

Implémentation de la Régression Logistique (From Scratch)

OULAD ALI Ayoub

22 octobre 2025

Table des matières

1	Introduction à la régression logistique	3
1.1	Estimation des paramètres	3
2	Préparation des Données	3
3	Maximum de vraisemblance et condition d'optimalité	3
3.1	Condition du maximum de vraisemblance	4
4	Fonctions Utilitaires et Gradient	4
5	Compte-rendu : Difficultés rencontrées lors de l'implémentation	5
5.1	Génération des données	5
5.2	Problèmes de convergence	5
6	Choix de la méthode d'optimisation	5
7	Visualisation de la sigmoïde	7
8	Implémentation de l'algorithme ADAM	7
9	Descente de gradient stochastique (SGD / Mini-batch)	9
9.1	Principe de la SGD	9
9.2	Difficultés rencontrées	9
9.3	Performances	10
10	Vérification des résultats	11
10.1	Matrice de confusion	11
11	Approche d'Optimisation par Moyennage	13
11.1	Principe de la Méthode	13
11.2	Justification Théorique et Convergence vers l'Argmin	13
11.3	Avantages et Inconvénients	14
12	Application à des datasets "réels"	16
12.1	Dataset Titanic	16
12.1.1	Test de la méthode par moyenne (sur Titanic)	16
12.1.2	Test de la descente de gradient simple (sur Titanic)	17
12.2	Dataset Diabète	18
13	Prolongement : Régression Logistique Multiclasse avec OvR (One vs Rest)	19
13.1	Principe	19
13.2	Implémentation et Test sur Iris	19
13.3	Avis sur OvR	21
A	Sources utilisées	21

1 Introduction à la régression logistique

La régression logistique est un modèle de classification utilisé pour prédire la probabilité qu'une observation appartienne à une classe binaire $Y \in \{0, 1\}$.

La régression logistique prédit une probabilité à l'aide de la **fonction sigmoïde** :

$$\mathbb{P}(Y = 1 | X) = \frac{1}{1 + \exp(-X^\top \beta - b)}$$

où :

- X est un vecteur de variables explicatives,
- β est le vecteur des coefficients,
- b est le biais (intercept).

1.1 Estimation des paramètres

Les paramètres β et b sont estimés par **maximisation de la vraisemblance**. Cela revient à minimiser la **fonction de perte logistique** (ou entropie croisée) :

$$L(\beta, b) = - \sum_{i=1}^n [y_i \log(p_i) + (1 - y_i) \log(1 - p_i)]$$

avec $p_i = \mathbb{P}(Y_i = 1 | X_i)$ la probabilité prédite.

2 Préparation des Données

Génération d'un jeu de données synthétiques pour la classification binaire.

```
1 import numpy as np
2 from sklearn.datasets import make_classification
3 import pandas as pd
4 from matplotlib import pyplot as plt
5 from numpy.linalg import norm
6 from sklearn.metrics import accuracy_score
7
8 # Je choisis 100k et 500k comme taille du jeu de données
9 # car empiriquement c'est celui que j'ai trouvé le meilleur
10 X, y = make_classification(n_samples=100000, n_features=5, n_classes=2,
11                           random_state=42)
12 df = pd.DataFrame(X, columns=['feature_1', 'feature_2', 'feature_3', 'feature_4',
13                     'feature_5'])
12 df['target'] = y
```

Listing 1 – Génération des données synthétiques

3 Maximum de vraisemblance et condition d'optimalité

En régression logistique, l'objectif est de trouver les paramètres (\mathbf{B}, b) qui maximisent la vraisemblance du modèle, c'est-à-dire la probabilité d'observer les données données X et y sous le modèle logistique. Cela revient à minimiser la **fonction d'entropie croisée** (ou **négative log-vraisemblance**) définie par :

$$\mathcal{L}(\mathbf{B}, b) = - \sum_{i=1}^n [y_i \log(p_i) + (1 - y_i) \log(1 - p_i)]$$

où $p_i = \sigma(\mathbf{x}_i^\top \mathbf{B} + b)$ est la probabilité prédite pour l'observation i par la fonction sigmoïde σ . La fonction d'entropie croisée est **convexe**, ce qui garantit l'existence d'un minimum global unique (l'argmin).

3.1 Condition du maximum de vraisemblance

Le **maximum de vraisemblance** (ou le minimum de l'entropie croisée) est atteint lorsque la dérivée de la fonction par rapport aux paramètres est nulle. Autrement dit :

$$\nabla_{\mathbf{B}} \mathcal{L}(\mathbf{B}, b) = \mathbf{0}, \quad \text{et} \quad \frac{\partial \mathcal{L}(\mathbf{B}, b)}{\partial b} = 0$$

Ce sont les **conditions du premier ordre**, caractéristiques des points stationnaires (ici, le minimum global unique). En pratique, les algorithmes d'optimisation itératifs comme la descente de gradient cherchent un point où le **gradient est proche de zéro** :

$$\|\nabla \mathcal{L}(\mathbf{B}, b)\| \approx 0$$

Cela constitue un critère d'arrêt courant pour l'optimisation.

4 Fonctions Utilitaires et Gradient

Implémentation des fonctions sigmoïde, prédiction, log-vraisemblance, entropie croisée et calcul des gradients.

```

1 def sigmoid(z):
2     # Ajout d'un clip pour eviter l'overflow/underflow
3     z = np.clip(z, -250, 250)
4     return 1 / (1 + np.exp(-z)) # Utilisation de 1 / (1 + exp(-z))
5
6 def f(B, b, X_data=X): # Renommer X pour eviter conflit de nom
7     z = X_data @ B + b
8     return sigmoid(z)
9
10 def log_vraisemblance(B, b, X_data=X, y_data=y, eps_clip=1e-15):
11     p = f(B, b, X_data)
12     p = np.clip(p, eps_clip, 1 - eps_clip) # evite log(0)
13     # Correction: log_vraisemblance est la somme, pas la negative somme
14     return (y_data * np.log(p) + (1 - y_data) * np.log(1 - p)).sum()
15
16 def entropie_croisee_binaire(B, b, X_data=X, y_data=y):
17     # entropie croisee = - log vraisemblance / n
18     n_samples = X_data.shape[0]
19     if n_samples == 0:
20         return 0
21     # Utiliser la moyenne pour la stabilite numerique et la comparaison
22     return -log_vraisemblance(B, b, X_data, y_data) / n_samples
23
24 def grad_B_entropie_croisee_binaire(B, b, X_data=X, y_data=y):
25     n_samples = X_data.shape[0]
26     if n_samples == 0:
27         return np.zeros_like(B)
28     # Gradient de l'entropie croisee MOYENNE
29     return (f(B, b, X_data) - y_data) @ X_data / n_samples
30
31 def grad_b_entropie_croisee_binaire(B, b, X_data=X, y_data=y):
32     n_samples = X_data.shape[0]
33     if n_samples == 0:
34         return 0
35     # Gradient de l'entropie croisee MOYENNE
36     return (f(B, b, X_data) - y_data).sum() / n_samples

```

Listing 2 – Fonctions de base et gradients

5 Compte-rendu : Difficultés rencontrées lors de l'implémentation

5.1 Génération des données

Les données ont été générées à l'aide de la fonction `make_classification` de `sklearn.datasets` avec 5 features et deux classes. Cette étape ne posait pas de difficulté particulière, mais a nécessité de bien comprendre la structure de la matrice X pour l'analyse suivante.

