

RÉPUBLIQUE ALGÉRIENNE DÉMOCRATIQUE ET POPULAIRE
Ministère de l'Enseignement Supérieur et de la Recherche Scientifique
Université de Mohamed El-Bachir El-Ibrahimi - BBA
Faculté des Mathématiques et Informatiques



MÉMOIRE

Présenté en vue de l'obtention du diplôme

Master en Informatique

Spécialité : Ingénierie de l'Informatique Décisionnelle

Thème

Revue Systématique de la Littérature sur les Méthodes d'Apprentissage Automatique pour l'Analyse des Big Data avec une Étude de Cas.

Présenté par :

- SAADI Imane
- BORDJI Zahra

Soutenu le : 20/06/2024, devant la commission d'examen suivante :

Mr. Salah MAACHE	Enseignant à l'université de BBA	Président
Mr. Makhlouf NAILI	Enseignant à l'université de BBA	Examineur
Mr. Oussama SENOUCI	Enseignant à l'université de BBA	Encadrant

Promotion 2023 / 2024

**“Je n’ai jamais rêvé de succès,
J’ai travaillé pour ça.”**

- Estee Lauder -

Résumé

Avec l'explosion du volume de données générées quotidiennement, le Big Data est devenu un enjeu majeur pour de nombreux domaines. L'importance du Big Data réside dans sa capacité à offrir des insights précieux et à faciliter la prise de décisions informées. Cependant, pour exploiter pleinement ce potentiel, il est essentiel d'utiliser des techniques d'apprentissage automatique qui permettent de traiter, analyser et extraire des informations pertinentes à partir de ces vastes ensembles de données.

Ce mémoire présente une revue systématique de la littérature sur les méthodes d'apprentissage automatique pour le traitement et l'analyse des Big Data, accompagnée d'une étude de cas. L'étude couvre diverses techniques d'apprentissage supervisé, non supervisé, semi-supervisé et profond, ainsi que leurs algorithmes, notamment les SVM, la régression, les arbres de décision, les réseaux de neurones convolutifs (CNN) et récurrents (RNN), ainsi que des techniques de clustering comme HDDC, SOM, FCM et k-means. Une méthodologie rigoureuse a été utilisée pour identifier et évaluer les études pertinentes. Dans l'étude de cas, l'algorithme k-means a été appliqué au jeu de données Iris, démontrant son efficacité pour identifier des patterns dans les données.

En conclusion, cette revue systématique a mis en évidence différentes techniques d'apprentissage automatique pour traiter les Big Data ainsi que leurs limites. Grâce à cette étude, nous avons identifié les problèmes actuels, ce qui permet d'explorer des pistes d'amélioration et de résolution de ces problèmes à l'avenir.

Mots clés : Big Data, Apprentissage Automatique, Revue Systématique, K-means, SVM, CNN, RNN, Clustering.

Abstract

With the explosion of data volume generated daily, Big Data has become a major concern across various domains. The significance of Big Data lies in its ability to provide valuable insights and facilitate informed decision-making. However, to fully harness this potential, it is essential to employ machine learning techniques that can process, analyze, and extract relevant information from these vast datasets.

This thesis presents a systematic literature review on machine learning methods for Big Data processing and analysis, accompanied by a case study. The study covers various supervised, unsupervised, semi-supervised, and deep learning techniques, along with their algorithms, including SVM, regression, decision trees, convolutional neural networks (CNN), recurrent neural networks (RNN), and clustering techniques such as HDDC, SOM, FCM, and k-means. A rigorous methodology was employed to identify and evaluate relevant studies. In the case study, the k-means algorithm was applied to the Iris dataset, demonstrating its effectiveness in identifying patterns within the data.

In conclusion, this systematic review has highlighted different machine learning techniques for addressing Big Data challenges and their limitations. Through this study, current issues have been identified, paving the way for exploring avenues for improvement and resolution of these issues in the future.

Keywords : Big Data, Machine Learning, Systematic Review, K-means, SVM, CNN, RNN, Clustering.

ملخص

مع انفجار حجم البيانات التي يتم إنشاؤها يوميًا، أصبحت البيانات الضخمة تحديًا رئيسيًا للعديد من المجالات. تكمن أهمية البيانات الضخمة في قدرتها على تقديم رؤى قيمة وتسهيل اتخاذ القرارات المستنيرة. ومع ذلك، للاستفادة الكاملة من هذا الإمكان، من الضروري استخدام تقنيات التعلم الآلي التي تمكن من معالجة وتحليل واستخراج المعلومات ذات الصلة من هذه المجموعات الضخمة من البيانات.

يقدم هذا البحث مراجعة منهجية للأدبيات حول طرق التعلم الآلي لمعالجة وتحليل البيانات الضخمة، مصحوبة بدراسة حالة. تغطي الدراسة تقنيات متنوعة للتعلم الموجه وغير الموجه ونصف الموجه والتعلم العميق، بالإضافة إلى خوارزمياتها مثل آلات الدعم الشعاعي، والانحدار، وأشجار القرار، والشبكات العصبية الالتفافية والعصبية المتكررة، وكذلك تقنيات التجميع مثل التحليل العنقودي والخرائط الذاتية التنظيم والتحلل العنقودي الضبابي والتجميع المتعدد. تم استخدام منهجية صارمة لتحديد وتقييم الدراسات ذات الصلة. في دراسة الحالة، تم تطبيق خوارزمية التجميع المتعدد على مجموعة بيانات "إيريس"، مما يثبت فعاليتها في تحديد الأنماط في البيانات.

في الختام، سلطت هذه المراجعة المنهجية الضوء على تقنيات التعلم الآلي المختلفة لمعالجة البيانات الضخمة وكذلك حدودها. بفضل هذه الدراسة، تم تحديد المشاكل الحالية، مما يسمح باستكشاف سبل التحسين وحل هذه المشاكل في المستقبل.

كلمات مفتاحية: البيانات الضخمة، التعلم الآلي، مراجعة منهجية، ك-المتوسطات، آلات المتجهات الداعمة، الشبكات العصبية الالتفافية، الشبكات العصبية التكرارية، التجميع

Remerciements

Nous remercions tout d'abord Dieu de nous avoir donné la force nécessaire pour achever notre travail. Nous exprimons nos sincères remerciements et notre gratitude à notre encadrant, le Dr. Senouci Oussama, pour tous ses efforts, exprimés dans les termes les plus élogieux.

En guise de reconnaissance, nous tenons à exprimer nos sincères remerciements à toutes les personnes qui ont contribué de près ou de loin au bon déroulement de nos stages de fin d'étude et à l'élaboration de ce modeste travail. Nous exprimons notre profonde gratitude envers tous les enseignants pour la qualité de leur enseignement, leurs conseils et leur soutien indéfectible envers tous les étudiants. Nous tenons également à remercier l'ensemble du personnel pour leur patience, leurs conseils avisés et leur intérêt manifeste pour nos travaux. Dans l'impossibilité de citer tous les noms, nos sincères remerciements vont à tous ceux et celles qui, de près ou de loin, ont permis par leurs conseils et leurs compétences la réalisation de ce mémoire.

Enfin, nous ne voulons pas oublier de remercier tout le corps professoral de l'Université Mohamed El Bachir El Ibrahimi pour le travail considérable qu'il accomplit afin de créer les conditions les plus favorables pour le déroulement de nos études. Merci à tous.

Dédicace

Ce mémoire est dédiée :

A mes chers parents, Nacir et Fairouz, qui ont été mes piliers et mes guides tout au long de ce parcours, je dédie ce projet avec un amour infini. Votre soutien inébranlable et vos encouragements constants ont été la lumière qui a éclairé mon chemin vers la réussite.

A ma sœur bien-aimée, Mounia, dont le soutien inconditionnel et l'amour sans faille ont été ma force et mon inspiration. C'est grâce à toi que je suis ici aujourd'hui, prête à franchir cette étape importante de ma vie. A mon frère cher, Fouad, je dédie ce travail avec tout mon cœur, reconnaissant pour vos sacrifices, vos conseils et votre amour qui m'ont permis d'atteindre ce moment. Vous êtes mes héros, et je vous aime plus que tout.

★ Zahra ★

Dédicace

Je suis honoré d'exprimer ma profonde gratitude à Dieu pour m'avoir accordé la force et le courage d'accomplir cet humble travail.

Je dédie ce travail à mes parents, Karima et Djamel. Votre amour inconditionnel, vos innombrables sacrifices et votre foi inébranlable en mes capacités ont été la base solide sur laquelle j'ai construit ma vie. Merci de m'avoir enseigné que la persévérance et le dévouement sont les clés du succès.

J'envoie également mon hommage à toute ma famille et mes amis. À mes grands-parents, Mama Houriya et Grand-père Brahim, votre sagesse et votre héritage ont enrichi mon esprit et mon âme. Je vous suis reconnaissant pour votre amour éternel et vos bénédictions silencieuses qui m'ont toujours accompagné.

Je voudrais également dédier ce travail à mes sœurs Fayza, Chaima, Ikhlassa, Ritadje et Ghofrane. Votre soutien indéfectible et vos encouragements constants ont été une source de force et de motivation. Je vous en serai éternellement reconnaissant.

À mes amies Chaima et Asmaa, pour les moments de joie partagés, les défis que nous avons relevés ensemble et les sourires partagés même dans les moments les plus sombres.

À mon fiancé, votre patience sans fin et votre soutien indéfectible ont été un phare dans ma vie. Votre présence a rendu chaque défi plus facile à surmonter et chaque victoire plus douce.

Et à tous ceux qui ont contribué directement ou indirectement à ce voyage, Votre présence, vos paroles et vos actions ont façonné mon parcours. Ce travail est autant le vôtre que le mien et je vous en suis profondément reconnaissant.

★ *Imane* ★

Table des matières

Abstract	ii
Liste des figures	xiii
Abréviations	xiv
Introduction Générale	1
1 Généralités sur le Contexte	4
1.1 Introduction	4
1.2 Big Data et son analyse	5
1.2.1 Définition du Big Data	5
1.2.2 Évolution historique du Big Data	6
1.2.3 Caractéristiques du Big Data	7
1.2.4 Structuration du Big Data	8
1.2.5 Traitement du Big Data	10
1.3 Apprentissage Automatique (Machine Learning)	12
1.3.1 Définition de l'Apprentissage Automatique	12
1.3.2 Évolution de l'Apprentissage Automatique	13
1.3.3 Types de l'Apprentissage Automatique	14
1.3.4 Applications de l'Apprentissage Automatique dans le traitement des Big Data	18
1.3.5 Outils utilisés dans le traitement des Big Data	20
1.4 Défis et limites actuels	22
1.4.1 Défis rencontrés dans l'analyse des Big Data	23
1.5 Conclusion	24

2	État de l'art	26
2.1	Introduction	26
2.1.1	Contexte de l'étude	26
2.1.2	Importance de l'application de l'apprentissage automatique dans l'analyse des Big Data	27
2.2	Application de l'Apprentissage Automatique dans l'Analyse des Big Data	27
2.2.1	Fondements de l'Apprentissage Automatique	28
2.2.2	Techniques d'Apprentissage Automatique pour l'Analyse des Big Data	31
2.2.3	Algorithmes Populaires et leur Pertinence	37
2.2.4	Défis et Opportunités	38
2.2.5	Gestion du Volume et de la Variété des données	38
2.2.6	Scalabilité des Algorithmes d'Apprentissage Automatique	39
2.3	Revue des Études Classiques	39
2.3.1	Comparaison des Études Classiques	45
2.3.2	Limitations des Études Classiques	46
2.4	Motivation pour l'utilisation de la Revue Systématique de la Littérature (SLR)	47
2.4.1	Avantages de la SLR par rapport aux enquêtes Classiques	47
2.4.2	Importance de Surmonter les Limitations Identifiées	47
2.5	Conclusion	48
3	Revue Systématique de la Littérature	49
3.1	Introduction	49
3.2	Méthodologie d'étude SLR	50
3.3	Motivation et Contributions	51
3.4	Méthodologie de Révision	51
3.4.1	Questions de Recherche	51
3.4.2	Méthodologie de Recherche	52
3.5	Classification des Solutions d'Apprentissage Automatique pour l'Ana- lyse des Big Data	57
3.5.1	Les Algorithmes d'Apprentissage Automatique de Classification pour le Traitement des Big Data	58

3.6	Analyse de Résultats et Discussion	71
3.6.1	Réponses au Questions de Recherche :	72
3.6.2	Considérations et Directions Futures	79
3.7	Conclusion	81
4	Application et Évaluation d'une Méthode d'Apprentissage Automatique pour l'Analyse du Big Data : Étude de Cas	82
4.1	Introduction	82
4.1.1	Contexte et objectifs de l'implémentation dans l'étude de cas . .	83
4.1.2	Importance de Colab et Python dans le traitement du Big Data	83
4.2	Configuration de l'Environnement de Développement	84
4.2.1	Configuration de Google Colab	84
4.3	Description de l'Étude de Cas et Préparation des Données	85
4.3.1	Présentation de l'étude de cas	85
4.3.2	Dataset	85
4.3.3	Nettoyage et pré-traitement des données	86
4.3.4	Analyse exploratoire des données	86
4.4	Sélection et Adaptation de la Méthode d'Apprentissage Automatique .	87
4.4.1	Justification du choix de l'algorithme	87
4.4.2	Adaptation de l'algorithme aux spécificités de l'étude de cas . .	87
4.5	Implémentation de la Méthode	88
4.5.1	Détails de l'implémentation de l'algorithme	88
4.5.2	Méthodes d'évaluation utilisées	93
4.6	Résultats	95
4.6.1	Visualisation des Clusters	101
4.6.2	Évaluation des Clusters	104
4.7	Comparaison avec "Agglomerative" Clustering	104
4.8	Discussion	106
4.8.1	Analyse critique des performances du modèle	107
4.8.2	Implications des résultats pour l'étude de cas	107
4.9	Conclusion	108
4.9.1	Résumé des contributions principales et de l'impact sur les objectifs de l'étude de cas	109

4.9.2	Réflexion sur l'importance de l'approche choisie pour le traitement et l'analyse de Big Data	109
	Conclusion Générale	111
	Bibliographie	113

