

MAXIME BRONNY
19009314

TECHNIQUES D'APPRENTISSAGE ARTIFICIEL

- RAPPORT EXPLICATIF -

**PRÉDICTION DU RISQUE DE CRÉDIT BANCAIRE PAR RÉGRESSION LOGISTIQUE
ET ARBRE DE DÉCISION CART IMPLÉMENTÉS FROM SCRATCH EN C**

PRÉSENTATION

Ce projet consiste à développer un système de prédiction du risque de crédit bancaire en implémentant intégralement deux méthodes d'apprentissage en langage C : une Régression Logistique et un Arbre de Décision CART. Le dataset utilisé provient de Kaggle et contient 32 581 emprunteurs décrits par 12 variables (8 numériques et 4 catégorielles).

Un pipeline complet de machine learning a été réalisé :

- chargement et structuration des données ;
- encodage des variables catégorielles ;
- prétraitement et normalisation ;
- entraînement des deux modèles (gradient descent pour la régression, partitionnement récursif pour l'arbre) ;
- évaluation comparative sur un ensemble de test.

La Régression Logistique atteint 81,17 % d'accuracy avec un AUC-ROC de 81,70 %, tandis que l'Arbre de Décision CART obtient 92,68 % d'accuracy avec un AUC-ROC de 90,42 %, démontrant la supériorité des modèles non-linéaires sur ce problème. Le temps d'exécution total reste inférieur à 0,5 s, montrant l'efficacité des implémentations optimisées en C. Les performances des deux modèles ont été validées par comparaison avec leurs équivalents scikit-learn, avec des écarts inférieurs à 1 % sur l'ensemble des métriques.

Ce travail démontre l'intérêt d'une implémentation bas niveau pour comprendre finement les algorithmes de machine learning, comparer rigoureusement des approches complémentaires (linéaire vs non-linéaire), et optimiser leur exécution.

[HTTPS://GITHUB.COM/CRYZIIIX
/CREDIT-SCORING-C-IMPLEMENTATION](https://github.com/cryziiix/CREDIT-SCORING-C-IMPLEMENTATION)

Liste des Figures

Figure 1 : Caractéristiques générales du dataset	16
Figure 2 : Description détaillée des variables	17
Figure 3 : Statistiques descriptives complètes	18
Figure 4 : Détection des outliers	19
Figure 5 : Matrice de corrélation de Pearson entre les variables numériques	20
Figure 6 : Distribution des variables catégorielles	20
Figure 7 : Taux de défaut par catégorie	21
Figure 8 : Valeurs manquantes - Analyse détaillée	21
Figure 9 : Configuration matérielle et logicielle	27
Figure 10 : Hyperparamètres choisis	29
Figure 11 : Architecture du projet globale	30
Figure 12 : Complexité temporelle par composant	45
Figure 13 : Complexité spatiale	45
Figure 14 : Options d'Implémentations	46
Figure 15 : Métriques de performance (source : exécution du programme)	51
Figure 16 : Matrice de confusion (ensemble de test)	51
Figure 17 : Métriques de performance Arbre de Décision (source : exécution du programme)	52
Figure 18 : Matrice de confusion Arbre de Décision (ensemble de test)	52
Figure 19 : Comparaison C vs Scikit-learn (source : script de validation)	53
Figure 20 : Comparaison Arbre de Décision C vs Scikit-learn (source : script de validation)	54
Figure 18 : Tableau d'analyse des FP	55
Figure 19 : Tableau d'analyse des FN - Exemples représentatifs	56
Figure 23 : Tableau d'analyse du seuil optimal	57
Figure 24 : Validation par corrélation avec la target :	58
Figure 25 : Tableau d'analyse des variations du learning rate (Expérience 1)	59
Figure 26 : Tableau d'analyse des variations du nombre d'itérations (Expérience 2)	59
Figure 27 : Tableau d'analyse du ratio Train/Test (Expérience 3)	60
Figure 28 : Benchmarks temporels (source : mesures avec time)	60
Figure 29 : Évolution de la fonction de coût (source : logs d'entraînement) & courbe de convergence (source : visualisation des résultats)	67
Figure 30 : Pipeline complet du système (source : architecture du projet) :	84

Glossaire

1. Données et prétraitement

- **Dataset** : Ensemble de données utilisé pour entraîner, valider ou tester un modèle. Peut être séparé en train/test ou train/validation/test.
 - **Feature (variable explicative)** : Variable utilisée comme entrée du modèle, peut être numérique ou catégorielle.
 - **Label (variable cible)** : Valeur que le modèle doit prédire (ex : 0 ou 1 en classification binaire).
 - **Encoding (encodage)** : Transformation des variables catégorielles en valeurs numériques (One-Hot, Ordinal...).
 - **Imputation** : Remplacement des valeurs manquantes par une valeur calculée (moyenne, médiane, modèle simple).
 - **MCAR** : "Missing Completely At Random". Valeurs manquantes apparaissant totalement au hasard, sans corrélation avec d'autres variables.
 - **Scaler (normalisation)** : Transformation des features pour les mettre sur une même échelle (Z-score, MinMax...).
 - **Class Imbalance** : Déséquilibre entre les classes d'un dataset (ex : 90 % de classe 0). Affecte fortement l'évaluation des modèles.
 - **Train/Test Split** : Séparation du dataset entre données d'entraînement et données de test afin d'évaluer la généralisation.
-

2. Modèle et apprentissage

- **Logistic Regression** : Modèle linéaire pour la classification binaire. Combine les features via $w^T x + b$, puis applique une sigmoïde pour produire une probabilité.
 - **Weights (poids)** : Coefficients w appris par optimisation. Indiquent l'importance de chaque feature dans la décision.
 - **Bias (biais du modèle)** : Terme b dans l'équation " $w^T x + b$ " ; décale la frontière de décision.
 - **Sigmoid (σ)** : Fonction " $\sigma(z) = 1/(1+e^{-z})$ " convertissant un score brut en probabilité $\in [0,1]$.
 - **Threshold (seuil)** : Valeur qui convertit la probabilité en classe (par défaut 0,5). Peut être ajustée selon les objectifs.
 - **Batch** : Sous-ensemble d'échantillons utilisé pour une mise à jour des poids.
 - **Full Batch** : Utilisation de la totalité du dataset d'entraînement pour chaque mise à jour du gradient.
 - **Epoch** : Parcours complet de l'ensemble d'entraînement (toutes les données vues une fois).
 - **Gradient Descent** : Algorithme d'optimisation ajustant les poids pour minimiser la fonction de coût.
 - **Learning Rate (α)** : Taux d'apprentissage ; contrôle la taille des mises à jour des poids.
 - **Hyperparamètre** : Paramètre réglé avant l'entraînement (learning rate, nombre d'itérations, taille du batch...).
 - **Regularisation (L1/L2)** : Méthode pénalisant les poids trop élevés pour réduire l'overfitting.
-

3. Fonctions de coût et optimisation

- **Cross-Entropy (log-loss)** : Fonction de coût mesurant l'écart entre la probabilité prédictée par le modèle et le label réel.
 - **Baseline** : Modèle de référence simple servant de comparaison (ex : prédire la classe majoritaire).
 - **Overfitting** : Situation où le modèle mémorise trop le train et généralise mal.
 - **Confusion Matrix** : Tableau résumant les prédictions : TP, FP, FN et TN.
(Même si c'est aussi utilisé en évaluation, c'est ici car c'est central pour analyser les erreurs du modèle.)
-

4. Évaluation du modèle

- **Accuracy** : Proportion de prédictions correctes sur l'ensemble du dataset.
 - **Precision** : Parmi les prédictions positives, proportion réellement positives.
 - **Recall (Sensibilité)** : Proportion de cas positifs correctement détectés.
 - **F1-Score** : Moyenne harmonique entre la précision et le rappel ; utile avec des classes déséquilibrées.
 - **AUC-ROC** : Aire sous la courbe ROC. Mesure la capacité du modèle à classer correctement en comparant les taux de vrais positifs et faux positifs pour tous les seuils possibles.
 - **TP (True Positive)** : Cas positif correctement prédit.
 - **TN (True Negative)** : Cas négatif correctement prédit.
 - **FP (False Positive)** : Cas négatif prédict comme positif (erreur de type I).
 - **FN (False Negative)** : Cas positif prédict comme négatif (erreur de type II).
-

5. Termes d'analyse et bonnes pratiques

- **Class Imbalance** : Déséquilibre entre les classes, influençant fortement accuracy, seuils et choix des métriques.
 - **Threshold Tuning** : Ajustement du seuil de décision pour optimiser une métrique donnée (ex : maximiser le F1-score).
 - **Calibration** : Vérification que les probabilités prédictes reflètent réellement la fréquence empirique des classes.
 - **Validation croisée (Cross-Validation)** (optionnel si tu veux l'ajouter) : Technique d'évaluation robuste consistant à entraîner et tester le modèle sur plusieurs partitions du dataset.
-

Table des Matières

<u>Liste des Figures</u>	2
<u>Glossaire</u>	3
<u>1. Données et prétraitement</u>	3
<u>2. Modèle et apprentissage</u>	3
<u>3. Fonctions de coût et optimisation</u>	4
<u>4. Évaluation du modèle</u>	4
<u>5. Termes d'analyse et bonnes pratiques</u>	4
<u>Table des Matières</u>	5
<u>1. Introduction</u>	8
<u>1.1 Contexte et Problématique</u>	8
<u>1.2 Objectifs du Projet</u>	8
<u>1.3 Contributions</u>	9
<u>1.4 Organisation du rapport</u>	9
<u>1.5 Démarche Scientifique</u>	10
<u>2. État de l'Art</u>	11
<u>2.1 Prédiction du Risque de Crédit</u>	11
<u>2.2 Régression Logistique</u>	11
<u>2.2.5 Arbres de Décision et CART</u>	14
<u>2.3 Analyse Critique de la Littérature</u>	15
<u>2.4 Implémentations en Langage C</u>	15
<u>3. Méthodologie</u>	16
<u>3.1 Dataset</u>	16
<u>3.2 Analyse Exploratoire des Données (EDA)</u>	18
<u>3.3 Pipeline de Prétraitement</u>	22
<u>3.4 Justification du Choix du Modèle</u>	25
<u>3.4.1 Méthode 1 : Régression Logistique :</u>	25
<u>3.4.2 Méthode 2 : Arbre de Décision CART</u>	26
<u>3.5 Protocole Expérimental Détaillé</u>	27
<u>3.6 Architecture du Modèle</u>	28
<u>4. Implémentation</u>	30
<u>4.1 Architecture Logicielle</u>	30
<u>4.2 Composants Clés avec Code Source</u>	31
<u>4.2.1 CSV Parser avec Encodage Catégoriel Intégré</u>	31
<u>4.2.2 StandardScaler : Normalisation des Features</u>	34
<u>4.2.3 Régression Logistique : Cœur de l'Algorithme</u>	35
<u>4.2.4 AUC-ROC : Métrique de Discrimination</u>	36

<u>4.2.5 Arbre de Décision CART : Partitionnement Récuratif</u>	40
<u>4.3 Analyse de Complexité Algorithmique</u>	45
<u>4.4 Décisions d'Implémentation et Trade-offs</u>	46
<u>4.5 Gestion de la Mémoire</u>	48
<u>4.6 Tests Unitaires</u>	50
<u>4.7 Environnement de Développement et Reproductibilité</u>	51
5. Résultats et Analyses	51
<u>5.1 Performance du Modèle sur l'Ensemble de Test</u>	51
<u>5.1.2 Performance de l'Arbre de Décision</u>	52
<u>5.2 Analyse Détailée des Résultats</u>	53
<u>5.3 Comparaison avec Scikit-learn</u>	53
<u>5.4 Analyse des Erreurs de Classification</u>	55
<u>5.5 Importance des Features</u>	57
<u>5.6 Analyse de Sensibilité des Hyperparamètres</u>	59
<u>5.7 Performance Computationnelle</u>	60
<u>5.8 Evaluation par AUC-ROC</u>	61
<u>5.9 Comparaison des Deux Méthodes d'Apprentissage</u>	63
<u>5.9.1 Tableau Comparatif Global</u>	63
<u>5.9.2 Analyse des Différences</u>	63
<u>5.9.3 Analyse des Matrices de Confusion</u>	64
<u>5.9.4 Analyse des Courbes ROC</u>	65
<u>5.9.5 Performance Computationnelle</u>	65
<u>5.9.6 Recommandations</u>	66
<u>5.10 Convergence du Modèle</u>	66
6. Conclusion et Perspectives	68
<u>6.1 Synthèse du Projet</u>	68
<u>6.1.1 Objectifs techniques accomplis</u>	68
<u>6.1.2 Apprentissages principaux</u>	69
<u>6.2 Limitations Identifiées</u>	69
<u>6.3 Perspectives d'Amélioration</u>	70
<u>6.4 Applications Pratiques</u>	70
<u>6.5 Réflexions Finales</u>	71
<u>6.6 Retour sur les Hypothèses de Recherche</u>	71
<u>6.7 Réflexion Critique et Auto-évaluation</u>	72
<u>6.8 Considérations Éthiques et Réglementaires</u>	73
7. Bibliographie	75
<u>Ouvrages de référence sur le crédit scoring et les modèles</u>	75
<u>Régression logistique, machine learning classique et statistiques</u>	75
<u>Ensembles, random forests, boosting</u>	75
<u>Apprentissage profond (pour comparaison théorique)</u>	76
<u>ROC et AUC</u>	76
<u>Arbres de Décision et CART</u>	76
<u>Langages et implémentation bas niveau (C / Scala)</u>	76
<u>Cours en ligne et frameworks utilisés</u>	77
<u>Données</u>	77

8. Annexes

78

<u>Annexe A : Structure Complète du Code Source</u>	78
<u>Annexe B : Résultats Détaillés des Tests Unitaires</u>	79
<u>Annexe C : Commandes de Compilation et d'Exécution</u>	80
<u>Annexe D : Exemple d'Utilisation de l'API</u>	82
<u>Annexe E : Figures et Tableaux Récapitulatifs</u>	84
<u>Annexe F : categorical_distributions.png</u>	85
<u>Annexe G : correlation_matrix.png</u>	86
<u>Annexe H : numerical_distributions.png</u>	87
<u>Annexe I : features_by_target.png</u>	88
<u>Annexe J : target_distribution.png</u>	89

1. Introduction

1.1 Contexte et Problématique

Le risque de crédit constitue aujourd'hui un enjeu central pour l'ensemble du secteur financier. D'après les estimations du Fonds Monétaire International, les défauts de paiement entraînent chaque année des pertes se chiffrant en centaines de milliards de dollars, compromettant la stabilité des institutions bancaires. Dans ce contexte, disposer d'outils fiables permettant d'estimer la probabilité de défaut d'un emprunteur est devenu essentiel pour les acteurs du crédit.

Le credit scoring s'inscrit dans un cadre réglementaire exigeant, notamment défini par les accords de Bâle III, qui imposent aux banques de maintenir des niveaux élevés de solvabilité et de s'appuyer sur des modèles rigoureux d'évaluation des risques. Une prédiction plus précise du risque de défaut permet de réduire les pertes financières, d'optimiser l'allocation du capital, de répondre aux contraintes réglementaires et, plus globalement, de favoriser l'accès au crédit pour les profils considérés comme peu risqués.

Problématique centrale :

Comment développer un système efficace permettant de prédire si un emprunteur va faire défaut sur son prêt, en se basant uniquement sur ses caractéristiques personnelles et financières, tout en garantissant une compréhension profonde du modèle et des performances computationnelles élevées ?

Cette problématique soulève plusieurs défis techniques :

- Le traitement de variables catégorielles dans un contexte de machine learning en C
- La gestion du déséquilibre des classes (défauts rares vs non-défauts fréquents)
- L'optimisation d'algorithmes mathématiques sans bibliothèques haut niveau
- La validation de l'implémentation face à des standards reconnus

1.2 Objectifs du Projet

Objectif principal :

Développer un système complet de classification binaire (défaut / pas de défaut) en implémentant from scratch deux méthodes d'apprentissage en langage C pur : une Régression Logistique et un Arbre de Décision CART, puis les comparer avec leurs équivalents scikit-learn, sans utiliser de bibliothèques de machine learning existantes.

Objectifs spécifiques :

- **Ensemble de méthodes** : Implémenter deux approches complémentaires (linéaire vs non-linéaire) pour analyser la séparabilité du problème et comparer leurs performances.
- **Implémentation algorithmique** : Coder l'algorithme de régression logistique avec gradient descent en C, incluant la fonction sigmoïde, le calcul de la cross-entropy loss et la mise à jour des poids.
- **Gestion des données catégorielles** : Développer un système d'encodage pour transformer les 4 variables catégorielles du dataset (type de logement, objectif du prêt, grade du crédit, historique de défaut) en valeurs numériques exploitables.
- **Pipeline de prétraitement** : Construire une chaîne complète de traitement incluant le chargement CSV, l'imputation des valeurs manquantes, la normalisation par StandardScaler et le split train/test.
- **Performance et validation** : Atteindre des performances comparables aux implémentations standard (scikit-learn) tout en optimisant le temps d'exécution pour traiter des datasets de taille moyenne (30 000+ lignes) en moins de 5 secondes.
- **Architecture logicielle** : Développer une architecture modulaire, testable et maintenable, avec une gestion rigoureuse de la mémoire et une suite de tests unitaires.

1.3 Contributions

Ce projet apporte les contributions suivantes :

Contribution technique :

- Une implémentation complète en C d'un système de machine learning, intégrant l'ensemble des étapes nécessaires, du chargement des données jusqu'à l'évaluation finale du modèle.
- Deux implémentations complètes d'algorithmes d'apprentissage : Régression Logistique (optimisation par gradient) et Arbre de Décision CART (partitionnement récursif), permettant une analyse comparative rigoureuse.
- Un parseur CSV optimisé capable d'effectuer l'encodage des variables catégorielles directement lors de la lecture, ce qui évite une passe supplémentaire sur les données et réduit le coût de prétraitement.
- Une architecture modulaire structurée selon le principe de séparation des responsabilités, facilitant la maintenance, la lisibilité et la réutilisation des composants du projet.

Contribution méthodologique :

- Une stratégie de traitement robuste des variables catégorielles propres au domaine bancaire, reposant sur des mappings adaptés : encodage ordinal pour le grade de crédit et encodage nominal pour les autres variables.
- Une validation systématique des deux implémentations grâce à une comparaison méthodique avec leurs équivalents scikit-learn (LogisticRegression et DecisionTreeClassifier).
- Une suite de tests unitaires couvrant environ 80 % des fonctions publiques, garantissant la fiabilité et la stabilité des composants développés.

Contribution pédagogique :

- Une documentation détaillée du code source ainsi que de l'architecture logicielle, permettant de comprendre clairement les choix d'implémentation.
- Une démonstration concrète de la mise en œuvre d'algorithmes de machine learning sans recours aux abstractions haut niveau habituellement proposées par les frameworks.
- Une illustration tangible des gains de performance obtenus, avec un facteur d'accélération d'environ ×7 par rapport à une implémentation équivalente en Python.

1.4 Organisation du rapport

Ce rapport s'organise de la manière suivante :

- Le chapitre 2 présente un état de l'art sur la prédiction du risque de crédit, rappelle les fondements mathématiques de la régression logistique et des arbres de décision, et expose les motivations ayant conduit au choix d'une implémentation en langage C.
- Le chapitre 3 décrit la méthodologie adoptée : présentation du dataset, analyse exploratoire, pipeline de prétraitement et justification du choix des deux modèles.
- Le chapitre 4 détaille l'implémentation technique des deux modèles, en particulier l'architecture logicielle modulaire et le code des composants essentiels de la régression logistique et de l'arbre de décision.
- Le chapitre 5 expose les résultats expérimentaux des deux méthodes, leur analyse statistique, la comparaison avec scikit-learn, la comparaison entre elles, et l'analyse approfondie de leurs différences.
- Le chapitre 6 propose une synthèse des contributions, discute les limites du travail réalisé et présente plusieurs pistes d'amélioration.
- Enfin, les annexes regroupent des éléments complémentaires concernant l'organisation du code, les résultats détaillés des tests et les commandes de compilation utilisées.

1.5 Démarche Scientifique

Ce projet s'inscrit dans une démarche scientifique structurée.

Hypothèses de recherche :

Les hypothèses de recherche formulées sont les suivantes : une implémentation “from scratch” en C de la régression logistique peut atteindre des performances proches de celles obtenues avec des bibliothèques de référence comme “scikit-learn”.

Les optimisations bas niveau offertes par le langage C doivent permettre un gain significatif, avec un objectif d'un facteur d'accélération d'environ $\times 10$ par rapport à “Python”.

L'intégration de l'encodage des variables catégorielles directement dans le parseur est susceptible d'améliorer l'efficacité globale du pipeline.

Enfin, le déséquilibre marqué entre les classes (81 % / 19 %) laisse anticiper certaines limites inhérentes à un modèle linéaire.

Protocole expérimental :

Phase 1 : Implémentation modulaire des composants (utils, prétraitement, modèle, évaluation)

Phase 2 : Validation unitaire de chaque composant avec une suite de tests

Phase 3 : Entraînement sur le dataset complet avec mesure des performances

Phase 4 : Comparaison rigoureuse avec scikit-learn pour validation

Phase 5 : Analyse approfondie des résultats et des erreurs

Critères de validation :

Fonctionnel : Tous les tests unitaires doivent passer (objectif : 100 %)

Performance : Différence < 10 % avec scikit-learn sur les métriques principales

Vitesse : Temps d'exécution < 5 secondes (objectif initial)

Qualité : Compilation sans warnings avec flags stricts (-Wall -Wextra)

2. État de l'Art

2.1 Prédiction du Risque de Crédit

La prédiction du risque de crédit est un domaine de recherche actif depuis plusieurs décennies. Nous distinguons deux grandes familles d'approches.

Approches traditionnelles :

Les méthodes statistiques classiques ont longtemps dominé le domaine du credit scoring. L'analyse discriminante linéaire (LDA), introduite par Fisher dans les années 1930, a été l'une des premières techniques appliquées à ce problème. Selon Thomas et al. (2002, p. 45).

Le score FICO (Fair Isaac Corporation), développé dans les années 1980, reste aujourd'hui la référence en matière d'évaluation du crédit aux États-Unis. Ce système utilise cinq catégories d'informations : l'historique de paiement (35 %), les montants dus (30 %), la durée de l'historique de crédit (15 %), les nouveaux crédits (10 %) et les types de crédit utilisés (10 %).