5.2 Problèmes de convergence

L'implémentation initiale d'une descente de gradient naïve a rencontré plusieurs obstacles :

- **Mauvais choix du pas ρ** : Des valeurs trop grandes ont fait diverger l'entropie croisée ; des valeurs trop petites ont ralenti la convergence ou causé des erreurs numériques ($\log(0)$).
- **Critère d'arrêt mal défini** : L'arrêt basé uniquement sur une tolérance ε sur la fonction coût ne suffisait pas. Un nouveau critère basé sur la variation absolue de l'entropie croisée et la norme du gradient a été utilisé :

$$\text{tol} = \max(|\mathcal{L}_k - \mathcal{L}_{k-1}|, \|\nabla \mathcal{L}_k\|^2)$$

où :

- \mathcal{L}_k est la valeur de l'entropie croisée à l'itération k ,
- $|\mathcal{L}_k - \mathcal{L}_{k-1}|$ mesure l'évolution de la fonction de coût,
- $\|\nabla \mathcal{L}_k\|^2 = \|\nabla_B \mathcal{L}_k\|^2 + \left(\frac{\partial \mathcal{L}_k}{\partial b}\right)^2$ est la norme au carré du gradient total.

On continue les itérations tant que $\text{tol} > \varepsilon$, avec ε un seuil de tolérance fixé (par exemple $\varepsilon = 10^{-2}$).

- **Gradients non nuls à l'optimum** : Les gradients de B et b restaient parfois élevés, suggérant que la convergence était encore incomplète ou piégée dans un minimum local (ou que la tolérance ϵ était trop grande).

6 Choix de la méthode d'optimisation

Dans ce projet, nous avons choisi d'utiliser une **méthode de descente de gradient** pour minimiser la fonction d'entropie croisée, qui correspond à la négative de la log-vraisemblance dans le cadre de la régression logistique.

Ce choix est motivé par notre **familiarité avec les méthodes de gradient**, développée dans le cadre du **cours d'optimisation**. La descente de gradient est une méthode simple à mettre en œuvre, tout en étant suffisamment robuste pour résoudre efficacement des problèmes d'optimisation différentiable. Elle permet d'approximer le minimum de la fonction objectif en suivant la direction du gradient opposé.

```
1 def descente_gradient(X_data=X, y_data=y, NitMax=100000, rho=1e-4, eps=1/100, B=None, b=None):
2     # Initialisation si B et b ne sont pas fournis
3     if B is None:
4         B = np.ones(X_data.shape[1])
5     if b is None:
6         b = 1.0 # Utiliser un flottant pour b
7
8     iter = 0
9     gradB = grad_B_entropie_croisee_binaire(B, b, X_data, y_data)
10    gradb = grad_b_entropie_croisee_binaire(B, b, X_data, y_data)
11    B_list = []
12    b_list = []
```

```

13     entropie_croisee_binaire_liste = [entropie_croisee_binaire(B, b, X_data,
14                                         y_data)]
15     norme_sq = norm(gradB)**2 + gradb**2 # Norme au carre
16     # Attention: la tolerance initiale doit etre basee sur la norme et non la
17     # valeur de la fonction
18     tol = norme_sq # Initialisation avec la norme du gradient au carre
19     print(f"Iteration 0: Entropie={entropie_croisee_binaire_liste[0]:.4f}, Norme^2
20           Grad={norme_sq:.4f}, Tol={tol:.4f}")
21
22
23     while iter < NitMax and tol > eps:
24         B_prev = B.copy() # Garder une copie pour verifier changement
25         b_prev = b
26
27         B = B - rho * gradB
28         b = b - rho * gradb
29
30         gradB = grad_B_entropie_croisee_binaire(B, b, X_data, y_data)
31         gradb = grad_b_entropie_croisee_binaire(B, b, X_data, y_data)
32
33         B_list.append(B)
34         b_list.append(b)
35         current_entropie = entropie_croisee_binaire(B, b, X_data, y_data)
36         entropie_croisee_binaire_liste.append(current_entropie)
37
38         norme_sq = norm(gradB)**2 + gradb**2 # Norme au carre
39         # La tolerance est le max entre la variation d'entropie et la norme du
40         # gradient au carre
41         tol = max(abs(entropie_croisee_binaire_liste[-2] -
42                     entropie_croisee_binaire_liste[-1]), norme_sq)
43
44         iter += 1
45         if iter % 1000 == 0: # Afficher de temps en temps
46             # Vefier si les parametres changent encore significativement
47             param_change_B = norm(B - B_prev) / (norm(B) + 1e-8)
48             param_change_b = abs(b - b_prev) / (abs(b) + 1e-8)
49             print(f"Iteration {iter}: Entropie={current_entropie:.4f}, Norme^2 Grad
50                   ={norme_sq:.4f}, Tol={tol:.4f}, dParam B={param_change_B:.2e}, dParam b={
51                   param_change_b:.2e}")
52
53
54     B_list = np.array(B_list)
55     b_list = np.array(b_list)
56     entropie_croisee_binaire_liste = np.array(entropie_croisee_binaire_liste)
57     cvg = tol <= eps
58     print(f"Fin: Iter={iter}, Entropie={entropie_croisee_binaire_liste[-1]:.4f},
59           Norme^2 Grad={norme_sq:.4f}, Tol={tol:.4f}, Converged={cvg}")
60
61     return B, b, iter, cvg, B_list, b_list, entropie_croisee_binaire_liste
62
63 # Exemple d'appel (peut prendre du temps)
64 # B_res, b_res, iterations, converged, _, _, cost_history = descente_gradient(
65 #     NitMax=10000, rho=1e-2, eps=1e-4)
66 # print(f"B = {B_res}")
67 # print(f"b = {b_res}")
68 # print(f"Iterations = {iterations}")
69 # print(f"Converged = {converged}")
70 # plt.plot(cost_history)
71 # plt.title("Evolution de l'entropie croisee")
72 # plt.xlabel("Iteration")
73 # plt.ylabel("Entropie Croisee Moyenne")
74 # plt.show()

```

Listing 3 – Implémentation de la descente de gradient

7 Visualisation de la sigmoïde

- **Objectif** : tracer $P(y = 1 | X)$ en fonction d'une seule variable x_i , en fixant les autres.
- **Problème initial** : la fonction sigmoïde était mal tracée (courbes plates ou inversées) car on utilisait directement `sigmoid(X)` au lieu de :

$$f(X) = \sigma(X \cdot B + b)$$

- **Solution** : On a reconstruit un X_{fixed} où seule la feature d'intérêt varie, et les autres sont fixées à la moyenne.

```
1 # Assurez-vous que B et b sont definis (resultats de descente_gradient)
2 # B_res, b_res, _, _, _, _ = descente_gradient() # Execution reelle si
   necessaire
3
4 # Supposons que B_res et b_res sont disponibles
5 B_example = np.array([0.598, 0.011, 1.104, 0.895, 1.116]) # Exemple de valeurs
6 b_example = 0.008
7
8 x_vals = np.linspace(X[:, 0].min(), X[:, 0].max(), 200)
9 X_fixed = np.tile(X.mean(axis=0), (len(x_vals), 1))
10 X_fixed[:, 0] = x_vals
11
12 def predict_proba_fixed(X_data, B_param, b_param):
13     # Assurez-vous que f est accessible ou redefinissez la prediction ici
14     z = X_data @ B_param + b_param
15     return sigmoid(z)
16
17 p = predict_proba_fixed(X_fixed, B_example, b_example)
18
19 # plt.figure() # Creer une nouvelle figure si besoin
20 # plt.plot(x_vals, p, label='Sigmoide (feature 1)', color='blue')
21 # Pour scatter, il faut utiliser les donnees originales X et y
22 # Afficher un sous-echantillon pour eviter la surcharge
23 # sample_indices = np.random.choice(X.shape[0], 500, replace=False)
24 # plt.scatter(X[sample_indices, 0], y[sample_indices], alpha=0.2, color='orange',
   ', label='y reel (echantillon)', s=10)
25
26 # plt.xlabel('feature_1')
27 # plt.ylabel('P(y=1)')
28 # plt.title("Sigmoide en fonction de la feature 1")
29 # plt.grid(True)
30 # plt.legend()
31 # plt.show() # Decommenter pour afficher le graphique
```