Table des figures

1.1	Caractéristiques du Big data.	7
1.2	Étapes du traitement des Big Data [1].	10
1.3	Types de l'Apprentissage Automatique.	14
1.4	Apprentissage Supervisé.	15
1.5	Apprentissage non Supervisé.	15
1.6	Exemple d'un Réseau de Neurones.	17
1.7	Outils utilisés dans le traitement des Big Data	21
2.1	Techniques de l'Apprentissage Automatique.	28
2.2	SVM	32
2.3	La configuration formelle d'un neurone.	33
2.4	Architecture d'un Réseau de Neurones Convolutionnels	36
3.1	Méthodologie d'étude Systématique.	50
3.2	Classification des Algorithmes de d'Apprentissage Automatique.	57
3.3	Nombre d'articles étudiés en fonction de l'année de publication.	72
4.1	Mesures des sépales et des pétales de différentes espèces d'iris	95
4.2	Distribution des caractéristiques des données Iris	96
4.3	Relation entre les caractéristiques	97
4.4	Matrice de corrélation des caractéristiques	98
4.5	Caractéristiques moyennes des iris dans chaque cluster	99
4.6	Évaluation des performances du clustering K-Means	100
4.7	Analyse de la courbe d'inertie pour le clustering K-Means.	100
4.8	la distribution des clusters dans chaque classe	101
4.9	Visualiser les clusters avec K-Means (2D)	102
4.10	Visualisation 3D des clusters avec K-Means dans le jeu de données Iris	103

4.11 Comparaison entre K-Means et AgglomerativeClustering	105
4.12 Comparaison des résultats entre K-Means et AgglomerativeClustering .	106

Abréviations

- **SGBD** : Système de Gestion de Base de Données
- **BD** : Big Data
- **ADD** : Analyse des Données
- **ML** : Machine Learning
- **IA** : Intelligence Artificielle
- **IoT** : Internet des Objets
- **GFS** : Google File System
- **SQL** : Structured Query Language
- **NoSQL** : Not Only SQL
- **RN** : Réseaux de Neurones
- **NLP** : Natural Language Processing
- **SLR** : Revue Systématique de la Littérature
- **SVM** : Support Vector Machine
- **MMH** : Hyperplan Marginal Maximal
- **DT** : Arbres de Décision
- **RNA** : Réseau Neuronal Artificiel
- **ACP** : Analyse en Composantes Principales
- **CNN** : Réseaux de Neurones Convolutionnels
- **RNN** : Réseaux de Neurones Récurents
- **MLP** : Perceptron Multicouche
- **MBD** : Massives Données Mobiles

- **MFFR** : Multifeature Fusion Retrieval
- **DPWSS** : Support par Processus Pondéré par la Sensibilité Différentielle
- **LASSO** : Least Absolute Shrinkage and Selection Operator
- **CART** : Classification and Regression Trees
- **FSOM** : Algorithme de Fuzzy Self-Organizing Map
- **HDCC** : High-Dimensional Data Clustering
- **EM** : Expectation-Maximization
- **FCM** : Fuzzy C-Means
- **GMM** : Modèle de Mélange Gaussien
- **GTC** : Graph-Based Temporal Classification
- **SimPLE** : Simulated Proximal Learning with Error
- **LSTM** : Long Short-Term Memory
- **GRU** : Gated Recurrent Unit
- **BiLSTM** : Bidirectional Long Short-Term Memory
- **CRF** : Conditional Random Field
- **SAEs** : Stacked Autoencoders
- **GANs** : Réseaux Génératifs Antagonistes

Introduction Générale

Dans le cadre de notre projet de fin d'études en Master 2 Ingénierie d'Informatique Décisionnelle, nous avons choisi d'explorer un domaine crucial dans le contexte de l'analyse des données contemporaine : **la revue systématique de la littérature sur les méthodes d'apprentissage automatique pour le traitement et l'analyse des Big Data**, accompagnée d'une étude de cas. Notre travail vise à fournir une compréhension approfondie des différentes approches existantes dans ce domaine en pleine expansion.

Cette étude examine en profondeur les différentes approches existantes dans le domaine de l'apprentissage automatique appliqué aux Big Data, mettant en lumière leurs avantages, leurs limites, et leurs domaines d'application spécifiques. En outre, notre travail cherche à appliquer ces méthodes à un cas concret afin d'évaluer leur efficacité et leur pertinence dans un contexte réel. Cette introduction présente le contexte et la problématique de notre mémoire, ainsi que les objectifs et les contributions de notre travail. Enfin, nous détaillons l'organisation de notre rapport.

Problématique & Motivations

Face à la croissance exponentielle des données et aux limites des méthodes traditionnelles de revue de littérature, il est crucial de mener une revue systématique et approfondie des méthodes d'apprentissage automatique dédiées au traitement des Big Data. Cette revue doit permettre de comprendre les différentes approches existantes, leurs avantages, leurs limites, et leurs applications spécifiques. De plus, il est important d'évaluer l'efficacité et la pertinence de ces méthodes dans des études de cas spécifiques

afin de fournir des recommandations pour leur utilisation future.

La fulgurante évolution technologique a provoqué une prolifération massive de données, ce qui a renforcé la nécessité d'utiliser efficacement les méthodes d'apprentissage automatique pour traiter et extraire des informations pertinentes des Big Data. Cette étude vise à examiner et à évaluer les capacités des méthodes d'apprentissage automatique dans ce contexte, en fournissant une analyse approfondie des approches existantes et en appliquant une seule méthode à un cas concret pour en évaluer l'efficacité et l'impact réel.

Objectifs & Contributions

Notre étude vise à explorer en profondeur les diverses méthodes d'apprentissage automatique utilisées pour le traitement et l'analyse des Big Data. À travers une revue systématique de la littérature, nous examinons un large éventail d'approches, allant des techniques classiques aux méthodes les plus récentes, afin de fournir un aperçu complet du paysage de l'apprentissage automatique dans ce domaine. En analysant les publications existantes, nous identifions les forces et les faiblesses de chaque méthode. De plus, notre recherche s'engage à aller au-delà de la simple exploration théorique en appliquant une méthode à un cas d'étude concret.

Pour ce faire, nous avons choisi l'algorithme k-Means, largement utilisé dans le domaine du clustering, et nous l'avons testé sur l'ensemble de données Iris, un ensemble de données classique en apprentissage automatique. Cette évaluation nous permettra de mieux comprendre l'efficacité de l'algorithme k-means dans un contexte de traitement des big data, basée sur une étude de cas spécifique. Bien que notre étude se concentre sur une analyse approfondie des méthodes d'apprentissage automatique, notre évaluation pratique fournit des informations précieuses sur la pertinence de ces méthodes pour relever les défis posés par les Big Data.

Organisation du mémoire

Le mémoire est organisé en quatre chapitres comme suit :

— **Chapitre 1 : Généralités sur le Contexte**

Ce chapitre a pour objectif de fournir une présentation détaillée des caractéristiques des 5V du Big Data, ainsi que de la structuration des données. Il aborde également les différentes étapes du traitement des Big Data et les types d'apprentissage automatique utilisés.

— **Chapitre 2 : État d'art**

Ce chapitre vise à réaliser une étude comparative des méthodes classiques, en mettant en évidence leurs limites, ainsi que des différents algorithmes d'apprentissage automatique en fonction de leur scalabilité. Cette analyse permettra de mieux comprendre les avantages et les inconvénients de chaque approche dans le contexte du traitement des Big Data.

— **Chapitre 3 : Revue systématique de la littérature**

Ce chapitre présente une revue systématique de la littérature sur les méthodes d'apprentissage automatique pour le traitement et l'analyse des Big Data. Cette étude vise à explorer et à synthétiser les approches actuelles, en mettant en lumière les techniques les plus efficaces et les défis associés à la gestion des vastes ensembles de données.

— **Chapitre 4 : Application et Évaluation d'une Méthode d'Apprentissage Automatique pour l'Analyse du Big Data : Étude de Cas**

Ce chapitre présente une étude de cas sur l'application de l'algorithme k-means au jeu de données Iris. L'objectif est de démontrer l'efficacité de cet algorithme de clustering pour identifier des motifs et des structures au sein des données.

Chapitre 1

Généralités sur le Contexte

1.1 Introduction

Le premier chapitre explore le monde du Big Data, en mettant en évidence son omniprésence dans notre quotidien et son rôle essentiel dans l'analyse de données. Cette analyse continue, alimentée par un flux constant de données, est cruciale dans divers secteurs, influençant l'avenir de manière significative. Des exemples concrets montrent l'impact du Big Data, notamment dans la conception de véhicules autonomes, le développement de médicaments efficaces et l'amélioration des processus décisionnels grâce à l'intelligence artificielle. Le chapitre définit ensuite le Big Data, en soulignant sa grande envergure et les défis qu'il pose à l'analyse traditionnelle. L'évolution historique du Big Data, des années 1950 à aujourd'hui, est examinée, mettant en lumière les avancées technologiques et les changements majeurs dans la collecte et le stockage des données.

Les caractéristiques principales du Big Data, représentées par les 5V (Volume, Vitesse, Variété, Véracité, et Valeur), sont expliquées en détail. Le chapitre distingue également les données structurées, semi-structurées et non structurées, offrant ainsi une vue d'ensemble sur la diversité des informations traitées dans le domaine du Big Data. Enfin, le processus de traitement du Big Data est décrit comme une série d'étapes essentielles, depuis la définition des objectifs jusqu'à l'intégrité et la sécurité des données. Une transition vers l'apprentissage automatique (Machine Learning) est introduite, soulignant son rôle crucial dans l'analyse des mégadonnées. Ce premier chapitre établit ainsi les bases nécessaires pour explorer en profondeur les aspects complexes et capti-

vants du Big Data et de son analyse.

1.2 Big Data et son analyse

Le flux continu de données alimente constamment notre monde, soumettant chaque aspect à une analyse approfondie. L'analyse des données se révèle cruciale dans tous les domaines, permettant d'extraire des significations profondes des données collectées et ouvrant ainsi la voie à un avenir remarquable. Des exemples concrets de cette influence se manifestent dans la conception de véhicules autonomes sécurisés, le développement de médicaments à l'efficacité optimale, et l'amélioration de nos processus décisionnels grâce à l'intelligence artificielle. Bien que l'acronyme de l'analyse des données (ADD) puisse différer de celui du Big Data, il reste essentiel pour donner un sens cohérent à toutes les informations que nous recueillons.

1.2.1 Définition du Big Data

Avant d'explorer les implications du Big Data, il est primordial de définir ce terme de plus en plus répandu. En 2011, l'auteur K. Crawford a caractérisé le Big Data dans le domaine scientifique comme un ensemble de données d'une envergure considérable, parfois comparable à la taille du peta ou de l'exabyte. Cette ampleur rend ardue, voire impossible, l'analyse et l'exploitation des données avec les outils traditionnels, nécessitant ainsi l'utilisation de superordinateurs et de structures réseau complexes.

Il est crucial de souligner que ce qui rend le Big Data distinct n'est pas simplement sa taille, mais surtout la capacité à exploiter efficacement cette masse de données. D'autres experts ont étendu la compréhension du Big Data en intégrant les notions de vitesse et de variété. En référence au travail de Beyer [2], une définition éclairante émerge : *"Le Big Data représente une quantité importante d'informations générées à une vitesse élevée, présentant une grande variété, nécessitant ainsi de nouvelles méthodes de traitement pour permettre une prise de décision améliorée, la découverte d'informations pertinentes et l'optimisation des processus"*.

1.2.2 Évolution historique du Big Data

Dans cette section consacrée à l'historique du Big Data, nous observons une évolution significative depuis les années 1950 jusqu'à nos jours [3]. Bien que le terme « Big Data » soit devenu courant récemment, les bases du processus de collecte et de stockage des données remontent aux années 1950, avec l'utilisation des premiers ordinateurs commerciaux. Jusqu'aux années 1990, la progression des données était lente en raison du coût élevé des ordinateurs, des supports de stockage, et de la limitation des sources de génération de données.

Au cours de cette période, les données étaient principalement structurées pour répondre aux besoins des systèmes d'information opérationnels. Cependant, des avancées technologiques, telles que l'émergence des bases de données parallèles dans les années 1980, ont montré une capacité significative à traiter et stocker progressivement des données. L'avènement du World Wide Web au début des années 1990 a entraîné une explosion de données, marquant le début de trois générations majeures de Big Data.

La première génération, Big Data 1.0 (1994–2004), a été caractérisée par le commerce électronique et le développement de techniques d'exploration pour analyser les activités en ligne. Google a joué un rôle clé en introduisant des modèles de programmation tels que GFS et MapReduce pour gérer les données à l'échelle de l'Internet. L'analyse du contenu Web était divisée en utilisabilité, structure Web et contenu Web, avec l'utilisation de techniques telles que la recherche d'information et le traitement du langage naturel.

La deuxième génération, Big Data 2.0 (2005–2014), a été influencée par le Web 2.0 et les médias sociaux. Les grandes entreprises ont lancé des projets de Big Data, et l'analyse des réseaux sociaux est devenue populaire, utilisant des données non structurées pour comprendre les sentiments et opinions des utilisateurs. Les outils d'analyse ont adopté des services basés sur le cloud avec un coût flexible.

La troisième génération, Big Data 3.0 (2015–présent), englobe les données des deux générations précédentes. Les applications Internet des objets (IoT) ont contribué massivement avec des données variées telles que des images, de l'audio et de la vidéo. L'IoT a créé un environnement technologique alimenté par des données générées par des dispositifs et capteurs uniques, favorisant le partage et la collaboration autonome sur les réseaux.

Le processus de gestion des données dans le Big Data suit généralement un cycle de vie comprenant la gestion des données (acquisition, extraction, nettoyage, intégration et agrégation), suivi de l'analyse des données (modélisation, analyse et interprétation).

1.2.3 Caractéristiques du Big Data

Les caractéristiques du Big Data, résumées par les « 5 V du Big Data », mettent en évidence les défis et les opportunités liés à la gestion, l'analyse et l'exploitation de vastes volumes de données dans divers domaines. La Figure 1.1 illustre ces caractéristiques.

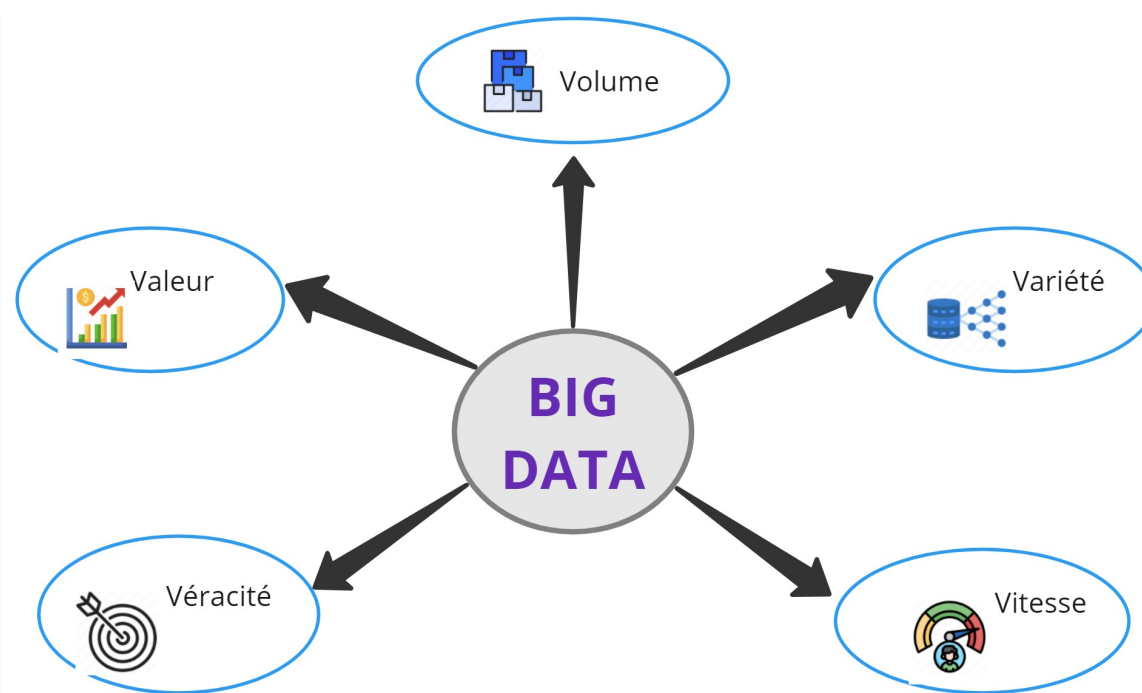


FIGURE 1.1 – Caractéristiques du Big data.

