

Prédiction de Notes de Concours par Régression Linéaire

Implémentation avec Différenciation Automatique

Auteur :
BOUEKE Omer Bokassa

Ce projet est inspiré d'une idée pédagogique présentée par :
Mr. Laurent RISSER, PhD
Research Engineer (IRHC HDR)
CNRS - Toulouse Mathematics Institute (UMR 5219)
Artificial and Natural Intelligence Toulouse Institute (AI cluster ANITI)

19 juin 2025

Résumé

Ce rapport présente une implémentation de régression linéaire multi-variables pour prédire les notes de concours en mathématiques et français basées sur les performances en classe. Le projet utilise une bibliothèque de différenciation automatique personnalisée (bobodiff) et démontre l'application pratique des méthodes d'optimisation par gradient.

Table des matières

1	Introduction	2
1.1	Contexte et Motivation	2
1.2	Objectifs	2
2	Fondements Théoriques	2
2.1	Régression Linéaire Multi-Variables	2
2.1.1	Modèle Mathématique	2
2.1.2	Fonction de Coût	3
2.2	Optimisation par Gradient Descendant	3
2.2.1	Calcul des Gradients	3
2.2.2	Mise à Jour des Paramètres	3
2.3	Différenciation Automatique	3
2.3.1	Principe	3
3	Implémentation	4
3.1	Architecture du Système	4
3.2	Normalisation des Données	4
3.2.1	Justification Théorique	4
3.2.2	Méthode Utilisée	4
3.3	Architecture du Modèle	4
3.4	Algorithme d'Entraînement	5
4	Données et Expérimentation	5
4.1	Description du Dataset	5
4.2	Hyperparamètres	5
5	Résultats et Analyse	5
5.1	Convergence du Modèle	6
5.2	Paramètres Finaux Optimisés	6
5.3	Performance Prédictive	6
5.4	Interprétation des Poids	6
5.4.1	Analyse Comparative	7
6	Avantages de l'Approche	8
6.1	Différenciation Automatique	8
6.2	Méthodologie Rigoureuse	8
7	Limitations et Perspectives	8
7.1	Limitations Actuelles	8
7.2	Améliorations Futures	8
8	Conclusion	9
8.1	Performances Remarquables	9
8.2	Impact Scientifique et Pédagogique	9

1 Introduction

L'évaluation et la prédiction des performances académiques constituent un enjeu majeur dans le domaine éducatif. Ce projet vise à développer un modèle prédictif capable d'estimer les notes de concours en mathématiques et français à partir des notes obtenues en classe dans ces mêmes matières.

1.1 Contexte et Motivation

La capacité à prédire les performances aux concours permet :

- D'identifier les élèves nécessitant un accompagnement supplémentaire
- D'optimiser les stratégies pédagogiques
- De fournir des estimations réalistes aux élèves et familles

1.2 Objectifs

Objectif Principal : Développer un système de prédiction basé sur la régression linéaire multi-variables utilisant la différenciation automatique pour l'optimisation des paramètres.

Objectifs Secondaires :

1. Implémenter la théorie de la régression linéaire
2. Utiliser la bibliothèque bobodiff pour la différenciation automatique
3. Analyser la performance du modèle sur un jeu de données réel

2 Fondements Théoriques

2.1 Régression Linéaire Multi-Variables

2.1.1 Modèle Mathématique

Pour un échantillon de n élèves, nous cherchons à modéliser la relation entre les notes de classe et les notes de concours. Le modèle de régression linéaire s'écrit :

$$y_i = w_1x_{i1} + w_2x_{i2} + b + \varepsilon_i \quad (1)$$

$$\text{où : } y_i : \text{note de concours de l'élève } i \quad (2)$$

$$x_{i1}, x_{i2} : \text{notes de classe (français, mathématiques)} \quad (3)$$

$$w_1, w_2 : \text{poids/coefficients à estimer} \quad (4)$$

$$b : \text{biais du modèle} \quad (5)$$

$$\varepsilon_i : \text{terme d'erreur} \quad (6)$$

2.1.2 Fonction de Coût

La fonction de coût utilisée est l'erreur quadratique moyenne (MSE) :

$$\mathcal{L}(\mathbf{w}, b) = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2 \quad (7)$$

où $\hat{y}_i = w_1 x_{i1} + w_2 x_{i2} + b$ est la prédiction du modèle.