Listing 4 – Visualisation de la sigmoïde par rapport à une feature

8 Implémentation de l'algorithme ADAM

L'algorithme ADAM a été testé pour améliorer la descente classique.

- **Erreurs initiales** : les coefficients de moment (β_1, β_2) étaient initialement nuls, neutralisant les avantages de l'algorithme.
- **Gestion des dimensions** : m et v ont été correctement définis comme des vecteurs de même forme que le gradient g .
- **Amélioration constatée** : ADAM a permis une convergence plus fluide, mais un bon réglage de α, β_1, β_2 restait crucial.

```
1 def ADAM(X_data=X, y_data=y, NitMax=100000, alpha=0.01, eps=1e-4, beta1=0.9,
   beta2=0.999, B=None, b=None):
2     # Initialisation si B et b ne sont pas fournis
```

```

3     if B is None:
4         B = np.zeros(X_data.shape[1])
5     if b is None:
6         b = 0.0
7
8     iter = 0
9     gradB = grad_B_entropie_croisee_binaire(B, b, X_data, y_data)
10    gradb = grad_b_entropie_croisee_binaire(B, b, X_data, y_data)
11
12    m_B = np.zeros(gradB.shape)
13    v_B = np.zeros(gradB.shape)
14    m_b, v_b = 0.0, 0.0
15
16    B_list = []
17    b_list = []
18    entropie_croisee_binaire_liste = [entropie_croisee_binaire(B, b, X_data,
19                                y_data)]
20
21    norme_sq = norm(gradB)**2 + gradb**2
22    tol = norme_sq # Initialisation avec la norme du gradient au carre
23    print(f"Iteration 0: Entropie={entropie_croisee_binaire_liste[0]:.4f}, Norme^2
24      Grad={norme_sq:.4f}, Tol={tol:.4f}")
25
26    while iter < NitMax and tol > eps:
27        B_prev = B.copy() # Pour debug
28        b_prev = b
29
30        m_B = beta1 * m_B + (1 - beta1) * gradB
31        v_B = beta2 * v_B + (1 - beta2) * gradB**2
32        m_b = beta1 * m_b + (1 - beta1) * gradb
33        v_b = beta2 * v_b + (1 - beta2) * gradb**2
34
35        # Correction du biais
36        m_B_hat = m_B / (1 - beta1***(iter + 1))
37        v_B_hat = v_B / (1 - beta2***(iter + 1))
38        m_b_hat = m_b / (1 - beta1***(iter + 1))
39        v_b_hat = v_b / (1 - beta2***(iter + 1))
40
41        # Mise a jour
42        B = B - alpha * m_B_hat / (np.sqrt(v_B_hat) + 1e-8)
43        b = b - alpha * m_b_hat / (np.sqrt(v_b_hat) + 1e-8)
44
45        B_list.append(B)
46        b_list.append(b)
47        current_entropie = entropie_croisee_binaire(B, b, X_data, y_data)
48        # Vérifier si l'entropie augmente (problème de taux d'apprentissage)
49        if len(entropie_croisee_binaire_liste) > 1 and current_entropie >
50            entropie_croisee_binaire_liste[-1] + 1e-6:
51            print(f"Attention: l'entropie augmente à l'itération {iter+1}")
52            # Option: réduire alpha, arrêter, ou revenir en arrière
53            # break
54            entropie_croisee_binaire_liste.append(current_entropie)
55
56    gradB = grad_B_entropie_croisee_binaire(B, b, X_data, y_data)
57    gradb = grad_b_entropie_croisee_binaire(B, b, X_data, y_data)
58
59    norme_sq = norm(gradB)**2 + gradb**2
60    tol = max(abs(entropie_croisee_binaire_liste[-2] -
61                  entropie_croisee_binaire_liste[-1]), norme_sq)
62
63    iter += 1

```

```

62     if iter % 100 == 0: # Afficher moins souvent pour ADAM
63         param_change_B = norm(B - B_prev) / (norm(B) + 1e-8)
64         param_change_b = abs(b - b_prev) / (abs(b) + 1e-8)
65         print(f"Iteration {iter}: Entropie={current_entropie:.4f}, Norme^2 Grad
66             ={norme_sq:.4f}, Tol={tol:.4f}, dParam B={param_change_B:.2e}, dParam b={param_change_b:.2e}")
67
68     B_list = np.array(B_list)
69     b_list = np.array(b_list)
70     entropie_croisee_binaire_liste = np.array(entropie_croisee_binaire_liste)
71     cvg = tol <= eps
72     print(f"Fin ADAM: Iter={iter}, Entropie={entropie_croisee_binaire_liste[-1]:.4
73         f}, Norme^2 Grad={norme_sq:.4f}, Tol={tol:.4f}, Converged={cvg}")
74
75     return B, b, iter, cvg, B_list, b_list, entropie_croisee_binaire_liste
76
77 # Exemple d'appel ADAM (peut etre plus rapide)
78 # B_adam, b_adam, iter_adam, cvg_adam, _, _, cost_adam = ADAM(alpha=0.01, eps=1e
79 # -6, NitMax=5000)
80 # print(f"ADAM B = {B_adam}")
81 # print(f"ADAM b = {b_adam}")
82 # plt.figure()
83 # plt.title("Evolution de l'entropie croisee (ADAM)")
84 # plt.xlabel("Iteration")
85 # plt.ylabel("Entropie Croisee Moyenne")
86 # plt.show()

```

Listing 5 – Implémentation de l'optimiseur ADAM

9 Descente de gradient stochastique (SGD / Mini-batch)

Nous avons initialement tenté d'implémenter une descente de gradient stochastique (SGD) pour optimiser notre modèle de régression logistique.

9.1 Principe de la SGD

La SGD (ou sa variante mini-batch) met à jour les paramètres du modèle en utilisant une seule observation (SGD) ou un petit sous-ensemble (mini-batch) à chaque itération, contrairement à la descente de gradient classique qui utilise l'ensemble des données.

Avantages :

- Convergence potentiellement plus rapide en début d'optimisation.
- Possibilité de gérer des jeux de données très volumineux (qui ne tiennent pas en mémoire).
- Peut échapper plus facilement aux minima locaux.