1. **Volume** : Le volume de données fait référence à la quantité de données générées par des entreprises ou des individus, une caractéristique souvent associée au Big Data. Dans tous les secteurs d'activité, les entreprises sont confrontées au défi de gérer la croissance constante du volume de données générées quotidiennement. La norme actuelle inclut des catalogues dépassant les 10 millions de produits, devenant ainsi plus courants que rares. Certains clients, qui gèrent non seulement des produits mais aussi leur propre base de clients, peuvent facilement accumuler des volumes dépassant le téraoctet de données [4].
2. **Vélocité (Vitesse)** : La rapidité se réfère à la fréquence à laquelle les données

sont produites, capturées et partagées. En raison des avancées technologiques récentes, tant les consommateurs que les entreprises génèrent désormais des quantités de données significatives en des laps de temps beaucoup plus courts [4].

3. **Variété** : La variété se rapporte à la diversité des formats de données. Le format traditionnel est celui de la base de données relationnelle, où l'information est stockée selon un schéma rigide et organisé, comme dans un tableau. Cette forme de données est qualifiée de « structurée ». Cependant, de nos jours, plus de 80 % (certains analystes avancent même 95 à 99 %!) des données sont considérées comme "non-structurées", c'est-à-dire qu'elles ne peuvent pas être facilement classées . Parmi ces données figurent le texte, Les courriels, les photos, les vidéos, les enregistrements vocaux, les messages, etc. Le Big Data offre la capacité de regrouper toutes ces données et de les analyser [5].
4. **Véracité** : La véracité aborde la question de la fiabilité limitée et du désordre inhérent à la donnée. Souvent, les données présentent un manque de qualité et de précision, les rendant difficiles à contrôler. L'une des missions du Big Data est d'apporter une certaine organisation à ce chaos, non pas en structurant les données, mais plutôt en organisant leur accès et en permettant l'association d'analyses correspondante aux besoins des utilisateurs [5].
5. **Valeur** : Du point de vue des entreprises, le « V » la plus cruciale est la valeur. Dans l'idéal, le «Big Data» devrait apporter une «Valeur» significative. Les équipes dédiées à l'analyse et à la recherche doivent prendre en compte, concevoir, et déterminer l'ampleur de cette valeur. Dans le contexte des activités commerciales, la valeur figure parmi les premières propriétés discutées et une estimation préliminaire de la valeur est souvent projetée dès le début d'un projet lié au «Big Data». En facilitant le développement de l'infrastructure nécessaire, le "Big Data" joue un rôle essentiel dans la création du socle sur lequel peuvent être implantés l'apprentissage automatique et l'intelligence artificielle [6].

1.2.4 Structuration du Big Data

Une structure de données englobe une collection de valeurs de données, de leurs relations, ainsi que des fonctions ou opérations applicables à ces données. Elle constitue un moyen d'organiser et de stocker des données dans un ordinateur, facilitant ainsi leur

consultation et leur modification de manière efficace.

Dans le contexte du Big Data, les données collectées, stockées et traitées peuvent provenir de divers domaines et être générées par différentes sources de données hétérogènes, engendrant ainsi une masse de données de types variés, qu'elles soient structurées, non structurées, ou semi-structurées :

1.2.4.1 Données Structurées

Les données structurées font référence à des données présentant un format et une longueur définis, facilement stockables, analysables et hautement organisées. En d'autres termes, ces données sont arrangées dans une structure reconnaissable pour permettre des requêtes efficaces et la récupération d'informations à des fins d'organisation.

Un exemple concret de données structurées est une base de données relationnelle utilisant le langage de requête structuré (SQL). Cette base de données contient des nombres organisés, des dates, des groupes de mots, et des nombres appelés chaînes/texte. Du fait de la structure transparente de la base de données, elle peut être explorée à l'aide d'algorithmes de recherche simples et directs, classés par type de données dans le contenu réel [7].

1.2.4.2 Données Semi-Structurées

En dépit de l'appellation parfois trompeuse de "données non structurées", un document texte peut être perçu, selon divers angles, comme un objet doté de structure. Certains documents renferment des éléments typographiques et de mise en page qui agissent comme des balises "flexibles" permettant d'identifier des composants cruciaux du document. Si certains documents sont qualifiés de "faiblement structurés" en raison du faible nombre d'indices typographiques ou de mises en page révélateurs de la structure, d'autres, à l'opposé, sont parfois catégorisés comme "semi-structurés" en raison de la présence d'éléments de mise en forme plus prononcés [8].

1.2.4.3 Données non Structurées

En général, les données non structurées ne sont pas organisées de manière tabulaire (comme dans une base de données relationnelle) et ne suivent pas un schéma prédéfini.

Elles peuvent provenir de sources telles que les réseaux sociaux, les flux de médias, les documents texte, les courriels, les images, les vidéos, etc. Ces données ne sont pas faciles à traiter avec des méthodes traditionnelles, car elles peuvent manquer de cohérence dans leur format et leur organisation [9].

1.2.5 Traitement du Big Data

Le traitement du Big Data comprend généralement les étapes suivantes, illustrées dans la Figure 1.2 :

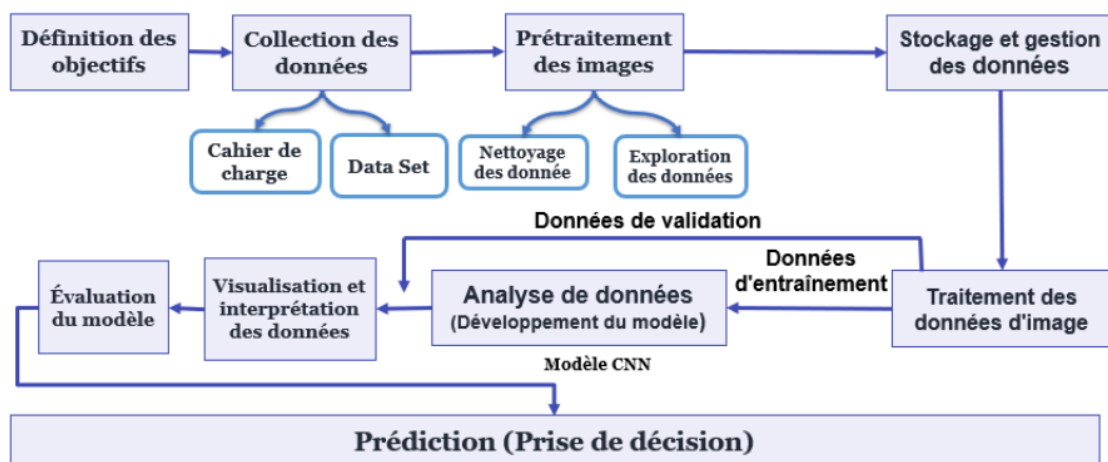


FIGURE 1.2 – Étapes du traitement des Big Data [1].

1. Définition des objectifs de traitement :

La première étape pour traiter les Big Data implique l'identification des objectifs à atteindre. Quelles informations ou questions aimeriez-vous extraire des données ? Il est crucial de définir ces objectifs de manière précise pour guider les étapes suivantes du processus [1].

2. Collecte et acquisition de données :

La première étape de la gestion des Big Data est la collecte de données. Cette phase implique l'extraction d'informations à partir de diverses sources, telles que les bases de données internes, les capteurs connectés, les réseaux sociaux, les fichiers journaux et les données GPS, entre autres. Ces données peuvent revêtir des formes structurées, semi-structurées ou non-structurées. Il est essentiel de collecter une quantité suffisante de données pour obtenir des informations pertinentes,

tout en évitant une collecte excessive qui pourrait rendre le processus d'analyse coûteux et complexe [1].

3. **Nettoyage et préparation des données :**

La préparation des données est essentiel en raison de la diversité des sources, nécessitant une gestion efficace pour éviter le gaspillage de stockage. Les étapes clés incluent l'intégration pour une vue cohérente des données, le nettoyage pour identifier et corriger les erreurs, et la réduction de la redondance pour améliorer la qualité des données. Ces méthodes visent à réduire les coûts de stockage tout en améliorant la précision de l'analyse. Cependant, il est crucial d'équilibrer les avantages et les coûts, car certaines techniques peuvent entraîner une augmentation de la charge de calcul [10, 11].

4. **Stockage et gestion des données :**

Le stockage des Big Data englobe l'ensemble des processus liés à la conservation, la gestion et l'accès des vastes ensembles de données. Il implique divers systèmes tels que le stockage massif et distribué, ainsi que des mécanismes spécifiques comme les bases de données NoSQL. Au-delà de l'entreposage physique, le stockage inclut la conception de systèmes pour garantir un accès efficace, une gestion cohérente et la résolution des défis liés aux données massives [11].

5. **Analyse des données :**

L'analyse des données consiste à utiliser des méthodes statistiques appropriées pour extraire des informations significatives à partir d'ensembles de données massifs. Elle vise à concentrer, extraire et affiner des données utiles cachées, identifiant ainsi les tendances et les lois inhérentes. Cette approche joue un rôle crucial dans le développement national, la compréhension des besoins clients et la prévision des tendances du marché. Bien que des méthodes d'analyse traditionnelles soient toujours pertinentes, l'analyse des mégadonnées se démarque en tant que technique spécifique pour traiter ces vastes ensembles de données [12].

6. **Visualisation et interprétation des données :**

La visualisation et l'interprétation des données constituent les étapes suivantes dans la gestion des Big Data. Pendant cette étape, les résultats de l'analyse sont présentés sous forme de graphiques, de tableaux, et de rapports pour une compréhension accrue. Des outils de visualisation tels que des graphiques, des

tableaux croisés dynamiques et des cartes peuvent être mobilisés pour représenter les données de manière claire et concise [1].

7. Intégrité et sécurité des données :

La phase finale de la gestion des Big Data implique d'assurer l'intégrité et la sécurité des données. Il est primordial de garantir la qualité, l'exactitude et la confidentialité des données en utilisant des techniques de cryptage, d'authentification et de sauvegarde afin de les préserver contre les risques de sécurité tels que les violations de données, les cyberattaques et les erreurs humaines. La mise en place de politiques de sécurité des données est essentielle pour régir l'accès aux données, les mots de passe et les autorisations. Il est également crucial de surveiller régulièrement les données pour détecter les anomalies et les incohérences [1].

1.3 Apprentissage Automatique (Machine Learning)

L'apprentissage automatique, également connu sous le nom de Machine Learning (ML), est une branche de l'intelligence artificielle visant à doter un ordinateur de la capacité d'apprendre à partir de données. Contrairement à une programmation classique, où les instructions sont explicitement définies par un humain, l'apprentissage automatique permet à un système informatique d'acquérir des connaissances, de reconnaître des schémas et de prendre des décisions de manière autonome, tout en nécessitant une intervention humaine minimale. Les algorithmes d'apprentissage automatique sont utilisés dans divers domaines tels que la reconnaissance d'images et de la parole, le traitement du langage naturel, ainsi que dans le processus de prise de décisions [13].

1.3.1 Définition de l'Apprentissage Automatique

L'apprentissage automatique est la discipline qui permet aux ordinateurs de fonctionner sans nécessiter une programmation explicite. En 1959, Arthur Samuel a formulé sa définition en tant que "domaine d'étude conférant aux ordinateurs la capacité d'apprendre sans être explicitement programmés". En 1997, Tom Mitchell a défini l'apprentissage automatique comme l'étude d'algorithmes informatiques conçus pour effectuer des tâches sans programmation explicite.

En 2020, l'apprentissage automatique est caractérisé comme une approche informatique visant à développer des modèles capables de résoudre des problèmes complexes à partir de données. Cette discipline a acquis une utilisation croissante dans divers aspects de notre quotidien, allant de la classification d'images à la prévention des maladies, en passant par la détection de cyber-attaques dans le domaine de la cyber-sécurité et dans la nouvelle ère industrielle. Son influence est significative dans notre vie quotidienne, et l'intégration de ces algorithmes vise à améliorer notre quotidien en offrant des services et des applications capables de prendre des décisions autonomes optimales [14, 15].

1.3.2 Évolution de l'Apprentissage Automatique

Dans cette section dédiée à l'historique de l'apprentissage automatique, nous plongerons dans un passé riche et étendu, influencé de manière significative par les progrès majeurs dans les domaines de l'intelligence artificielle et de la statistique [1].

Dès les premières années de la décennie 1950, des chercheurs éminents tels que Marvin Minsky et John McCarthy ont jeté les bases des concepts préliminaires de l'apprentissage automatique. Les années 1960 ont marqué l'introduction d'algorithmes de régression linéaire et d'arbres de décision, ouvrant la voie à la capacité de faire des prédictions à partir de données. L'émergence des algorithmes d'apprentissage supervisé, comme les réseaux de neurones, a caractérisé les années 1980, offrant des solutions aux problèmes complexes.

Au début des années 1990, le développement des algorithmes d'apprentissage non supervisé, notamment les algorithmes de regroupement (clustering), a permis d'explorer les structures cachées des données. Le développement des algorithmes d'apprentissage profond, exploitant des réseaux de neurones profonds, a connu un développement significatif dans les années 2000, permettant de traiter de manière plus efficace des données complexes telles que les images et les textes. La décennie des années 2010 a été marquée par l'émergence des algorithmes d'apprentissage par renforcement, autorisant les machines à s'adapter à leur environnement en apprenant à partir de retours d'information. Plus récemment, les algorithmes d'apprentissage automatique distribué, en temps réel, ont gagné en popularité, en particulier dans la prise de décision et l'optimisation, résolvant ainsi des problèmes complexes en temps réel.