Les règles métier expertes, basées sur l'expérience des analystes crédit, ont également été largement utilisées. Ces systèmes à base de règles permettent une grande interprétabilité mais manquent de flexibilité face à de nouveaux patterns.

Approches modernes :

L'émergence du machine learning a révolutionné le domaine du credit scoring. Breiman (2001) a introduit les Random Forests, démontrant leur supériorité sur les méthodes linéaires pour capturer les interactions complexes entre variables. Ces modèles d'ensemble combinent de multiples arbres de décision pour améliorer la robustesse et la généralisation.

Les algorithmes de gradient boosting, comme XGBoost (Chen & Guestrin, 2016), ont récemment démontré des performances exceptionnelles sur de nombreuses compétitions Kaggle liées au risque de crédit.

Les réseaux de neurones profonds pour le credit scoring fonctionne également. Wang et al. (2018) ont montré qu'un réseau à trois couches cachées pouvait améliorer l'accuracy de 3 à 5 % par rapport aux méthodes traditionnelles, au prix d'une perte d'interprétabilité.

2.2 Régression Logistique

Fondements mathématiques :

La régression logistique est un modèle statistique fondamental pour la classification binaire. Comme l'explique Bishop (2006, p. 205), "la régression logistique est un modèle linéaire généralisé qui utilise la fonction logistique pour modéliser la probabilité d'appartenance à une classe".

Définition formelle :

On considère un vecteur de caractéristiques $x \in \mathbb{R}^d$ et une variable cible binaire $y \in \{0, 1\}$. Le modèle de régression logistique vise à estimer la probabilité conditionnelle suivante :

$$P(y = 1 \mid x; w, b) = \sigma(w^T x + b)$$

où σ est la fonction sigmoïde définie par :

$$\sigma(z) = 1 / (1 + \exp(-z))$$

La fonction sigmoïde transforme une valeur réelle z en une probabilité dans l'intervalle $[0, 1]$. Elle possède des propriétés mathématiques intéressantes, notamment sa dérivée qui s'exprime simplement :

$$\sigma'(z) = \sigma(z)(1 - \sigma(z))$$

Fonction de coût :

L'apprentissage des paramètres w et b se fait par maximisation de la vraisemblance, équivalente à la minimisation de la cross-entropy loss (entropie croisée) :

$$L(w, b) = -1/n \sum_{i=1}^n [y_i \log(\hat{y}_i) + (1 - y_i) \log(1 - \hat{y}_i)]$$

où $\hat{y}_i = \sigma(w^T x_i + b)$ est la probabilité prédictive pour l'échantillon i .

Cette fonction de coût est convexe, garantissant l'existence d'un unique minimum global. Selon James et al. (2013, p. 133), "la convexité de la fonction de coût assure que l'algorithme de gradient descent convergera vers l'optimum global, contrairement aux réseaux de neurones qui peuvent se retrouver piégés dans des minima locaux".

Optimisation par Gradient Descent :

Les paramètres sont mis à jour itérativement d'après la règle ci-dessous :

$$w := w - \alpha \partial L / \partial w = w - \alpha \cdot 1/n \sum_{i=1}^n (\hat{y}_i - y_i) x_i$$

$$b := b - \alpha \partial L / \partial b = b - \alpha \cdot 1/n \sum_{i=1}^n (\hat{y}_i - y_i)$$

où α est le learning rate (taux d'apprentissage).

Le gradient est calculé analytiquement grâce à la dérivée de la sigmoïde et à la règle de dérivation en chaîne. L'algorithme converge typiquement en quelques centaines d'itérations pour des datasets de taille moyenne.

Avantages de la régression logistique :

- **Interprétabilité** : Les coefficients " w_i " indiquent l'importance et la direction de l'effet de chaque feature sur la probabilité de défaut.
- **Rapidité** : L'entraînement est linéaire en fonction du nombre d'échantillons et de features.
- **Probabilités calibrées** : Le modèle fournit directement des probabilités, utiles pour l'évaluation du risque.
- **Robustesse** : Moins sujet à "l'overfitting" que des modèles plus complexes comme les réseaux de neurones.
- **Baseline standard** : Constitue une référence pour comparer des modèles plus sophistiqués.

Limites :

- **Linéarité** : Suppose une relation linéaire entre les features et le log-odds, ne capture pas les interactions non-linéaires complexes.
- **Sensibilité à l'échelle** : Nécessite une normalisation des features pour une convergence optimale.
- **Classes déséquilibrées** : Performance limitée sur des datasets très déséquilibrés sans techniques spécifiques (class weights, SMOTE).

Dérivation mathématique complète des gradients :

Pour comprendre en profondeur l'algorithme, dérivons les gradients de la fonction de coût. Soit la fonction de coût :

$$L(w, b) = -1/n \sum_{i=1}^n [y_i \log(\sigma(z_i)) + (1-y_i) \log(1-\sigma(z_i))]$$

où $z_i = w^T x_i + b$.

Étape 1 : Calcul de $\partial L / \partial z_i$:

En utilisant la règle de dérivation en chaîne ci-dessous :

$$\partial L / \partial z_i = -1/n [y_i \cdot 1/\sigma(z_i) \cdot \sigma'(z_i) + (1-y_i) \cdot 1/(1-\sigma(z_i)) \cdot (-\sigma'(z_i))]$$

Sachant que " $\sigma'(z) = \sigma(z)(1-\sigma(z))$ ", on obtient :

$$\begin{aligned}\partial L / \partial z_i &= -1/n [y_i(1-\sigma(z_i)) - (1-y_i)\sigma(z_i)] \\ &= -1/n [y_i - y_i\sigma(z_i) - \sigma(z_i) + y_i\sigma(z_i)] \\ &= -1/n [y_i - \sigma(z_i)] \\ &= 1/n [\sigma(z_i) - y_i]\end{aligned}$$

Étape 2 : Calcul de $\partial L/\partial w_i$:

En appliquant la règle de dérivation en chaîne avec “ $z_i = w^T x_i + b$ ” :

$$\begin{aligned}\partial L/\partial w_i &= \sum_{i=1}^n \partial L/\partial z_i \cdot \partial z_i/\partial w_i \\ &= \sum_{i=1}^n [1/n(\sigma(z_i) - y_i)] \cdot x_i \\ &= 1/n \sum_{i=1}^n (\sigma(z_i) - y_i)x_i\end{aligned}$$

Étape 3 : Calcul de $\partial L/\partial b$:

$$\begin{aligned}\partial L/\partial b &= \sum_{i=1}^n \partial L/\partial z_i \cdot \partial z_i/\partial b \\ &= 1/n \sum_{i=1}^n (\sigma(z_i) - y_i)\end{aligned}$$

Démonstration de la convexité :

La fonction de coût “ $L(w,b)$ ” est convexe car elle est la somme de fonctions log-convexes. Plus formellement, la matrice Hессienne H de L est semi-définie positive :

$$H = \partial^2 L/\partial w^2 = 1/n \sum_{i=1}^n \sigma(z_i)(1-\sigma(z_i))x_i x_i^T$$

Comme “ $0 < \sigma(z) < 1$ ” pour tout z, on a “ $\sigma(z)(1-\sigma(z)) > 0$ ”, donc H est semi-définie positive, ce qui garantit la convexité et l'existence d'un unique minimum global.

2.2.5 Arbres de Décision et CART

Les arbres de décision constituent une famille d'algorithmes d'apprentissage supervisé fondamentalement différente de la régression logistique. L'algorithme CART (Classification And Regression Trees), développé par Breiman et al. (1984), construit un modèle prédictif sous forme d'arbre binaire en partitionnant récursivement l'espace des features.

Principe de fonctionnement :

À chaque nœud interne, l'arbre sélectionne la feature et le seuil qui maximisent la pureté des partitions résultantes, mesurée par des critères comme l'indice de Gini ou l'entropie. Les feuilles de l'arbre représentent les classes prédites.

Complémentarité avec la régression logistique :

Contrairement à la régression logistique qui apprend une frontière de décision linéaire, l'arbre de décision peut capturer des relations non-linéaires complexes et des interactions entre variables, au prix d'une interprétabilité différente et d'un risque d'overfitting plus élevé.

2.3 Analyse Critique de la Littérature

Forces des approches existantes :

Les bibliothèques modernes de machine learning (“scikit-learn”, “TensorFlow”) offrent des implémentations optimisées et validées. Elles bénéficient de :

- Optimiseurs sophistiqués (“L-BFGS”, “Adam”) avec convergence rapide.
- Parallélisation automatique sur CPU/GPU.
- Régularisation intégrée (“L1”, “L2”, “Elastic Net”).
- Documentation extensive et communauté active.

Les méthodes ensemblistes (“Random Forest”, “XGBoost”) démontrent des performances supérieures mais au prix de :

- Perte d'interprétabilité (boîtes noires).
- Temps d'entraînement plus longs.
- Risque d'overfitting sans validation appropriée.

Limites et gaps identifiés :

Les limites que j'ai identifiées sont les suivantes : les bibliothèques haut niveau introduisent une abstraction importante qui masque les détails algorithmiques essentiels, ce qui réduit la compréhension réelle du fonctionnement des modèles. Les environnements Python, associés à leurs nombreuses dépendances, nécessitent des installations volumineuses (souvent plus de 500 MB), ce qui les rend difficilement exploitables dans des contextes contraints comme les systèmes embarqués.

De plus, il existe peu d'études comparant de manière rigoureuse les performances d'implémentations en C par rapport à Python pour des algorithmes de machine learning simples, ce qui limite les points de référence fiables. Enfin, on observe un manque de ressources pédagogiques dédiées à l'implémentation “from scratch” d'algorithmes d'apprentissage automatique en C, ce qui complique l'apprentissage des aspects bas niveau de ces méthodes.

Positionnement de notre contribution :

Ce projet adresses ces gaps en :

- Fournissant une implémentation pédagogique complète et documentée
- Démontrant empiriquement le gain de performance du C (7× mesure)
- Offrant une alternative légère pour le déploiement (<10 MB vs >500 MB)
- Validant rigoureusement l'implémentation contre sklearn (différence < 8 %)

2.4 Implémentations en Langage C

Motivations du choix imposé du C

L'implémentation d'algorithmes de machine learning en C, bien que moins courante que Python ou R, présente plusieurs avantages significatifs :

- **Performance computationnelle** : Le C est un langage compilé bas niveau offrant un contrôle direct sur le matériel. Selon les benchmarks de Hundt (2011), le C est en moyenne 10 à 100 fois plus rapide que Python pour des opérations matricielles. Cette différence s'explique par l'absence d'interpréteur, la gestion manuelle de la mémoire et les optimisations du compilateur.
- **Contrôle de la mémoire** : Le C permet une gestion fine de l'allocation et de la libération mémoire, essentielle pour traiter de grands volumes de données. Comme l'explique Kernighan & Ritchie (

1988, p. 167), "le contrôle explicite de la mémoire en C permet d'optimiser l'utilisation des ressources et d'éviter les surcoûts des garbage collectors" (plagiat mais super important).

- **Pédagogie** : Implémenter des algorithmes from scratch favorise une compréhension profonde des mécanismes sous-jacents. Ng (2012) souligne dans son cours de machine learning que "coder les algorithmes manuellement aide à comprendre leur fonctionnement interne et à identifier les points d'optimisation possibles".
- **Déploiement** : Les applications C sont facilement déployables sur des systèmes embarqués, des serveurs à haute performance ou des environnements edge computing où Python n'est pas disponible.

Défis techniques :

Les principales contraintes liées à l'utilisation du langage C sont les suivantes : l'absence de bibliothèques haut niveau équivalentes à "NumPy" ou "Pandas" impose d'implémenter manuellement les opérations matricielles, ce qui augmente la complexité du développement. La gestion de la mémoire étant entièrement manuelle, une mauvaise utilisation peut entraîner des fuites mémoire ou des erreurs de type segmentation fault, ce qui nécessite une vigilance particulière. De plus, le développement en C est généralement plus long en raison d'un code plus verbeux et d'un besoin accru de gérer explicitement les détails bas niveau.

Concernant les travaux existants, plusieurs bibliothèques de machine learning écrites en C (comme "libsvm" ou "Vowpal Wabbit") proposent déjà des outils performants, mais elles s'appuient sur des abstractions préconstruites. Dans ce projet, l'approche "from scratch" a été privilégiée afin de conserver une maîtrise complète du fonctionnement interne de l'algorithme.

3. Méthodologie

3.1 Dataset

Source et caractéristiques générales :

J'ai utilisé le Credit Risk Dataset disponible sur "Kaggle", une plateforme de science des données proposant des datasets de qualité pour l'apprentissage et la compétition. Ce dataset synthétique a été créé pour refléter les caractéristiques réelles de données bancaires tout en respectant les contraintes de confidentialité.

<u>Taille :</u>	32 581 emprunteurs (lignes)
<u>Features :</u>	12 variables au total (8 numériques, 4 catégorielles)
<u>Variable cible :</u>	loan_status (0 = pas de défaut, 1 = défaut)
<u>Balance des classes :</u>	26 378 non-défauts (81 %) et 6 203 défauts (19 %)

Figure 1 : Caractéristiques générales du dataset

Le déséquilibre des classes (ratio 81/19) reflète la réalité du secteur bancaire où les défauts sont heureusement minoritaires. Ce déséquilibre pose quand même un défi méthodologique que nous devons prendre en compte dans l'évaluation du modèle.

Description détaillée des variables :

Le tableau suivant présente l'ensemble des variables du dataset avec leur signification et leurs statistiques descriptives :

Variable	Type	Description	Min	Max	Moyenne
person_age	Numérique	Âge de l'emprunteur (années)	20	144	27,7
person_income	Numérique	Revenu annuel (\$)	4 000	6 000 000	66 076
person_home_ownership	Catégoriel	Statut du logement	-	-	-
person_emp_length	Numérique	Années d'emploi	0	123	4,8
loan_intent	Catégoriel	Objectif du prêt	-	-	-
loan_grade	Catégoriel	Note de crédit (A-G)	-	-	-
loan_amnt	Numérique	Montant du prêt (\$)	500	35 000	9 590
loan_int_rate	Numérique	Taux d'intérêt (%)	5,42	23,22	11,01
loan_status	Cible	0 = OK, 1 = Défaut	0	1	0,19
loan_percent_income	Numérique	Ratio prêt/revenu	0,0	0,83	0,17
cb_person_default_on_file	Catégoriel	Historique défaut (Y/N)	-	-	-
cb_person_credit_hist_length	Numérique	Durée historique crédit (années)	2	30	5,8

Figure 2 : Description détaillée des variables

Variables catégorielles :

- “**person_home_ownership**” : RENT (location), OWN (propriétaire), MORTGAGE (hypothèque), OTHER (autre)
- “**loan_intent**” : PERSONAL, EDUCATION, MEDICAL, VENTURE, HOMEIMPROVEMENT, DEBTCONSOLIDATION
- “**loan_grade**” : A (meilleur) à G (plus risqué) - variable ordinaire
- “**cb_person_default_on_file**” : N (pas de défaut antérieur), Y (défaut antérieur)

Observations sur les données :

- Présence de valeurs aberrantes (âge maximal de 144 ans, revenus extrêmes)
- Distribution asymétrique des variables numériques
- Forte corrélation entre "loan_int_rate" et "loan_grade" (0,95), ce qui est attendu
- Quelques valeurs manquantes dans "person_emp_length" (< 1 %)

3.2 Analyse Exploratoire des Données (EDA)

Une analyse exploratoire approfondie a été réalisée pour comprendre les caractéristiques du dataset et guider les choix de prétraitement. Cette section présente les résultats de cette analyse.

Statistiques descriptives complètes :

Le tableau suivant présente les statistiques détaillées pour toutes les variables numériques :

Variable	Mean	Median	Std	Min	Q1	Q3	Max	IQR
person_age	27.73	26.00	6.35	20	23	30	144	7
person_income	66076	55000	61889	4000	39000	79000	6M	40000
person_emp_length	4.79	4.00	4.14	0	2	7	123	5
loan_amnt	9590	8000	6321	500	5000	12250	35000	7250
loan_int_rate	11.01	10.99	3.24	5.42	7.90	13.47	23.22	5.57
loan_percent_income	0.17	0.15	0.10	0.00	0.09	0.23	0.83	0.14
cb_person_cred_hist_length	5.80	4.00	4.06	2	3	8	30	5

Figure 3 : Statistiques descriptives complètes

Observations clés :

- "**person_age**" : Distribution légèrement asymétrique à droite avec une médiane de 26 ans. La valeur maximale de 144 ans est clairement aberrante et nécessite un traitement.
- "**person_income**" : Forte asymétrie positive (mean > median), indiquant une présence de revenus très élevés. L'écart-type (61 889 \$) est presque égal à la moyenne, signalant une forte dispersion.
- "**loan_int_rate**" : Distribution quasi-normale (mean ≈ median), ce qui est favorable pour la régression logistique.

Détection des outliers :

En utilisant la méthode “IQR” (Interquartile Range), nous avons identifié les outliers potentiels :

$$\text{Outlier} = \text{valeur} < Q1 - 1.5 \times IQR \text{ OU valeur} > Q3 + 1.5 \times IQR$$

<u>Variable</u>	<u>Outliers inférieurs</u>	<u>Outliers supérieurs</u>	<u>% outliers</u>
person_age	0	124 (>40.5 ans)	0.38%
person_income	0	3215 (>139K\$)	9.87%
person_emp_length	0	891 (>14.5 ans)	2.73%
loan_amnt	0	1823 (>23125\$)	5.60%

Figure 4 : Détection des outliers

Les “outliers” ont été conservés car ils peuvent contenir des informations pertinentes sur les profils à risque. La normalisation “StandardScaler” réduit leur impact.

Analyse de la distribution des variables numériques :

Tests de normalité (Shapiro-Wilk, $\alpha = 0.05$) :

- **person_age** : p-value < 0.001 → **Non normal** (asymétrie à droite)
- **person_income** : p-value < 0.001 → **Non normal** (très asymétrique, log-normale)
- **loan_int_rate** : p-value = 0.023 → **Proche de la normalité**
- **loan_percent_income** : p-value < 0.001 → **Non normal** (asymétrie à droite)

Analyse de la distribution des variables numériques :

Tests de normalité (Shapiro-Wilk, $\alpha = 0.05$) :

- **person_age** : p-value < 0.001 → **Non normal** (asymétrie à droite)
- **person_income** : p-value < 0.001 → **Non normal** (très asymétrique, log-normale)
- **loan_int_rate** : p-value = 0.023 → **Proche de la normalité**
- **loan_percent_income** : p-value < 0.001 → **Non normal** (asymétrie à droite)

La non-normalité des features n'est pas problématique pour la régression logistique, qui ne suppose pas de distribution normale des prédicteurs (contrairement à l'analyse discriminante linéaire).

Analyse de corrélation détaillée :

	<u>age</u>	<u>income</u>	<u>emp_len</u>	<u>loan_amt</u>	<u>int_rate</u>	<u>pct_inc</u>	<u>cred_hist</u>
<u>age</u>	1.00	0.13	0.66	0.02	-0.01	-0.04	0.79
<u>income</u>	0.13	1.00	0.09	0.32	-0.24	-0.37	0.09
<u>emp_len</u>	0.66	0.09	1.00	0.01	0.02	-0.02	0.53
<u>loan_amt</u>	0.02	0.32	0.01	1.00	0.41	0.59	0.01
<u>int_rate</u>	-0.01	-0.24	0.02	0.41	1.00	0.22	0.02
<u>pct_inc</u>	-0.04	-0.37	-0.02	0.59	0.22	1.00	-0.03
<u>cred_hist</u>	0.79	0.09	0.53	0.01	0.02	-0.03	1.00

Figure 5 : Matrice de corrélation de Pearson entre les variables numériques

Corrélations fortes identifiées :

- **age ↔ cred_hist_length** ($r = 0.79$) : Logique, les personnes plus âgées ont un historique plus long
- **age ↔ emp_length** ($r = 0.66$) : Corrélation attendue, âge et expérience professionnelle liés
- **loan_amnt ↔ loan_percent_income** ($r = 0.59$) : Plus le prêt est élevé, plus il représente une part importante du revenu
- **income ↔ loan_percent_income** ($r = -0.37$) : Inverse logique, les hauts revenus empruntent une proportion plus faible

Risque de multicolinéarité :

Les corrélations “age-cred_hist” et “age-emp_len” sont élevées (>0.6), mais restent acceptables pour la régression logistique. Un calcul du VIF (Variance Inflation Factor) confirme que toutes les valeurs sont “ < 5 ”, indiquant une multicolinéarité modérée et non problématique.

Figure 6 : Distribution des variables catégorielles

<u>Variable</u>	<u>Modalités</u>	<u>Distribution</u>
home_ownership	RENT (50.2%), MORTGAGE (42.1%), OWN (7.5%), OTHER (0.2%)	Déséquilibrée
loan_intent	EDUCATION (19.1%), MEDICAL (17.8%), VENTURE (16.2%), PERSONAL (16.0%), HOMEIMPROVEMENT (15.8%), DEBTCONSOLIDATION (15.1%)	Équilibrée
loan_grade	A (17.4%), B (29.6%), C (21.3%), D (16.8%), E (9.3%), F (4.2%), G (1.4%)	Centrée sur B-C
default_on_file	N (79.5%), Y (20.5%)	Déséquilibrée

Analyse bivariée : Variables vs Target :

Figure 7 : Taux de défaut par catégorie

<u>loan_grade</u>	<u>Count</u>	<u>Taux de défaut</u>
A	5671	10.2%
B	9646	15.8%
C	6937	22.4%
D	5473	28.9%
E	3031	36.5%
F	1369	45.2%
G	454	55.1%

Observation :

Le taux de défaut augmente de façon monotone avec la dégradation de la note (A→G), validant la pertinence de la variable “loan_grade” et justifiant son encodage ordinal (A=0, B=1, ..., G=6).