Inconvénients :

- Convergence plus "bruyante" et moins stable vers la fin.
- Nécessité de bien ajuster le taux d'apprentissage (souvent avec une décroissance).
- Critère d'arrêt plus complexe à définir (le gradient sur un batch n'est pas représentatif du gradient global).

9.2 Difficultés rencontrées

L'implémentation de la SGD pure (batch_size=1) s'est avérée trop coûteuse en temps de calcul pour notre jeu de données. Chaque itération nécessitait le calcul du gradient pour une seule observation. La variante mini-batch est un compromis.

9.3 Performances

Les performances finales de la SGD/mini-batch peuvent être bonnes, mais nécessitent un réglage attentif des hyperparamètres (taille du batch, taux d'apprentissage, critère d'arrêt).

```

1 def descente_gradient_mini_batch(X_data=X, y_data=y, NitMax=1000, rho=1e-2, eps
2     =1e-4, batch_size=128, B=None, b=None):
3
4     if B is None:
5         B = np.zeros(X_data.shape[1])
6     if b is None:
7         b = 0.0
8
9     iter = 0
10    n_samples = X_data.shape[0]
11    if n_samples == 0:
12        print("Erreur: Dataset vide")
13        return B, b, 0, False # Retourner False pour convergence
14
15    # Pour Mini-batch, on n'utilise pas la variation de l'entropie comme critere d
16    # 'arret principal
17    # car elle est trop bruitee. On se base sur NitMax ou potentiellement une
18    # validation periodique.
19    # Le calcul de tol initial est juste indicatif sur le premier batch.
20
21    # Faire le premier calcul hors boucle pour initialiser entropie_liste
22    indices = np.random.choice(n_samples, batch_size, replace=False)
23    X_batch = X_data[indices]
24    y_batch = y_data[indices]
25
26    gradB = grad_B_entropie_croisee_binaire(B, b, X_batch, y_batch)
27    gradb = grad_b_entropie_croisee_binaire(B, b, X_batch, y_batch)
28    norme_sq_batch = norm(gradB)**2 + gradb**2
29    current_entropie_batch = entropie_croisee_binaire(B, b, X_batch, y_batch)
30    # entropie_croisee_binaire_liste = [current_entropie_batch] # On ne stocke pas
31    # forcement
32
33    print(f"MiniBatch Iter 0: Entropie Batch={current_entropie_batch:.4f}, Norme^2
34          Grad Batch={norme_sq_batch:.4f}")
35
36    while iter < NitMax: # On s'arrete principalement sur NitMax pour SGD/
37        MiniBatch
38
39        indices = np.random.choice(n_samples, batch_size, replace=False)
40        X_batch = X_data[indices]
41        y_batch = y_data[indices]
42
43        gradB = grad_B_entropie_croisee_binaire(B, b, X_batch, y_batch)
44        gradb = grad_b_entropie_croisee_binaire(B, b, X_batch, y_batch)
45
46        B = B - rho * gradB
47        b = b - rho * gradb
48
49        # Optionnel: Affichage periodique et verification convergence sur le set
50        # complet
51        if (iter + 1) % 100 == 0:
52            current_entropie_full = entropie_croisee_binaire(B, b, X_data, y_data)
53
54            gradB_full = grad_B_entropie_croisee_binaire(B, b, X_data, y_data)
55            gradb_full = grad_b_entropie_croisee_binaire(B, b, X_data, y_data)
56            norme_sq_full = norm(gradB_full)**2 + gradb_full**2
57            print(f"MiniBatch Iter {iter+1}: Entropie Full={current_entropie_full
58                  :.4f}, Norme^2 Grad Full={norme_sq_full:.4f}")
59            # Critere d'arret base sur la norme du gradient sur le set complet
60            if norme_sq_full < eps:

```

```

51         print(f"Convergence atteinte a l'iteration {iter+1} (norme
52             gradient)")
53             break
54
55     iter += 1
56
57 # Calcul final sur tout le dataset pour statut final
58 final_entropie_full = entropie_croisee_binaire(B, b, X_data, y_data)
59 gradB_full = grad_B_entropie_croisee_binaire(B, b, X_data, y_data)
60 gradb_full = grad_b_entropie_croisee_binaire(B, b, X_data, y_data)
61 norme_sq_full = norm(gradB_full)**2 + gradb_full**2
62 converged = norme_sq_full <= eps # Verification a la fin
63 print(f"Fin MiniBatch: Iter={iter}, Entropie Full={final_entropie_full:.4f},
64     Norme^2 Grad Full={norme_sq_full:.4f}, Converged={converged}")
65
66 return B, b, iter, converged # On retourne les derniers B, b
67
68 # B_mb, b_mb, iter_mb, cvg_mb = descente_gradient_mini_batch(NitMax=2000, rho=1e
# -2, eps=1e-4, batch_size=64)
69 # print(f"MiniBatch B = {B_mb}")
70 # print(f"MiniBatch b = {b_mb}")

```

Listing 6 – Implémentation de la descente de gradient Mini-Batch

Remarque : La méthode nécessite souvent un grand nombre d’itérations, et la convergence est plus oscillante.

10 Vérification des résultats

Afin de vérifier la robustesse des résultats, nous générerons un vecteur `y_pred` où chaque élément indique si la probabilité prédictive par la sigmoïde est supérieure à 0.5. Ainsi, chaque élément de `y_pred` sera égal à 1 si la probabilité est supérieure ou égale à 0.5, et à 0 sinon.

Cela peut être formulé de la manière suivante :

$$y_{\text{pred},i} = \begin{cases} 1 & \text{si } \hat{P}(y_i = 1 | X_i) \geq 0.5 \\ 0 & \text{sinon} \end{cases}$$

Cela nous permet de transformer les probabilités continues en prédictions binaires, que nous pourrons ensuite comparer avec les vraies valeurs y pour évaluer la performance du modèle.

10.1 Matrice de confusion

Une des méthodes les plus courantes pour mesurer sa performance est la **matrice de confusion**. Elle fournit une vue d’ensemble de la précision d’un modèle de classification en comparant les prédictions faites par le modèle avec les véritables étiquettes de classe.

Elle contient quatre éléments principaux :

- **Vrai Positif (TP)** : Le modèle a correctement prédit la classe positive (1).
- **Faux Positif (FP)** : Le modèle a prédit la classe positive alors que l’étiquette réelle est négative (0). (Erreur de Type I)
- **Vrai Négatif (TN)** : Le modèle a correctement prédit la classe négative (0).
- **Faux Négatif (FN)** : Le modèle a prédit la classe négative alors que l’étiquette réelle est positive (1). (Erreur de Type II)

		Prédit	
		Positif (1)	Négatif (0)
Réel	Positif (1)	TP	FN
	Négatif (0)	FP	TN

À partir de cette matrice, nous pouvons calculer plusieurs indicateurs :

1. **Exactitude (Accuracy)** : Proportion de prédictions correctes.

$$\text{Accuracy} = \frac{TP + TN}{TP + TN + FP + FN}$$

2. **Précision (Precision)** : Proportion des prédictions positives qui sont réellement positives.

$$\text{Precision} = \frac{TP}{TP + FP}$$

3. **Rappel (Recall) ou Sensibilité** : Proportion des vrais positifs détectés par le modèle.

$$\text{Recall} = \frac{TP}{TP + FN}$$

4. **Score F1** : Moyenne harmonique entre la précision et le rappel.

