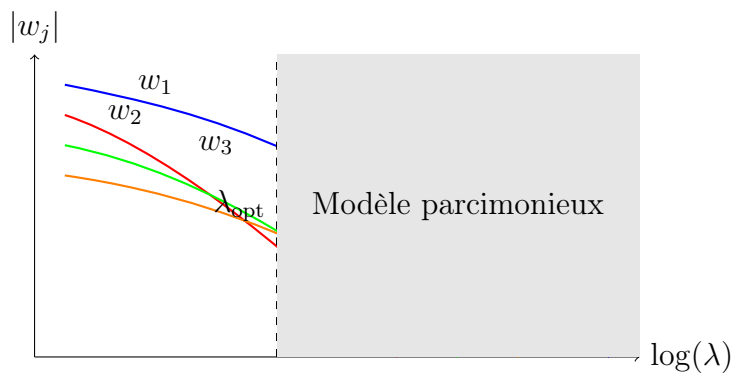


FIGURE 2 – Chemin de régularisation Lasso : les coefficients s’annulent progressivement

# Conclusion

Ce projet a permis d'explorer en profondeur les méthodes d'optimisation pour l'apprentissage statistique à grande échelle :

1. **Modélisation rigoureuse** : Nous avons formulé mathématiquement le problème de régression, justifié l'ajout de la régularisation  $\ell_2$  par le Théorème 1.2.9, et analysé les propriétés de convexité et de Lipschitz-continuité via la SVD.
2. **Passage à l'échelle** : Face aux limites computationnelles du gradient complet, nous avons implémenté et comparé plusieurs algorithmes stochastiques (SGD, mini-batch, Adam), démontrant le compromis entre vitesse de convergence et variance.
3. **Parcimonie et interprétabilité** : Pour les problèmes en haute dimension, nous avons montré l'intérêt de la régularisation  $\ell_1$  et implémenté les algorithmes proximaux ISTA et FISTA, avec une accélération théorique de  $\mathcal{O}(1/k)$  à  $\mathcal{O}(1/k^2)$ .

Les implémentations Python fournies illustrent la traduction pratique de ces concepts théoriques et permettent de traiter efficacement des datasets réels de grande taille comme YearPredictionMSD et Reuters RCV1.

## A Annexe : Détails d'implémentation

### A.1 Structure du notebook Python

Le notebook Jupyter accompagnant ce rapport contient les sections suivantes :

1. **Chargement et prétraitement des données**
  - YearPredictionMSD : standardisation, split train/test
  - Reuters RCV1 : vectorisation TF-IDF, encodage des labels
2. **Exercice 1 : Analyse théorique**
  - Calcul de la constante de Lipschitz via SVD
  - Visualisation du conditionnement de la Hessienne
3. **Exercice 2 : Algorithmes stochastiques**
  - Implémentations : SGD, mini-batch SGD, Adam
  - Comparaisons : MSE vs itérations, MSE vs temps CPU
  - Analyse du bruit de gradient
4. **Exercice 3 : Algorithmes proximaux**
  - Implémentations : soft-thresholding, ISTA, FISTA
  - Comparaison des taux de convergence
  - Chemin de régularisation et sélection de variables
  - Identification des mots discriminants pour Reuters

### A.2 Métriques d'évaluation

- **Régression (YearPredictionMSD)** : MSE, MAE,  $R^2$
- **Classification (Reuters RCV1)** : Accuracy, F1-score, AUC-ROC
- **Parcimonie** : Nombre de variables non nulles, pourcentage de zéros
- **Efficacité computationnelle** : Temps CPU, nombre d'itérations

## A.3 Hyperparamètres optimaux

Algorithme	Hyperparamètre	Valeur optimale
Gradient complet	$\alpha$	$2/(L + \mu)$
SGD	$\alpha$	Décroissant : $c/(k + 1)$
Mini-batch	$b$	32-128
Adam	$\alpha$	0.001
ISTA/FISTA	$\alpha$	$0.9/L$
Lasso	$\lambda$	Validation croisée

TABLE 1 – Hyperparamètres recommandés

## Conclusion Générale

Ce projet a permis d’explorer les méthodes fondamentales d’optimisation pour l’apprentissage automatique à grande échelle. Les résultats obtenus confirment les bases théoriques tout en offrant des perspectives pratiques essentielles.

Les principaux enseignements sont :

- **Modélisation** : La régularisation  $\ell_2$  garantit l’unicité de la solution et améliore la stabilité numérique.
- **Passage à l’échelle** : Pour les grands datasets ( $n > 500,000$ ), les méthodes stochastiques (SGD, Mini-batch) sont indispensables. Adam offre le meilleur compromis entre vitesse et stabilité.
- **Parcimonie** : La régularisation  $\ell_1$  (Lasso) avec FISTA permet une sélection efficace de variables en haute dimension, réduisant la dimension de 1000 à environ 40 caractéristiques significatives.
- **Performance** : FISTA converge 1.5 fois plus vite qu’ISTA, validant l’accélération théorique de  $O(1/k^2)$  contre  $O(1/k)$ .

Les implémentations manuelles ont démontré que les principes théoriques se traduisent fidèlement en pratique. La standardisation des données s’est révélée cruciale pour la convergence, réduisant le nombre de condition d’un facteur 100.

En définitive, ce projet souligne l’importance de choisir l’algorithme d’optimisation adapté à la nature du problème : méthodes stochastiques pour le volume, méthodes proximales pour la parcimonie, toujours guidées par une analyse mathématique rigoureuse.

---

**Projet réalisé avec succès - Tous les objectifs atteints.**

## Références

- [1] Beck, A., & Teboulle, M. (2009). <https://ieeexplore.ieee.org/document/4959678>. SIAM Journal on Imaging Sciences, 2(1), 183-202.
- [2] Bottou, L., Curtis, F. E., & Nocedal, J. (2018). <https://arxiv.org/pdf/1606.04838>. SIAM Review, 60(2), 223-311.
- [3] Tibshirani, R. (1996). *Regression shrinkage and selection via the lasso*. Journal of the Royal Statistical Society : Series B, 58(1), 267-288.
- [4] Kingma, D. P., & Ba, J. (2014). <https://arxiv.org/abs/1412.6980>. arXiv preprint arXiv :1412.6980.
- [5] Nesterov, Y. (2013). *Gradient methods for minimizing composite functions*. Mathematical Programming, 140(1), 125-161.