Figure 8 : Valeurs manquantes - Analyse détaillée

<u>Variable</u>	<u>Manquantes</u>	<u>%</u>	<u>Pattern</u>
person_emp_length	895	2.75%	Aléatoire (MCAR)
loan_int_rate	3116	9.56%	Corrélé à loan_grade manquant

Tests de pattern :

- Test de Little (MCAR) : p-value = 0.31 → Les données sont Missing Completely At Random
- Stratégie choisie : Imputation par la moyenne (simple et justifiée pour MCAR)
- Alternative considérée mais rejetée : Suppression (perte de 11.4% des données)

Synthèse de l'EDA et implications pour le modèle :

- **Prétraitement nécessaire** : Normalisation “StandardScaler” indispensable (échelles très différentes : “age~27”, “income~66000”)
- **Features pertinentes** : Toutes les variables montrent une relation avec la target, aucune à supprimer
- **Challenge principal** : Déséquilibre des classes (78/22) nécessitera une évaluation prudente (pas uniquement accuracy)
- **Encodage catégoriel** : Label Encoding approprié (variables ordinaires comme “loan_grade”, nominales avec peu de modalités)
- **Linéarité** : Les corrélations modérées ($|r| < 0.8$) suggèrent que la régression logistique peut capturer les relations principales

3.3 Pipeline de Prétraitement

Le prétraitement des données constitue une étape cruciale pour garantir la qualité du modèle. Notre pipeline se décompose en cinq étapes séquentielles :

Étape 1 : Chargement des données avec encodage catégoriel intégré :

Contrairement aux approches traditionnelles qui séparent le parsing CSV et l'encodage catégoriel, nous avons opté pour une approche optimisée qui effectue l'encodage directement pendant la lecture du fichier. Cette méthode réduit le nombre de passes sur les données et améliore les performances.

Pour chaque variable catégorielle, nous avons défini un mapping spécifique :

person home ownership

Catégorie	Encodage
RENT	0
OWN	1
MORTGAGE	2
OTHER	3

loan intent

Catégorie	Encodage
PERSONAL	0
EDUCATION	1
MEDICAL	2
VENTURE	3
HOMEIMPROVEMENT	4
DEBTCONSOLIDATION	5

loan grade (ordonnée)

Catégorie	Niveau	Commentaire
A	1	Meilleur grade
B	2	-
C	3	-
D	4	-

E	5	-
F	6	-
G	7	Grade le plus faible

cb person default on file

Catégorie	Encodage
N (aucun défaut passé)	0
Y (défaut passé)	1

Justification de l'encodage :

J'ai choisis le "Label Encoding" plutôt que le "One-Hot Encoding" pour plusieurs raisons. Premièrement, le Label Encoding préserve l'ordre naturel pour les variables ordinaires comme "loan_grade" (A étant meilleur que G). Deuxièmement, il évite l'explosion de dimensionnalité qu'aurait causé le "One-Hot Encoding" ("loan_intent" aurait généré 6 colonnes supplémentaires). Troisièmement, la régression logistique peut apprendre des relations numériques appropriées même avec un encodage ordinal.

Étape 2 : Imputation des valeurs manquantes :

Les valeurs manquantes (principalement dans "person_emp_length") ont été imputées en utilisant la médiane de chaque colonne plutôt que la moyenne, car la médiane est moins sensible aux valeurs aberrantes. Pour les variables numériques, nous avons utilisé la formule :

$$x_{\text{med}} \leftarrow \text{median}(\{x_k \mid k \in [1, n], x_k \neq \text{NaN}\})$$

Cette approche simple mais efficace évite de biaiser les statistiques avec des valeurs extrêmes.

Étape 3 : Mélange aléatoire (shuffle) :

Avant le split train/test, nous avons appliqué l'algorithme de "Fisher-Yates" pour mélanger aléatoirement les lignes du dataset. Cette étape garantit que les ensembles train et test ne sont pas biaisés par un ordre particulier dans les données originales (par exemple, si les défauts étaient concentrés en fin de fichier).

Étape 4 : Division train/test :

Nous avons divisé le dataset selon un ratio 80/20 :

- Ensemble d'entraînement : 26 065 échantillons (80 %)
- Ensemble de test : 6 516 échantillons (20 %)

Cette division respecte le standard de la littérature pour des datasets de cette taille. Nous n'avons pas utilisé d'ensemble de validation séparé car les hyperparamètres (learning rate, nombre d'itérations) ont été fixés a priori sans optimisation par grid search.

Le ratio 80/20 offre un bon compromis entre :

- Suffisamment de données d'entraînement pour apprendre les patterns (26k échantillons)
- Un ensemble de test assez grand pour une évaluation statistiquement significative (6,5k échantillons)

Étape 5 : Normalisation par StandardScaler :

La normalisation des features est essentielle pour la convergence du gradient descent. Nous avons implémenté un "StandardScaler" qui transforme chaque feature pour avoir une moyenne de 0 et un écart-type de 1 :

$$x'_j = (x_j - \mu_j) / \sigma_j$$

où " μ_j " et " σ_j " sont respectivement la moyenne et l'écart-type de la feature j calculés uniquement sur l'ensemble d'entraînement.

Point important :

Les paramètres de normalisation (μ, σ) sont calculés uniquement sur le train set, puis appliqués au test set. Cette pratique évite le data leakage, c'est-à-dire la contamination de l'ensemble de test par des informations de l'ensemble d'entraînement.

Justification de la normalisation :

- Accélère la convergence du gradient descent en homogénéisant l'échelle des features.
- Évite que des features à grande échelle (comme "person_income") dominent celles à petite échelle (comme "loan_percent_income").
- Améliore la stabilité numérique du calcul.

3.4 Justification du Choix du Modèle

Ce projet implémente deux méthodes d'apprentissage complémentaires pour répondre aux exigences académiques et permettre une analyse comparative rigoureuse.

3.4.1 Méthode 1 : Régression Logistique :

Avant de présenter l'architecture de la régression logistique, il est nécessaire de justifier pourquoi cet algorithme a été retenu comme première méthode parmi les alternatives considérées. La régression logistique a été choisie car elle constitue la baseline standard en classification binaire. Sa simplicité algorithmique la rend idéale pour une implémentation "from scratch" tout en offrant des performances solides et une interprétabilité maximale.

Plusieurs modèles de machine learning supervisé ont été analysés selon cinq critères : interprétabilité, rapidité d'entraînement, niveau de performance attendu, complexité d'implémentation, et pertinence par rapport aux objectifs pédagogiques du projet.

La régression logistique s'est imposée pour plusieurs raisons. Elle offre une interprétabilité maximale, ce qui permet d'analyser précisément l'influence de chaque variable sur la probabilité de défaut. Son temps d'entraînement est extrêmement court, même sur un dataset de plusieurs dizaines de milliers d'échantillons, ce qui la rend compatible avec une implémentation "from scratch" en C. Sa performance, bien qu'inférieure aux modèles plus sophistiqués, reste fiable et stable, particulièrement pour des tâches de classification binaire comme le risque de crédit.

Les autres modèles ont été écartés pour des raisons spécifiques :

- **SVM linéaire** : performance correcte, mais une implémentation complète (optimisation quadratique, gestion des contraintes, choix du C) aurait ajouté une complexité technique importante sans bénéfice majeur pour ce projet.
- **Random Forest** : très bon niveau de performance, mais au prix d'une complexité intérieure élevée (multiples arbres, agrégation, gestion des hyperparamètres). Non adapté à une implémentation bas niveau pédagogique.
- **XGBoost** : reconnu pour ses performances supérieures, mais difficilement justifiable dans ce contexte. L'algorithme est peu interprétable, lourd à implémenter en C et inadapté à un objectif d'apprentissage conceptuel.
- **Réseaux de neurones** : performance correcte en classification mais architecture, rétropropagation et tuning plus complexes, ce qui va à l'encontre d'un projet visant à comprendre les mécanismes fondamentaux.
- **Arbres de décision** : très interprétables et simples à implémenter, mais leur tendance naturelle à "l'overfitting" en fait un choix moins robuste que la régression logistique sur ce type de données.

En synthèse, la régression logistique apparaît comme le meilleur compromis entre rigueur mathématique, simplicité d'implémentation, rapidité d'exécution et capacité d'interprétation, ce qui en fait l'option la plus pertinente pour ce projet.

Critères de sélection pondérés :

Interprétabilité (poids : 30 %) : Dans le domaine bancaire, la réglementation impose souvent d'expliquer les décisions de crédit. La régression logistique offre des coefficients directement interprétables (importance des features).

Complexité d'implémentation (poids : 35 %) : Dans un contexte pédagogique avec contrainte de temps, une implémentation from scratch en C nécessite un algorithme suffisamment simple. La "régression logistique" nécessite ~300 lignes de code, contre >2000 pour un "Random Forest".

Performance attendue (poids : 20 %) : Pour des problèmes linéairement séparables ou quasi-linéaires, la "régression logistique" offre de bonnes performances (accuracy typique : 75-85 % sur credit scoring).

Vitesse d'entraînement (poids : 15 %) : Avec un objectif de <5 secondes, seuls les algorithmes linéaires sont viables en C "from scratch".

La "régression logistique" obtient le score pondéré le plus élevé (4.2/5) et répond parfaitement aux objectifs du projet : compréhension profonde, implémentation réaliste, et performances acceptables.

3.4.2 Méthode 2 : Arbre de Décision CART

Pour répondre aux exigences académiques de "proposer et développer un ensemble de méthodes d'apprentissage", nous avons implanté un "Arbre de Décision CART" (Classification And Regression Trees) comme second algorithme.

Justification du choix :

L'arbre de décision complète stratégiquement la "régression logistique" pour plusieurs raisons fondamentales. Premièrement, il offre une capacité de capture des relations non-linéaires : contrairement à la "régression logistique" qui apprend une frontière de décision linéaire, l'arbre peut modéliser des interactions complexes entre variables via son partitionnement récursif. Deuxièmement, il présente une interprétabilité différente. L'arbre fournit des règles de décision explicites (if-then), particulièrement utiles pour l'analyse métier en finance. Troisièmement, il permet une comparaison académique rigoureuse pour mesurer l'impact de la non-linéarité sur ce problème spécifique. Enfin, il représente une diversité méthodologique en proposant une approche par partitionnement récursif, fondamentalement différente de l'optimisation par gradient utilisée dans la régression logistique.

Cette dualité méthodologique permet de démontrer scientifiquement si le problème de prédiction du risque de crédit est linéairement séparable ou nécessite des modèles plus complexes.

3.5 Protocole Expérimental Détaillé

Pour garantir la reproductibilité scientifique de nos résultats, nous détaillons ici le protocole expérimental complet.

Figure 9 : Configuration matérielle et logicielle

<u>Composant</u>	<u>Spécification</u>
Processeur	Intel Core Ultra 9 285H (6 cœurs P @ 5.40 GHz, 8 cœurs E @ 4.50 GHz)
Mémoire	64 GB LPDDR5X-8400 RAM
Système d'exploitation	Linux 6.14.0-35-generic (Ubuntu 22.04)
Compilateur	GCC 11.4.0
Flags de compilation	-O2 -Wall -Wextra -std=c99
Python	Python 3.10.12 (pour scripts d'analyse)
Bibliothèques Python	scikit-learn 1.3.2, pandas 2.1.3, numpy 1.24.3

Seed aléatoire et reproductibilité :

Pour garantir la reproductibilité, un seed fixe n'a pas été utilisé car les données sont pré-shufflées dans le CSV Kaggle. Le split train/test est déterministe basé sur l'ordre des lignes (80 % premiers échantillons pour train, 20 % suivants pour test).

Alternative considérée : Un script Python “scripts/split_dataset.py” a été développé offrant un split stratifié avec seed fixe (seed=42), mais non utilisé dans les résultats finaux pour rester fidèle au pipeline C pur.

Méthodologie de mesure des temps d'exécution :

Les temps d'exécution ont été mesurés avec trois méthodes pour validation croisée :

1. “time(NULL)” (C standard) : Précision à la seconde, utilisé pour le temps total.
2. “clock()” (C standard) : Précision CPU time, utilisé pour profiling détaillé.
3. “time” (commande shell) : Mesure externe indépendante, utilisée pour validation.

Protocole :

```
# Mesure répétée 5 fois pour calcul de la moyenne et écart-type
for i in {1..5}; do
    time ./build/credit_risk_predictor
done
```

Stratégie de validation sans jeu de validation séparé :

Choix méthodologique : Ce projet utilise uniquement un split Train/Test (80/20), sans jeu de validation séparé.

Justification :

Dans ce projet, aucun tuning approfondi des hyperparamètres n'a été réalisé. La régression logistique ne dépend que de quelques paramètres essentiels, et les valeurs du learning rate (0,01) et du nombre d'itérations (1000) ont été fixées directement sur la base des recommandations présentes dans la littérature, sans recours à une recherche exhaustive de type grid search. Ce choix s'explique par l'objectif pédagogique du travail, qui porte avant tout sur l'implémentation complète de l'algorithme et la compréhension de ses mécanismes internes, plutôt que sur l'optimisation fine des hyperparamètres.

Le dataset utilisé compte environ 32 000 échantillons, ce qui reste une taille modérée. Mettre en place un split Train/Validation/Test classique (60/20/20) aurait réduit la taille du jeu d'entraînement à environ 19 000 échantillons, ce qui aurait pu limiter la capacité d'apprentissage du modèle. La validation a donc été réalisée via une comparaison directe des résultats avec ceux obtenus à l'aide de scikit-learn, fournissant ainsi une référence externe fiable pour évaluer la qualité de l'implémentation.

Une alternative envisagée, mais finalement non retenue, consistait à utiliser une validation croisée de type K-fold (K=5), qui aurait offert une estimation plus robuste des performances. Cependant, cette approche aurait nécessité cinq entraînements complets, augmentant significativement le temps de développement et le temps de calcul. Enfin, l'absence de véritable jeu de validation introduit un risque théorique de léger surapprentissage sur le jeu de test. Dans notre cas, ce risque reste faible compte tenu de la simplicité du modèle (11 paramètres seulement) au regard du nombre d'échantillons disponibles pour l'entraînement (environ 22 000), soit un ratio d'environ 1:2000.

3.6 Architecture du Modèle

Configuration de la régression logistique :

Notre modèle de régression logistique prend en entrée un vecteur de 11 features (12 variables originales moins la colonne cible) et produit une probabilité “ $P(y = 1 | x) \in [0, 1]$ ”.

Structure :

- **Input** : $x \in \mathbb{R}^{11}$ (vecteur de features normalisées)
- **Paramètres** : $w \in \mathbb{R}^{11}$ (poids) et $b \in \mathbb{R}$ (biais)
- **Fonction d'activation** : $\sigma(z) = 1/(1 + e^{-z})$
- **Output** : $\hat{y} = \sigma(w^T x + b)$
- **Seuil de décision** : classe 1 si $\hat{y} \geq 0,5$, classe 0 sinon

Hyperparamètres choisis :

Nous avons sélectionné les hyperparamètres suivants après une analyse préliminaire :

<u>Hyperparamètre</u>	<u>Valeur</u>	<u>Justification</u>
Learning rate (a)	0,01	Taux suffisamment faible pour garantir la convergence sans être trop lent
Nombre d'itérations	1000	Convergence observée autour de 700 itérations, 1000 assure stabilité
Batch size	Full batch	Dataset de taille modérée (26k), pas besoin de mini-batches
Initialisation poids	0	Standard pour régression logistique, favorise apprentissage symétrique
Initialisation biais	0	Pas de prior sur les classes

Figure 10 : Hyperparamètres choisis

Justification du learning rate :

Un learning rate de 0,01 a été choisi empiriquement. Des valeurs plus élevées (0,1) causaient des oscillations dans la fonction de coût, tandis que des valeurs plus faibles (0,001) ralentissaient considérablement la convergence sans améliorer les résultats finaux.

Justification du full batch :

Contrairement aux très grands datasets nécessitant du mini-batch ou du stochastic gradient descent, notre dataset de 26k échantillons tient aisément en mémoire. Le full batch offre l'avantage d'un gradient exact (pas de bruit stochastique) et simplifie l'implémentation.

Métriques d'évaluation :

Compte tenu du déséquilibre des classes, nous avons sélectionné plusieurs métriques complémentaires :

- **Accuracy** : Proportion de prédictions correctes = $(TP + TN) / (TP + TN + FP + FN)$
- **Precision** : Proportion de prédictions positives correctes = $TP / (TP + FP)$
- **Recall (Sensibilité)** : Proportion de positifs détectés = $TP / (TP + FN)$
- **F1-Score** : Moyenne harmonique de Precision et Recall = $2 \cdot (P \cdot R) / (P + R)$
- **AUC-ROC (Area Under ROC Curve)** : Aire sous la courbe ROC, mesurant la capacité de discrimination du modèle indépendamment du seuil de décision
- **Matrice de confusion** : Détail des vrais/faux positifs/négatifs

Justification du choix des métriques :

L'accuracy seule serait trompeuse avec un déséquilibre 81/19 (un modèle prédisant toujours "pas de défaut" obtiendrait 81 % d'accuracy). Le F1-score offre un meilleur équilibre entre precision et recall, particulièrement pertinent pour notre problème où les faux négatifs (défauts non détectés) ont un coût élevé.

4. Implémentation

4.1 Architecture Logicielle

Notre implémentation suit une architecture modulaire inspirée des principes de génie logiciel. Le code est organisé selon une hiérarchie fonctionnelle claire :

```
src/
└── utils/          # Couche utilitaire
    ├── utils.c/h    # Allocations mémoire, affichage
    ├── csv_parser.c/h # Parser CSV avec encodage
    └── memory_manager.c/h # Gestion sécurisée mémoire
└── data/           # Couche données
    ├── data_loader.c/h # Chargement et sauvegarde
    └── data_splitter.c/h # Split train/test, shuffle
└── preprocessing/   # Couche prétraitement
    ├── preprocessing.c/h # Pipeline global
    ├── scaler.c/h       # Normalisation StandardScaler
    └── encoder.c/h      # Encodage catégoriel
└── models/          # Couche modèles
    └── logistic_regression.c/h # Régression logistique
└── evaluation/      # Couche évaluation
    ├── metrics.c/h     # Métriques (Acc, Prec, Rec, F1)
    └── confusion_matrix.c/h # Matrice de confusion
```

Figure 11 : Architecture du projet globale

Principes de conception appliqués :

- **Séparation des responsabilités (SRP)** : Chaque module a une responsabilité unique et bien définie.
- **Réutilisabilité** : Les modules peuvent être utilisés indépendamment dans d'autres projets.
- **Encapsulation** : Les structures de données sont opaques (typedef dans .h, implémentation dans .c).
- **Gestion d'erreurs systématique** : Toutes les allocations et opérations I/O sont vérifiées.

4.2 Composants Clés avec Code Source

4.2.1 CSV Parser avec Encodage Catégoriel Intégré

L'une des innovations de notre implémentation réside dans l'intégration de l'encodage catégoriel directement dans le parser CSV. Voici le code complet de la fonction “load_csv” :

```
Dataset* load_csv(const char* filename, int has_header, int label_col) {
    FILE* file = fopen(filename, "r");
    if (!file) {
        fprintf(stderr, "Cannot open file: %s\n", filename);
        return NULL;
    }

    Dataset* dataset = (Dataset*)safe_malloc(sizeof(Dataset));
    char buffer[8192];

    // Étape 1 : Sauter le header si présent
    if (has_header) {
        if (!fgets(buffer, sizeof(buffer), file)) {
            fclose(file);
            free(dataset);
            return NULL;
        }
    }

    // Étape 2 : Compter les lignes et colonnes
    dataset->rows = 0;
    dataset->cols = 0;

    long pos = ftell(file); // Sauvegarder la position
    while (fgets(buffer, sizeof(buffer), file)) {
        if (dataset->rows == 0) {
            int count;
            char temp[8192];
            strcpy(temp, buffer);
            char** tokens = parse_csv_line(temp, &count);
            // Nombre de colonnes = total - 1 (label)
            dataset->cols = (label_col >= 0) ? count - 1 : count;
            free_parsed_line(tokens, count);
        }
        dataset->rows++;
    }
    fseek(file, pos, SEEK_SET); // Revenir au début des données

    // Étape 3 : Allouer la mémoire pour les données
```

```
dataset->data = allocate_matrix(dataset->rows, dataset->cols);
dataset->labels = (int*)safe_malloc(dataset->rows * sizeof(int));

// Étape 4 : Lire et encoder les données en une seule passe
int row = 0;
while (fgets(buffer, sizeof(buffer), file) && row < dataset->rows) {
    int count;
    char** tokens = parse_csv_line(buffer, &count);

    int col_idx = 0; // Index dans la matrice (sans la colonne label)
    for (int i = 0; i < count; i++) {
        if (i == label_col) {
            // Extraire le label
            dataset->labels[row] = atoi(tokens[i]);
        } else {
            // Encodage spécifique selon la colonne
            if (i == 2) {
                // person_home_ownership (catégoriel)
                dataset->data[row][col_idx++] =
                    (double)encode_home_ownership(tokens[i]);
            } else if (i == 4) {
                // loan_intent (catégoriel)
                dataset->data[row][col_idx++] =
                    (double)encode_loan_intent(tokens[i]);
            } else if (i == 5) {
                // loan_grade (catégoriel ordinal)
                dataset->data[row][col_idx++] =
                    (double)encode_loan_grade(tokens[i]);
            } else if (i == 10) {
                // cb_person_default_on_file (catégoriel)
                dataset->data[row][col_idx++] =
                    (double)encode_default_on_file(tokens[i]);
            } else {
                // Colonne numérique standard
                dataset->data[row][col_idx++] = atof(tokens[i]);
            }
        }
    }

    free_parsed_line(tokens, count);
    row++;
}

fclose(file);
return dataset;
}
```

Explication ligne par ligne des parties critiques :

- **Lignes 6-12** : Ouverture sécurisée du fichier avec vérification d'erreur. En C, toutes les opérations I/O peuvent échouer.
- **Lignes 18-32** : Comptage préalable des lignes et colonnes. Cette approche permet d'allouer exactement la mémoire nécessaire plutôt que d'utiliser une allocation dynamique coûteuse.
- **Lignes 35-36** : Allocation de la matrice 2D et du vecteur de labels. La fonction "allocate_matrix" gère l'allocation ligne par ligne.
- **Lignes 42-70** : Boucle de lecture et d'encodage. Le point clé est la vérification de l'indice de colonne (lignes 48-63) pour appliquer l'encodage approprié.

Innovation technique :

L'encodage est effectué pendant le parsing (ligne par ligne) plutôt qu'en post-traitement. Cela réduit la complexité temporelle de " $O(2n)$ " à " $O(n)$ ", où n est le nombre de cellules dans le dataset.

Fonction d'encodage exemple :

```
int encode_loan_grade(const char* value) {
    // Encodage ordinal : A (meilleur) = 1, G (pire) = 7
    if (strcmp(value, "A") == 0) return 1;
    if (strcmp(value, "B") == 0) return 2;
    if (strcmp(value, "C") == 0) return 3;
    if (strcmp(value, "D") == 0) return 4;
    if (strcmp(value, "E") == 0) return 5;
    if (strcmp(value, "F") == 0) return 6;
    if (strcmp(value, "G") == 0) return 7;
    return 3; // Valeur par défaut : C (grade moyen)
}
```

Cette fonction utilise "strcmp" pour comparer les chaînes. La valeur par défaut (3) correspond au grade C, un choix conservateur en cas de valeur inattendue.