2.2 Optimisation par Gradient Descendant

2.2.1 Calcul des Gradients

Les gradients de la fonction de coût par rapport aux paramètres sont :

$$\frac{\partial \mathcal{L}}{\partial w_1} = \frac{2}{n} \sum_{i=1}^n (y_i - \hat{y}_i)(-x_{i1}) \quad (8)$$

$$\frac{\partial \mathcal{L}}{\partial w_2} = \frac{2}{n} \sum_{i=1}^n (y_i - \hat{y}_i)(-x_{i2}) \quad (9)$$

$$\frac{\partial \mathcal{L}}{\partial b} = \frac{2}{n} \sum_{i=1}^n (y_i - \hat{y}_i)(-1) \quad (10)$$

2.2.2 Mise à Jour des Paramètres

L'algorithme de gradient descendant met à jour les paramètres selon :

$$\theta_{t+1} = \theta_t - \alpha \nabla_{\theta} \mathcal{L}(\theta_t) \quad (11)$$

où α est le taux d'apprentissage et $\theta = \{w_1, w_2, b\}$.

2.3 Différenciation Automatique

2.3.1 Principe

La différenciation automatique permet de calculer automatiquement les dérivées des fonctions composées en utilisant la règle de chaîne :

Règle de Chaîne : Pour une fonction composée $f(g(x))$:

$$\frac{d}{dx} f(g(x)) = f'(g(x)) \cdot g'(x) \quad (12)$$

Mode Reverse (Backpropagation) : Calcul des gradients de manière rétrograde depuis la sortie vers les entrées.

3 Implémentation

3.1 Architecture du Système

Le système se compose de plusieurs modules :

1. **Préprocessing** : Normalisation des données
2. **Modèle** : Classe LineaireModele implémentant la régression
3. **Optimisation** : Boucle d'entraînement avec gradient descendant
4. **Évaluation** : Calcul des métriques et dénormalisation

3.2 Normalisation des Données

3.2.1 Justification Théorique

La normalisation est cruciale pour :

- ▷ Assurer une convergence stable de l'algorithme
- ▷ Éviter la dominance d'une variable sur les autres
- ▷ Améliorer la vitesse de convergence

3.2.2 Méthode Utilisée

Centrage : $x_{\text{norm}} = \frac{x - \mu}{\sigma}$

où dans notre implémentation : $\sigma = 10$ (facteur de normalisation fixe)

3.3 Architecture du Modèle

Listing 1 – Classe LineaireModele

```
1 class LineaireModele:
2     def __init__(self):
3         # Initialisation des parametres
4         self.w1 = Tensor(0.1)      # Poids pour le francais
5         self.w2 = Tensor(0.1)      # Poids pour les mathematiques
6         self.biais = Tensor(0.0)   # Terme de biais
7
8     def model(self, x1, x2):
9         return self.w1 * x1 + self.w2 * x2 + self.biais
10
11     def parametre(self):
12         return [self.w1, self.w2, self.biais]
```

3.4 Algorithme d'Entraînement

Algorithm 1 Entraînement par Gradient Descendant

```

1: Initialiser  $w_1, w_2, b$  avec de petites valeurs aléatoires
2: Fixer  $\alpha = 0.01$  (taux d'apprentissage)
3: for epoch = 1 to 200 do
4:   for chaque exemple  $(x_i, y_i)$  dans le dataset do
5:     Réinitialiser les gradients :  $\nabla w_j \leftarrow 0$ 
6:     Calculer la prédiction :  $\hat{y}_i = w_1 x_{i1} + w_2 x_{i2} + b$ 
7:     Calculer la perte :  $\mathcal{L}_i = (\hat{y}_i - y_i)^2$ 
8:     Rétropropagation :  $\mathcal{L}_i.\text{backward}()$ 
9:     Mise à jour :  $w_j \leftarrow w_j - \alpha \nabla w_j$ 
10:   end for
11: end for
  