$$\text{F1-Score} = 2 \times \frac{\text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}}$$

```

1 # Supposons que B_res, b_res sont les resultats de descente_gradient
2 # et B_mb, b_mb sont les resultats de descente_gradient_mini_batch
3
4 # Utiliser de nouvelles donnees de test
5 X_new, y_new = make_classification(n_samples=1000, n_features=5, n_classes=2,
6     random_state=123)
7
8 # Predictions avec la descente classique (exemple)
9 B_gd_example = np.array([0.512, 0.011, 1.016, 1.026, 1.014]) # Remplacer par vos
10    vraies valeurs
11 b_gd_example = 0.008
12
13 y_pred_gd = f(B_gd_example, b_gd_example, X_new)
14 y_pred_gd_class = (y_pred_gd >= 0.5).astype(int)
15 accuracy_gd = accuracy_score(y_new, y_pred_gd_class)
16 print(f"Precision (Accuracy) pour la descente classique: {accuracy_gd:.3f}")
17
18 # Predictions avec Mini-Batch (exemple)
19 B_mb_example = np.array([0.369, 0.070, 0.853, 1.376, 0.716]) # Remplacer par vos
20    vraies valeurs
21 b_mb_example = 0.057
22
23 y_pred_mb = f(B_mb_example, b_mb_example, X_new)
24 y_pred_mb_class = (y_pred_mb >= 0.5).astype(int)
25 accuracy_mb = accuracy_score(y_new, y_pred_mb_class)
26 print(f"Precision (Accuracy) pour la descente mini-batch: {accuracy_mb:.3f}")
27
28 # Fonction Matrice de Confusion
29 def matrice_de_confusion(y_true, y_pred):
30     # Classes sont 0 et 1
31     TP = np.sum((y_true == 1) & (y_pred == 1))
32     FN = np.sum((y_true == 1) & (y_pred == 0))
33     FP = np.sum((y_true == 0) & (y_pred == 1))
34     TN = np.sum((y_true == 0) & (y_pred == 0))
35     return np.array([[TP, FN], [FP, TN]])
36
37 cm_gd = matrice_de_confusion(y_new, y_pred_gd_class)
38 print("\nMatrice de confusion (Descente Gradient):")
39 print(cm_gd)
40
41 TP_gd = cm_gd[0, 0]
42 FN_gd = cm_gd[0, 1]
```

```

40 FP_gd = cm_gd[1, 0]
41 TN_gd = cm_gd[1, 1]
42
43 precision_gd = TP_gd / (TP_gd + FP_gd) if (TP_gd + FP_gd) != 0 else 0
44 recall_gd = TP_gd / (TP_gd + FN_gd) if (TP_gd + FN_gd) != 0 else 0
45 f1_gd = 2 * (precision_gd * recall_gd) / (precision_gd + recall_gd) if (
    precision_gd + recall_gd) != 0 else 0
46
47 print(f"Precision (Precision GD): {precision_gd:.3f}")
48 print(f"Rappel (Recall GD): {recall_gd:.3f}")
49 print(f"Score F1 (GD): {f1_gd:.3f}")
50
51 cm_mb = matrice_de_confusion(y_new, y_pred_mb_class)
52 print("\nMatrice de confusion (Mini-Batch):")
53 print(cm_mb)
54 # Calculer Precision, Recall, F1 pour Mini-Batch de maniere similaire si besoin

```

Listing 7 – Calcul de la précision et de la matrice de confusion

11 Approche d’Optimisation par Moyennage

Cette section détaille une approche d’optimisation alternative explorée dans ce projet, issue d’une réflexion initiée lors d’un travail antérieur sur la "Régression Linéaire et Descente de Gradient from scratch". L’idée est d’exploiter la convexité du problème et les propriétés statistiques pour approcher l’optimum global.

11.1 Principe de la Méthode

L’approche consiste à :

1. Diviser le jeu de données d’entraînement $D = \{(X_i, y_i)\}_{i=1}^N$ en M sous-ensembles (splits) D_1, \dots, D_M . Idéalement, ces splits sont disjoints ou tirés aléatoirement.
2. Pour chaque sous-ensemble D_j , effectuer une optimisation complète (par exemple, par descente de gradient) pour minimiser l’entropie croisée $\mathcal{L}_j(B, b)$ calculée sur D_j . Cela donne un jeu de paramètres optimaux (\hat{B}_j, \hat{b}_j) pour ce split.

$$(\hat{B}_j, \hat{b}_j) = \arg \min_{B, b} \mathcal{L}_j(B, b)$$

3. Calculer les paramètres finaux (\bar{B}, \bar{b}) en moyennant les paramètres obtenus sur tous les splits :

$$\bar{B} = \frac{1}{M} \sum_{j=1}^M \hat{B}_j \quad \text{et} \quad \bar{b} = \frac{1}{M} \sum_{j=1}^M \hat{b}_j$$

11.2 Justification Théorique et Convergence vers l’Argmin

Pourquoi cette méthode est-elle susceptible de converger vers l’argmin global (B^*, b^*) ?

1. **Convexité** : Comme mentionné, l’entropie croisée $\mathcal{L}(B, b)$ est convexe. Son minimum global (B^*, b^*) est unique.

2. **Estimateurs Locaux** : Chaque (\hat{B}_j, \hat{b}_j) est le minimum de $\mathcal{L}_j(B, b)$. Si le sous-ensemble D_j est suffisamment grand et représentatif de la distribution globale des données, alors \mathcal{L}_j est une bonne approximation de \mathcal{L} , et son minimum (\hat{B}_j, \hat{b}_j) sera un estimateur (bruité) du minimum global (B^*, b^*) .

3. **Loi des Grands Nombres** : Considérons (\hat{B}_j, \hat{b}_j) comme une variable aléatoire dont l’espérance est (B^*, b^*) (ou proche de celle-ci si les splits sont représentatifs). En moyennant un grand nombre M de ces estimateurs (idéalement issus de splits indépendants ou faiblement

corrélés), la **Loi des Grands Nombres** stipule que la moyenne (\bar{B}, \bar{b}) converge en probabilité vers l'espérance de ces estimateurs :

$$\bar{B} \xrightarrow[M \rightarrow \infty]{P} \mathbb{E}[\hat{B}_j] \approx B^*$$

$$\bar{b} \xrightarrow[M \rightarrow \infty]{P} \mathbb{E}[\hat{b}_j] \approx b^*$$

Cette convergence est d'autant plus vraie que le **jeu de données initial est grand**, permettant de créer de nombreux sous-ensembles D_j qui restent suffisamment grands pour que (\hat{B}_j, \hat{b}_j) soit un estimateur raisonnable de (B^*, b^*) .

4. Résultats Empiriques : Les expériences menées (voir code ci-dessous) montrent que les paramètres (\bar{B}, \bar{b}) obtenus par cette méthode donnent une bonne précision (*Accuracy* $\approx 85\%$ sur le jeu de données synthétique), ce qui conforte l'idée que (\bar{B}, \bar{b}) est une bonne approximation de l'argmin global (B^*, b^*) .

11.3 Avantages et Inconvénients

Avantages :

- **Parallélisation :** Les M optimisations peuvent être effectuées en parallèle, réduisant potentiellement le temps total d'exécution sur des architectures multi-coeurs ou distribuées.
- **Robustesse :** Le moyennage réduit la variance des estimateurs individuels due à l'échantillonnage spécifique de chaque split, rendant le résultat final potentiellement plus stable et robuste aux découpages aléatoires.