En résumé, l'historique de l'apprentissage automatique illustre la transformation de ce

domaine en un pilier essentiel de l'intelligence artificielle, avec une évolution constante des algorithmes pour traiter des données de plus en plus complexes et résoudre des problèmes de plus en plus difficiles.

1.3.3 Types de l'Apprentissage Automatique

Grâce au succès de l'apprentissage automatique, diverses méthodes ont vu le jour pour exploiter son potentiel. La Figure 1.3 illustre les différents types d'apprentissage automatique.

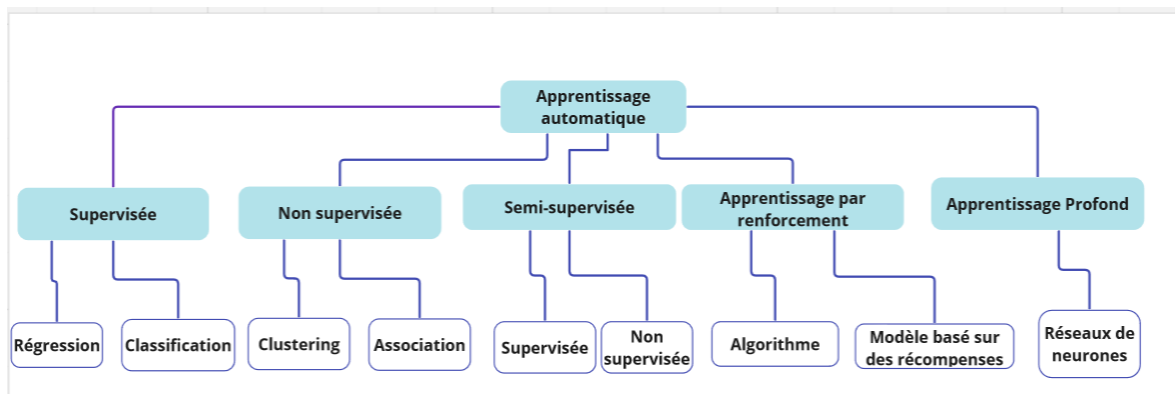


FIGURE 1.3 – Types de l'Apprentissage Automatique.

1.3.3.1 Apprentissage Supervisé

L'apprentissage supervisé est une méthode d'apprentissage automatique dans laquelle un algorithme informatique est formé à partir d'exemples de données étiquetées afin d'accomplir une tâche spécifique. Dans un système d'apprentissage supervisé, les algorithmes utilisent les données d'entraînement pour apprendre à faire des prédictions sur de nouvelles données. Les algorithmes apprennent en comparant leurs prédictions avec des étiquettes correctes associées aux données d'entraînement, puis en ajustant leurs modèles en conséquence.

Cette approche est largement utilisée pour diverses tâches, telles que la classification de données, la régression linéaire, la reconnaissance d'images, la reconnaissance de la parole, etc. L'apprentissage supervisé trouve une application étendue dans les domaines de la science des données, de l'intelligence artificielle et de l'apprentissage automatique [1]. La Figure 1.4 présente une illustration des différents types d'apprentissage supervisé.

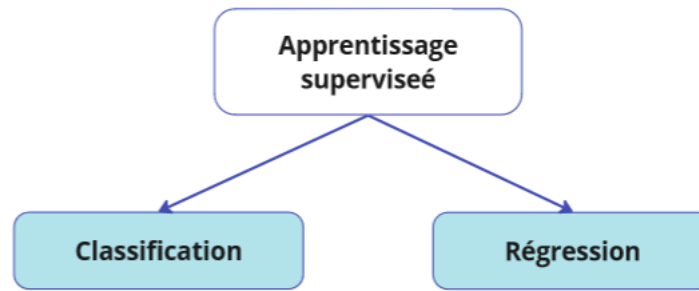


FIGURE 1.4 – Apprentissage Supervisé.

1.3.3.2 Apprentissage non Supervisé

L'apprentissage non supervisé est une branche de l'apprentissage automatique caractérisée par l'analyse et le regroupement de données non étiquetées. Ces algorithmes permettent de découvrir des modèles cachés ou des regroupements de données sans nécessiter d'intervention humaine. Dans un système d'apprentissage non supervisé, l'algorithme est libre de détecter des motifs et des modèles dans les données sans guidance explicite. Les algorithmes peuvent employer des techniques comme l'association, le clustering et la création de modèles afin de révéler des structures cachées dans les données [16]. La Figure 1.5 illustre différents types d'apprentissage supervisé.

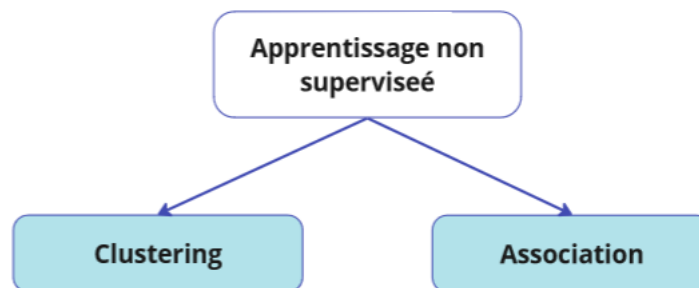


FIGURE 1.5 – Apprentissage non Supervisé.

1.3.3.3 Apprentissage Semi-Supervisé

L'apprentissage semi-supervisé représente une catégorie intermédiaire au sein du domaine de l'apprentissage automatique, se situant entre l'apprentissage supervisé et l'apprentissage non supervisé. Cette approche combine habilement des données étiquetées et non étiquetées dans le processus d'entraînement d'un algorithme. Ce processus démarre par le marquage, où certaines données sont étiquetées avec les bonnes

réponses (valeurs cibles). Cette phase d'apprentissage à propos de la création d'un modèle pondéré capable de prédire les réponses pour des données similaires qui n'ont pas encore été étiquetées. Ainsi, l'apprentissage semi-supervisé exploite simultanément des données étiquetées et non étiquetées pour adapter le modèle. Les méthodes d'apprentissage semi-supervisé étendent les techniques de l'apprentissage choisi, soit non-supervisé ou supervisé, afin d'inclure des informations supplémentaires typiques de l'autre paradigme d'apprentissage [17].

1.3.3.4 Apprentissage par Renforcement

L'apprentissage par renforcement est une méthode de l'apprentissage automatique qui permettant de former des modèles d'intelligence artificielle. Dans ce cadre, l'agent IA ou l'algorithme acquiert des stratégies de manière autonome, visant à apprendre, au fil d'expériences successives, les actions optimales pour résoudre des problèmes spécifiques. L'apprentissage par renforcement s'inscrit dans la branche de l'apprentissage machine où un agent, qu'il s'agisse d'une machine ou d'un logiciel, développe des compétences décisionnelles en interagissant avec un environnement dynamique. Contrairement à l'apprentissage supervisé, où l'agent est formé sur des données étiquetées, l'apprentissage par renforcement repose sur un processus d'essais et d'erreurs. Au fil du temps, l'agent accumule l'expérience nécessaire pour élaborer une stratégie décisionnelle optimale en vue d'optimiser les récompenses. Les algorithmes d'apprentissage par renforcement trouvent des applications variées, notamment dans la gestion des robots, l'automatisation des usines, l'optimisation des chaînes de livraison, les réseaux électriques intelligents, et bien d'autres domaines [18].

1.3.3.5 Apprentissage Profond

L'apprentissage profond, une branche de l'apprentissage automatique, lui-même étant un sous-ensemble de l'intelligence artificielle, a la capacité d'analyser des données non structurées telles que des images, des vidéos, du texte, etc. Cette technique repose sur des réseaux de neurones composés de plusieurs couches interconnectées, similaire au fonctionnement du cerveau humain.

Ces réseaux de neurones sont constitués de couches de nœuds interconnectés, également appelés "neurones", conçus pour traiter et analyser de grandes quantités de données

d'entrée. Chaque neurone dans un réseau d'apprentissage profond reçoit des entrées des neurones de la couche précédente, utilisant ces entrées pour effectuer une opération mathématique appelée fonction d'activation, produisant ainsi une sortie. Cette sortie est ensuite transmise à la couche suivante de neurones, le processus se répétant jusqu'à la production de la sortie finale. Plus le nombre de couches est élevé, plus l'apprentissage est profond, permettant de résoudre des problèmes très complexes.

Les algorithmes d'apprentissage profond peuvent être formés de manière supervisée ou non supervisée, en fonction du type de problème à résoudre. L'apprentissage supervisé implique de fournir à l'algorithme des données d'entraînement étiquetées, où la sortie correcte ou "étiquette" est connue pour chaque entrée. L'algorithme apprend ensuite à faire des prévisions sur de nouvelles données en se basant sur les modèles identifiés lors de l'entraînement. D'autre part, l'apprentissage non supervisé implique de fournir à l'algorithme des données non étiquetées, et l'algorithme doit identifier des modèles ou des caractéristiques dans les données par lui-même.

Le traitement du langage naturel offre à une machine la possibilité de comprendre et interpréter le langage humain, permettant ainsi d'obtenir des résultats de pointe dans de nombreux domaines [19]. La Figure 1.6 montre un exemple de réseau de neurones.

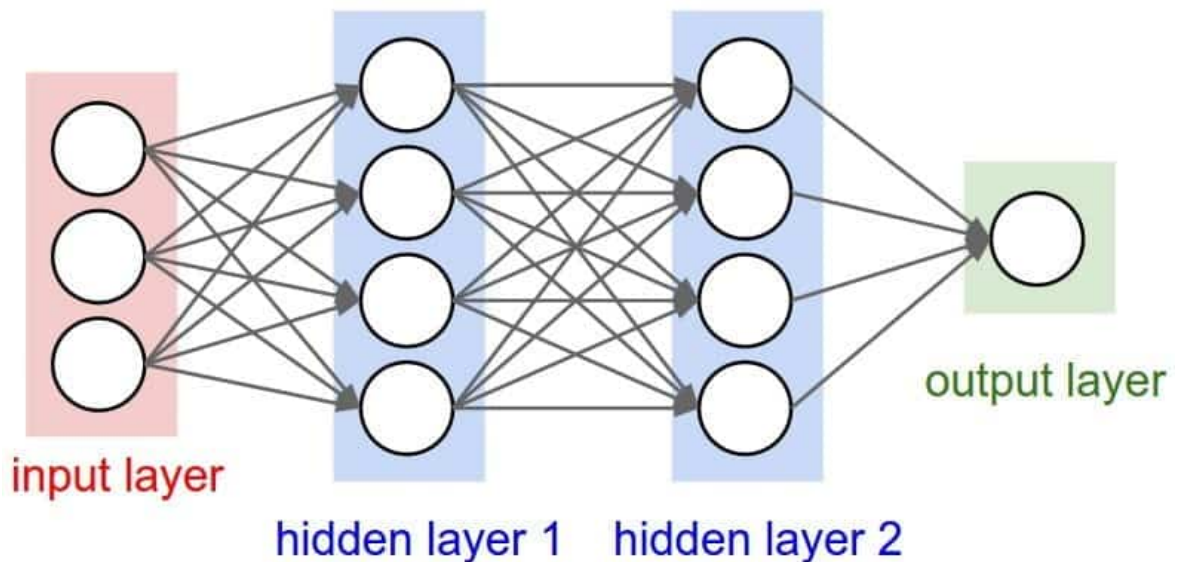


FIGURE 1.6 – Exemple d'un Réseau de Neurones.

1.3.4 Applications de l'Apprentissage Automatique dans le traitement des Big Data

L'avènement de l'apprentissage automatique a profondément transformé le paysage de l'analyse des Big Data, offrant une méthode efficace pour explorer et interpréter d'importantes volumes de données. En exploitant les algorithmes d'apprentissage automatique, les entreprises sont désormais en mesure d'extraire des informations cruciales de leurs vastes ensembles de données. Cette avancée technologique permet une compréhension approfondie des données, facilitant ainsi la prise de décisions éclairées et la découverte de tendances significatives au sein de ces ensembles massifs, voici quelque applications de l'apprentissage automatique dans l'analyse du Big Data [20].

1. **Analyse prédictive** : L'analyse prédictive, émergeant comme une application majeure de l'apprentissage automatique dans le domaine de l'analyse du Big Data. En utilisant des algorithmes d'apprentissage automatique, les entreprises peuvent explorer les données historiques afin d'identifier des modèles et de formuler des prédictions sur des événements futurs. Cette approche repose sur la capture des relations entre les variables explicatives et les variables prédites à partir d'occurrences passées. Elle permet, par exemple, d'anticiper des tendances de comportement, que ce soit dans le passé, le présent ou le futur. Toutefois, la précision et la pertinence des résultats dépendront fortement du niveau d'analyse des données et de la qualité des hypothèses, soulignant ainsi l'importance cruciale de ces facteurs dans le processus d'analyse prédictive [21].
2. **Détection de fraude** : Une application cruciale de l'apprentissage automatique dans le domaine de l'analyse du Big Data est la détection des fraudes, notamment celle liée aux cartes de crédit l'un des types de fraude les plus étudiés. Les algorithmes d'apprentissage automatique peuvent repérer des modèles indiquant une activité frauduleuse. Les fraudeurs utilisent diverses méthodes, telles que l'association non autorisée avec des commerces, le phishing par e-mail, la création de faux sites Web bancaires. Les techniques d'apprentissage automatique sont essentielles pour développer des systèmes automatisés de détection de fraude, mais la construction de modèles généralisés présente des défis tels que le déséquilibre de classe, la présence d'échantillons non étiquetés et la nécessité de

scalabilité. Pour relever ces défis, les systèmes de la détection des fraudes doivent intégrer des approches supervisées et semi-supervisées, adaptées au contexte du Big Data. L'apprentissage automatique est crucial pour analyser efficacement les données massives et renforcer la sécurité financière contre la fraude [22].

3. Systèmes de recommandation : Les systèmes de recommandation représentent une application répandue de l'apprentissage automatique dans le domaine de l'analyse du Big Data. Ces systèmes utilisent des algorithmes d'apprentissage automatique pour analyser le comportement des utilisateurs et générer des recommandations personnalisées, comme cela est illustré dans le cas de services de streaming tels que Netflix. Les étapes du système de recommandation commencent par l'obtention de données à partir de diverses sources éducatives, avec une intégration et une suppression des incohérences. Ensuite, une sélection d'attributs est effectuée pour réduire le volume de données. Le prétraitement des données comprend le nettoyage des valeurs manquantes et la gestion des données bruyantes à l'aide de techniques de lissage. La transformation des données peut impliquer la discrétisation ou la génération de hiérarchies. La phase suivante consiste à utiliser des techniques de fouille de données pour générer des recommandations. Enfin, les techniques de recommandation sont catégorisées en filtrage collaboratif, filtrage basé sur le contenu, système basé sur la connaissance et systèmes hybrides [23].