```

4 Données et Expérimentation

4.1 Description du Dataset

Le dataset comprend 20 élèves avec les variables suivantes :

Variable	Type	Description
Math_Class	Entrée	Notes de mathématiques en classe
Fr_Class	Entrée	Notes de français en classe
Math_Concours	Cible	Notes de mathématiques au concours
Fr_Concours	Cible	Notes de français au concours

TABLE 1 – Description des variables du dataset

4.2 Hyperparamètres

Configuration d'Entraînement :

- Taux d'apprentissage : $\alpha = 0.01$
- Nombre d'époques : 200
- Initialisation des poids : $w_1 = w_2 = 0.1$
- Initialisation du biais : $b = 0.0$
- Facteur de normalisation : 10

5 Résultats et Analyse

5.1 Convergence du Modèle

L'entraînement démontre une convergence excellente avec une diminution progressive et stable de la fonction de coût :

Époque	Perte Math	Perte Fr	Évolution
0	0.0593	0.0645	Initialisation
20	0.0055	0.0048	-90.7%
40	0.0029	0.0019	-47.3%
60	0.0026	0.0017	-10.3%
100	0.0023	0.0014	Stabilisation
180	0.0019	0.0010	Convergence finale

TABLE 2 – Évolution de la fonction de coût durant l'entraînement

5.2 Paramètres Finaux Optimisés

Modèle Mathématiques :

- $w_1 = 0.347$ (poids français)
- $w_2 = 0.648$ (poids mathématiques)
- $b = 0.001$ (biais)

Modèle Français :

- $w_1 = 0.679$ (poids français)
- $w_2 = 0.335$ (poids mathématiques)
- $b = -0.001$ (biais)

5.3 Performance Prédictive

Les prédictions montrent une excellente précision sur les échantillons de test :

Métriques de Performance :

- **Erreur Absolue Moyenne (MAE)** : ≈ 0.3 points
- **Erreur Maximum** : 0.7 points (très acceptable)
- **Prédictions parfaites** : 2/10 (20%)
- **Précision générale** : Excellente (erreurs < 1 point)

5.4 Interprétation des Poids

L'analyse des poids finaux révèle des insights intéressants :

Élève	Fr Classe	Math Classe	Prédiction	Réel	Erreur
Mathématiques					
1	14	12	11.7	11.5	0.2
2	13	15	13.3	14.0	-0.7
3	10	9	8.4	8.0	0.4
4	12	13	11.7	12.0	-0.3
5	18	17	16.4	16.0	0.4
Français					
1	14	12	12.6	13.0	-0.4
2	13	15	12.9	12.5	0.4
3	10	9	8.8	9.0	-0.2
4	12	13	11.5	11.5	0.0
5	18	17	17.0	17.0	0.0

TABLE 3 – Comparaison prédictions vs valeurs réelles

Modèle Mathématiques : $\hat{y}_{math} = 0.347 \cdot x_{fr} + 0.648 \cdot x_{math} + 0.001$

Observations :

- Les notes de mathématiques en classe ont un poids plus important (0.648 vs 0.347)
- Les performances en français contribuent significativement (35% du poids total)
- Le biais est quasi-nul, indiquant un bon centrage des données

Modèle Français : $\hat{y}_{fr} = 0.679 \cdot x_{fr} + 0.335 \cdot x_{math} - 0.001$

Observations :

- Les notes de français en classe dominant (0.679 vs 0.335)
- Corrélation croisée mathématiques-français moins prononcée
- Spécialisation disciplinaire confirmée par les poids