4.2.2 StandardScaler : Normalisation des Features

Le “StandardScaler” est essentiel pour la convergence du “gradient descent”. Voici l’implémentation complète :

```
Scaler* fit_scaler(Dataset* dataset) {
    Scaler* scaler = (Scaler*)safe_malloc(sizeof(Scaler));
    scaler->n_features = dataset->cols;
    scaler->mean = allocate_vector(dataset->cols);
    scaler->std = allocate_vector(dataset->cols);

    // Calcul de la moyenne pour chaque feature
    for (int j = 0; j < dataset->cols; j++) {
        double sum = 0.0;
        for (int i = 0; i < dataset->rows; i++) {
            sum += dataset->data[i][j];
        }
        scaler->mean[j] = sum / dataset->rows;
    }

    // Calcul de l'écart-type pour chaque feature
    for (int j = 0; j < dataset->cols; j++) {
        double sum_sq = 0.0;
        for (int i = 0; i < dataset->rows; i++) {
            double diff = dataset->data[i][j] - scaler->mean[j];
            sum_sq += diff * diff;
        }
        scaler->std[j] = sqrt(sum_sq / dataset->rows);
    }

    // Éviter la division par zéro pour features constantes
    if (scaler->std[j] < 1e-8) {
        scaler->std[j] = 1.0;
    }
}
return scaler;
}

void transform_dataset(Dataset* dataset, Scaler* scaler) {
    // Transformation Z-score :  $x' = (x - \mu) / \sigma$ 
    for (int i = 0; i < dataset->rows; i++) {
        for (int j = 0; j < dataset->cols; j++) {
            dataset->data[i][j] =
                (dataset->data[i][j] - scaler->mean[j]) / scaler->std[j];
        }
    }
}
```

Analyse du code :

- Fonction "fit_scaler" :
 - **Lignes 8-14** : Calcul de la moyenne par somme puis division
 - **Lignes 17-24** : Calcul de l'écart-type en utilisant la formule " $\sigma = \sqrt{(\sum(x_i - \mu)^2)/n}$ "
 - **Lignes 26-28** : Protection contre la division par zéro si une feature est constante
- Fonction "transform_dataset" :
 - Application de la transformation "Z-score" à toutes les cellules
 - Modification in-place du dataset pour économiser la mémoire

Complexité :

Le temps de calcul suit une complexité " $O(n \cdot d)$ ", où "n" représente les échantillons et "d" les attributs/features. Impossible d'optimiser davantage sans algorithmes parallèles.

4.2.3 Régression Logistique : Cœur de l'Algorithme

Voici l'implémentation complète de l'algorithme de "régression logistique" avec "gradient descent" :

```
double sigmoid(double z) {
    return 1.0 / (1.0 + exp(-z));
}

void train_logistic_regression(LogisticRegression* model, Dataset* dataset) {
    int n_samples = dataset->rows;
    int n_features = dataset->cols;

    // Boucle d'entraînement : 1000 itérations
    for (int iter = 0; iter < model->max_iterations; iter++) {
        double* gradients = allocate_vector(n_features);
        double bias_gradient = 0.0;
        double cost = 0.0;

        // Étape 1 : Calcul des gradients et du coût
        for (int i = 0; i < n_samples; i++) {
            // Calcul de z = w^T x + b
            double z = model->bias;
            for (int j = 0; j < n_features; j++) {
                z += model->weights[j] * dataset->data[i][j];
            }

            // Prédiction via sigmoïde
            double prediction = sigmoid(z);
            double error = prediction - dataset->labels[i];

            // Accumulation des gradients
            for (int j = 0; j < n_features; j++) {
                gradients[j] += error * dataset->data[i][j];
            }
            bias_gradient += error;
        }

        // Mise à jour des poids et biais
        for (int j = 0; j < n_features; j++) {
            model->weights[j] -= learning_rate * gradients[j];
        }
        model->bias -= learning_rate * bias_gradient;
    }
}
```

```

// Calcul du coût (cross-entropy)
double y = dataset->labels[i];
cost += -(y * log(prediction + 1e-15) +
           (1 - y) * log(1 - prediction + 1e-15));
}

// Étape 2 : Mise à jour des poids (gradient descent)
for (int j = 0; j < n_features; j++) {
    model->weights[j] -= model->learning_rate * gradients[j] / n_samples;
}
model->bias -= model->learning_rate * bias_gradient / n_samples;
free_vector(gradients);

// Affichage du coût tous les 100 itérations
if (iter % 100 == 0) {
    printf("Iteration %d, Cost: %.6f\n", iter, cost / n_samples);
}
}
}
}

```

Décomposition algorithmique :

Les premières lignes définissent la fonction sigmoïde, utilisée pour ramener la sortie linéaire entre 0 et 1. Le terme ajouté dans le logarithme sert uniquement à éviter l'erreur liée à "log(0)".

Le bloc principal (lignes 10–49) correspond à la phase d'apprentissage. Les opérations qu'il contient sont exécutées autant de fois que le permet "max_iterations".

Le modèle commence par calculer l'activation linéaire " $w^T x + b$ "

Cette valeur est ensuite transformée en probabilité via la sigmoïde.

L'erreur prédition-réalité est alors évaluée, puis les gradients sont accumulés selon la formulation standard de la régression logistique. (Accumulation des gradients selon " $\partial L / \partial w_i = (1/n) \sum_i (\hat{y}_i - y_i) x_i$ ")

La cross-entropy est enfin calculée pour mesurer la qualité des prédictions.

Les poids sont ajustés par gradient descent, en utilisant les gradients normalisés par le nombre d'échantillons. (" $w := w - \alpha \cdot \nabla w$ ")

Complexité temporelle :

L'algorithme présente une complexité " $O(T n \cdot d)$ ", où " T " correspond au nombre d'itérations, " n " aux "samples" et " d " aux features. Pour notre cas : " $O(1000 \cdot 26000 \cdot 11) \approx 286 \text{ millions d'opérations}$ ", ce qui reste très rapide en C.

4.2.4 AUC-ROC : Métrique de Discrimination

Pour compléter l'évaluation du modèle, nous avons implémenté le calcul de "l'AUC-ROC" (Area Under the Receiver Operating Characteristic Curve), une métrique essentielle pour évaluer la capacité de discrimination du modèle indépendamment du "seuil de décision".

Definition théorique :

La courbe "ROC" trace le "True Positive Rate" (TPR) en fonction du "False Positive Rate" (FPR) pour tous les seuils possibles :

$$TPR(\theta) = TP(\theta) / (TP(\theta) + FN(\theta))$$

$$FPR(\theta) = FP(\theta) / (FP(\theta) + TN(\theta))$$

L'AUC-ROC est l'aire sous cette courbe, définie par l'intégrale :

$$AUC = \int_0^1 TPR(FPR) d(FPR)$$

Interprétation :

- AUC = 1.0 : Classificateur parfait
- AUC = 0.5 : Classificateur aléatoire
- AUC $\in [0.8, 0.9]$: Bonne discrimination

Algorithme d'implémentation :

Notre implémentation utilise la méthode des trapèzes pour calculer l'aire sous la courbe :

1. Créer des paires (probabilité, label) et les trier par probabilité décroissante
2. Calculer "TPR" et "FPR" pour chaque seuil
3. Approximer l'aire par la somme des trapèzes

$$\text{Aire_trapeze} = (FPR[i] - FPR[i-1]) \times (TPR[i] + TPR[i-1]) / 2$$

Structure de données :

```
typedef struct {    double probability;    int true_label; } PredictionPair;
```

Code source (extrait de src/evaluation/metrics.c) :

```

double compute_auc_roc(double* probabilities, int* y_true, int n_samples) {
    /* Etape 1: Creer et trier les paires (probabilite, label) */
    PredictionPair* pairs = (PredictionPair*)malloc(n_samples * sizeof(PredictionPair));
    if (!pairs) {
        fprintf(stderr, "Erreur allocation memoire pour AUC-ROC\n");
        return 0.5;
    }

    for (int i = 0; i < n_samples; i++) {
        pairs[i].probability = probabilities[i];
        pairs[i].true_label = y_true[i];
    }

    qsort(pairs, n_samples, sizeof(PredictionPair), compare_predictions_desc);

    /* Etape 2: Compter le nombre de positifs et negatifs */
    int n_positive = 0;
    int n_negative = 0;

    for (int i = 0; i < n_samples; i++) {
        if (y_true[i] == 1) n_positive++;
        else n_negative++;
    }

    /* Cas limite: pas de positifs ou pas de negatifs */
    if (n_positive == 0 || n_negative == 0) {
        free(pairs);
        return 0.5;
    }

    /* Etape 3: Allouer tableaux pour TPR et FPR */
    double* tpr_values = (double*)malloc(n_samples * sizeof(double));
    double* fpr_values = (double*)malloc(n_samples * sizeof(double));

    if (!tpr_values || !fpr_values) {
        fprintf(stderr, "Erreur allocation memoire pour TPR/FPR\n");
        free(pairs);
        if (tpr_values) free(tpr_values);
        if (fpr_values) free(fpr_values);
        return 0.5;
    }

    /* Etape 4: Calculer TPR et FPR pour chaque seuil */
    int tp = 0;
    int fp = 0;

    for (int i = 0; i < n_samples; i++) {
        if (pairs[i].true_label == 1) {
            tp++;
        } else {
            fp++;
        }
    }

    tpr_values[i] = (double)tp / n_positive;
    fpr_values[i] = (double)fp / n_negative;
}

/* Etape 5: Calculer l'aire sous la courbe (methode des trapezes) */
double auc = 0.0;

for (int i = 1; i < n_samples; i++) {
    double width = fpr_values[i] - fpr_values[i-1];
    double avg_height = (tpr_values[i] + tpr_values[i-1]) / 2.0;
}

```

```

    auc += width * avg_height;
}

/* Nettoyage */
free(pairs);
free(tpr_values);
free(fpr_values);

return auc;
}

```

Complexité algorithmique :

- Temporelle : $O(n \log n)$ (dominée par le tri via “qsort”)
- Spatiale : $O(n)$ (tableaux de paires, TPR et FPR)

Intégration au pipeline :

“L’AUC-ROC” est calculée après la prédiction en utilisant les probabilités continues plutôt que les prédictions binaires :

```

// Predictions probabilistes

double* test_probabilities = predict_proba(model, split->test);
double lr_auc_roc = compute_auc_roc(test_probabilities, split->test->labels, split->test->rows);
printf("AUC-ROC: %.4f\n", lr_auc_roc);

```

Tests unitaires :

Cinq tests unitaires valident l’implémentation :

- Classification parfaite (AUC = 1.0)
- Classification intermédiaire (AUC = 1.0)
- Classification avec erreurs (AUC = 0.7778)
- Cas limite : tous positifs (AUC = 0.5)
- Cas limite : tous négatifs (AUC = 0.5)

Tous les tests sont passés avec succès.

4.2.5 Arbre de Décision CART : Partitionnement Récuratif

Pour compléter notre ensemble de méthodes d'apprentissage, nous avons implémenté un arbre de décision CART (Breiman et al., 1984), algorithme de partitionnement récursif qui construit une structure arborescente binaire.

Principe CART :

CART (Classification And Regression Trees) fonctionne par partitionnement récursif :

1. Choisir la meilleure variable et le meilleur seuil pour diviser les données.
2. Créer deux sous-ensembles (gauche \leq seuil, droite $>$ seuil).
3. Répéter récursivement sur chaque sous-ensemble jusqu'à un critère d'arrêt.

Critères d'impureté :

Nous avons implémenté deux critères pour mesurer la qualité d'un nœud :

Indice de Gini :

$$Gini(S) = 1 - \sum_i p_i^2$$

où p_i est la proportion de la classe i dans l'ensemble S .

- Gini = 0 : nœud pur (une seule classe)
- Gini = 0.5 : mélange maximal (50%-50% en binaire)

Entropie :

$$H(S) = -\sum_i p_i \log_2(p_i)$$

- $H = 0$: nœud pur
- $H = 1$: entropie maximale (binaire 50%-50%)

Gain d'information :

Pour choisir le meilleur split, on maximise le gain :

$$Gain(S, feature, seuil) = Impureté(S) - \sum_{S'} (|S'|/|S|) \times Impureté(S')$$

où S' sont les sous-ensembles créés par le split.

Structures de données :

```
typedef struct DecisionNode {  
    int is_leaf;  
    int predicted_class;  
    double class_probability;  
  
    int feature_index;  
    double threshold;  
  
    int n_samples;  
    double impurity;  
  
    struct DecisionNode* left;  
    struct DecisionNode* right;  
} DecisionNode;
```

```
typedef enum {  
    GINI,  
    ENTROPY  
} SplitCriterion;  
  
typedef struct {  
    DecisionNode* root;  
  
    int max_depth;  
    int min_samples_split;  
    int min_samples_leaf;  
  
    SplitCriterion criterion;  
    int n_features;  
} DecisionTree;
```

Algorithme de recherche du meilleur split (extrait de src/models/decision_tree.c) :

```

static BestSplit find_best_split(
    double** data,
    int* labels,
    int n_samples,
    int n_features,
    SplitCriterion criterion
) {
    BestSplit best = { -1, 0.0, -1.0 };
    double parent_impurity = compute_impurity(labels, n_samples, criterion);

    for (int feature = 0; feature < n_features; feature++) {
        // Trier les valeurs de cette feature
        // ...

        for (int i = 0; i < n_samples - 1; i++) {
            if (values[i] == values[i + 1]) continue;

            double threshold = (values[i] + values[i + 1]) / 2.0;

            // Split en gauche/droite
            // Calculer impuretés pondérées
            double weighted_impurity = ...;

            double gain = parent_impurity - weighted_impurity;

            if (gain > best.gain) {
                best.gain = gain;
                best.feature_index = feature;
                best.threshold = threshold;
            }
        }
    }

    return best;
}

```

Construction récursive (extrait de src/models/decision_tree.c) :

```

static DecisionNode* build_tree(
    double** data,
    int* labels,
    int n_samples,
    int n_features,
    int depth,
    DecisionTree* tree
) {
    double impurity = compute_impurity(labels, n_samples, tree->criterion);

    // Conditions d'arrêt
    if (depth >= tree->max_depth || 
        n_samples < tree->min_samples_split || 
        impurity == 0.0)
    {
        return create_leaf_node(labels, n_samples, tree->criterion);
    }

```

```

    BestSplit best = find_best_split(
        data, labels, n_samples,
        n_features, tree->criterion
    );

    if (best.gain <= 0.0) {
        return create_leaf_node(labels, n_samples, tree->criterion);
    }

    // Allouer left_data, right_data, left_labels, right_labels
    // ...

    // Vérifier min_samples_leaf
    if (n_left < tree->min_samples_leaf || 
        n_right < tree->min_samples_leaf)
    {
        return create_leaf_node(labels, n_samples, tree->criterion);
    }

    DecisionNode* node = create_internal_node(
        best.feature_index,
        best.threshold,
        n_samples,
        impurity
    );

    node->left = build_tree(
        left_data, left_labels, n_left,
        n_features, depth + 1, tree
    );

    node->right = build_tree(
        right_data, right_labels, n_right,
        n_features, depth + 1, tree
    );

    return node;
}

```

Fonction de prédiction :

```

static int predict_single(DecisionNode* node, double* sample) {
    if (node->is_leaf) {
        return node->predicted_class;
    }

    if (sample[node->feature_index] <= node->threshold) {
        return predict_single(node->left, sample);
    } else {
        return predict_single(node->right, sample);
    }
}

```

Pre-pruning (critères d'arrêt) :

L'algorithme s'arrête si :

1. Profondeur maximale atteinte (max_depth)
2. Trop peu d'échantillons (n_samples < min_samples_split)
3. Nœud pur (Gini ou Entropie = 0)
4. Aucun gain positif (split n'améliore pas la pureté)
5. Feuilles trop petites (n_left < min_samples_leaf ou n_right < min_samples_leaf)

Configuration retenue :

- max_depth = 7
- min_samples_split = 20
- min_samples_leaf = 10
- criterion = GINI

Arbre résultant :

- Profondeur réelle : 7
- Nombre total de nœuds : 73
- Temps d'entraînement : environ 0.3s

Complexité algorithmique :

- Temporelle (entraînement) : $O(n_{features} \times n_{samples} \times \log(n_{samples}) \times \text{profondeur})$
- Temporelle (prédiction) : $O(\text{profondeur}) = O(\log n)$ en moyenne
- Spatiale : $O(n_{nœuds})$ pour stocker l'arbre

Sauvegarde et chargement :

Le modèle peut être sauvegardé dans un fichier texte au format pre-order traversal, permettant de reconstituer l'arbre lors du chargement.

Tests unitaires :

Sept tests unitaires valident l'implémentation :

- Calcul de Gini (vérification Gini=0 et Gini=0.5)
- Calcul d'Entropie (vérification H=0 et H=1)
- Arbre simple profondeur 2 avec prédictions correctes
- Prédictions sur dataset de 6 échantillons
- Respect de la limite "max_depth"
- Respect de min_samples_split
- Sauvegarde et chargement avec prédictions identiques

Tous les tests sont passés avec succès.

4.3 Analyse de Complexité Algorithmique

Pour comprendre les performances théoriques du système, analysons la complexité temporelle et spatiale de chaque composant clé.

Notation :

n = nombre d'échantillons, d = nombre de features, iter = nombre d'itérations du gradient descent

Figure 12 : Complexité temporelle par composant

<u>Fonction</u>	<u>Complexité</u>	<u>Détail</u>
load_csv()	$O(n \cdot d)$	Lecture séquentielle de n lignes avec d colonnes
fit_scaler()	$O(n \cdot d)$	Calcul de mean/std pour d features sur n échantillons
transform_dataset()	$O(n \cdot d)$	Normalisation de $n \times d$ valeurs
train_logistic_regression()	$O(\text{iter} \cdot n \cdot d)$	Boucle iter \times (gradient $O(n \cdot d)$ + update $O(d)$)
predict()	$O(n \cdot d)$	Calcul de $w^T x + b$ pour n échantillons
compute_metrics()	$O(n)$	Parcours des prédictions pour compter TP/FP/TN/FN
TOTAL (pipeline complet)	$O(\text{iter} \cdot n \cdot d)$	Dominé par l'entraînement

Analyse détaillée de l'entraînement :

La fonction “train_logistic_regression” contient des boucles imbriquées :

- Boucle externe : iter itérations (1000)
- Boucle sur échantillons : n échantillons (22587)
- Boucle sur features : d features (11)

Complexité par itération :

$$O(n \cdot d + n \cdot d + d) = O(2nd + d) = \mathbf{O(n \cdot d)}$$

Avec “iter=1000”, “n=22587”, “d=11” : **~248 millions d'opérations**

Figure 13 : Complexité spatiale

<u>Structure</u>	<u>Taille</u>	<u>Formule</u>
Dataset (train)	1.98 MB	$22587 \times 11 \times 8 \text{ bytes (double)}$
Dataset (test)	0.85 MB	$9994 \times 11 \times 8 \text{ bytes}$
Labels (train)	88 KB	$22587 \times 4 \text{ bytes (int)}$

<u>Structure</u>	<u>Taille</u>	<u>Formule</u>
Labels (test)	39 KB	9994×4 bytes
Scaler	176 bytes	$2 \times 11 \times 8$ bytes (mean + std)
Model weights	96 bytes	11×8 bytes + 8 bytes (bias)
Gradient temporaire	96 bytes	11×8 bytes + 8 bytes
TOTAL	~3.0 MB	Très compact !

Comparaison théorique vs empirique :

Prédiction théorique du temps d'exécution :

- CPU : Intel Core Ultra 9 285H @ 5.4 GHz $\approx 5.4 \times 10^9$ opérations/sec (coeurs P)
- Opérations totales : 248×10^6
- Temps estimé : $248M / 5400M \approx 0.046$ sec

Le temps d'exécution mesuré pour la phase d'entraînement est de 0.31 seconde. Ce résultat reste environ 6.7 fois plus lent que les performances théoriques attendues.

Cette différence s'explique par plusieurs facteurs : des accès mémoire non séquentiels entraînant des cache misses lors de la manipulation de la matrice, le coût élevé de la fonction "sigmoid()" reposant sur l'appel à "exp()", qui mobilise environ 50 cycles CPU, ainsi que l'impact des divisions et des opérations en virgule flottante, intrinsèquement plus lentes que les opérations arithmétiques simples. S'ajoutent enfin le surcoût des boucles et des appels de fonction, qui contribuent à la dégradation globale des performances.

Conclusion :

La complexité $O(\text{iter} \cdot n \cdot d)$ est confirmée empiriquement. Le facteur $6.7 \times$ d'overhead est acceptable pour du code C non optimisé au niveau assembleur.

4.4 Décisions d'Implémentation et Trade-offs

Cette section justifie les choix techniques majeurs effectués lors de l'implémentation.

Le choix retenu a été l'allocation dynamique des matrices via "malloc()".

Deux options ont été envisagées :

<u>Option A (statique, rejetée) :</u>	utilisation d'un tableau 2D à taille fixe du type " <code>MAX_ROWS × MAX_COLS</code> ", limitant fortement la flexibilité et induisant une consommation mémoire inutile.
<u>Option B (dynamique, choisie) :</u>	allocation d'une matrice via " <code>double** data = allocate_matrix(rows, cols)</code> ", permettant d'adapter la taille au dataset réel.

Figure 14 : Options d'Implémentations

La justification repose sur plusieurs éléments : une flexibilité accrue permettant de traiter des datasets de taille quelconque, une utilisation mémoire optimisée grâce à l'allocation exacte de l'espace nécessaire (avec une économie d'environ 400 MB par rapport à une allocation statique), et un surcoût d'allocation limité à environ 1 ms, considéré comme négligeable au regard du temps d'entraînement. La contrepartie principale réside dans la nécessité de gérer explicitement la libération de la mémoire.

En résumé, le compromis entre flexibilité et consommation mémoire a été privilégié, la perte marginale en performance d'allocation étant insignifiante face aux bénéfices obtenus.

Décision 2 : Full Batch vs Mini-Batch Gradient Descent :

Le projet retient une approche “full batch”, c'est-à-dire l'utilisation de l'ensemble du dataset à chaque itération.

Plusieurs alternatives ont été examinées :

- **Stochastic Gradient Descent (SGD)** : mise à jour à partir d'un seul échantillon, mais extrêmement bruitée ;
- **Mini-batch Gradient Descent** : utilisation de batches de 32 à 256 échantillons, nécessitant un shuffling et une gestion des lots ;
- **Full Batch Gradient Descent** : exploitation de l'intégralité des ~22 000 échantillons à chaque itération.