4. Analyse des sentiments :

L'analyse de sentiment vise à analyser et résumer les opinions exprimées dans les vastes volumes de données générées par les utilisateurs. Elle cherche à déterminer automatiquement si un texte généré par l'utilisateur exprime une opinion positive, négative ou neutre à l'égard d'une entité telle qu'un produit, une personne, un sujet ou un événement. Les données proviennent de diverses sources, notamment la publication ou le partage de données sur les sites de médias sociaux, des vidéos, des films audio, etc. Cette masse de données est appelée big data, et elle peut se présenter sous forme structurée, semi-structurée ou non-structurée, adaptée à l'analyse de sentiment. Différentes approches sont utilisées pour l'analyse de sentiment, comprenant l'approche basée sur les lexiques, l'approche basée sur l'apprentissage automatique, et l'approche hybride. L'approche basée sur l'apprentis-

sage automatique fait appel à des algorithmes d'apprentissage supervisé ou non supervisé pour classer les données. Elle implique l'utilisation de deux ensembles de documents : un ensemble d'entraînement et un ensemble de test. Un classificateur d'apprentissage supervisé se forme à partir de l'ensemble d'entraînement en se basant sur les attributs distinctifs du texte, puis évalue sa performance sur l'ensemble de tests. Plusieurs algorithmes d'apprentissage automatique, tels que Maximum Entropy, Naive Bayes, et Support Vector Machines, sont couramment utilisés pour la classification des tweets.

Les applications de surveillance des médias sociaux et les entreprises dépendent toutes de l'analyse de sentiment et de l'apprentissage automatique pour les aider à obtenir des informations sur les mentions, les marques et les produits [24, 25].

5. **Traitement du langage naturel** : Le traitement du langage naturel (NLP) se positionne comme une application de l'apprentissage automatique dans l'analyse des Big Data. Cette application tire parti des techniques avancées de traitement du langage naturel pour faciliter l'interrogation des bases de données en utilisant des langues naturelles telles que l'anglais.

Cette application permet aux utilisateurs d'interroger les vastes bases de données en utilisant des langues naturelles. À travers cette approche novatrice, l'objectif est de rendre l'analyse de données plus accessible, même pour ceux qui ne sont pas familiers avec les requêtes SQL complexes. Cette évolution s'inscrit dans le contexte global de l'essor des Big Data, offrant des solutions plus conviviales et efficaces pour explorer et comprendre ces vastes ensembles d'informations [26].

1.3.5 Outils utilisés dans le traitement des Big Data

L'apprentissage automatique est devenu un aspect important de l'analyse du Big Data. Face à l'augmentation constante du volume et de la complexité des données, l'analyse manuelle est devenue pratiquement irréalisable. L'apprentissage automatique offre une solution à ce défi en utilisant des algorithmes et des modèles statistiques pour assimiler les données et formuler des prédictions ou des décisions pertinentes.

Cependant, pour mettre en œuvre l'apprentissage automatique dans l'analyse du Big Data, certains outils et technologies sont nécessaires. Nous explorerons certains des outils et technologies les plus couramment utilisés pour l'apprentissage automatique

dans l'analyse du Big Data. La synthèse de ces outils est représentée dans la Figure 1.7.

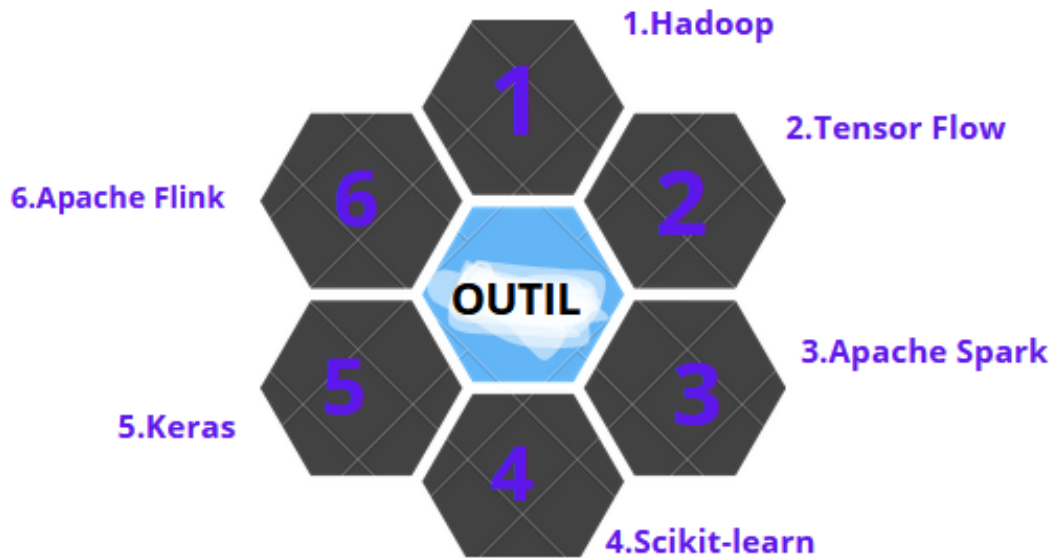


FIGURE 1.7 – Outils utilisés dans le traitement des Big Data

1. **Hadoop** : Hadoop émerge comme un framework open source largement adopté pour le stockage et le traitement de vastes ensembles de données. Il offre un système de fichiers distribué qui autorise le traitement simultané des données sur plusieurs nœuds. Hadoop présente également une compatibilité avec diverses bibliothèques d'apprentissage automatique, telles que Mahout et Spark MLlib. Ces bibliothèques mettent à disposition un ensemble d'algorithmes destinés à l'exploration de données, au clustering, à la classification, et à la régression [20].
2. **TensorFlow** : TensorFlow représente une bibliothèque open source conçue par Google pour la création et l'entraînement de modèles d'apprentissage automatique. Cette bibliothèque offre une architecture souple, facilitant le déploiement sur différentes plates-formes telles que les processeurs, les GPU et les dispositifs mobiles. En plus de sa flexibilité, TensorFlow met à disposition une gamme variée d'algorithmes couvrant l'apprentissage en profondeur, l'apprentissage par renforcement, et le traitement du langage naturel [20].
3. **Apache Spark** : Apache Spark est un framework open source qui fournit un moteur d'analyse unifié pour le traitement de grands ensembles de données. Il prend en charge divers langages de programmation, tels que Java, Scala, et Py-

thon. Spark fournit également un ensemble de bibliothèques d'apprentissage automatique telles que MLlib et GraphX. Ces bibliothèques fournissent un ensemble d'algorithmes pour la classification, la régression, le clustering, et le traitement graphique [20].

4. **Scikit-learn** : Scikit-learn est une bibliothèque d'apprentissage automatique pour Python qui fournit un ensemble d'algorithmes pour l'exploration de données, la classification, la régression, et le clustering. Il fournit également des outils pour la sélection, le prétraitement, et l'évaluation de modèles. Scikit-learn est largement utilisé dans l'industrie et le monde universitaire en raison de sa facilité d'utilisation et de sa flexibilité [20].
5. **Keras** : Keras est une bibliothèque de réseaux neuronaux de haut niveau pour Python qui fournit un moyen simple et efficace de créer et de former des modèles d'apprentissage en profondeur. Il prend en charge divers backends, tels que TensorFlow, Theano, et CNTK. Keras fournit également un ensemble de modèles pré-entraînés pour la classification d'images et de textes, la détection d'objets et le traitement du langage naturel [20].
6. **Apache Flink** : Apache Flink est un framework open source qui fournit un moteur de traitement de flux distribué pour l'analyse en temps réel. Il prend en charge divers langages de programmation, tels que Java, Scala, et Python. Flink fournit également un ensemble de bibliothèques d'apprentissage automatique telles que FlinkML. Ces bibliothèques fournissent un ensemble d'algorithmes de classification, de régression, et de clustering [20].

1.4 Défis et limites actuels

L'apprentissage automatique, en tant que technologie innovante, a joué un rôle révolutionnaire qui a permis de découvrir des modèles et des informations cachés au sein du Big Data. Malgré ses avancées significatives et son potentiel considérable, cette approche rencontre toutefois divers défis et limites lorsqu'il s'agit de traiter le Big Data.

1.4.1 Défis rencontrés dans l'analyse des Big Data

Dans cette section, nous explorerons certains défis de l'apprentissage automatique dans l'analyse du Big Data.

1. **Qualité des données :** L'intégration de l'apprentissage automatique dans l'analyse du Big Data est confrontée à un défi significatif lié à la qualité des données. Les algorithmes d'apprentissage automatique dépendent de vastes ensembles de données pour s'entraîner efficacement. Toutefois, la présence de données de mauvaise qualité peut avoir des effets négatifs sur la précision des résultats générés par l'algorithme. Les problèmes de qualité des données peuvent découler de diverses sources, telles que des erreurs de saisie, des données incomplètes ou des formats de données incohérents.

Afin de relever ce défi, il est impératif pour les entreprises de garantir la propreté, la cohérence, et l'exactitude de leurs données. Cela implique une attention particulière aux processus de collecte, de stockage, et de maintenance des données, visant à éliminer toute source potentielle d'erreur [20].

2. **Intégration des données :** L'intégration des données représente un autre défi significatif dans l'application de l'apprentissage automatique à l'analyse du Big Data. Les entreprises stockent souvent leurs données dans différents formats et emplacements. L'intégration de ces données pour créer une source unique de vérité peut être un processus complexe et long. Les algorithmes d'apprentissage automatique nécessitent des données issues de diverses sources pour un entraînement efficace. Afin de surmonter cet obstacle, les entreprises doivent investir dans des outils d'intégration de données capables d'automatiser le processus d'unification des données provenant de sources multiples.

Cette automatisation faciliterait non seulement le processus d'intégration mais également la création d'un ensemble de données cohérent, essentiel pour l'entraînement et la performance des algorithmes d'apprentissage automatique [20].

3. **Confidentialité des données :** La confidentialité des données émerge comme une préoccupation cruciale dans la mise en place de l'apprentissage automatique pour l'analyse du Big Data. Les algorithmes d'apprentissage automatique exigent l'accès à des informations sensibles, telles que les données clientes, financières, et personnelles. Il est impératif que les entreprises veillent à mettre en place des me-

sures de sécurité adéquates pour prévenir tout accès non autorisé à ces données. Ces mesures englobent l'adoption de techniques avancées de cryptage, de systèmes de contrôle d'accès, et de méthodes de masquage des données. Ainsi, les organisations doivent investir de manière significative dans ces dispositifs de sécurité afin d'assurer la confidentialité des données, garantissant ainsi une utilisation appropriée des informations sensibles et répondant aux exigences réglementaires en matière de protection des données [20].

4. **Pénurie de talents** : La mise en place de l'apprentissage automatique dans l'analyse du Big Data requiert des compétences spécialisées et des connaissances approfondies. Toutefois, le marché connaît actuellement une pénurie de data scientists qualifiés et d'experts en apprentissage automatique. Cela rend difficile pour les entreprises la recherche des talents adéquats pour intégrer l'apprentissage automatique dans leurs projets d'analyse de Big data. Pour surmonter ce défi, les entreprises doivent consacrer des ressources à des programmes de formation visant à améliorer les compétences de leur personnel existant ou collaborer avec des prestataires externes possédant l'expertise nécessaire [20].
5. **Sélection de l'algorithme** : Le choix approprié de l'algorithme d'apprentissage automatique revêt une importance cruciale pour la réussite d'un projet d'analyse de Big Data. Cependant, il existe de nombreux algorithmes disponibles, chacun ayant ses propres forces et faiblesses. Les entreprises doivent minutieusement évaluer leurs besoins avant de sélectionner l'algorithme le plus adapté. Par exemple, si l'objectif est de prédire le taux de désabonnement des clients, il se peut qu'un algorithme d'arbre de décision soit plus efficace qu'un algorithme de réseau neuronal. Il est donc essentiel de faire un choix éclairé en fonction des spécificités du projet et des objectifs visés [20].

1.5 Conclusion

En conclusion, nous avons exploré le panorama complet du big data et de son analyse, en soulignant son évolution rapide. Les caractéristiques fondamentales du big data 5V, la structuration et la gestion de ces données massives sont cruciales pour extraire des informations utiles. Grâce aux avancées de l'apprentissage automatique,

les big data les plus complexes sont désormais traités efficacement. Ces avancées ont permis des applications diverses comme l'analyse prédictive, la détection de fraude, les systèmes de recommandation, l'analyse des sentiments, et le traitement du langage naturel, démontrant l'impact de l'apprentissage automatique dans la résolution de problèmes complexes liés au big data. Les outils de traitement du big data sont essentiels pour une gestion et une analyse efficace, facilitant la prise de décision et l'innovation continue. En résumé, l'intersection du big data et de l'apprentissage automatique crée un paysage dynamique qui façonne l'avenir de la science des données, avec de nombreux défis et opportunités à comprendre pour rester à la pointe de l'innovation.

Dans le prochain chapitre, intitulé "État de l'art", nous examinerons les différentes études classiques ainsi que des exemples d'algorithmes d'apprentissage automatique.

Chapitre 2

État de l'art

2.1 Introduction

Dans ce deuxième chapitre, nous explorerons en détail l'application de l'apprentissage automatique dans l'analyse des Big Data. Cela met en lumière l'importance cruciale de cette approche dynamique face à la prolifération de volumes massifs et variés de données. La structure du chapitre comprend une présentation des fondements de l'apprentissage automatique, des techniques spécifiques adaptées à l'analyse des Big Data, des algorithmes populaires, des cas d'utilisation spécifiques, des défis rencontrés, ainsi qu'une revue des études classiques avec une motivation pour la Revue Systématique de la Littérature (SLR). L'accent est mis sur la nécessité de surmonter les limitations identifiées pour une compréhension approfondie de ce domaine en constante évolution.

2.1.1 Contexte de l'étude

Le chapitre explore la fusion significative de la technologie et des opportunités dans le domaine des données. Il met en avant l'importance des vastes ensembles de données générés rapidement à partir de diverses sources, offrant des occasions uniques d'extraire des informations exploitables grâce à des techniques d'analyse avancées, notamment l'apprentissage automatique. Le chapitre présente diverses techniques d'apprentissage automatique et examine des cas d'utilisation spécifiques aux Big Data, tout en explorant les défis et opportunités de l'application de ces méthodes. Une revue des études classiques est également fournie, comparant les revues académiques spécialisées dans les études classiques. Enfin, le chapitre aborde les limitations des études classiques et sou-

ligne l'importance des revues systématiques de la littérature pour une compréhension approfondie des applications de l'apprentissage automatique dans l'analyse des Big Data. Dans l'ensemble, il offre une vue d'ensemble complète des développements et des enjeux dans ce domaine dynamique.