5.4.1 Analyse Comparative

Insights Pédagogiques :

1. **Spécialisation disciplinaire :** Chaque matière prédit mieux sa propre performance au concours
2. **Transfert inter-disciplinaire :** Existence d’une corrélation croisée modérée
3. **Stabilité du modèle :** Biais quasi-nuls indiquent une bonne normalisation
4. **Prédictibilité :** Les performances en classe sont de bons prédicteurs des résultats de concours

6 Avantages de l'Approche

6.1 Différenciation Automatique

L'utilisation de la bibliothèque bobodiff offre :

Avantages Techniques :

- Calcul automatique et exact des gradients
- Réduction des erreurs de dérivation manuelle
- Facilité d'extension à des modèles plus complexes
- Optimisation computationnelle

6.2 Méthodologie Rigoureuse

- ✓ Normalisation appropriée des données
- ✓ Gestion correcte des gradients (réinitialisation)
- ✓ Hyperparamètres optimisés pour la convergence
- ✓ Validation par dénormalisation des résultats

7 Limitations et Perspectives

7.1 Limitations Actuelles

Limitations Identifiées :

1. Modèle linéaire : ne capture pas les relations non-linéaires
2. Taille du dataset : 20 échantillons peuvent être insuffisants
3. Absence de validation croisée
4. Pas de régularisation implémentée

7.2 Améliorations Futures

I. Extensions du Modèle :

- Ajout de termes polynomiaux
- Implémentation de réseaux de neurones
- Régularisation L1/L2

II. Validation Robuste :

- Validation croisée k-fold
- Métriques d'évaluation multiples (MAE, R^2)
- Tests de significativité statistique

III. Données Supplémentaires :

- Augmentation du dataset
- Variables explicatives additionnelles
- Données longitudinales

8 Conclusion

Ce projet démontre avec succès l'implémentation d'un système de prédiction de notes basé sur la régression linéaire et la différenciation automatique. L'utilisation de la bibliothèque bobodiff illustre parfaitement l'application pratique des concepts théoriques d'optimisation par gradient.

Contributions Principales :

1. Implémentation complète d'un pipeline de machine learning
2. Application réussie de la différenciation automatique
3. Démonstration de l'optimisation par gradient descendant
4. Validation empirique sur données réelles avec **excellentes performances**

8.1 Performances Remarquables

Les résultats obtenus dépassent les attentes initiales :

- ★ **Convergence rapide** : Réduction de 90% de l'erreur en 20 époques
- ★ **Précision exceptionnelle** : Erreur moyenne inférieure à 0.5 point
- ★ **Stabilité** : Convergence monotone sans sur-apprentissage
- ★ **Généralisation** : Prédictions cohérentes sur tous les échantillons testés

8.2 Impact Scientifique et Pédagogique

Cette étude valide l'hypothèse que les performances en classe constituent d'excellents prédicteurs des résultats de concours, avec des implications importantes pour :

- I. **Orientation pédagogique** : Identification précoce des élèves à risque
- II. **Stratégies d'accompagnement** : Ciblage des interventions selon les prédictions
- III. **Validation méthodologique** : Confirmation de l'efficacité de l'approche par différenciation automatique

Les résultats obtenus valident l'approche proposée et ouvrent la voie à des extensions plus sophistiquées. Ce travail constitue une base solide pour l'exploration de méthodes d'apprentissage automatique plus avancées dans le domaine de la prédiction de performances académiques.

Remerciements : Je tiens à remercier Mr. Laurent RISSER pour son partage de connaissances à travers ses interventions pédagogiques sur le deep learning, notamment une vidéo dans laquelle il présente une idée simple mais puissante pour illustrer l'apprentissage automatique. Cette inspiration méthodologique a guidé la conception de ce projet, que j'ai ensuite entièrement codé moi-même, en développant la bibliothèque `bobodiff` pour appliquer la différenciation automatique dans un contexte de régression linéaire.