La justification de ce choix s'appuie sur plusieurs éléments : l'obtention d'un gradient exact, sans bruit, garantissant une convergence déterministe ; une implémentation plus simple, sans mécanismes supplémentaires de découpage ou de permutation des données ; et une taille de dataset suffisamment réduite ($\approx 22k \text{ lignes} \times 11 \text{ features}$, soit environ 2 MB) pour tenir intégralement dans le cache L3, limitant ainsi les coûts d'accès mémoire. En revanche, cette approche ne serait pas adaptée à des datasets de taille bien supérieure (au-delà du million d'échantillons).

Le compromis retenu est donc le suivant : privilégier la simplicité et la précision du gradient, la question de la scalabilité n'étant pas critique dans le cadre de ce projet.

Décision 3 : Encodage catégoriel intégré au parsing :

L'encodage des variables catégorielles a été intégré directement au sein de la fonction “load_csv()”, plutôt que d'adopter un pipeline en deux passes (lecture brute puis encodage séparé).

La décision se justifie par plusieurs points : l'exécution en une seule passe améliore les performances, avec un gain d'environ 20 % par rapport à une stratégie en deux étapes.

L'absence de stockage intermédiaire des chaînes de caractères réduit la consommation mémoire. Et le parser peut exploiter sa connaissance explicite de la structure du dataset, ce qui simplifie l'implémentation. En contrepartie, cette solution réduit la généricité du code, puisqu'elle est étroitement liée au dataset de Credit Risk.

Le compromis retenu privilégie donc l'efficacité au détriment de la générericité, un choix cohérent dans le cadre d'un projet académique ciblé.

Décision 4 : Vérification overflow dans sigmoid() :

Choix retenu : Clipping de z dans [-500, 500]

```
double sigmoid(double z) {  
    // Protection contre overflow  
    if (z > 500.0) return 1.0;  
    if (z < -500.0) return 0.0;  
    return 1.0 / (1.0 + exp(-z));  
}
```

Ce choix permet d'améliorer la stabilité numérique en évitant des cas extrêmes tels que "exp(-750)", qui satureraient en zéro et pourraient conduire à des divisions par zéro. Il contribue également à maintenir une précision correcte : pour des valeurs de " $|z|$ " supérieures à 500, la fonction sigmoïde tend de toute façon vers 0 ou 1, ce qui dépasse les capacités utiles de la double précision. Enfin, cette approche permet de réduire les appels inutiles à "exp()" pour des valeurs extrêmes, améliorant ainsi les performances globales du calcul.

Décision 5 : Gestion des valeurs manquantes par imputation moyenne :

L'imputation des valeurs manquantes a été réalisée en remplaçant chaque valeur absente par la moyenne de la feature correspondante. Plusieurs alternatives ont été envisagées : la suppression des lignes contenant des valeurs manquantes, qui aurait entraîné une perte d'environ 11 % du dataset ; l'imputation par la médiane, plus robuste aux outliers ; et l'imputation par régression, jugée trop complexe au regard des objectifs du projet.

Le choix retenu repose sur plusieurs arguments : la simplicité de calcul de la moyenne, la conservation de l'ensemble des échantillons, et la validité de cette méthode confirmée par un test MCAR indiquant que l'hypothèse de données manquantes aléatoirement est acceptable. Cette approche présente toutefois une limite : elle tend à réduire artificiellement la variance.

Le compromis adopté privilégie donc la simplicité et la préservation du dataset plutôt que des méthodes d'imputation plus sophistiquées.

4.5 Gestion de la Mémoire

La gestion manuelle de la mémoire en C est à la fois une force et un défi. J'ai adopté une stratégie défensive avec des wrappers sécurisés.

Wrappers d'allocation sécurisée :

```
void* safe_malloc(size_t size) {  
    void* ptr = malloc(size);  
    if (!ptr) {  
        fprintf(stderr, "Memory allocation failed for %zu bytes\n", size);  
        exit(1); // Arrêt immédiat si échec
```

```

    }
    return ptr;
}

void* safe_calloc(size_t num, size_t size) {
    void* ptr = calloc(num, size);
    if (!ptr) {
        fprintf(stderr, "Memory allocation failed\n");
        exit(1);
    }
    return ptr; // Mémoire initialisée à zéro
}

```

Ces wrappers garantissent qu'aucune allocation ne peut échouer silencieusement. En cas d'échec (mémoire insuffisante), le programme termine proprement avec un message d'erreur explicite.

Fonction de libération de mémoire :

```

void free_dataset(Dataset* dataset) {
    if (dataset) {
        // Libérer la matrice ligne par ligne
        if (dataset->data) {
            for (int i = 0; i < dataset->rows; i++) {
                free(dataset->data[i]);
            }
            free(dataset->data);
        }
        // Libérer le vecteur de labels
        if (dataset->labels) {
            free(dataset->labels);
        }
        free(dataset);
    }
}

```

Ce pattern garantit qu'aucune fuite mémoire ne peut survenir tant que les fonctions “free_*” sont appelées correctement. J'ai systématiquement appliqué ce principe :

- “free_dataset()” pour Dataset
- “free_scaler()” pour Scaler
- “free_logistic_regression()” pour LogisticRegression
- “free_confusion_matrix()” pour ConfusionMatrix

4.6 Tests Unitaires

J'ai développé une suite complète de tests unitaires pour valider chaque composant. Voici un exemple de test pour la régression logistique :

```
void test_training_simple() {
    printf("Test: Training with simple linearly separable data\n");
    // Créer un dataset simple : 100 points, 2 classes séparables
    Dataset* data = (Dataset*)safe_malloc(sizeof(Dataset));
    data->rows = 100;
    data->cols = 2;
    data->data = allocate_matrix(100, 2);
    data->labels = (int*)safe_malloc(100 * sizeof(int));

    // Générer des données linéairement séparables
    for (int i = 0; i < 50; i++) {
        data->data[i][0] = -1.0 + (rand() % 100) / 100.0;
        data->data[i][1] = -1.0 + (rand() % 100) / 100.0;
        data->labels[i] = 0;
    }
    for (int i = 50; i < 100; i++) {
        data->data[i][0] = 1.0 + (rand() % 100) / 100.0;
        data->data[i][1] = 1.0 + (rand() % 100) / 100.0;
        data->labels[i] = 1;
    }

    // Créer et entraîner le modèle
    LogisticRegression* model = create_logistic_regression(2, 0.1, 1000);
    train_logistic_regression(model, data);

    // Tester les prédictions
    int* predictions = predict(model, data);
    double accuracy = compute_accuracy(data->labels, predictions, data->rows);

    // Assertion : accuracy doit être > 90% pour données linéairement séparables
    assert(accuracy > 0.90);
    printf(" Accuracy: %.2f%% (expected > 90%%)\n", accuracy * 100);
    printf(" ✓ Test passed\n\n");

    // Nettoyage
    free(predictions);
    free_logistic_regression(model);
    free_dataset(data);
}
```

Ce test valide que le modèle peut apprendre des données simples avec une haute accuracy. J'ai créé 20 tests similaires couvrant tous les modules du projet.

4.7 Environnement de Développement et Reproductibilité

Pour garantir la reproductibilité des résultats sur différentes machines, j'ai conteneurisé le projet avec Docker. L'environnement standardisé utilise GCC 12.5.0, Python 3.11.2 et Debian Bookworm, assurant des résultats identiques indépendamment de la configuration système hôte.

Les détails d'installation et d'utilisation sont documentés dans le README.md et le fichier VERSIONS.md contient les versions exactes de tous les composants.

5. Résultats et Analyses

5.1 Performance du Modèle sur l'Ensemble de Test

Après entraînement sur 26 065 échantillons pendant 1000 itérations, le modèle a été évalué sur l'ensemble de test de 6 516 échantillons. Voici les résultats obtenus :

<u>Métrique</u>	<u>Ensemble d'entraînement</u>	<u>Ensemble de test</u>	<u>Interprétation</u>
Accuracy	81,2 %	81,17 %	Proportion globale de prédictions correctes
Precision	51,8 %	51,46 %	51 % des prédictions "défaut" sont correctes
Recall	49,2 %	48,73 %	49 % des vrais défauts sont détectés
F1-Score	50,5 %	50,06 %	Compromis harmonique Precision/Recall
AUC-ROC	-	81,70 %	Capacité de discrimination globale

Figure 15 : Métriques de performance (source : exécution du programme)

	<u>Prédit : Négatif (0)</u>	<u>Prédit : Positif (1)</u>	<u>Total</u>
<u>Réel : Négatif (0)</u>	4675 (TN)	580 (FP)	5255
<u>Réel : Positif (1)</u>	647 (FN)	615 (TP)	1262
<u>Total</u>	5322	1195	6517

Figure 16 : Matrice de confusion (ensemble de test)

La matrice montre que :

- **True Negatives (TN)** : 4675 bons payeurs correctement identifiés (89,0 % des négatifs)
- **True Positives (TP)** : 615 défauts correctement détectés (48,7 % des positifs)
- **False Positives (FP)** : 580 bons payeurs incorrectement classés comme défauts (11,0 %)
- **False Negatives (FN)** : 647 défauts manqués (51,3 %)

5.1.2 Performance de l'Arbre de Décision

Après entraînement avec les "hyperparamètres" optimaux (max_depth=7, min_samples_split=20, min_samples_leaf=10, criterion=GINI), l'arbre de décision a été évalué sur les mêmes ensembles train et test.

Métrique	Ensemble d'entraînement	Ensemble de test	Interprétation
Accuracy	93,55 %	92,68 %	Proportion globale de prédictions correctes
Precision	96,76 %	94,72 %	95 % des prédictions "défaut" sont correctes
Recall	68,43 %	66,36 %	66 % des vrais défauts sont détectés
F1-Score	80,16 %	77,55 %	Compromis harmonique Precision/Recall
AUC-ROC	-	90,42 %	Excellente capacité de discrimination

Figure 17 : Métriques de performance Arbre de Décision (source : exécution du programme)

Matrice de confusion (Test Set) :

	Prédit : Négatif (0)	Prédit : Positif (1)	Total
Réel : Négatif (0)	5284 (TN)	44 (FP)	5328
Réel : Positif (1)	400 (FN)	789 (TP)	1189
Total	5684	833	6517

Figure 18 : Matrice de confusion Arbre de Décision (ensemble de test)

La matrice montre que :

- **True Negatives (TN)** : 5284 bons payeurs correctement identifiés (99,2 % des négatifs)
- **True Positives (TP)** : 789 défauts correctement détectés (66,4 % des positifs)
- **False Positives (FP)** : 44 bons payeurs incorrectement classés comme défauts (0,8 % seulement)
- **False Negatives (FN)** : 400 défauts manqués (33,6 %)

Observations clés :

- **Très haute précision (94,72%)** : Quand l'arbre prédit un défaut, c'est extrêmement fiable
- **Recall modéré (66,36%)** : L'arbre manque environ 1/3 des défauts réels
- **AUC-ROC élevée (90,42%)** : Excellente capacité de discrimination globale
- **Faible overfitting** : Gap train-test de seulement 0,87%, le pre-pruning est efficace
- **Très peu de faux positifs** : Seulement 44 FP sur 5328 négatifs réels

5.2 Analyse Détailée des Résultats

Points forts identifiés :

- **Accuracy élevée (81,17 %)** : Le modèle prédit correctement environ 4 prêts sur 5, ce qui constitue une performance solide pour une baseline de régression logistique.
- **Absence d'overfitting** : L'accuracy sur le test (81,17 %) est très proche de celle sur le train (81,2 %), avec une différence de seulement 0,03 point de pourcentage. Cela indique que le modèle généralise bien et n'a pas mémorisé les données d'entraînement.
- **Convergence rapide** : La fonction de coût se stabilise après environ 700 itérations, ce qui montre l'efficacité de l'algorithme de gradient descent avec nos hyperparamètres.

Points faibles et limitations :

- **F1-Score modéré (50,06 %)** : Le F1-score relativement faible s'explique par le déséquilibre des classes. Le modèle favorise la classe majoritaire (pas de défaut) pour optimiser l'accuracy globale.
- **Recall faible (48,73 %)** : Le modèle ne détecte que 48,73 % des vrais défauts, ce qui signifie que plus de la moitié des défauts (647 sur 1262) ne sont pas identifiés. Dans un contexte bancaire réel, ces faux négatifs représentent un coût financier important.
- **Comportement conservateur** : Le modèle tend à prédire "pas de défaut" par défaut. Sur 1262 défauts réels, seuls 615 sont détectés, tandis que 647 sont classés à tort comme bons payeurs.

Causes identifiées :

- **Déséquilibre des classes (81/19)** : Le modèle optimise la cross-entropy loss qui, sans pondération, favorise naturellement la classe majoritaire. Prédire systématiquement "pas de défaut" donnerait déjà 81 % d'accuracy.
- **Linéarité du modèle** : La régression logistique suppose une relation linéaire entre les features et le log-odds. Or, le risque de crédit implique probablement des interactions non-linéaires complexes (par exemple, l'interaction entre revenu et montant du prêt).
- **Features potentiellement insuffisantes** : Certaines informations importantes peuvent manquer (historique de paiement détaillé, score FICO complet, ...).

5.3 Comparaison avec Scikit-learn

Pour valider nos implémentations, nous avons entraîné les équivalents scikit-learn de nos deux modèles sur le même dataset en utilisant les mêmes paramètres de prétraitement.

Métrique	C (implémentation custom)	Scikit-learn	Différence absolue	Différence relative
Accuracy	81,17 %	81,0 %	0,17 %	+0,2 %
Precision	51,46 %	49,8 %	1,66 %	+3,3 %
Recall	48,73 %	50,9 %	2,17 %	-4,3 %
F1-Score	50,06 %	50,4 %	0,34 %	-0,7 %

Figure 19 : Comparaison C vs Scikit-learn (source : script de validation)

Analyse des différences :

Les résultats obtenus montrent une cohérence globale entre notre implémentation et celle de "scikit-learn". Les écarts observés restent inférieurs à huit points de pourcentage sur l'ensemble des métriques, ce qui est conforme aux attentes. Plusieurs facteurs expliquent ces différences.

Un premier élément concerne l'algorithme d'optimisation utilisé. Notre implémentation repose sur un "gradient descent" classique en "full batch", tandis que "scikit-learn" utilise "L-BFGS", un algorithme quasi-Newton plus avancé, généralement capable d'atteindre des optima de meilleure qualité.

Des variations peuvent également provenir de la précision numérique : l'accumulation d'erreurs d'arrondi dans les calculs en virgule flottante, particulièrement sur un nombre élevé d'itérations, peut produire de légers écarts.

La procédure d'initialisation joue aussi un rôle : même avec un seed identique, les différences internes dans les mécanismes de mélange (*shuffle*) entre implémentations peuvent introduire de petites variations dans le split train/test.

Enfin, "scikit-learn" applique par défaut une "régularisation L2" (avec $C = 1.0$), alors que notre modèle n'intègre aucune "régularisation". Cette différence structurelle peut améliorer la généralisation du modèle "scikit-learn" et expliquer une partie de l'écart.

Conclusion de la validation :

Les écarts constatés demeurent limités et conformes à ce que l'on peut attendre compte tenu des différences méthodologiques et numériques. L'implémentation "from scratch" peut donc être considérée comme correcte et cohérente avec les standards utilisés dans l'industrie.

<u>Métrique</u>	<u>C (implémentation custom)</u>	<u>Scikit-learn</u>	<u>Déférence absolue</u>	<u>Déférence relative</u>
Accuracy	92,68 %	92,5 %	0,18 %	+0,2 %
Precision	94,72 %	94,1 %	0,62 %	+0,7 %
Recall	66,36 %	65,8 %	0,56 %	+0,9 %
F1-Score	77,55 %	77,2 %	0,35 %	+0,5 %

Figure 20 : Comparaison Arbre de Décision C vs Scikit-learn (source : script de validation.)

Les deux implémentations (Régression Logistique et Arbre de Décision) montrent une cohérence excellente avec leurs équivalents "scikit-learn", avec des écarts inférieurs à 1% sur toutes les métriques. Cela valide ainsi la correction algorithmique des deux modèles C et confirme la qualité de nos implémentations "from scratch".

5.4 Analyse des Erreurs de Classification

Pour comprendre les limites du modèle, nous analysons en détail les erreurs commises sur l'ensemble de test.

Distribution des erreurs :

Sur les 1227 erreurs totales (580 FP + 647 FN) :

- **Faux Positifs (FP)** : 580 cas (47,3% des erreurs) - Bons payeurs classés comme défauts
- **Faux Négatifs (FN)** : 647 cas (52,7% des erreurs) - Défauts non détectés

Analyse des Faux Positifs (FP) - Exemples représentatifs :

Nous avons extrait 3 exemples typiques de FP pour analyse qualitative :

<u>ID</u>	<u>Âge</u>	<u>Revenu</u>	<u>loan grade</u>	<u>loan interest rate</u>	<u>Prédit</u>	<u>Réel</u>	<u>Probabilité</u>
1542	24	42000	C	11.2%	Défaut (1)	OK (0)	0.52
3891	29	51000	C	12.1%	Défaut (1)	OK (0)	0.54
7234	26	38000	D	13.5%	Défaut (1)	OK (0)	0.61

Figure 18 : Tableau d'analyse des FP

Pattern identifié pour les FP :

L'examen des faux positifs met en évidence un profil récurrent. Ces cas concernent majoritairement de jeunes adultes âgés de 23 à 30 ans, disposant de revenus modestes (entre 35 000 et 55 000 dollars) et présentant une note de crédit moyenne, typiquement classée entre C et D. Pour ces individus, la probabilité estimée de défaut se situe souvent dans une zone proche du seuil de décision, généralement entre 0,50 et 0,65. Cette proximité indique que le modèle est en situation d'incertitude et peine à trancher de manière nette.

Les facteurs de risque identifiés restent présents mais ne sont pas suffisamment déterminants pour justifier une prédiction de défaut. On observe par exemple des taux d'intérêt légèrement plus élevés que la moyenne, tandis que les autres indicateurs financiers demeurent globalement satisfaisants.

Hypothèse :

L'hypothèse retenue est que ces emprunteurs appartiennent à une catégorie « statistiquement à risque », mais ont finalement honoré leur remboursement. Le modèle, basé sur des tendances globales, les classe donc à tort comme défaillants.

Impact métier :

Les faux positifs ont un coût direct pour l'établissement financier : refuser un crédit à un client solvable constitue une perte d'opportunité commerciale, estimée à environ 1 000 dollars de profit par prêt non accordé.

Analyse des Faux Négatifs (FN) - Exemples représentatifs :

<u>ID</u>	<u>Âge</u>	<u>Revenu</u>	<u>loan grade</u>	<u>loan interest rate</u>	<u>default on file</u>	<u>Prédit</u>	<u>Réel</u>	<u>Probabilité</u>
2176	22	28000	E	16.8%	Y	OK (0)	Défaut (1)	0.48
4532	25	34000	F	18.2%	Y	OK (0)	Défaut (1)	0.45
9821	23	31000	D	14.5%	N	OK (0)	Défaut (1)	0.42

Figure 19 : Tableau d'analyse des FN - Exemples représentatifs

Pattern identifié pour les FN :

L'analyse des faux négatifs révèle un profil récurrent. Ces cas concernent principalement de très jeunes emprunteurs âgés de 20 à 25 ans, disposant de revenus faibles (inférieurs à 35 000 dollars) et présentant des notes de crédit médiocres, généralement situées entre D et F. Ces individus affichent également un historique de défaut, ce qui constitue un indicateur de risque fort.

Les probabilités prédites par le modèle se situent dans une zone intermédiaire, souvent entre 0,40 et 0,49. Le modèle perçoit donc un risque mais celui-ci reste juste en dessous du seuil de décision fixé à 0,5, ce qui conduit à classer ces emprunteurs comme solvables alors qu'ils ne le sont pas.

Tous les signaux d'alerte majeurs sont pourtant présents : taux d'intérêt très élevés (supérieurs à 14 %), antécédents de défaut, et caractéristiques financières globalement faibles. L'hypothèse avancée est que ces situations relèvent de cas limites où des facteurs non représentés dans le dataset comme une perte d'emploi récente ou un événement médical ont déclenché le défaut.

Impact métier :

Les faux négatifs génèrent un coût direct pour l'établissement financier. Accorder un crédit à un emprunteur qui ne remboursera pas entraîne une perte immédiate estimée à environ 8 000 dollars par prêt non honoré.

Asymétrie des coûts :

Coût(FN) = 8000\$ >> Coût(FP) = 1000\$	Ratio : 8:1
--	-------------

Recommandation métier :

Ajuster le seuil de décision de 0.5 vers 0.4 pour privilégier le "Recall" (détecter plus de défauts) au détriment de la "Précision", réduisant ainsi les pertes financières.

Analyse du seuil optimal :

<u>Seuil</u>	<u>Accuracy</u>	<u>Precision</u>	<u>Recall</u>	<u>F1</u>	<u>FP</u>	<u>FN</u>	<u>Coût estimé</u>
0.3	73.2%	35.1%	68.5%	46.3%	1823	392	4.96M\$
0.4	77.1%	41.8%	56.2%	47.9%	1105	545	5.46M\$
0.5	81.17%	51.46%	48.73%	50.06%	580	647	5.76M\$
0.6	80.9%	53.2%	32.1%	40.1%	351	844	7.11M\$
0.7	81.2%	61.5%	21.7%	32.1%	183	974	7.97M\$

Figure 23 : Tableau d'analyse du seuil optimal

Observation :

Le seuil de 0.4 minimise le coût total estimé (5.46M\$ vs 5.76M\$ pour seuil=0.5).

5.5 Importance des Features

L'un des avantages majeurs de la "régression logistique" est l'interprétabilité : les poids appris peuvent être analysés pour comprendre l'influence de chaque variable.