2.1.2 Importance de l'application de l'apprentissage automatique dans l'analyse des Big Data

L'application de l'apprentissage automatique dans l'analyse des Big Data présente des opportunités significatives dans divers domaines. Les ensembles de données massifs et variés permettent l'extraction d'information exploitable grâce à des techniques avancées d'analyse. L'adaptabilité de l'apprentissage automatique aux données disponibles est cruciale dans un contexte de diversité des sources de données. Cette technologie polyvalente trouve des applications dans la prise de décision, la prévision, et divers secteurs tels que la santé, la science, l'ingénierie, les affaires, et la finance.

Ces techniques d'apprentissage automatique offrent l'opportunité de traiter les Big Data et de résoudre des problèmes complexes. L'apprentissage automatique contribue également à la gestion du volume massif et de la variété des données, mettant en avant des solutions telles que la parallélisation des algorithmes et l'apprentissage en ligne. La scalabilité des algorithmes d'apprentissage automatique est soulignée comme essentielle pour traiter efficacement les vastes ensembles de données caractéristiques des Big Data, offrant des avantages tels que le traitement rapide et la capacité à s'adapter à l'augmentation des volumes de données.

En résumé, l'application de l'apprentissage automatique dans l'analyse des Big Data joue un rôle clé dans l'exploitation efficace des données massives et complexes, offrant des solutions innovantes et des opportunités prometteuses.

2.2 Application de l'Apprentissage Automatique dans l'Analyse des Big Data

Représente une convergence significative de technologie et d'opportunités dans le domaine moderne des données. Les ensembles de données massifs, rapides à générer, et

provenant de diverses sources offrent des perspectives uniques pour extraire des informations exploitables grâce à des techniques d'analyse avancées, notamment l'apprentissage automatique. Cette approche, en tant que sous-discipline de l'intelligence artificielle, permet aux systèmes informatiques d'apprendre à exécuter des tâches spécifiques de manière automatique en s'adaptant aux données disponibles. Les applications de l'apprentissage automatique sont vastes, allant de la prise de décision à la prévision, et sa pertinence s'étend à des domaines variés tels que la santé, la science, l'ingénierie, les affaires et la finance. En examinant les différents types de tâches en apprentissage automatique, comme illustré par la Figure 2.1, il devient clair que cette technologie est essentielle pour exploiter efficacement les Big Data et offre des opportunités prometteuses pour résoudre des problèmes complexes [27].

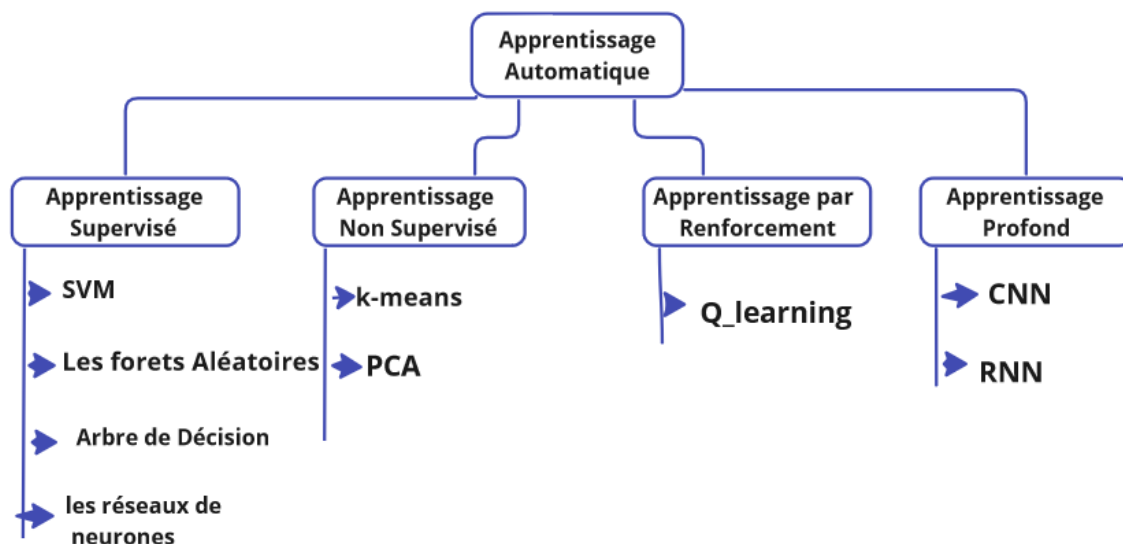


FIGURE 2.1 – Techniques de l'Apprentissage Automatique.

2.2.1 Fondements de l'Apprentissage Automatique

Dans l'analyse des big data utilisant des méthodes d'application automatique, voici quelques-unes des techniques couramment utilisées [28] :

1. **Apprentissage Supervisé** : Cette approche consiste à entraîner un modèle sur un ensemble de données étiquetées, où le modèle apprend à prédire des étiquettes pour de nouvelles données. Les algorithmes populaires incluent les machines à vecteurs de support (SVM), les arbres de décision (DT), les forêts aléatoires, et les réseaux de neurones.

2. **Apprentissage Non Supervisé** : Contrairement à l'apprentissage supervisé, cette approche ne nécessite pas d'étiquettes pour les données d'entraînement. Les algorithmes d'apprentissage non supervisés sont utilisés pour la segmentation, la réduction de dimensionnalité, et la détection d'anomalies. Les techniques incluent le clustering (comme le k-means) et l'analyse en composantes principales (PCA).
3. **Apprentissage par Renforcement** : Cette approche consiste à apprendre à prendre des décisions séquentielles en interagissant avec un environnement. Les algorithmes d'apprentissage par renforcement sont utilisés pour l'optimisation et le contrôle, et comprennent des techniques telles que les processus décisionnels markoviens (MDP) et les algorithmes Q-learning.
4. **Apprentissage Profond** : Cette catégorie utilise des réseaux de neurones artificiels profonds pour modéliser et comprendre des données complexes. Elle est particulièrement efficace pour le traitement d'images, de vidéos, de texte et de parole. Les algorithmes les plus populaires sont les réseaux de neurones convolutionnels (CNN) et les réseaux de neurones récurrents (RNN).

2.2.1.1 Domaines d'Application

Dans cette section, nous abordons les domaines d'application de chaque catégorie de l'apprentissage automatique.

1. **Domaines d'application pour l'apprentissage Supervisée** : L'apprentissage automatique supervisé est utilisé pour résoudre des problèmes de classification et de régression en enseignant à un ordinateur à reconnaître des schémas dans des ensembles de données étiquetés. Les données annotées sont utilisées pour former un modèle, qui est ensuite utilisé pour prédire de nouvelles données non annotées. Les algorithmes populaires incluent la régression linéaire, la régression logistique, les arbres de décision et la machine à vecteurs de support (SVM). L'objectif de l'apprentissage automatique supervisé est de produire des modèles capables de généraliser à de nouvelles données et de prendre des décisions précises en fonction des informations fournies [29].
2. **Domaines d'application pour l'apprentissage Non Supervisée** : L'apprentissage non supervisé se distingue de l'apprentissage supervisé par le fait que les algorithmes apprennent des caractéristiques des données sans avoir de réponses

correctes ou d'enseignant pour les guider. Lorsque de nouvelles données sont introduites, ces algorithmes utilisent les caractéristiques qu'ils ont précédemment apprises pour identifier les structures ou les classes des données. Cette méthode est principalement utilisée pour regrouper des données similaires et pour réduire la dimensionnalité des données en identifiant les caractéristiques les plus importantes [30].

3. Domaines d'application pour l'Apprentissage Profond :

Apprentissage profond est une approche d'apprentissage automatique qui utilise des réseaux de neurones artificiels profonds pour apprendre des représentations hiérarchiques des données. Ces réseaux de neurones profonds sont capables d'apprendre des caractéristiques complexes à partir des données en utilisant plusieurs couches de traitement non linéaire. L'apprentissage en profondeur est souvent utilisé pour des tâches telles que la reconnaissance d'images, la reconnaissance vocale, la traduction automatique, et d'autres domaines où des modèles complexes et abstraits sont nécessaires pour traiter des données massives et non structurées. [31].

4. Domaines d'application pour l'Apprentissage Par Renforcement : L'apprentissage par renforcement est une méthode d'apprentissage automatique qui permet à un agent ou à un algorithme d'acquérir des stratégies de manière autonome en interagissant avec un environnement dynamique. Contrairement à l'apprentissage supervisé, où l'agent est formé sur des données étiquetées, l'apprentissage par renforcement repose sur un processus d'essais et d'erreurs. Au fil du temps, l'agent accumule de l'expérience pour élaborer une stratégie décisionnelle optimale en vue d'optimiser les récompenses. Les algorithmes d'apprentissage par renforcement, tels que les processus décisionnels markoviens (MDP) et les algorithmes Q-learning, sont utilisés pour l'optimisation et le contrôle dans des domaines variés tels que la gestion des robots, l'automatisation des usines, l'optimisation des chaînes de livraison, des réseaux électriques intelligents, et bien d'autres [32].

2.2.1.2 Importance dans le Contexte des Big Data

Les techniques de l'apprentissage automatique jouent un rôle crucial dans l'analyse des Big Data en permettant le traitement de vastes ensembles de données pour en extraire des informations précieuses. Par exemple, elles sont utilisées avec succès pour classifier des maladies telles que le cancer du sein en utilisant des approches comme les arbres de décision, le gain d'information et les ensembles de gènes. De plus, l'utilisation de techniques méta-heuristiques combinées à des algorithmes comme la régression logistique a permis d'améliorer la prédiction de la gravité des maladies et la prise de décision clinique. Ces avancées illustrent l'importance des techniques de l'apprentissage automatique dans l'analyse des Big Data, notamment dans des domaines critiques comme la santé [33].

2.2.2 Techniques d'Apprentissage Automatique pour l'Analyse des Big Data

Dans cette section, nous explorons comment les techniques d'apprentissage automatique analysent les Big Data. Ces méthodes détectent des modèles complexes et extraient des informations cruciales des vastes ensembles de données, améliorant ainsi la prise de décision et divers processus dans de nombreux secteurs.

2.2.2.1 Apprentissage Supervisé :

Parmi les algorithmes les plus utilisés en apprentissage automatique supervisé, on trouve [29] :

1. Machines à Vecteurs de Support :

Les SVM, ou machines à vecteurs de support, sont des algorithmes d'apprentissage supervisés utilisés pour la classification et la régression. Leur caractéristique principale est l'utilisation d'une capacité de noyau pour créer un hyperplan séparant les classes dans les données, ajusté itérativement pour minimiser l'erreur. L'objectif est de trouver un hyperplan marginal maximal (HMM) en suivant des étapes itératives de création et de sélection d'hyperplans optimaux. On distingue les SVM simples, adaptées à la régression linéaire et à la classification, des SVM à noyau, offrant plus de flexibilité pour les données non linéaires en ajustant un

hyperplan dans un espace de dimension supérieure. Leur utilisation répandue en apprentissage automatique s'explique par leur capacité à découvrir des relations complexes sans nécessiter de modifications importantes, particulièrement aux ensembles de données plus petits avec des dizaines à des milliers de fonctionnalités, où ils fournissent généralement des résultats précis..

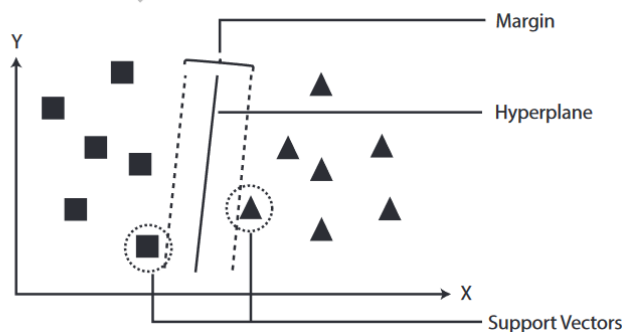


FIGURE 2.2 – SVM
[29]

Les principaux termes associés aux SVM sont illustrés dans la figure 2.2 suivante :
 Support Vecteur : Ils sont les points de données les plus proches de l'hyperplan et identifient à définir la ligne de séparation. Hyperplane : Il s'agit du plan de décision qui sépare les différentes classes dans un ensemble de données. Marge : Elle représente l'écart entre deux lignes sur des points de données de différentes classes, calculé comme la distance perpendiculaire entre la ligne et le vecteur de support.

2. Arbres de Décision (DT) :

Les arbres de décision sont des structures qui dépendent des valeurs des éléments. Ils utilisent la stratégie du gain d'information pour déterminer quel élément dans l'ensemble de données fournit les meilleures données, ce qui en fait le nœud racine, et ainsi de suite jusqu'à ce que chaque cas de l'ensemble de données soit organisé. Chaque branche de l'arbre de décision représente un élément de l'ensemble de données. Ils sont largement utilisés pour la classification.

Dans une analyse d'arbre de décision, celui-ci est utilisé pour représenter visuellement et décrire le processus de décision. Comme son nom l'indique, il utilise une représentation arborescente des choix. Les modèles d'arbre sont la variable objective qui peut prendre un ensemble discret de valeurs, appelé arbres de clas-

sification. Dans ce modèle, les feuilles représentent les étiquettes de classe, tandis que les combinaisons de caractéristiques menant à ces étiquettes sont représentées par les branches de l'arbre.

3. Algorithme de Réseau Neuron Artificiel :

Les réseaux de neurones artificiels (RNA) sont des programmes informatiques qui s'inspirent du fonctionnement biologique du cerveau humain. Ces réseaux sont conçus avec l'idée de regrouper des neurones artificiels en couches, imitant ainsi la manière dont le cerveau traite l'information, afin de convertir des entrées en sorties significatives. Chaque neurone artificiel, en tant qu'unité fondamentale, reçoit des données d'entrée provenant des neurones précédents, chacune de ces entrées étant pondérée par un poids (w) représentant la force de la connexion. À travers une fonction d'activation et en utilisant la somme de ces poids, chaque unité produit des sorties significatives par le biais de sa seule issue [34]. La Figure 2.3 illustre la structure formelle d'un neurone de l'ANN.

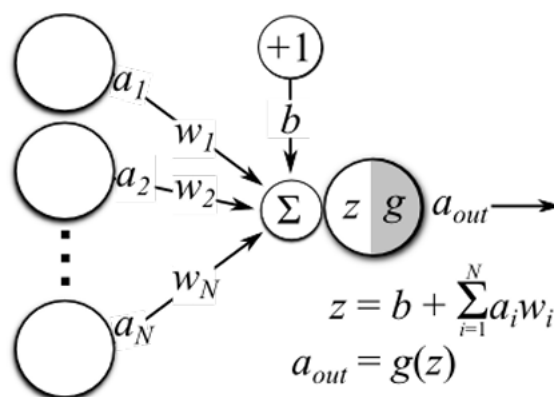


FIGURE 2.3 – La configuration formelle d'un neurone. [34]

où :

- a_i : représente les entrées du neurone.
- b : chaque neurone dispose d'une entrée supplémentaire de valeur 1, avec un poids de connexion spécifique b . Cela assure qu'il y aura toujours une activation dans le neurone, même si toutes les autres entrées sont nulles ($a_i = 0$).
- w_i : les poids associés aux entrées.