Inconvénients :

- **Dépendance à la taille des données :** La méthode est plus pertinente lorsque le dataset initial est **grand**, afin que chaque sous-ensemble soit suffisamment représentatif. Sur de petits datasets, les estimateurs locaux peuvent être très bruités et leur moyenne moins fiable.
- **Efficacité en calcul séquentiel :** Sans parallélisation, effectuer M optimisations complètes peut être moins efficace que des méthodes itératives globales comme ADAM ou SGD.
- **Pas un algorithme standard :** Cette approche spécifique n'est pas un algorithme d'optimisation standard nommé dans la littérature, bien qu'elle s'inspire de concepts connus.

```

1 # Assurez-vous que X et y sont definis (grand dataset initial)
2 n_splits = 500
3 if X.shape[0] < n_splits:
4     print("Attention: Moins d'échantillons que de splits demandés.")
5     n_splits = max(1, int(X.shape[0]/10)) # S'assurer d'avoir au moins 10 points
6     par split si possible
7
8 X_list = np.array_split(X, n_splits)
9 y_list = np.array_split(y, n_splits)
10
11 # Fonction pour estimer sur un seul split
12 def estimateur_split(X_sub, y_sub):
13     # Utiliser la fonction descente_gradient definie precedemment
14     # Ajuster les hyperparametres si necessaire pour une convergence raisonnable
15     # sur les petits sets
16     B_est, b_est, _, cvg, _, _ = descente_gradient(X_data=X_sub, y_data=y_sub,
17         NitMax=2000, rho=5e-3, eps=1e-3, B=np.zeros(X_sub.shape[1]), b=0.0) # Params
18         ajustes pour splits plus petits
19     if not cvg:
20         print("Attention: l'optimisation sur un split n'a pas atteint la tolerance
21             .")
22     return B_est, b_est

```

```

18
19 Est = []
20 print(f"Estimation sur {n_splits} sous-échantillons...")
21 successful_estimations = 0
22 for i, (X_split, y_split) in enumerate(zip(X_list, y_list)): # Renommer pour
    clarte
23     if X_split.shape[0] > 1: # Vérifier qu'il y a assez de données
        # print(f"Processing split {i+1}/{n_splits} (taille {X_split.shape[0]})")
        ... # Décommenter pour verbose
24     try:
25         B, b = estimateur_split(X_split, y_split) # Appeler la fonction d'
    estimation
26         if np.isfinite(B).all() and np.isfinite(b): # Vérifier si résultats
    valides
27             Est.append((B, b))
28             successful_estimations += 1
29         else:
30             print(f"Résultats non finis sur split {i+1}, ignore.")
31     except Exception as e:
32         print(f"Erreur sur le split {i+1}: {e}")
33     else:
34         print(f"Split {i+1} trop petit ({X_split.shape[0]} échantillons), ignore
    .")
35 print(f"Nombre d'estimations réussies: {successful_estimations}/{n_splits}")
36
37 # Calculer la moyenne des paramètres estimés
38 if Est: # S'assurer qu'il y a des résultats valides
39     B_final_mean = np.mean([B for B, b in Est], axis=0)
40     b_final_mean = np.mean([b for B, b in Est], axis=0)
41
42     print("\nMoyenne des B estimés sur les splits :")
43     print(B_final_mean)
44     print("Moyenne des b estimés sur les splits :")
45     print(b_final_mean)
46
47 # Évaluer la performance du modèle moyen sur de nouvelles données
48 X_new_eval, y_new_eval = make_classification(n_samples=1000, n_features=5,
    n_informative=5, n_redundant=0, n_classes=2, random_state=42) # Utiliser un
    nom différent
49
50     y_pred_mean = f(B_final_mean, b_final_mean, X_new_eval) # f doit être
    accessible globalement ou redéfinie
51     y_pred_mean_class = (y_pred_mean >= 0.5).astype(int)
52     accuracy_mean = accuracy_score(y_new_eval, y_pred_mean_class)
53     print(f"\nPrecision (Accuracy) du modèle moyen sur nouveau test set: {
    accuracy_mean:.3f}")
54
55 # Comparer avec la descente de gradient sur tout le dataset
56 print("\nComparaison avec Descente de Gradient classique (exécution rapide)
    ...")
57 B_full_comp, b_full_comp, _, _, _, _, _ = descente_gradient(X_data=X, y_data
    =y, NitMax=5000, rho=1e-3, eps=1e-4) # Exécution plus courte pour
    comparaison
58     y_pred_full_comp = f(B_full_comp, b_full_comp, X_new_eval)
59     y_pred_full_comp_class = (y_pred_full_comp >= 0.5).astype(int)
60     accuracy_full_comp = accuracy_score(y_new_eval, y_pred_full_comp_class)
61     print(f"Precision (Accuracy) Descente Gradient classique sur nouveau test
    set: {accuracy_full_comp:.3f}")
62
63
64
65 else:
66     print("Aucune estimation réussie sur les splits.")

```

Listing 8 – Estimation par moyenne sur sous-échantillons

Remarque : L'efficacité et la précision de cette méthode par moyennage dépendent fortement de la taille globale des données, du nombre de splits, et de la bonne convergence de l'optimisation sur chaque split.

12 Application à des datasets "réels"

12.1 Dataset Titanic

On teste l'algorithme sur un dataset concernant le Titanic où la valeur à prédire est 1 (survécu) ou 0 (décédé), en utilisant les features `pclass`, `sex`, `age`.

```

1 # Import nécessaire si non fait avant
2 from sklearn.datasets import fetch_openml
3 from sklearn.model_selection import train_test_split
4 from sklearn.preprocessing import StandardScaler
5
6 titanic = fetch_openml(name='titanic', version=1, as_frame=True, parser='auto')
7 # Specifier parser
8 df_titanic = titanic.frame
9
10 # Garder un sous-ensemble simple et propre, enlever les lignes avec NaN
11 df_titanic = df_titanic[['pclass', 'sex', 'age', 'survived']].dropna()
12
13 # Encoder la variable 'sex' en 0/1
14 df_titanic['sex'] = df_titanic['sex'].map({'male': 0, 'female': 1})
15 # Convertir 'survived' en entier 0/1
16 df_titanic['survived'] = df_titanic['survived'].astype(int)
17
18 # 80% train, 20% test
19 df_train, df_test = train_test_split(df_titanic, test_size=0.2, random_state=42,
20                                     stratify=df_titanic['survived'])
21
22 # Separer features et labels
23 X_train_titanic = df_train[['pclass', 'sex', 'age']].values
24 y_train_titanic = df_train['survived'].values
25 X_test_titanic = df_test[['pclass', 'sex', 'age']].values
26 y_test_titanic = df_test['survived'].values
27
28 # Standardiser les données (important pour la descente de gradient)
29 scaler_titanic = StandardScaler()
30 X_train_titanic_scaled = scaler_titanic.fit_transform(X_train_titanic)
31 X_test_titanic_scaled = scaler_titanic.transform(X_test_titanic)
```