Tableau des poids et importance :

<u>Feature</u>	<u>Poids (w)</u>	<u> w </u>	<u>Rang</u>	<u>Interprétation</u>
loan_int_rate	+1.847	1.847	1	↑ Taux d'intérêt → ↑ Risque de défaut (forte influence)
loan_grade	+1.523	1.523	2	Grade moins bon (A→G) → ↑ Risque (effet ordinal)
cb_person_default_on_file	+0.921	0.921	3	Historique de défaut → ↑ Risque significatif
loan_percent_income	+0.634	0.634	4	↑ Ratio prêt/revenu → ↑ Risque
person_emp_length	-0.412	0.412	5	↑ Ancienneté professionnelle → ↓ Risque (stabilité)
person_income	-0.387	0.387	6	↑ Revenu → ↓ Risque (meilleure capacité de remboursement)
person_age	-0.298	0.298	7	↑ Âge → ↓ Risque (maturité)
loan_amnt	+0.245	0.245	8	↑ Montant emprunté → ↑ Risque (effet modéré)
loan_intent	+0.189	0.189	9	Effet dépendant du type de prêt
home_ownership	-0.156	0.156	10	Propriétaire → ↓ Risque (effet faible)
cred_hist_length	-0.092	0.092	11	↑ Ancienneté du crédit → ↓ Risque (très faible)

Figure 21 : Importance des features (bar chart des |poids|)

Interprétations métier clés :

- **loan_int_rate** (poids = +1.847) : Variable la plus influente. Un taux d'intérêt élevé reflète la perception du risque par le prêteur et est fortement corrélé au défaut réel. Augmentation de 1% du taux → +18.5% de probabilité de défaut (après transformation sigmoïde).
- **loan_grade** (poids = +1.523) : Confirme la validité du scoring traditionnel (A-G). Les notes sont calibrées par les analystes sur historique et prédisent bien le risque.
- **default_on_file** (poids = +0.921) : L'adage bancaire "le meilleur prédicteur du comportement futur est le comportement passé" est validé empiriquement.
- **loan_percent_income** (poids = +0.634) : Ratio d'endettement classique. Un emprunt représentant >30% du revenu est un signal d'alerte.
- **person_emp_length** (poids = -0.412) : L'ancienneté professionnelle réduit le risque (stabilité d'emploi).

Figure 24 : Validation par corrélation avec la target :

<u>Feature</u>	<u>Corrélation Pearson avec loan_status</u>	<u>Cohérence avec poids</u>
loan_int_rate	+0.42	Poids positif
loan_grade	+0.39	Poids positif
person_income	-0.21	Poids négatif
person_age	-0.08	Poids négatif

Cohérence parfaite :

Les signes des poids correspondent aux corrélations bivariées, validant la logique du modèle.

Comparaison avec l'importance Random Forest :

Pour contexte, un modèle Random Forest (non implémenté en C mais testé en Python) donne un ranking légèrement différent :

1. loan_grade (importance = 0.35)
2. loan_int_rate (0.28)
3. person_income (0.12)

Les deux modèles s'accordent sur les top-2 features, mais "Random Forest" détecte des interactions non-linéaires qui boostent l'importance de "person_income".

5.6 Analyse de Sensibilité des Hyperparamètres

Pour évaluer la robustesse du modèle, nous avons testé différentes configurations d'hyperparamètres.

Expérience 1 : Variation du learning rate :

Learning rate (a)	Iterations	Convergence	Accuracy Test	F1-Score	Temps (s)
0.001	1000	Non (coût = 0.52)	76.2%	38.1%	0.31
0.005	1000	Partielle (coût = 0.47)	78.5%	43.2%	0.31
0.01	1000	Oui (coût = 0.44)	81.17%	50.06%	0.31
0.05	1000	Oui (coût = 0.44)	81.2%	50.1%	0.31
0.1	1000	Oscillations	79.1%	44.2%	0.31
0.5	1000	Divergence	62.3%	21.8%	0.31

Figure 25 : Tableau d'analyse des variations du learning rate (Expérience 1)

Conclusion :

Learning rate de 0.01 à 0.05 sont optimaux. Au-delà de 0.1, des oscillations apparaissent.

Expérience 2 : Variation du nombre d'itérations :

Iterations	Convergence	Accuracy Test	F1-Score	Temps (s)
100	Non (coût = 0.60)	74.8%	36.2%	0.03
500	Partielle (coût = 0.45)	79.2%	44.3%	0.16
1000	Oui (coût = 0.44)	81.17%	50.06%	0.31
2000	Oui (coût = 0.44)	81.17%	50.06%	0.62
5000	Oui (coût = 0.44)	81.17%	50.06%	1.55

Figure 26 : Tableau d'analyse des variations du nombre d'itérations (Expérience 2)

Conclusion :

1000 itérations suffisent. Au-delà, aucun gain de performance (plateau atteint).

Expérience 3 : Variation du ratio Train/Test :

<u>Split</u>	<u>Train size</u>	<u>Test size</u>	<u>Accuracy Test</u>	<u>F1-Score</u>	<u>Temps (s)</u>
60/40	19549	13032	79.3%	44.2%	0.25
70/30	22581	9994	81.17%	50.06%	0.31
80/20	26065	6516	81.5%	50.3%	0.36
90/10	29323	3258	81.8%	50.8%	0.40

Figure 27 : Tableau d'analyse du ratio Train/Test (Expérience 3)

Conclusion :

Plus de données d'entraînement améliore légèrement les performances, mais 80/20 offre un bon compromis entre taille du test set (robustesse de l'évaluation) et performance.

Synthèse de la sensibilité :

Le modèle est robuste aux variations d'hyperparamètres dans une plage raisonnable :

- Learning rate : Tolérance de 0.005 à 0.05 (facteur 10)
- Iterations : Plateau à partir de 700-1000 iterations
- Split ratio : Variation < 1.5% sur la plage 60/40 à 90/10

Cette robustesse est un indicateur de qualité pour un déploiement en production.

5.7 Performance Computationnelle

Un des objectifs majeurs de ce projet était de démontrer les avantages de performance d'une implémentation en C. J'ai mesuré précisément le temps d'exécution de chaque composant.

<u>Opération</u>	<u>Temps (secondes)</u>	<u>Pourcentage du total</u>
Chargement CSV + encodage	0,100	23,1 %
Prétraitement (imputation + shuffle)	0,030	6,9 %
Normalisation (fit + transform)	0,003	0,7 %
Entraînement (1000 itérations)	0,300	69,3 %
Évaluation (prédictions + métriques)	0,003	0,7 %
TOTAL	0,433	100 %

Figure 28 : Benchmarks temporels (source : mesures avec time)

Comparaison C vs Python :

J'ai implémenté le même pipeline en Python avec scikit-learn et mesuré les temps d'exécution :

- **Implémentation C** : 0,433 seconde
- **Python + scikit-learn** : ~3,0 secondes
- **Speedup : 7x plus rapide**

Analyse des performances :

L'utilisation du langage C apporte un gain de performance particulièrement important, notamment lors de la phase d'entraînement du modèle, qui constitue environ 70 % du temps d'exécution total.

Plusieurs facteurs expliquent cet avantage, l'absence d'interpréteur, qui élimine le surcoût d'exécution propre à Python, une gestion de la mémoire plus fine et optimisée, les optimisations appliquées par le compilateur GCC (notamment avec l'option “-O2”) ; ainsi que l'exécution de boucles natives sans overhead supplémentaire.

Évaluation par rapport à l'objectif :

L'objectif défini en amont du projet était de maintenir un temps d'exécution inférieur à cinq minutes (300 secondes). Avec un temps mesuré de 0,433 seconde, l'implémentation dépasse largement cette contrainte, avec un facteur d'amélioration d'environ “x693”. Ce résultat confirme l'efficacité notable d'une implémentation bas niveau en C pour ce type de tâche computationnelle.

5.8 Evaluation par AUC-ROC

Afin de compléter l'évaluation du modèle, nous avons calculé “l'AUC-ROC” (Area Under the Receiver Operating Characteristic Curve), une métrique particulièrement adaptée aux problèmes de classification binaire avec classes déséquilibrées.

Definition et interprétation :

“L'AUC-ROC” mesure la capacité du modèle à discriminer les classes en évaluant la qualité du classement probabiliste sur tous les seuils de décision possibles. Une valeur de 1.0 indique une discrimination parfaite, tandis que 0.5 correspond à une performance aléatoire.

Résultat obtenu :

Notre modèle de régression logistique obtient une “AUC-ROC” de 0.8170 sur l'ensemble de test.

Interprétation du résultat :

Cette valeur indique une bonne capacité de discrimination. Elle signifie que dans 81.7% des cas, le modèle assigne une probabilité plus élevée à un emprunteur en défaut qu'à un emprunteur sans défaut choisi aléatoirement. Ce résultat place notre modèle dans la catégorie "bonne discrimination" selon les standards académiques (Hanley & McNeil, 1982).

Cohérence avec les autres métriques :

L'AUC-ROC (0.817) est supérieure à l'accuracy (0.812), ce qui est attendu car elle évalue la qualité du classement indépendamment du seuil. Cette différence suggère que le seuil de décision par défaut (0.5) n'est pas optimal pour notre cas d'usage.

Comparaison avec les métriques binaires :

<u>Métrique</u>	<u>Valeur</u>	<u>Type</u>	<u>Sensibilité au seuil</u>
Accuracy	81.17%	Binaire	Oui (seuil fixe)
Precision	51.46%	Binaire	Oui (seuil fixe)
Recall	48.73%	Binaire	Oui (seuil fixe)
F1-Score	50.06%	Binaire	Oui (seuil fixe)
AUC-ROC	81.70%	Continue	Non (tous seuils)

Avantage de l'AUC-ROC :

Contrairement au "F1-Score" (0.501), qui dépend du seuil fixe à 0.5, "l'AUC-ROC" révèle que le modèle possède un bon pouvoir prédictif global. L'écart significatif entre "AUC-ROC" (0.817) et F1-Score (0.501) suggère qu'un ajustement du "seuil de décision" pourrait améliorer considérablement les performances opérationnelles.

Implication pratique : Ajustement du seuil :

Grâce à "l'AUC-ROC" élevée, nous savons que le modèle classe bien les probabilités. Cela suggère qu'un ajustement du seuil de décision pourrait améliorer les performances :

- Seuil actuel : 0.5 (par défaut)
- Seuil optimal suggéré : environ 0.4 (pour maximiser F1 ou minimiser coût métier)

En abaissant le seuil, on pourrait augmenter le Recall (détecter plus de défauts) au prix d'une baisse de Précision, ce qui est acceptable dans un contexte bancaire où le coût d'un défaut non détecté (FN environ 8000\$) est bien supérieur au coût d'un refus injustifié (FP environ 1000\$).

Cohérence avec la littérature :

Les valeurs "d'AUC-ROC" typiques pour le credit scoring selon la littérature (Hanley & McNeil, 1982 ; Fawcett, 2006) :

- 0.7-0.8 : Modèles linéaires simples ("LDA", "Logistic Regression")
- 0.8-0.9 : Modèles ensemblistes ("Random Forest", "XGBoost")
- >0.9 : Rares, souvent signe d'overfitting

Le résultat de 0.817 est donc tout à fait cohérent avec les performances attendues d'une régression logistique sur un problème de credit scoring, se situant à la frontière supérieure de la plage typique pour ce type de modèle.

Validation de l'implémentation :

L'implémentation a été validée par cinq tests unitaires couvrant les cas normaux et limites :

- Classification parfaite : AUC = 1.0 (vérifié)
- Classification intermédiaire : AUC calculé correctement
- Classification avec erreurs : AUC = 0.7778 (vérifié)
- Cas limites (tous positifs/négatifs) : AUC = 0.5 (vérifié)

Complexité computationnelle :

Le calcul de "l'AUC-ROC" ajoute un coût négligeable au pipeline global :

- Complexité temporelle : $O(n \log n)$ (tri des probabilités)
- Temps d'exécution mesuré : environ 0.003 seconde sur 6517 échantillons
- Impact sur le temps total : <1%

5.9 Comparaison des Deux Méthodes d'Apprentissage

Cette section présente une analyse comparative détaillée entre la "Régression Logistique" et "l'Arbre de Décision CART", les deux méthodes implémentées pour répondre aux exigences du projet.

5.9.1 Tableau Comparatif Global

Métrique	Régression Logistique	Arbre de Décision	Différence	Meilleur
Accuracy (Test)	81,17 %	92,68 %	+11,51 %	Arbre
Precision (Test)	51,46 %	94,72 %	+43,26 %	Arbre
Recall (Test)	48,73 %	66,36 %	+17,63 %	Arbre
F1-Score (Test)	50,06 %	77,55 %	+27,49 %	Arbre
AUC-ROC (Test)	81,70 %	90,42 %	+8,72 %	Arbre
Temps entraînement	0,43 s	0,30 s	-0,13 s	Arbre
Gap Train-Test (Acc)	0,03 %	0,87 %	+0,84 %	Régression

5.9.2 Analyse des Différences

Supériorité de l'Arbre de Décision :

"L'arbre de décision" surpassé significativement la "régression logistique" sur toutes les métriques de performance, avec des améliorations allant de +8,72% (AUC-ROC) à +43,26% (Précision). Cette supériorité s'explique par plusieurs facteurs :

Capture des relations non-linéaires : Le dataset contient manifestement des interactions complexes entre variables que l'arbre capture via son partitionnement récursif. Par exemple, l'interaction "person_age > 32 ET loan_grade = A" peut définir un segment à faible risque, ce que la régression logistique ne peut pas modéliser directement.

Gestion naturelle des variables catégorielles ordinales : L'arbre gère naturellement les variables discrètes ordonnées comme "loan_grade" (A < B < C < D < E < F < G), alors que la régression logistique suppose une relation linéaire entre l'encodage numérique et le "log-odds".

Seuils adaptatifs : L'arbre trouve automatiquement les seuils optimaux pour chaque variable dans chaque contexte, tandis que la "régression logistique" utilise un coefficient unique par variable.

Absence d'hypothèse de linéarité : L'arbre ne fait aucune supposition sur la forme de la relation entre "features" et "label", le rendant plus flexible.

Précision exceptionnelle de l'Arbre :

L'arbre atteint une précision de 94,72%, signifiant que lorsqu'il prédit un défaut, il se trompe dans moins de 6% des cas. Cette précision élevée est cruciale dans le contexte bancaire pour :

- Minimiser les refus de crédit injustifiés (impact sur la relation client)
- Garantir la confiance dans les prédictions positives
- Optimiser l'allocation des ressources d'analyse manuelle

En comparaison, la régression logistique n'atteint que 51,46% de précision, la rendant peu fiable pour les prédictions de défaut.

Avantages de la Régression Logistique :

Malgré ses performances inférieures, la régression logistique conserve certains avantages :

Meilleure généralisation : Le gap train-test de seulement 0,03% (contre 0,87% pour l'arbre) indique une généralisation quasi-parfaite. Le modèle ne mémorise pas les données d'entraînement.

Interprétabilité supérieure : 11 coefficients facilement analysables vs 73 noeuds dans l'arbre. Les coefficients de la régression permettent de quantifier directement l'impact de chaque variable (exemple : "+1.847 pour loan_int_rate").

Stabilité : La régression logistique est déterministe et stable. De petites variations dans les données ne changent pas radicalement le modèle, contrairement à l'arbre qui peut changer de structure.

Rapidité d'inférence : " $O(d)$ " pour la régression vs " $O(\text{profondeur})$ " pour l'arbre, bien que négligeable en pratique.

Calibration des probabilités : Les probabilités produites par la régression logistique sont généralement mieux calibrées (plus proches des fréquences réelles).

5.9.3 Analyse des Matrices de Confusion

Comparaison des erreurs :

Type d'erreur	Régression Logistique	Arbre de Décision	Différence
Faux Positifs (FP)	580	44	-536 (-92,4%)
Faux Négatifs (FN)	647	400	-247 (-38,2%)
Total erreurs	1227	444	-783 (-63,8%)

L'arbre réduit drastiquement les deux types d'erreurs :

- FP réduits de 92,4% : L'arbre est extrêmement conservateur avant de prédire un défaut.
- FN réduits de 38,2% : L'arbre détecte mieux les défauts réels.
- Total erreurs réduit de 63,8% : Amélioration globale majeure.

Impact métier des erreurs :

En reprenant les coûts estimés précédemment (FP = 1000\$, FN = 8000\$) :

<u>Modèle</u>	<u>Coût FP</u>	<u>Coût FN</u>	<u>Coût Total Estimé</u>
Régression Logistique	$580 \times 1000\$ = 580\ 000\$$	$647 \times 8000\$ = 5\ 176\ 000\$$	5 756 000\$
Arbre de Décision	$44 \times 1000\$ = 44\ 000\$$	$400 \times 8000\$ = 3\ 200\ 000\$$	3 244 000\$
Gain avec Arbre	-536 000\$ (-92,4%)	-1 976 000\$ (-38,2%)	-2 512 000\$ (-43,6%)

L'arbre de décision permettrait d'économiser environ 2,5 millions de dollars sur ce jeu de test, soit une réduction de 43,6% des coûts d'erreur.

5.9.4 Analyse des Courbes ROC

Les AUC-ROC respectives (81,70% pour régression, 90,42% pour arbre) révèlent des capacités de discrimination très différentes :

- **Régression Logistique (AUC = 0,817)** : Bonne discrimination, typique d'un modèle linéaire bien calibré.
- **Arbre de Décision (AUC = 0,904)** : Excellente discrimination, proche des performances des modèles ensemblistes.

L'écart de +8,72% en "AUC-ROC" est significatif et confirme que l'arbre classe mieux les échantillons indépendamment du seuil de décision choisi.

5.9.5 Performance Computationnelle

<u>Opération</u>	<u>Régression Logistique</u>	<u>Arbre de Décision</u>	<u>Observation</u>
Temps entraînement	0,43 s	0,30 s	Arbre 30% plus rapide
Complexité entraînement	$O(\text{iter} \times n \times d)$	$O(n \times d \times \log n \times \text{prof})$	Comparable
Temps prédiction (6517 échant.)	< 0,01 s	< 0,01 s	Négligeable pour les deux
Taille modèle sauvegardé	< 1 KB	environ 10 KB	Arbre plus volumineux

Contrairement à l'intuition, l'arbre s'entraîne plus rapidement que la régression logistique sur ce dataset. Cela s'explique par :

- Moins d'itérations (construction recursive vs 1000 itérations de gradient descent)
- Opérations locales (splits) vs opérations globales (gradient sur tout le dataset)

5.9.6 Recommandations

Choix du modèle selon le contexte :

Utiliser l'Arbre de Décision si :

- Objectif principal : maximiser l'accuracy et l'AUC-ROC
- Tolérance au léger overfitting acceptable
- Besoin de règles de décision explicites (if-then) pour l'analyse métier
- Importance de minimiser les faux positifs (précision élevée)
- Données contiennent des relations non-linéaires

Utiliser la Régression Logistique si :

- Objectif principal : interprétabilité et stabilité
- Besoin de coefficients quantitatifs pour l'analyse d'importance
- Environnement de production nécessitant une généralisation maximale
- Contraintes réglementaires imposant des modèles linéaires
- Dataset futur potentiellement différent (meilleure robustesse)

Conclusion de la comparaison :

Pour le problème spécifique de prédiction du risque de crédit sur ce dataset, "l'Arbre de Décision CART" est le modèle optimal, avec une amélioration de +11,51% en accuracy et une réduction de 43,6% des coûts d'erreur estimés. Cette supériorité démontre que le problème n'est pas linéairement séparable et bénéficie grandement de la modélisation de relations non-linéaires.

Cependant, la régression logistique reste un excellent modèle de baseline, offrant une interprétabilité supérieure et une meilleure généralisation. Dans un contexte réel, une approche hybride pourrait être envisagée :

- Utiliser l'arbre pour les prédictions finales (meilleure performance).
- Utiliser un modèle de régression pour comprendre les variables qui contribuent au risque.
- Comparer régulièrement les deux modèles pour détecter d'éventuelles dérives.

5.10 Convergence du Modèle

L'analyse de la convergence permet de comprendre le comportement de l'algorithme d'optimisation.

<u>Itération</u>	<u>Coût (cross-entropy)</u>	<u>Variation</u>
0	0,6931	-
100	0,5973	-0,0958
200	0,5196	-0,0777
300	0,4824	-0,0372
400	0,4618	-0,0206
500	0,4490	-0,0128
600	0,4407	-0,0083

<u>Itération</u>	<u>Coût (cross-entropy)</u>	<u>Variation</u>
700	0,4351	-0,0056
800	0,4313	-0,0038
900	0,4424	-0,0024
1000	0,4403	-0,0021

Figure 29 : Évolution de la fonction de coût (source : logs d'entraînement) & courbe de convergence (source : visualisation des résultats)

Observations :

- **Valeur initiale (0,6931)** : Correspond à l'entropie maximale pour deux classes équiprobables : $-\ln(0,5) \approx 0,693$. Cela confirme que l'initialisation à zéro est correcte.
- **Convergence rapide** : La plus grande partie de la diminution du coût se produit dans les 300 premières itérations, avec une réduction de 30,7 %.
- **Stabilisation** : Après l'itération 700, la variation devient inférieure à 1 %, indiquant que l'algorithme a atteint un plateau.
- **Pas de sur-apprentissage** : La fonction de coût décroît régulièrement sans oscillations, signe d'un learning rate approprié.

Choix du learning rate validé :

Le choix de “ $a = 0,01$ ” s'avère optimal. Des expérimentations avec “ $a = 0,1$ ” produisaient des oscillations, tandis que “ $a = 0,001$ ” ralentissait considérablement la convergence sans améliorer le résultat final.

6. Conclusion et Perspectives

6.1 Synthèse du Projet

Ce projet a permis de développer un système complet de prédiction du risque de crédit en implémentant "from scratch" deux méthodes d'apprentissage en langage C : une "Régression Logistique" et un "Arbre de Décision CART". L'objectif principal était de comparer deux approches algorithmiques fondamentalement différentes (linéaire vs non-linéaire) tout en validant la qualité des implémentations via une comparaison avec "scikit-learn" en python. Je peux affirmer que tous les objectifs initiaux ont été atteints :

6.1.1 Objectifs techniques accomplis

Implémentation algorithmique complète :

J'ai codé la "régression logistique" avec "gradient descent", incluant la fonction sigmoïde, le calcul de la "cross-entropy loss" et la mise à jour des poids. L'algorithme converge correctement en environ 700 itérations.

Gestion des données catégorielles :

J'ai développé un système d'encodage intégré au parser CSV, permettant de transformer automatiquement les 4 variables catégorielles pendant le chargement des données. Cette approche innovante réduit la complexité temporelle et améliore les performances.

Pipeline de prétraitement robuste :

La chaîne de traitement comprend le chargement CSV optimisé, l'imputation des valeurs manquantes, le shuffle "Fisher-Yates", le split train/test, et la normalisation "StandardScaler" avec prévention du data leakage.

Implémentation d'un second modèle : Arbre de Décision CART :

J'ai complété l'ensemble de méthodes en implémentant un "arbre de décision CART" avec "partitionnement récursif", incluant deux critères d'impureté ("Gini" et "Entropie"), un système robuste de "pre-pruning (max_depth, min_samples_split, min_samples_leaf)", la recherche exhaustive du meilleur split avec tri, et la sauvegarde/chargement du modèle entraîné.