- z : la fonction d'intégration des entrées.
- g : fonction d'activation.
- a_{out} : la sortie du neurone.

2.2.2.2 Apprentissage Nom Supervisée :

Parmi les algorithmes les plus utilisés en apprentissage automatique nom supervisé, on trouve :

1. **K-means :**

L'algorithme K-means est un algorithme itératif de regroupement de données. Utilisant la distance comme mesure, il divise un ensemble de données en K classes en calculant la moyenne des distances pour chaque classe, à partir d'un ensemble initial de centroïdes. K-means est une méthode non supervisée visant à regrouper des éléments similaires en K groupes. Son objectif est de minimiser la variance intra-classe en réduisant la distance euclidienne entre les éléments d'une classe et son centre. Pour cela, il ajuste les centres itérativement en réassignant les éléments aux classes les plus proches. Ce processus se répète jusqu'à ce que les centres convergent, indiquant que l'algorithme a atteint sa solution [35].

2. **L'Analyse en composantes principales :**

L'analyse en composantes principales (ACP) constitue un algorithme d'apprentissage non supervisé fréquemment utilisé pour la réduction de la dimensionnalité dans le domaine de l'apprentissage automatique. Ce processus statistique opère une transformation orthogonale sur des observations caractérisées par des corrélations, les convertissant ainsi en un ensemble de caractéristiques linéairement non corrélées. Les résultats de cette transformation sont baptisés "composants principaux". L'ACP est largement adopté pour l'analyse exploratoire des données et la modélisation prédictive, se révélant être un outil efficace pour élaborer des modèles robustes en diminuant les variances au sein d'un ensemble de données. Ses applications concrètes englobent des domaines tels que le traitement d'images, les systèmes de recommandation de films, et l'optimisation de l'allocation de puissance dans divers canaux de communication. En tant que technique d'extraction de caractéristiques, l'ACP conserve les variables cruciales tout en éliminant celles de moindre importance. Les fondements mathématiques de l'algorithme reposent

sur des concepts tels que la variance, la covariance, les valeurs propres et les vecteurs propres [36].

2.2.2.3 Apprentissage par Renforcement :

Parmi les algorithmes les plus utilisés en apprentissage automatique par renforcement, on trouve :

1. **Algorithme de Q-Learning** : Le Q-learning est un algorithme d'apprentissage par renforcement qui cherche à déterminer la politique optimale pour maximiser la récompense cumulée d'un agent dans un environnement. L'algorithme utilise une table de qualité (Q-table) pour stocker les valeurs de qualité associées à chaque paire état-action. Chaque entrée $Q[s, a]$ représente la récompense totale anticipée par l'agent lorsqu'il démarre à l'état s , effectue l'action a , et suit une politique donnée.

Le principe de l'algorithme repose sur l'exploration de l'environnement par l'agent, où il choisit des actions de manière aléatoire (exploration) ou basées sur les valeurs de la Q-table (exploitation). L'algorithme ajuste les valeurs de la Q-table en utilisant une formule de mise à jour qui prend en compte la récompense obtenue, la valeur maximale anticipée dans le nouvel état, un taux d'apprentissage (α), et un facteur d'actualisation (γ). Cette mise à jour vise à propager les récompenses importantes vers les états permettant de les atteindre. Le processus d'apprentissage consiste en l'initialisation de la Q-table, la répétition d'un certain nombre d'épisodes, le choix d'actions, l'exécution des actions, l'observation des récompenses, et la mise à jour des valeurs de la Q-table. L'algorithme vise à converger vers la politique optimale pour prendre des décisions optimales à chaque pas de temps, en maximisant la récompense cumulée [37].

2.2.2.4 Apprentissage Profond :

Parmi les algorithmes les plus utilisés en apprentissage profond, on trouve :

1. **Réseaux de Neurones Convolutionnels (CNN)** :

Les réseaux de neurones convolutionnels (CNN) sont des modèles spécialisés dans le traitement d'images et sont actuellement les plus performants pour la classification d'images. Composés de deux parties distinctes, les CNN extraient d'abord

des caractéristiques des images à l'aide de filtres de convolution, générant des cartes de convolutions. Ensuite, ces caractéristiques sont combinées dans des canapés entièrement connectés pour la classification finale. Les CNN s'inspirent du fonctionnement du cortex visuel des vertébrés et produisent généralement une distribution de probabilité sur les différentes catégories d'images en sortie. Cependant, les perceptrons multicouches utilisés dans les CNN peuvent rencontrer des difficultés avec les images de grande taille en raison de la croissance exponentielle du nombre de connexions, ce qui peut poser des problèmes de scalabilité [31].

L'architecture d'un CNN se compose d'une séquence de couches de traitement, comprenant la couche de convolution (Convolution Layer) pour le traitement des données dans un champ récepteur, la couche de sou-échantillonnage (Subsampling) qui permet la compression de l'information en particulier la taille de l'image intermédiaire, la couche "entièrement connectée". Ces couches combinent les caractéristiques extraites par la partie convolutive pour effectuer la classification de l'image, et enfin, la couche neurone par catégorie : Chaque neurone dans la dernière couche représente une catégorie, et les valeurs de sortie sont normalisées pour anciennement une distribution de probabilité sur les différentes catégories (Neuron per Category) [31], comme illustré dans la Figure 2.4.

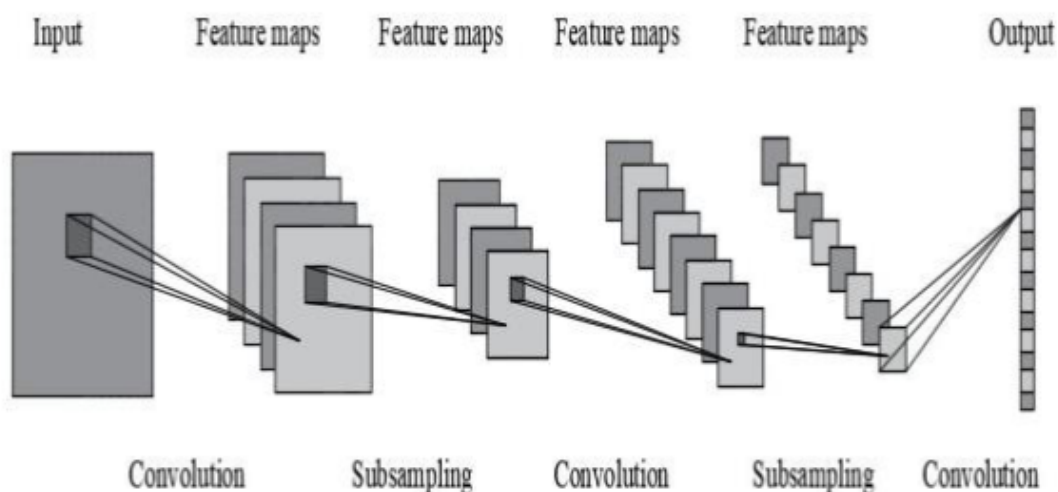


FIGURE 2.4 – Architecture d'un Réseau de Neurones Convolutionnels [32].

2. **Algorithme de Réseaux de Neurones Récurrents (RNN)** : sont des modèles d'apprentissage profond efficaces pour traiter les séries temporelles et le langage

naturel. Ils conservent les informations des entrées précédentes dans leur état interne, ce qui permet de capturer les relations entre des données successives. Pour surmonter le problème de disparition du gradient, des variantes comme les LSTM et les GRU ont été développées pour gérer les dépendances à long terme. Les RNN, notamment avec des architectures avancées comme RNN-BiLSTM-CRF, ont montré des performances remarquables dans des applications telles que la reconnaissance vocale et la traduction automatique, en extrayant et modélisant efficacement les caractéristiques des données séquentielles [38].

2.2.3 Algorithmes Populaires et leur Pertinence

Les algorithmes populaires tels que la régression logistique, la PLS, et les réseaux de neurones ont démontré des performances intéressantes dans le contexte industriel de la production d'électricité nucléaire. Cependant, il est crucial de ne pas tirer des conclusions hâtives sur l'efficacité générale de ces techniques et d'envisager l'évaluation de leur performance à travers d'autres algorithmes. La régression logistique est fréquemment utilisée pour estimer la probabilité d'une réponse binaire basée sur un ou plusieurs prédicteurs. Les réseaux de neurones, bien qu'ils produisent des résultats difficiles à interpréter, représentent une généralisation non linéaire du modèle de régression logistique.

De plus, les machines à vecteurs de support (SVM) sont employées pour construire un hyperplan séparant les données en deux classes en fonction des modalités possibles de la variable cible. Il est également essentiel de prendre en compte des techniques telles que la régression LASSO, qui effectue une sélection continue de sous-ensembles de prédicteurs en fonction d'un paramètre de réglage. Ces différentes approches offrent des possibilités variées pour la modélisation et la prédiction dans divers domaines, y compris la production d'électricité nucléaire. Il est recommandé d'explorer et de comparer plusieurs algorithmes afin de trouver celui qui convient le mieux à un problème spécifique [39].

2.2.3.1 Cas d'utilisation Spécifiques aux Big Data

Les cas d'utilisation spécifiques au Big Data incluent la classification distribuée et collective de données massives en utilisant un grand ensemble de classificateurs sur le

cloud pour améliorer la précision des résultats.

Cette approche se caractérise par son aspect totalement réparti, sans entité centrale, et permet une labellisation des données selon la classe ayant le plus fort poids, en définissant les résultats locaux des classificateurs et leurs facteurs de performance. Cette méthode vise à améliorer la véracité dans le traitement des Big Data en utilisant une approche collective et consensuelle au sein de chaque voisinage des nœuds du système [40].

2.2.4 Défis et Opportunités

Les défis dans l'application de l'apprentissage automatique à l'analyse des Big Data pour comprendre les opinions politiques résident dans la nécessité d'une analyse multi-échelle, la validité des données numériques par rapport aux enquêtes traditionnelles et l'établissement de liens entre les choix conceptuels et méthodologiques. Améliorer l'efficacité des prédictions, comprendre les structures en réseaux et explorer simultanément les échelles micro et macroscopiques sont des opportunités. Les équipes interdisciplinaires et le rôle du savoir expert sont cruciaux. L'exploration de méthodes multi-échelles et de catégories de sentiment offre des perspectives prometteuses pour surmonter ces défis dans l'analyse des opinions politiques à partir des Big Data [41].

2.2.5 Gestion du Volume et de la Variété des données

L'impératif de gérer de manière pédagogique le volume important et la diversité des données est souligné, mettant en exergue les défis posés par le Big Data, notamment la croissance exponentielle des données. Au-delà de la quantité, l'attention se porte sur la diversité des informations, fréquemment non structurées. Pour relever ces défis, en se focalisant sur l'apprentissage automatique à grande échelle, il met en avant la parallélisation des algorithmes et l'apprentissage en ligne. Évoquant l'utilisation de technologies distribuées telles que Spark, PySpark, Weka-MOA, et WekaSpark, il propose des solutions pour traiter ces données massives [42].

2.2.6 Scalabilité des Algorithmes d'Apprentissage Automatique

L'évolutivité des algorithmes d'apprentissage automatique est essentielle pour traiter efficacement les vastes ensembles de données caractéristiques de ce domaine. Cette capacité à évoluer de manière efficace offre des avantages majeurs, tels qu'un traitement rapide, la possibilité de parallélisme, et une adaptabilité à l'augmentation des volumes de données. Cependant, la conception complexe et les exigences en ressources, notamment en termes de puissance de calcul, représentent des défis. Cette scalabilité demeure un élément clé pour tirer pleinement partie de l'analyse des Big Data, jouant un rôle crucial dans la capacité à extraire des insights significatifs néanmoins de ces vastes ensembles de données [43].

2.3 Revue des Études Classiques

Dans cette section, nous examinerons certaines études classiques (enquêtes) qui passent en revue les travaux proposés dans la littérature sur Revue Systématique sur les méthodes d'apprentissage automatique pour le traitement et l'analyse des Big Data. L'objectif de cette étude est d'identifier les limitations de ce type d'étude et de les mettre en évidence.

Récemment, Jiyang Xie et al. [44] ont présenté une étude abordant les défis et les opportunités liés à l'analyse des données massives mobiles (MBD) en utilisant des méthodes d'apprentissage automatique. L'article met en lumière l'importance croissante des données massives mobiles en raison de la prolifération des appareils mobiles et des réseaux sans fil, caractérisant les MBD par cinq aspects principaux : le volume, la vitesse, la variété, la valeur, et la véracité. des données. Ces caractéristiques représentent des défis uniques en termes de prétraitement des données et de développement de méthodes d'analyse adaptées. L'article explore diverses méthodes d'analyse des données massives mobiles, telles que la division et la conquête, l'échantillonnage, les machines à vecteurs de support, les arbres de décision, les réseaux neuronaux et les machines d'apprentissage extrême, distinguées leur importance pour traiter efficacement les données massives et extraire des informations significatives. De plus, il examine des applications spécifiques, notamment la modélisation des canaux sans fil, l'analyse du comportement

en ligne et hors ligne des individus, et la reconnaissance de la parole dans l'Internet des véhicules. En outre, l'article souligne les défis futurs auxquels l'analyse des données massives mobiles est confrontée, tels que l'Internet mobile à grande échelle, les problèmes de surajustement et de sous-ajustement, les problèmes de généralisation, l'apprentissage croisé, des modalités et l'extension des dimensions du canal. Ces défis ont nécessité une réflexion approfondie et des avancées continues dans le domaine de l'apprentissage automatique pour exploiter pleinement le potentiel des données massives mobiles. En conclusion, l'article met en lumière l'importance croissante de l'analyse des données massives mobiles basée sur l'apprentissage automatique dans divers domaines tels que les communications sans fil.

Cet article met en évidence les limitations des études sur l'analyse des données massives mobiles basée sur l'apprentissage automatique. Ces limitations comprennent la nécessité de développer davantage les méthodes actuelles, malgré leurs performances satisfaisantes dans les tests de données réelles. Les défis associés à l'Internet mobile à grande échelle, aux problèmes de surajustement et de sous-ajustement, à la généralisation, à l'apprentissage croisé des modalités et à l'extension des dimensions du canal sont également évoqués. De plus, la charge importante imposée aux systèmes de transmission sans fil en raison de la croissance des appareils mobiles et de la vitesse élevée de l'Internet mobile nécessite des améliorations des technologies de communication sans fil. Ces limitations soulignent la nécessité continue de développer les méthodes d'analyse des données massives mobiles pour relever les défis actuels et futurs du domaine.