Listing 9 – Chargement et préparation du dataset Titanic

12.1.1 Test de la méthode par moyenne (sur Titanic)

```

1 n_splits_titanic = 50 # Moins de splits car dataset plus petit
2 X_list_titanic = np.array_split(X_train_titanic_scaled, n_splits_titanic)
3 y_list_titanic = np.array_split(y_train_titanic, n_splits_titanic)
4
5 Est_titanic = []
6 print(f"Estimation sur {n_splits_titanic} sous-échantillons (Titanic)... ")
7 successful_estimations_titanic = 0
8 for i, (X_sub, y_sub) in enumerate(zip(X_list_titanic, y_list_titanic)): #
9     Renommer pour clarte
10    if X_sub.shape[0] > 1:
11        # print(f"Processing split {i+1}/{n_splits_titanic} (Titanic)... ") #
12        Decommenter pour verbose
```

```

11     try:
12         # Utiliser directement descente_gradient en passant les donnees du split
13         B, b, _, cvg, _, _ = descente_gradient(X_data=X_sub, y_data=y_sub,
14         NitMax=5000, rho=1e-2, eps=1e-3) # Ajuster hyperparams
15         if np.isfinite(B).all() and np.isfinite(b):
16             Est_titanic.append((B, b))
17             successful_estimations_titanic += 1
18             # if not cvg:
19                 #     print(f"Warning: Split {i+1} (Titanic) n'a pas atteint la
20                 # tolerance.")
21             else:
22                 print(f"Resultats non finis sur split {i+1} (Titanic), ignore.")
23     except Exception as e:
24         print(f"Erreur sur le split {i+1} (Titanic): {e}")
25     print(f"Nombre d'estimations reussies (Titanic): {successful_estimations_titanic
26         }/{n_splits_titanic}}")
27
28 if Est_titanic:
29     B_titanic_mean = np.mean([B for B, b in Est_titanic], axis=0)
30     b_titanic_mean = np.mean([b for B, b in Est_titanic], axis=0)
31
32     print("\nMoyenne des B (Titanic) :")
33     print(B_titanic_mean)
34     print("Moyenne des b (Titanic) :")
35     print(b_titanic_mean)
36
37     # Evaluer sur le test set Titanic
38     y_pred_titanic_mean = f(B_titanic_mean, b_titanic_mean,
39     X_test_titanic_scaled)
40     y_pred_titanic_mean_class = (y_pred_titanic_mean >= 0.5).astype(int)
41     accuracy_titanic_mean = accuracy_score(y_test_titanic,
42     y_pred_titanic_mean_class)
43     print(f"\nPrecision (Accuracy) du modele moyen sur test set Titanic: {accuracy_titanic_mean:.3f}")
44
45     cm_titanic_mean = matrice_de_confusion(y_test_titanic,
46     y_pred_titanic_mean_class)
47     print("\nMatrice de confusion (Moyenne, Titanic):")
48     print(cm_titanic_mean)
49 else:
50     print("Aucune estimation reussie sur les splits (Titanic).")

```

Listing 10 – Estimation par moyenne sur Titanic

La méthode des moyennes peut moins bien fonctionner sur des datasets plus petits.

12.1.2 Test de la descente de gradient simple (sur Titanic)

```

1 print("\nDescente de gradient simple sur tout le dataset d'entrainement Titanic
      ...")
2 # Appeler descente_gradient directement avec les donnees Titanic
3 B_titanic_full, b_titanic_full, _, cvg_full_titanic, _, _, _ = descente_gradient(
4     (
5         X_data=X_train_titanic_scaled,
6         y_data=y_train_titanic,
7         NitMax=100000, rho=1e-3, eps=1e-6 # Tenter une tolerance plus faible
8     )
9
10    print("\nB (Titanic Full):")
11    print(B_titanic_full)

```

```

11 print("b (Titanic Full):")
12 print(b_titanic_full)
13 print(f"Convergence atteinte (Full Titanic): {cvg_full_titanic}")
14
15 # Evaluer sur le test set Titanic
16 y_pred_titanic_full = f(B_titanic_full, b_titanic_full, X_test_titanic_scaled)
17 y_pred_titanic_full_class = (y_pred_titanic_full >= 0.5).astype(int)
18 accuracy_titanic_full = accuracy_score(y_test_titanic, y_pred_titanic_full_class
    )
19 print(f"\nPrecision (Accuracy) du modèle full sur test set Titanic: {accuracy_titanic_full:.3f}")
20
21 cm_titanic_full = matrice_de_confusion(y_test_titanic, y_pred_titanic_full_class
    )
22 print("\nMatrice de confusion (Full, Titanic):")
23 print(cm_titanic_full)
24
25 TP_full = cm_titanic_full[0, 0]
26 FN_full = cm_titanic_full[0, 1]
27 FP_full = cm_titanic_full[1, 0]
28 TN_full = cm_titanic_full[1, 1]
29
30 precision_full = TP_full / (TP_full + FP_full) if (TP_full + FP_full) != 0 else 0
31 recall_full = TP_full / (TP_full + FN_full) if (TP_full + FN_full) != 0 else 0
32 f1_full = 2 * (precision_full * recall_full) / (precision_full + recall_full) if (precision_full + recall_full) != 0 else 0
33
34 print(f"Precision (Precision Full): {precision_full:.3f}")
35 print(f"Rappel (Recall Full): {recall_full:.3f}")
36 print(f"Score F1 (Full): {f1_full:.3f}")

```

Listing 11 – Descente de gradient simple sur Titanic

Les résultats montrent une précision raisonnable, qui pourrait être améliorée (par exemple avec ADAM, régularisation, etc.).

12.2 Dataset Diabète

On applique la même démarche sur un dataset de prédiction du diabète.

```

1 # Charger le dataset (Assurez-vous que le fichier diabetes2.csv est accessible)
2 try:
3     df_diabetes = pd.read_csv("diabetes2.csv") # Mettre le bon chemin si nécessaire
4 except FileNotFoundError:
5     print("Fichier diabetes2.csv non trouvé. Assurez-vous qu'il est dans le répertoire.")
6     # Créer un dataframe vide pour éviter les erreurs suivantes, ou arrêter.
7     df_diabetes = pd.DataFrame() # Initialiser comme vide
8
9 if not df_diabetes.empty:
10     # Séparer features (X) et target (y - 'Outcome')
11     X_diabetes = df_diabetes.drop('Outcome', axis=1).values
12     y_diabetes = df_diabetes['Outcome'].values
13
14     # Split train/test
15     X_train_db, X_test_db, y_train_db, y_test_db = train_test_split(
16         X_diabetes, y_diabetes, test_size=0.2, random_state=42, stratify=y_diabetes
17     )
18
19     # Standardiser
20     scaler_db = StandardScaler()

```

```

21 X_train_db_scaled = scaler_db.fit_transform(X_train_db)
22 X_test_db_scaled = scaler_db.transform(X_test_db)
23
24 # Entrainement avec la descente de gradient
25 print("\nDescente de gradient simple sur dataset Diabete...")
26 # Appeler descente_gradient avec les donnees Diabete
27 B_db_full, b_db_full, _, cvg_db_full, _, _, _ = descente_gradient(
28     X_data=X_train_db_scaled,
29     y_data=y_train_db,
30     NitMax=100000, rho=1e-3, eps=1e-6 # Tenter tolerance plus faible
31 )
32
33 print("\nB (Diabetes Full):")
34 print(B_db_full)
35 print("b (Diabetes Full):")
36 print(b_db_full)
37 print(f"Convergence atteinte (Full Diabetes): {cvg_db_full}")
38
39 # Evaluation
40 y_pred_db_full = f(B_db_full, b_db_full, X_test_db_scaled)
41 y_pred_db_full_class = (y_pred_db_full >= 0.5).astype(int)
42 accuracy_db_full = accuracy_score(y_test_db, y_pred_db_full_class)
43 print(f"\nPrecision (Accuracy) du modele full sur test set Diabete: {accuracy_db_full:.3f}")
44
45 cm_db_full = matrice_de_confusion(y_test_db, y_pred_db_full_class)
46 print("\nMatrice de confusion (Full, Diabete):")
47 print(cm_db_full)
48 else:
49     print("\nDataset Diabete non charge, section ignoree.")