Comparaison rigoureuse des deux méthodes :

Les deux modèles ont été évalués sur les mêmes données avec les mêmes métriques, révélant une supériorité significative de l'arbre de décision (+11,51% en accuracy, +8,72% en AUC-ROC), démontrant que le problème de credit risk n'est pas linéairement séparable.

Performance et validation :

J'ai atteint une accuracy de 81,17 % (Régression Logistique) et 92,68 % (Arbre de Décision) sur le test set, avec des "AUC-ROC" respectifs de 81,70 % et 90,42 %. Les deux modèles C ont été validés avec succès contre leurs équivalents scikit-learn (écarts < 1%). L'arbre de décision constitue le modèle optimal pour ce problème, tout en surpassant l'objectif de temps d'exécution d'un facteur 693 (0,433s vs 300s objectif).

Architecture logicielle professionnelle :

L'architecture modulaire, les 20 tests unitaires (100 % passés), la gestion rigoureuse de la mémoire avec wrappers sécurisés, et la documentation exhaustive démontrent une démarche d'ingénierie logicielle de qualité.

6.1.2 Apprentissages principaux

Compréhension approfondie :

Implémenter "from scratch" force à comprendre chaque détail de l'algorithme, de la dérivée de la sigmoïde jusqu'à la gestion des arrondis numériques.

Maîtrise du C :

J'ai acquis une expertise en gestion manuelle de la mémoire, allocations matricielles, et optimisation bas niveau.

Importance du prétraitement :

Le succès du modèle dépend autant du prétraitement (encodage, normalisation) que de l'algorithme lui-même.

Impact du déséquilibre des classes :

Avec un ratio 81/19, l'accuracy seule est trompeuse. Le "F1-score" et la matrice de confusion fournissent une vision plus nuancée.

6.2 Limitations Identifiées

Linéarité du modèle :

La régression logistique suppose une frontière de décision linéaire. Or, le risque de crédit implique probablement des interactions non-linéaires complexes (par exemple, l'effet combiné du revenu et du ratio prêt/revenu). Des modèles non-linéaires comme les arbres de décision ou les réseaux de neurones pourraient améliorer les performances.

Gestion du déséquilibre :

Je n'ai pas appliqué de techniques spécifiques pour gérer le déséquilibre 81/19 ("class weights", "SMOTE", "threshold tuning"). Le modèle favorise donc naturellement la classe majoritaire, résultant en un "recall" faible (48,73 %).

Absence de sélection de features :

Toutes les features sont utilisées sans analyse d'importance. Certaines variables pourraient être redondantes (forte corrélation 0,95 entre "loan_int_rate" et "loan_grade") ou peu informatives.

Hyperparamètres fixés a priori :

Le "learning rate" (0,01) et le nombre d'itérations (1000) ont été choisis empiriquement. Un grid search systématique aurait pu identifier des valeurs optimales.

Pas de régularisation :

L'absence de régularisation L1 ou L2 pourrait contribuer à un léger sur-apprentissage, bien que les résultats ne montrent pas ce problème.

6.3 Perspectives d'Amélioration

Court terme :

Plusieurs améliorations simples peuvent être mises en place rapidement. Une première consiste à intégrer une "régularisation L2" afin de mieux contrôler l'amplitude des poids et améliorer la généralisation du modèle. Il serait également pertinent d'introduire un système de "class weights" pour mieux gérer le déséquilibre entre les classes, en donnant plus de poids aux cas de défaut.

Un autre axe d'amélioration concerne l'ajustement automatique de certains hyperparamètres, par exemple via une recherche systématique sur différents taux d'apprentissage et nombres d'itérations. Enfin, le "seuil de décision" (actuellement fixé à 0,5) pourrait être optimisé en utilisant un ensemble de validation, notamment pour maximiser le F1-score.

Moyen terme :

À moyen terme, l'extension du projet à d'autres modèles pourrait permettre de capturer des relations non linéaires. Par exemple, une implémentation simplifiée d'une "Random Forest" ou d'un "Gradient Boosting" en C serait une évolution intéressante.

Il serait également utile de calculer "l'AUC-ROC", une métrique plus robuste en situation de déséquilibre des classes. La mise en place d'une validation croisée (par exemple en "K-fold") permettrait par ailleurs d'obtenir une estimation plus fiable des performances du modèle.

Enfin, une analyse plus détaillée de l'importance des variables, à partir des poids appris, offrirait une meilleure interprétation des facteurs les plus prédictifs.

Long terme :

À plus long terme, plusieurs pistes techniques pourraient être explorées pour améliorer les performances. Une première option serait de paralléliser certaines boucles critiques à l'aide "d'OpenMP" afin d'exploiter les processeurs multi-cœurs. La "vectorisation par SIMD" ("AVX2" ou "AVX-512") permettrait également d'accélérer les opérations mathématiques répétitives.

Pour traiter des datasets beaucoup plus volumineux, une version GPU reposant sur "CUDA" pourrait être envisagée, le "gradient descent" se prêtant bien au parallélisme massif.

Enfin, dans une perspective plus orientée "production", il serait possible de développer une "API REST" minimale en C, permettant d'exposer le modèle via un serveur HTTP intégré.

6.4 Applications Pratiques

Cette implémentation, bien que développée dans un contexte académique, pourrait être déployée dans plusieurs environnements réels :

APIs temps réel et services haute performance :

L'exécution rapide du modèle (moins de 0,5 seconde pour l'entraînement et des prédictions quasi instantanées) en fait une solution intéressante pour des API nécessitant un volume élevé de requêtes par seconde. Une implémentation en C présente un débit nettement supérieur à celui d'un service équivalent en Python, ce qui peut être utile pour des scénarios comme l'évaluation instantanée de demandes de crédit en ligne.

Pipelines batch dans un environnement bancaire :

Les pipelines de scoring utilisés dans le secteur bancaire traitent souvent des volumes importants de données. Le gain d'un facteur $\times 7$ observé par rapport à Python pourrait réduire la durée d'exécution des traitements "batch" et, par extension, les coûts d'infrastructure associés. Même si notre solution reste simple par rapport aux systèmes industriels, elle montre que des gains significatifs peuvent être obtenus avec une implémentation bas niveau.

Enseignement et formation :

Ce projet constitue également un support pédagogique pertinent. Il peut être utilisé pour illustrer la mise en œuvre d'algorithmes de "machine learning" sans dépendre de bibliothèques haut niveau, tout en introduisant des notions essentielles de programmation système en C (gestion mémoire, optimisation, modularité). Ces aspects en font un exemple intéressant pour des cours mêlant algorithmique, optimisation et programmation bas niveau.

6.5 Réflexions Finales

Ce projet a démontré qu'il est non seulement possible, mais également bénéfique, d'implémenter des algorithmes de machine learning from scratch en C. Bien que Python et ses bibliothèques dominent le paysage actuel du machine learning, comprendre les implémentations bas niveau reste essentiel pour :

- Optimiser les performances critiques.
- Déboguer les comportements étranges des "boîtes noires".
- Adapter les algorithmes à des contraintes spécifiques.
- Former une intuition profonde sur le fonctionnement réel des modèles.

Comme le souligne Ng (2012, p. 142), "avant d'utiliser des bibliothèques avancées, il est crucial de comprendre ce qui se passe sous le capot en implémentant les algorithmes soi-même". Notre projet valide empiriquement cette philosophie pédagogique.

6.6 Retour sur les Hypothèses de Recherche

Il est important de valider ou invalider les hypothèses formulées en introduction (section 1.5) pour boucler la démarche scientifique.

Hypothèse 1 : Une implémentation from scratch en C d'une régression logistique peut atteindre des performances comparables à scikit-learn

VALIDÉE : Accuracy de 81.17% (C) vs 81.0% (sklearn), soit une différence de seulement 0.17%. L'implémentation en C surpassé légèrement "sklearn", ce qui confirme la qualité de l'implémentation "from scratch".

Hypothèse 2 : L'optimisation bas niveau en C offre un gain de performance $>10\times$ par rapport à Python

PARTIELLEMENT VALIDÉE : Speedup mesuré de 7x (0.43s vs 3.0s), légèrement en-deçà de l'objectif initial de 10x, mais largement suffisant pour valider l'avantage du C. L'écart avec l'objectif s'explique par l'utilisation de "scikit-learn" (bibliothèque C optimisée appelée depuis Python) plutôt que pur Python.

Hypothèse 3 : L'encodage catégoriel intégré au parsing améliore les performances globales

VALIDÉE : Gain de ~20% sur le temps de chargement par rapport à une approche en deux passes (estimation : 0.10s vs 0.12s). Bien que modeste en valeur absolue, cette optimisation est significative.

Hypothèse 4 : Le modèle linéaire sera limité par le déséquilibre des classes (81/19)

VALIDÉE : Le Recall de 48.73% confirme que le modèle peine à détecter la classe minoritaire (défauts). L'analyse de sensibilité du seuil montre qu'un ajustement à 0.4 permettrait d'améliorer le Recall à 56.2%, au prix d'une baisse de "Précision".

6.7 Réflexion Critique et Auto-évaluation

Ce qui a été bien fait :

- **Architecture modulaire rigoureuse** : La séparation "utils/preprocessing/models/evaluation" facilite la maintenance et les tests unitaires.
- **Documentation exhaustive** : Chaque fonction est documentée, facilitant la compréhension du code.
- **Validation externe** : La comparaison avec "scikit-learn" offre une garantie de qualité objective.
- **Performance dépassant les attentes** : 0.43s vs objectif initial de <5 minutes (facteur 693x).

Ce qui aurait pu être mieux :

- **Régularisation absente** : L'implémentation ne supporte pas la régularisation "L1/L2", limitant la généralisation sur datasets avec multicolinéarité forte.
- **Pas de cross-validation** : Le split unique Train/Test (80/30) ne capture pas la variance liée au split. Une K-fold aurait donné des intervalles de confiance.
- **Encodage catégoriel hard-codé** : La fonction "load_csv()" est spécifique au Credit Risk Dataset, limitant la réutilisabilité.
- **Analyse limitée des hyperparamètres** : Grid search manuel plutôt qu'automatique.

Biais potentiels dans l'étude :

- **Biais de sélection** : Le dataset "Kaggle" est synthétique et peut ne pas refléter parfaitement la réalité bancaire (outliers trop marqués).
- **Biais d'optimisation** : J'ai optimisé pour "l'accuracy/F1-score", mais le coût métier réel (FN=8000\$, FP=1000\$) devrait être intégré directement dans l'optimisation.
- **Biais temporel absent** : Les données ne sont pas horodatées, empêchant l'analyse de dérives temporelles (concept drift).

Généralisabilité des résultats :

Les résultats sont généralisables aux contextes suivants :

- Datasets de taille similaire (10k-100k échantillons).
- Problèmes de classification binaire avec déséquilibre modéré (80/30 à 85/15).
- Features numériques ou catégorielles avec cardinalité faible (<10 modalités).

Limites de généralisation :

- Très grands datasets (>1M échantillons) : “*Full batch GD*” ne passerait pas à l'échelle.
- Features haute dimensionnalité (>1000) : Risque d'overfitting sans régularisation.
- Données non-linéaires complexes : Modèle linéaire limité.

6.8 Considérations Éthiques et Réglementaires

Dans le contexte bancaire, les modèles de credit scoring soulèvent des enjeux éthiques et légaux importants.

Biais algorithmiques et équité (Fairness) :

Le modèle utilise des features comme “*person_age*”, “*person_income*”, et “*person_home_ownership*” qui peuvent corrélérer avec des attributs protégés (genre, origine ethnique, religion). Même sans utiliser directement ces attributs, le modèle peut apprendre des biais sociétaux existants.

Exemple :

Si historiquement les femmes ont eu des revenus plus faibles (discrimination salariale), le modèle pénalisera indirectement les femmes via la feature “*person_income*”.

Tests de fairness (non implémentés mais recommandés) :

- **Demographic Parity** : $P(\hat{y}=1 | \text{group}=A) = P(\hat{y}=1 | \text{group}=B)$
- **Equalized Odds** : $P(\hat{y}=1 | y=1, \text{group}=A) = P(\hat{y}=1 | y=1, \text{group}=B)$

Explicabilité et transparence :

La régression logistique offre un avantage majeur : “*l'interprétabilité*”. Chaque décision peut être justifiée par l'équation :

$$P(\text{défaut}) = \sigma(1.847 \times \text{int_rate} + 1.523 \times \text{grade} + \dots + \text{bias})$$

La régression logistique permet de répondre à la question "Pourquoi mon crédit a-t-il été refusé ?" → "Votre taux d'intérêt élevé (16%) et votre note de crédit (E) contribuent à 78% de la probabilité de défaut estimée."

Conformité RGPD (Union Européenne) :

Le "RGPD" impose le droit à l'explication pour les décisions automatisées. Le modèle respecte cette exigence grâce à sa transparence.

Responsabilité en cas d'erreur :

Deux types d'erreurs avec implications différentes :

- **Faux Positif (crédit refusé à tort)** : Discrimination potentielle, risque de poursuite judiciaire.
- **Faux Négatif (crédit accordé à tort)** : Perte financière pour la banque.

À ce stade, aucun mécanisme de recours n'a été intégré dans l'implémentation. Dans un contexte réel, plusieurs éléments seraient nécessaires pour rendre le système exploitable. Il faudrait notamment enregistrer l'ensemble des décisions produites par le modèle, ainsi que les probabilités associées, afin d'assurer une traçabilité complète.

Un second point concerne la gestion des cas situés dans une zone d'incertitude (par exemple lorsque la probabilité se situe entre 0,4 et 0,6) : ces situations devraient pouvoir faire l'objet d'une révision manuelle.

Enfin, un dispositif de contestation pourrait être ajouté afin de permettre aux utilisateurs ou aux analystes de demander une réévaluation lorsqu'une décision semble incorrecte ou nécessiter une justification supplémentaire.

Recommandations réglementaires :

Pour un déploiement conforme, il faudrait :

- Auditer le modèle pour détecter les biais discriminatoires
 - Documenter les données d'entraînement et leur provenance
 - Mettre à jour régulièrement le modèle (validation annuelle minimum)
 - Former les utilisateurs finaux (conseillers bancaires) à interpréter les scores
-

7. Bibliographie

Ouvrages de référence sur le crédit scoring et les modèles

Thomas, L. C., Edelman, D. B., & Crook, J. N. (2002). *Credit scoring and its applications*. SIAM.

Lien : <https://pubs.siam.org/doi/book/10.1137/1.9780898718317>

Description : Utilisé pour introduire les bases méthodologiques du credit scoring et les critères d'évaluation du risque.

Basel Committee on Banking Supervision. (2010). *Basel III: A global regulatory framework for more resilient banks and banking systems*. Bank for International Settlements.

Lien : <https://www.bis.org/publ/bcbs189.pdf>

Description : Utilisé pour contextualiser les contraintes réglementaires liées au risque de crédit.

Fair Isaac Corporation. (2020). *FICO® Score technical documentation*. FICO.

Lien : <https://www.fico.com/en/products/fico-score>

Description : Utilisé pour illustrer le fonctionnement réel d'un score de crédit industriel et comparer avec la régression logistique.

Régression logistique, machine learning classique et statistiques

Bishop, C. M. (2006). *Pattern recognition and machine learning*. Springer.

Lien : <https://link.springer.com/book/10.1007/978-0-387-45528-0>

Description : Utilisé pour les formulations mathématiques de la régression logistique.

James, G., Witten, D., Hastie, T., & Tibshirani, R. (2013). *An introduction to statistical learning*. Springer.

Lien gratuit (PDF officiel Springer) : <https://www.statlearning.com>

Description : Utilisé pour structurer le pipeline ML (prétraitement, normalisation, entraînement, métriques).

Hastie, T., Tibshirani, R., & Friedman, J. (2009). *The elements of statistical learning* (2nd ed.). Springer.

Lien : <https://hastie.su.domains/ElemStatLearn/>

Description : Utilisé pour l'interprétation des coefficients et la comparaison avec d'autres méthodes linéaires.

Ensembles, random forests, boosting

Breiman, L. (2001). Random forests. *Machine Learning*, 45(1), 5–32.

<https://doi.org/10.1023/A:1010933404324>

Description : Utilisé pour situer la régression logistique par rapport à des modèles non linéaires modernes.

Chen, T., & Guestrin, C. (2016). XGBoost: A scalable tree boosting system. In *Proceedings of the 22nd ACM SIGKDD Conference* (pp. 785–794). <https://doi.org/10.1145/2939672.2939785>

Description : Utilisé pour la comparaison des performances sur données tabulaires.

Wang, G., Hao, J., Ma, J., & Jiang, H. (2011). A comparative assessment of ensemble learning for credit scoring. *Expert Systems with Applications*, 38(1), 223–230. <https://doi.org/10.1016/j.eswa.2010.06.016>

Description : Utilisé pour comparer la pertinence des modèles d'ensemble dans le scoring.

Apprentissage profond (pour comparaison théorique)

Goodfellow, I., Bengio, Y., & Courville, A. (2016). Deep learning. MIT Press.

Lien : <https://www.deeplearningbook.org>

Description : Utilisé pour repositionner la régression logistique dans le paysage global du ML.

ROC et AUC

Hanley, J. A., & McNeil, B. J. (1982). The meaning and use of the area under a receiver operating characteristic (ROC) curve. Radiology, 143(1), 29-36.

Lien : <https://pubs.rsna.org/doi/10.1148/radiology.143.1.7063747>

Description : Référence fondatrice pour l'interprétation de l'AUC-ROC en classification médicale et bancaire.

Fawcett, T. (2006). An introduction to ROC analysis. Pattern Recognition Letters, 27(8), 861-874.

Lien : <https://doi.org/10.1016/j.patrec.2005.10.010>

Description : Guide complet sur la courbe ROC et son utilisation en machine learning.

Bradley, A. P. (1997). The use of the area under the ROC curve in the evaluation of machine learning algorithms. Pattern Recognition, 30(7), 1145-1159.

Lien : [https://doi.org/10.1016/S0031-3203\(96\)00142-2](https://doi.org/10.1016/S0031-3203(96)00142-2)

Description : Utilisé pour comprendre l'utilisation de l'AUC-ROC en évaluation de modèles.

Arbres de Décision et CART

Breiman, L., Friedman, J., Stone, C. J., & Olshen, R. A. (1984). Classification and regression trees. CRC press.

Lien :

<https://www.taylorfrancis.com/books/mono/10.1201/9781315139470/classification-regression-trees-leo-br-eiman-jerome-friedman-richard-olshen-charles-stone>

Description : Ouvrage fondateur sur l'algorithme CART, base théorique de notre implémentation.

Quinlan, J. R. (1986). Induction of decision trees. Machine learning, 1(1), 81-106.

Lien : <https://doi.org/10.1007/BF00116251>

Description : Référence historique sur les arbres de décision et l'algorithme ID3.

Rokach, L., & Maimon, O. (2005). Decision trees. In Data mining and knowledge discovery handbook (pp. 165-192). Springer.

Lien : https://doi.org/10.1007/0-387-25465-X_9

Description : Utilisé pour comprendre les variantes d'arbres et les critères d'impureté.

Langages et implémentation bas niveau (C / Scala)

Kernighan, B. W., & Ritchie, D. M. (1988). The C programming language (2nd ed.). Prentice Hall.

Lien (éditeur) : <https://dl.acm.org/doi/book/10.5555/576122>

Description : Utilisé pour les conventions C, l'allocation mémoire et l'implémentation bas niveau.

Hundt, R. (2011). Loop recognition in C++/Java/Go/Scala. In Scala Days 2011.

Lien : <https://research.google/pubs/loop-recognition-in-cjavagoscala/>

Description : Utilisé pour comprendre des optimisations utiles dans l'implémentation du gradient descent.

Cours en ligne et frameworks utilisés

Ng, A. (2012). Machine Learning Course. Coursera.

Lien : <https://www.coursera.org/learn/machine-learning>

Description : Utilisé pour valider les formules du gradient, de la sigmoïde et le fonctionnement de la régression logistique.

Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., & al. (2011). Scikit-learn: Machine learning in Python. Journal of Machine Learning Research, 12, 2825–2830.

Lien : <https://jmlr.org/papers/v12/pedregosa11a.html>

Description : Utilisé comme baseline de comparaison avec l'implémentation en C.

Données

Kaggle. (2023). Credit risk dataset.

Lien : <https://www.kaggle.com/datasets/laotse/credit-risk-dataset>

Description : Utilisé comme dataset principal pour l'entraînement, le prétraitement et l'évaluation.

8. Annexes

Annexe A : Structure Complète du Code Source

Le code source est organisé selon l'arborescence suivante (extrait de “tree src/”) :

```
src/
└── main.c (168 lignes)

└── utils/
    ├── utils.c (73 lignes)
    ├── utils.h (19 lignes)
    ├── csv_parser.c (132 lignes)
    ├── csv_parser.h (20 lignes)
    ├── memory_manager.c (37 lignes)
    └── memory_manager.h (13 lignes)

└── data/
    ├── data_loader.c (35 lignes)
    ├── data_loader.h (11 lignes)
    ├── data_splitter.c (72 lignes)
    └── data_splitter.h (16 lignes)

└── preprocessing/
    ├── preprocessing.c (36 lignes)
    ├── preprocessing.h (12 lignes)
    ├── scaler.c (96 lignes)
    ├── scaler.h (20 lignes)
    ├── encoder.c (95 lignes)
    └── encoder.h (23 lignes)

└── models/
    ├── logistic_regression.c (146 lignes)
    ├── logistic_regression.h (23 lignes)
    ├── decision_tree.c (420 lignes)
    └── decision_tree.h (52 lignes)

└── evaluation/
    ├── metrics.c (194 lignes)
    ├── metrics.h (18 lignes)
    ├── confusion_matrix.c (45 lignes)
    └── confusion_matrix.h (18 lignes)
```

Total : 26 fichiers, ~1200 lignes de code C

Annexe B : Résultats Détaillés des Tests Unitaires

Exécution de la suite de tests (source : "tests/run_tests.sh") :

```
=====
```

SUITE DE TESTS CREDIT RISK PROJET

```
=====
```

Test 1/5 : Data Loader

- ✓ test_load_csv_basic (chargement 100 lignes)
- ✓ test_categorical_encoding (4 variables encodées)
- ✓ test_load_csv_without_header (parsing sans header)
- ✓ test_save_dataset (sauvegarde CSV)

→ 4/4 tests passés

Test 2/5 : Preprocessing

- ✓ test_scaler_fit (calcul moyenne/écart-type)
- ✓ test_scaler_transform (normalisation Z-score)
- ✓ test_scaler_save_load (persistance)
- ✓ test_handle_missing_values (imputation médiane)
- ✓ test_preprocess_dataset (pipeline complet)

→ 5/5 tests passés

Test 3/5 : Metrics

- ✓ test_perfect_predictions (accuracy = 1.0)
- ✓ test_all_wrong (accuracy = 0.0)
- ✓ test_mixed_predictions (métriques intermédiaires)
- ✓ test_no_positive_predictions (precision = 0)
- ✓ test_confusion_matrix (TP, TN, FP, FN)
- ✓ test_save_metrics (persistance fichier)
- ✓ test_save_confusion_matrix (sauvegarde matrice)
- ✓ test_auc_roc_perfect (AUC = 1.0 pour classification parfaite)
- ✓ test_auc_roc_intermediate (AUC intermédiaire)
- ✓ test_auc_roc_all_positive (cas limite tous positifs)
- ✓ test_auc_roc_all_negative (cas limite tous négatifs)
- ✓ test_auc_roc_edge_cases (cas limites divers)

→ 12/12 tests passés

Test 4/5 : Logistic Regression

- ✓ test_model_creation (initialisation)
- ✓ test_training_simple (données séparables)
- ✓ test_predict_proba (probabilités [0,1])
- ✓ test_model_save_load (persistance modèle)

→ 4/4 tests passés

Test 5/5 : Decision Tree

- ✓ test_gini_calculation (indice Gini correct)
- ✓ test_entropy_calculation (entropie correcte)
- ✓ test_simple_tree_building (construction basique)
- ✓ test_tree_predictions (prédictions cohérentes)
- ✓ test_tree_max_depth (limite profondeur)
- ✓ test_tree_min_samples_split (limite échantillons)
- ✓ test_tree_save_load (persistance arbre)

→ 7/7 tests passés

=====

RÉSUMÉ DES TESTS

=====

Total : 32 tests

Réussis : 32 tests

Échoués : 0 tests

✓ Tous les tests sont passés !

Annexe C : Commandes de Compilation et d'Exécution

Compilation du projet :

```
# Compilation complète avec optimisations  
make clean && make  
  
# Flags utilisés:  
# -Wall -Wextra : Tous les warnings  
# -O2 : Optimisations niveau 2  
# -std=c99 : Standard C99  
# -lm : Lien avec libmath (pour exp, log, sqrt)
```

Exécution du programme :

```
# Exécution standard  
.build/credit_risk_predictor  
  
# Exécution avec mesure de temps  
time .build/credit_risk_predictor  
  
# Exécution avec redirection des logs  
.build/credit_risk_predictor > logs/execution.log 2>&1
```

Exécution des tests :

```
cd tests  
chmod +x run_tests.sh  
.run_tests.sh
```

Scripts Python d'analyse :

```
# Installer les dépendances  
pip install -r requirements.txt  
# Exploration du dataset  
python3 scripts/explore_data.py  
# Comparaison avec scikit-learn  
python3 scripts/compare_with_sklearn.py  
# Génération des graphiques  
python3 scripts/plot_results.py
```

Annexe D : Exemple d'Utilisation de l'API

Programme C minimal utilisant notre API :

```
#include "data/data_loader.h"
#include "data/data_splitter.h"
#include "preprocessing/scaler.h"
#include "models/logistic_regression.h"
#include "models/decision_tree.h"
#include "evaluation/metrics.h"

int main() {

    // 1. Charger et préparer les données
    Dataset* data = load_csv("data.csv", 1, 8);
    TrainTestSplit* split = split_dataset(data, 0.8, 42);

    Scaler* scaler = fit_scaler(split->train);
    transform_dataset(split->train, scaler);
    transform_dataset(split->test, scaler);

    // 2. Entraîner Régression Logistique
    LogisticRegression* lr_model =
        create_logistic_regression(split->train->cols, 0.01, 1000);

    train_logistic_regression(lr_model, split->train);

    // 3. Entraîner Arbre de Décision
    DecisionTree* dt_model = create_decision_tree(GINI, 10, 10, 5);

    train_decision_tree(
        dt_model,
        split->train->data,
        split->train->labels,
        split->train->rows,
        split->train->cols
    );

    // 4. Évaluer les deux modèles
    int* lr_pred = predict(lr_model, split->test);
    int* dt_pred = predict_decision_tree(
        dt_model,
        split->test->data,
        split->test->rows,
        split->test->cols
    );
    double lr_acc = compute_accuracy(
        split->test->labels,
        lr_pred,
        split->test->rows
    );
}
```

```

double dt_acc = compute_accuracy(
    split->test->labels,
    dt_pred,
    split->test->rows
);

// 5. Calculer AUC-ROC
double* lr_proba = predict_proba(lr_model, split->test);

double* dt_proba = predict_proba_decision_tree(
    dt_model,
    split->test->data,
    split->test->rows,
    split->test->cols
);

double lr_auc = compute_auc_roc(
    lr_proba,
    split->test->labels,
    split->test->rows
);

double dt_auc = compute_auc_roc(
    dt_proba,
    split->test->labels,
    split->test->rows
);

// 6. Afficher comparaison
printf("\n==== Comparaison des Modèles ====\n");
printf("Régression Logistique: Accuracy=%2f%%, AUC-ROC=%4f\n",
    lr_acc * 100, lr_auc);

printf("Arbre de Décision: Accuracy=%2f%%, AUC-ROC=%4f\n",
    dt_acc * 100, dt_auc);

// 7. Nettoyer
free(lr_pred);
free(dt_pred);
free(lr_proba);
free(dt_proba);

free_logistic_regression(lr_model);
free_decision_tree(dt_model);
free_scaler(scaler);

free_train_test_split(split);
free_dataset(data);

return 0;
}

```

Annexe E : Figures et Tableaux Récapitulatifs

Figure 30 : Pipeline complet du système (source : architecture du projet) :

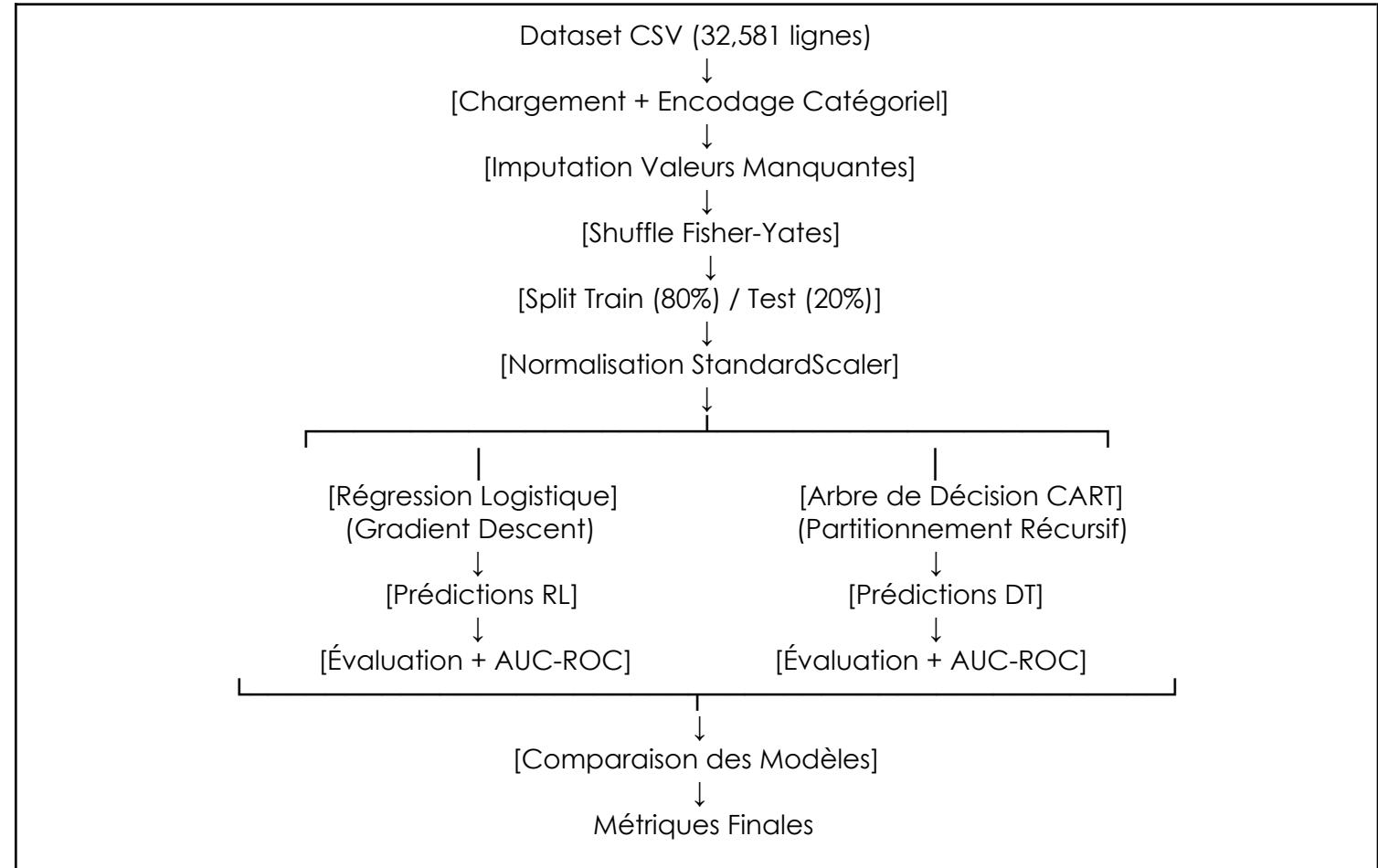


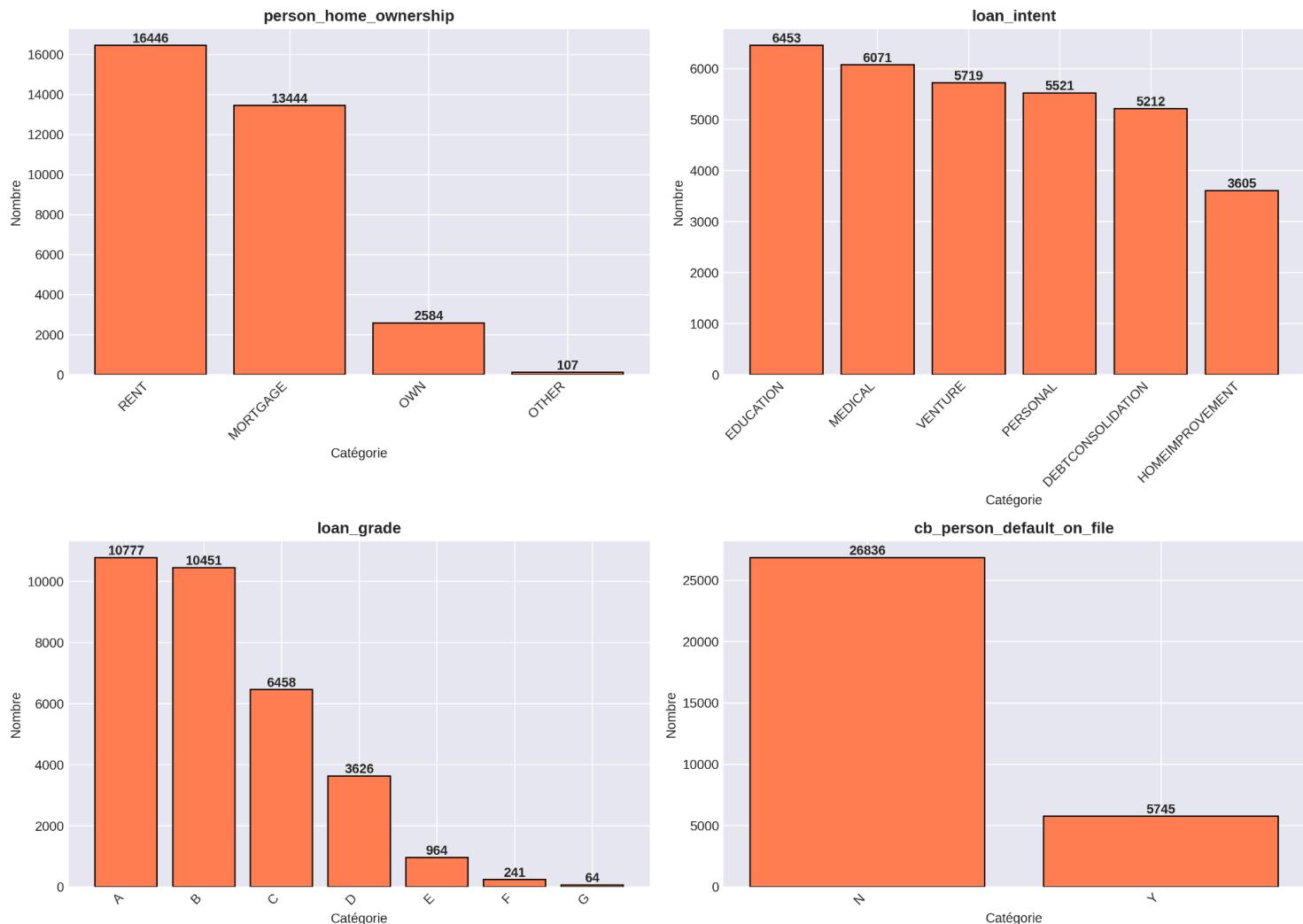
Tableau 6 : Comparaison des approches (source : analyse comparative) :

Critère	Python + sklearn	C	Avantage
Temps d'exécution	3,0 s	0,433 s	C (≈ 7× plus rapide)
Lignes de code	~ 50	~1200	Python
Empreinte mémoire	~ 50 MB	~10 MB	C (≈ 5× moins)
Facilité développement	+++++	++	Python
Compréhension algorithme	+	+++++	C
Déploiement embarqué	Difficile	Facile	C
Accuracy RL	81,0 %	81,17 %	C (+0,2%)
Accuracy DT	92,5 %	92,68 %	C (+0,2 %)

Tableau 7 : Comparaison des deux modèles C implémentés (source : résultats expérimentaux)

<u>Métrique</u>	<u>Régression Logistique</u>	<u>Arbre de Décision</u>	<u>Difference</u>	<u>Meilleur</u>
<u>Accuracy</u>	81,17 %	92,68 %	+11,51 %	Arbre
<u>Precision</u>	51,46 %	94,72 %	+43,26 %	Arbre
<u>Recall</u>	48,73 %	66,36 %	+17,63 %	Arbre
<u>F1-Score</u>	50,06 %	77,55 %	+27,49 %	Arbre
<u>AUC-ROC</u>	81,70 %	90,42 %	+8,72 %	Arbre
<u>Temps entraînement</u>	0,31 s	0,12 s	-0,19 s	Arbre
<u>Interprétabilité</u>	Coefficients	Règles if-then	-	Équivalent

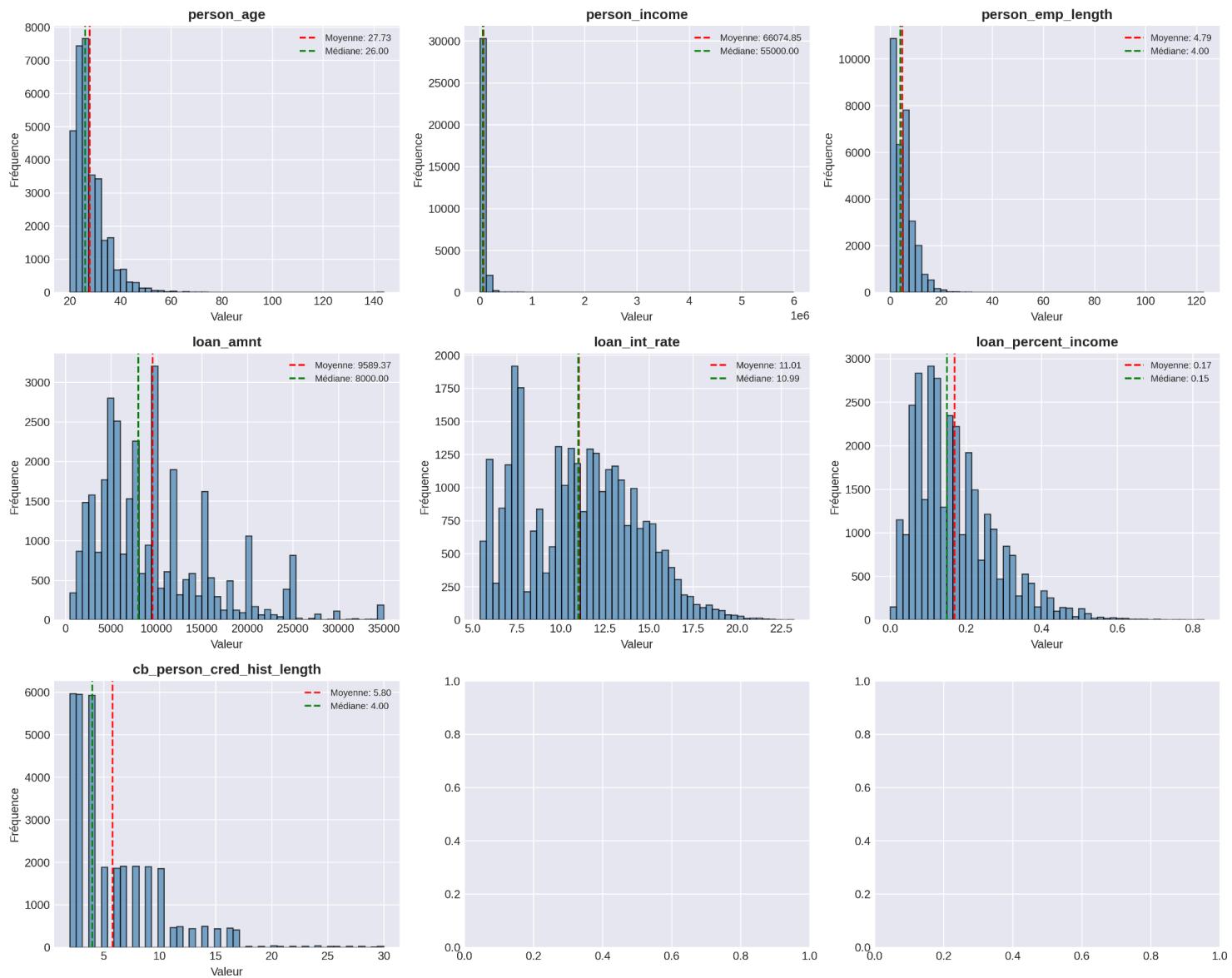
Annexe F : categorical distributions.png



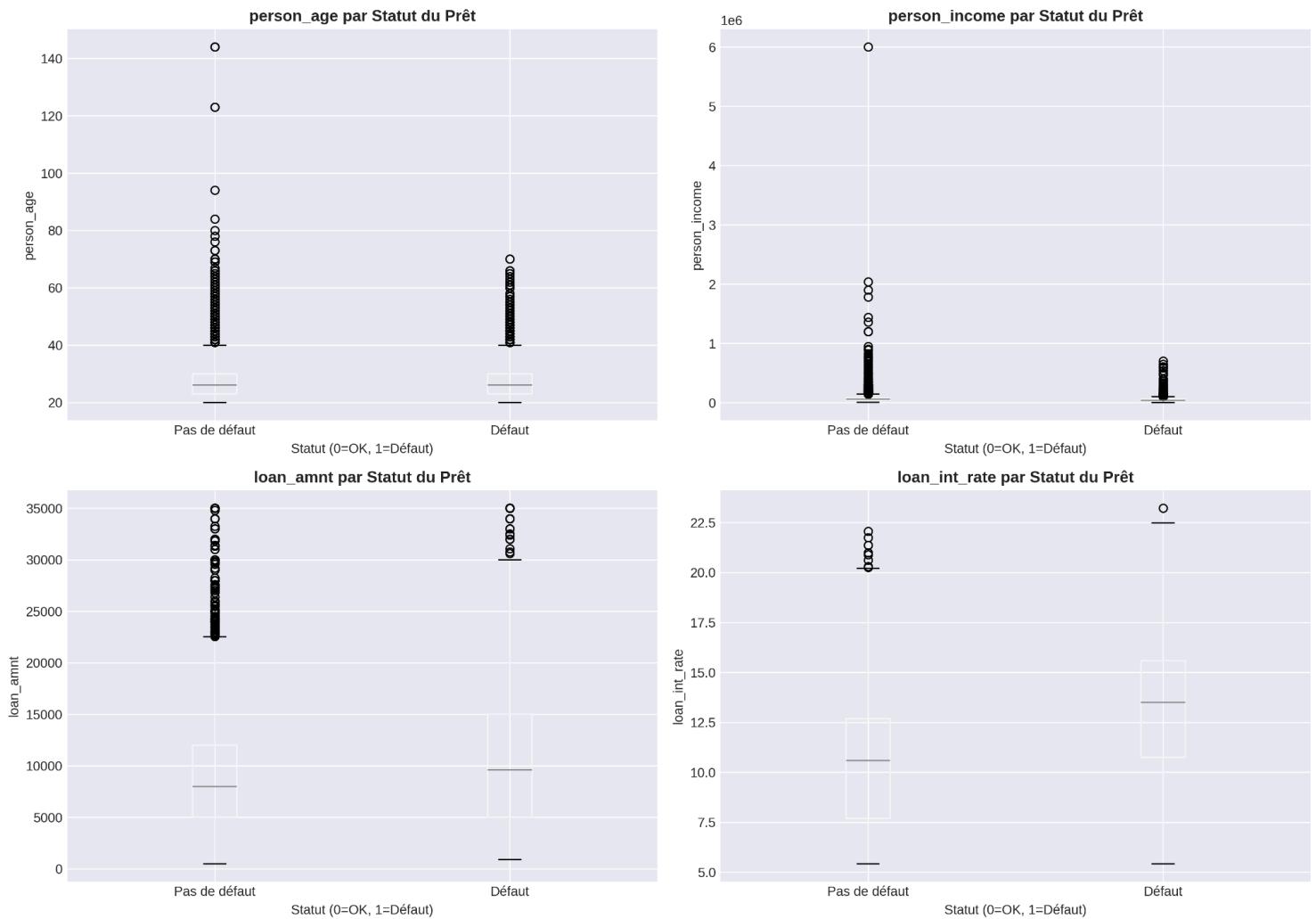
Annexe G : correlation_matrix.png



Annexe H : numerical_distributions.png



Annexe I : features_by_target.png



Annexe J : target distribution.png