```

Listing 12 – Chargement et préparation du dataset Diabète

Les hyperparamètres (pas, tolérance) ont été ajustés pour tenter d'améliorer la convergence.

13 Prolongement : Régression Logistique Multiclasse avec OvR (One vs Rest)

13.1 Principe

Pour étendre la régression logistique binaire à un problème de classification avec $K > 2$ classes, la méthode OvR (ou OvA - One vs All) consiste à entraîner K classificateurs binaires indépendants. Pour chaque classe $k \in \{0, \dots, K-1\}$:

- On considère les exemples de la classe k comme positifs (étiquette 1).
- On considère tous les autres exemples (des classes $j \neq k$) comme négatifs (étiquette 0).
- On entraîne un modèle de régression logistique binaire (\hat{B}_k, \hat{b}_k) sur ces nouvelles étiquettes.

Lors de la prédiction pour une nouvelle instance X_{new} :

- On calcule la probabilité de sortie de chaque classifieur k : $p_k = \sigma(X_{\text{new}} \cdot \hat{B}_k + \hat{b}_k)$.
- La classe prédite est celle qui correspond au classifieur ayant donné la plus haute probabilité :

$$\text{Prediction} = \arg \max_k p_k$$

13.2 Implémentation et Test sur Iris

- Simple à mettre en œuvre en réutilisant le modèle binaire.
- On teste sur le dataset Iris ($K = 3$ classes).

```

1 from sklearn.datasets import load_iris
2
3 iris = load_iris()
4 X_iris = iris.data
5 y_iris = iris.target
6
7 # Standardiser les features
8 scaler_iris = StandardScaler()
9 X_iris_scaled = scaler_iris.fit_transform(X_iris)
10
11 # Split train/test
12 X_train_iris, X_test_iris, y_train_iris, y_test_iris = train_test_split(
13     X_iris_scaled, y_iris, test_size=0.2, random_state=42, stratify=y_iris
14 )

```

Listing 13 – Chargement et préparation du dataset Iris

```

1 def OVA_train(X_train, y_train, num_classes):
2     B_liste = []
3     b_liste = []
4     for k in range(num_classes):
5         # Creer les etiquettes binaires pour la classe k
6         y_binaire_k = (y_train == k).astype(int)
7
8         print(f"\nEntrainement du classifieur pour la classe {k} vs Rest...")
9         # Entrainer le classifieur binaire k
10        # Appeler descente_gradient directement
11        # Utiliser une initialisation a zero pour B, b peut aider
12        B_k, b_k, _, cvg_k, _, _, _ = descente_gradient(
13            X_data=X_train,
14            y_data=y_binaire_k,
15            NitMax=50000, rho=1e-2, eps=1e-4, # Augmenter rho, reduire eps?
16            B=np.zeros(X_train.shape[1]), b=0.0
17        )
18
19        B_liste.append(B_k)
20        b_liste.append(b_k)
21        print(f"Parametres trouves pour classe {k}: Converged={cvg_k}")
22        # print(f"B={B_k}, b={b_k}") # Decommenter pour voir les params
23
24    # Retourner des arrays numpy standards
25    return np.array(B_liste), np.array(b_liste)
26
27 def OVA_predict(X_test, B_liste, b_liste):
28     num_samples = X_test.shape[0]
29     num_classes = len(B_liste)
30
31     probabilities = np.zeros((num_samples, num_classes))
32
33     for k in range(num_classes):
34         B_k = B_liste[k]
35         b_k = b_liste[k]
36         # Calculer les scores/probabilites pour le classifieur k
37         # Utiliser la fonction f globale (ou la rendre accessible)
38         # S'assurer que f utilise X_test comme data
39         probabilities[:, k] = f(B_k, b_k, X_data=X_test) # Specifier X_data
40
41     # Trouver la classe avec la probabilite maximale pour chaque echantillon
42     predictions = np.argmax(probabilities, axis=1)
43
44     return predictions
45
46 # Entrainement OvR sur Iris
47 num_classes_iris = len(np.unique(y_iris))

```

```

48 B_liste_iris, b_liste_iris = OVA_train(X_train_iris, y_train_iris,
    num_classes_iris)
49
50 # Predictions sur le test set Iris
51 y_pred_iris = OVA_predict(X_test_iris, B_liste_iris, b_liste_iris)
52
53 # Evaluation
54 accuracy_iris = accuracy_score(y_test_iris, y_pred_iris)
55 print(f"\nPrecision (Accuracy) OvR sur test set Iris: {accuracy_iris:.3f}")
56
57 # Affichage de la matrice de confusion pour multiclass
58 from sklearn.metrics import confusion_matrix
59 cm_iris = confusion_matrix(y_test_iris, y_pred_iris)
60 print("\nMatrice de confusion OvR (Iris):")
61 print(cm_iris)
62
63 # Visualisation (exemple pour les 2 premières features)
64 # plt.figure(figsize=(8, 6))
65 # plt.scatter(X_test_iris[:, 0], X_test_iris[:, 1], c=y_pred_iris, cmap='viridis',
66 #             marker='x', label='Predictions OvR')
67 # plt.scatter(X_train_iris[:, 0], X_train_iris[:, 1], c=y_train_iris, cmap='viridis',
68 #             marker='o', alpha=0.3, label='Entrainement')
69 # plt.xlabel(iris.feature_names[0])
70 # plt.ylabel(iris.feature_names[1])
71 # plt.title('Predictions OvR vs Entrainement (Iris - 2 premières features)')
72 # plt.legend()
73 # plt.show()

```

Listing 14 – Entraînement et prédiction OvR

13.3 Avis sur OvR

Avantages :

- Conceptuellement simple et facile à implémenter si un classifieur binaire existe.
- Modulaire : chaque classifieur est indépendant.

Inconvénients :

- **Lent :** Nécessite l'entraînement de K modèles, ce qui peut être long si K est grand ou si l'entraînement d'un modèle est coûteux.
- Peut souffrir de déséquilibre de classes lors de l'entraînement de chaque classifieur binaire (une classe contre toutes les autres).
- Les probabilités obtenues ne sont pas calibrées (leur somme pour une instance n'est pas nécessairement 1).

Conclusion : La méthode OvR est une approche simple pour la classification multiclass, mais peut être coûteuse en calcul. D'autres méthodes (comme Softmax ou OvO) existent.

A Sources utilisées

1. [Google ML Crash Course - Logistic Regression \(FR\)](#)
2. [DATatab - Régression logistique](#)
3. [SciML Lectures - INRIA \(chapitre sur la régression logistique\)](#) (Source la plus importante)
4. [GRETIS 2019 - Article PDF par Anas Barakat](#)
5. [Dremio - One-vs-all Classification](#)