Dans une étude récente, Walaa Medhat et al. [45] présenté dans l'article "Big data techniques et applications d'apprentissage profond analytique : A Survey" offrant une analyse approfondie des techniques de deep learning dans le domaine de l'analyse de données massives. Les auteurs mettent en avant l'importance croissante du deep learning dans le traitement de grands volumes de données non-structurées et soulignent son rôle crucial dans diverses applications du monde réel. Une taxonomie des techniques de deep learning est présentée, mettant en évidence les forces et les limites de chaque approche, ainsi que les orientations futures de recherche dans le domaine. Les auteurs discutent des revues de la littérature, des modèles de processus, des frameworks, et des

ensembles de données de référence utilisés dans le domaine du deep learning pour l'analyse de données massives. Ils abordent également les performances des différentes mises en œuvre et les défis rencontrés. L'article explore les applications du deep learning dans des domaines variés tels que la santé, la finance, l'éducation, et l'IoT, en utilisant des modèles tels que CNN, LSTM, RNN, et GAN pour des tâches d'analyse prédictive, de classification, de détection d'anomalies, et bien d'autres. Les auteurs soulignent les défis spécifiques rencontrés dans l'application du deep learning à l'analyse de données massives, tout en mettant en lumière les orientations futures de la recherche dans ce domaine en constante évolution. En outre, l'article aborde l'utilisation de modèles de deep learning spécifiques tels que LSTM, DBN, CNN, et RNN dans des domaines tels que la prévision de charge de travail, la prédiction de l'état de santé des moteurs, la reconnaissance d'images, les chatbots, la classification du cancer, la détection d'activités, la détection d'intrusions, le diagnostic de pannes, la prédiction de faillite, l'analyse du marché boursier, la prédiction du trafic, la communication sans fil, l'analyse de la santé, et la prédiction de l'énergie.

Les études sur le deep learning dans le domaine de l'analyse de données massives présentent plusieurs limites. Tout d'abord, certaines études ne fournissent pas d'explication détaillée de la mise en œuvre des techniques de deep learning dans l'analyse de données massives. De plus, il est à noter que certaines études ne clarifient pas la relation entre les types d'analyse de données massives et les domaines d'application. En outre, la structure des études présentées manque parfois d'une organisation systématique, et la méthode de sélection des articles n'est pas toujours claire. Enfin, certaines études ne fournissent pas d'évaluation analytique ou de taxonomie de l'intégration des données massives dans les applications IoT dans différents domaines. Ces limitations soulignent la nécessité d'une approche plus complète et structurée dans la recherche sur le deep learning pour l'analyse de données massives.

Amir H. Gandomi et al., [46] dans une étude récente, des chercheurs présentent un article abordant l'utilisation croissante des technologies d'apprentissage automatique pour l'analyse des mégadonnées, mettant en avant l'importance de l'intelligence artificielle et des méthodes d'apprentissage automatique, en particulier l'apprentissage profond, pour relever les défis posés par le volume, la variété et la vitesse des données dans le contexte des mégadonnées. L'analyse des sentiments est un des domaines ex-

plorés dans l'article, permettant de comprendre les opinions et les émotions des utilisateurs à partir de leurs interactions en ligne. Une étude a été menée sur l'analyse des sentiments sur Twitter pour évaluer les attitudes populaires avant, pendant, et après les élections, en comparant ces opinions aux résultats réels des élections. Les auteurs ont construit un ensemble de données à partir de l'API Twitter, pré-traité les données, extrait les caractéristiques pertinentes à l'aide de TF-IDF, puis utilisé le classificateur Naive Bayes pour recueillir les opinions publiques. Un autre domaine abordé est la détection d'intrusions, qui utilise des techniques d'apprentissage automatique pour identifier les comportements malveillants dans les réseaux informatiques, renforçant ainsi la sécurité des systèmes. Les auteurs soulignent que ces technologies permettent non seulement de traiter efficacement les données mais aussi de découvrir des modèles et des tendances cachées, offrant ainsi des perspectives précieuses pour la prise de décisions éclairées. En résumé, l'article met en lumière l'importance croissante des technologies d'apprentissage automatique pour l'analyse des mégadonnées, en soulignant les divers domaines d'application et en mettant en avant les avantages qu'elles offrent pour relever les défis actuels en matière de données. Ces technologies sont essentielles pour exploiter pleinement le potentiel des mégadonnées et pour permettre aux entreprises et aux chercheurs de tirer des enseignements précieux à partir de ces vastes ensembles de données.

L'une des limites des études sur l'analyse des sentiments réside dans la difficulté à interpréter correctement les nuances et les subtilités des émotions humaines à partir de données textuelles. Les modèles d'apprentissage automatique peuvent parfois avoir du mal à saisir le contexte et l'ironie, ce qui peut entraîner des erreurs d'interprétation des sentiments. De plus, les biais présents dans les données d'entraînement peuvent influencer les résultats de l'analyse des sentiments, conduisant à des conclusions partielles ou erronées. En ce qui concerne la détection d'intrusions, une limitation majeure réside dans la capacité des attaquants à contourner les systèmes de sécurité basés sur l'apprentissage automatique. Les cybercriminels peuvent utiliser des techniques sophistiquées pour tromper les modèles d'apprentissage automatique et éviter d'être détectés, rendant la détection d'intrusion plus difficile et moins fiable. Il est donc essentiel de prendre en compte ces limitations lors de l'utilisation des technologies d'apprentissage automatique pour l'analyse des mégadonnées et de continuer à développer des

approches et des techniques plus robustes pour surmonter ces défis.

Wei Li et al., [47], dans une étude récente, les auteurs présentent un article examinant l'utilisation croissante de l'apprentissage automatique (ML) et de l'Internet des objets (IoT) dans le secteur de la santé, mettant l'accent sur les systèmes de santé intelligents. L'article explore différentes recherches utilisant des algorithmes prédictifs et des réseaux neuronaux artificiels pour résoudre des problèmes spécifiques tels que l'analyse des données massives, la prédiction des dommages sismiques, l'alignement des ontologies biomédicales, la modélisation de l'impact structurel sous-marin, les schémas de communication dans les villes intelligentes et les réseaux d'estimation de l'âge. It talks about the problems and opportunities that come with using automatic learning techniques to analyze large amounts of data in the health field. For example, it talks about how to manage resources, keep them safe, make sure they can work with other systems, and use AI to analyze large amounts of data. These are all areas that need extra attention for AI to be used successfully in health care. De plus, l'article discute des systèmes et des approches utilisées dans l'apprentissage automatique pour les soins de santé intelligents, tels que les systèmes de recommandation, les systèmes de prédiction, l'agrégation de données, l'assistance à la vie, et l'analyse sécurisée. Ces systèmes exploitent l'IoT et les techniques d'apprentissage automatique pour améliorer la surveillance de la santé, la prédiction des maladies, l'assistance aux patients, et la sécurité des données dans les environnements de soins de santé. En conclusion, l'utilisation de l'apprentissage automatique dans l'IoT pour les soins de santé présente des opportunités pour améliorer les soins aux patients et les processus de prise de décision.

Certaines études sur l'apprentissage automatique et l'Internet des objets (IoT) dans les soins de santé peuvent être basées sur des modèles théoriques ou des simulations informatiques, ce qui pourrait ne pas refléter pleinement la complexité et la variabilité des environnements réels de soins de santé. De plus, la disponibilité et la qualité des données utilisées dans ces études peuvent varier, ce qui pourrait affecter la précision et la fiabilité des résultats obtenus. Par ailleurs, les études pourraient se concentrer sur des aspects spécifiques de l'apprentissage automatique et de l'IoT dans les soins de santé, laissant de côté d'autres aspects importants tels que la confidentialité des données, l'éthique ou l'acceptation par les utilisateurs. De plus, certaines recherches pourraient ne pas prendre en compte les contraintes budgétaires ou les ressources limitées dans

les environnements de soins de santé réels, ce qui pourrait limiter la mise en œuvre pratique des solutions proposées. Enfin, il est important de noter que les technologies évoluent rapidement, et certaines des études examinées pourraient ne pas prendre en compte les dernières avancées en matière d'apprentissage automatique et d'IoT, ce qui pourrait limiter la pertinence à long terme des résultats obtenus.

Amir Masoud Rahmani et al., [48] Dans une étude récente, nous présentons un article approfondi sur les mécanismes d'analyse de données volumineuses pilotées par l'intelligence artificielle, mettant l'accent sur la santé, l'agriculture, les médias sociaux, et d'autres domaines. L'article explore diverses techniques, dont l'apprentissage automatique, les méthodes basées sur les connaissances, les algorithmes de prise de décision, et les méthodes de recherche. Concernant l'analyse supervisée des données volumineuses, il met en avant des techniques telles que les classificateurs ensemblistes basés sur des algorithmes de forêt aléatoire, les machines à vecteurs de support, et d'autres. L'importance de l'évolutivité, de l'efficacité, de la précision et de la confidentialité est soulignée. De plus, une vue d'ensemble des mécanismes d'apprentissage automatique comme les réseaux de neurones récurrents, les réseaux neuronaux LSTM, et les réseaux neuronaux convolutifs sont proposés. L'article explore également des méthodes de recherche et d'optimisation, telles que les algorithmes évolutionnaires multi-objectifs, l'optimisation par essai de particules et l'approximation stochastique de perturbation simultanée, en mettant en évidence leurs avantages et inconvénients. Des études de cas concrets sont présentées dans des domaines spécifiques, mais l'examen incomplet d'articles non directement liés aux données volumineuses pourrait introduire un biais dans la compréhension globale des mécanismes d'analyse de données volumineuses.

Une autre limitation de l'étude réside dans le manque de comparaison technique par rapport aux méthodes proposées, ce qui pourrait restreindre la capacité des chercheurs à évaluer pleinement l'efficacité des différentes approches. De plus, le processus de sélection des articles n'est pas clairement évoqué, suscitant interrogations sur la représentativité de la littérature examinée. En outre, une catégorisation des articles en fonction de certains facteurs n'est pas fournie, compliquant la comparaison et la synthèse des résultats. Certains articles examinés ne présentent pas de métriques qualitatives pour évaluer les techniques d'analyse de données, limitant ainsi la compréhension de la qualité des méthodes étudiées. Enfin, l'absence d'une taxo-

nomie détaillée basée sur les techniques d'intelligence artificielle dans l'étude pourrait restreindre la compréhension des différentes approches utilisées dans l'analyse de données volumineuses. Ces limitations soulignent la nécessité de mener des recherches supplémentaires pour combler ces lacunes et améliorer la compréhension des mécanismes d'analyse de données volumineuses pilotées par l'intelligence artificielle.

2.3.1 Comparaison des Études Classiques

Cette étude se consacre à une analyse approfondie des revues académiques spécialisées dans les études classiques. En examinant les divergences et similitudes entre ces publications en matière de contenu, de méthodologie et d'approches éditoriales, elle aspire à fournir un aperçu éclairé aux chercheurs et étudiants s'investissant dans ce domaine captivant, comme présenté dans le Tableau 2.1.

TABLE 2.1 – Comparaison des Revues d'études classiques avec des approches d'Apprentissage Automatique

Référence de l'étude	Année	Type d'étude	Type d'apprentissage	Contribution Principale	Taxonomie
Jiyang Xie et al. [44]	2018	Classique	Apprentissage profond	Analyse des données massives mobiles	Caractéristiques des MBD, méthodes d'analyse
Walaa Medhat et al. [45]	2024	Classique	Apprentissage profond	Techniques de deep learning pour l'analyse de données massives	Taxonomie des techniques, applications, défis
Amir H. Gandomi et al. [46]	2022	Classique	Multi-type	Utilisation croissante de l'apprentissage automatique pour l'analyse des mégadonnées	Domaines d'application, technologies d'apprentissage automatique
Wei Li et al. [47]	2021	Classique	Multi-type	Utilisation croissante de l'apprentissage automatique et de l'IoT dans le secteur de la santé	Applications dans la santé, systèmes et approches utilisés
Amir Masoud Rahmani et al. [48]	2021	SLR	Apprentissage automatique Supervisée	Mécanismes d'analyse de données volumineuses pilotées par l'intelligence artificielle	Techniques d'analyse, méthodes de recherche et d'optimisation

2.3.2 Limitations des Études Classiques

Dans cette partie spécifique traitant des limites des études classiques, il est possible que les recherches traditionnelles sur un sujet donné présentent des contraintes par rapport aux études qui recourent à la méthode SLR (Revue Systématique de la Littérature) [49] :

- a) **Sélection biaisée des articles** : Les recherches classiques peuvent être limitées dans leur sélection d'articles, entraînant un biais potentiel dans les résultats et ne couvrant pas l'intégralité du domaine de recherche.
- b) **Manque de méthodologie claire** : Les recherches conventionnelles peuvent manquer d'une méthodologie explicite pour l'identification, la sélection et l'évaluation des articles de recherche pertinents, ce qui pourrait compromettre la fiabilité et l'objectivité des conclusions obtenues.
- c) **Manque d'exhaustivité** : Les recherches traditionnelles pourraient omettre d'inclure l'ensemble des travaux de recherche disponibles sur le sujet, engendrant ainsi une perspective partielle ou incomplète de l'état de l'art et des avancées récentes.
- d) **Absence d'évaluation de la qualité de la recherche** : Les recherches conventionnelles peuvent ne pas procéder à une évaluation de la qualité méthodologique des articles inclus, ce qui peut nuire à la validité des conclusions formulées.
- e) **Absence de synthèse systématique des résultats** : Les recherches classiques peuvent souffrir d'une absence de synthèse systématique et organisée des résultats des différentes études, compliquant ainsi la comparaison et l'identification de tendances ou de motifs communs.

En revanche, les études recourant à la méthode SLR sont élaborées de manière à dépasser ces limitations. Elles adoptent une méthodologie rigoureuse pour identifier, sélectionner et évaluer de manière systématique les articles pertinents, offrant ainsi une analyse approfondie et impartiale des résultats de recherche. Ceci permet d'obtenir une vision globale et précise de l'état de l'art dans un domaine spécifique.

2.4 Motivation pour l'utilisation de la Revue Systématique de la Littérature (SLR)

La Revue Systématique de la Littérature (SLR) est un outil essentiel pour synthétiser les connaissances existantes sur un sujet donné. En utilisant une SLR, les chercheurs peuvent identifier les tendances, les facteurs de succès, les risques potentiels, et les bonnes pratiques liées à l'agilité dans les projets bancaires. Cela permet d'obtenir une vue d'ensemble complète et objective du domaine, aidant ainsi les praticiens et les décideurs à prendre des décisions éclairées basées sur des preuves solides.

2.4.1 Avantages de la SLR par rapport aux enquêtes Classiques

La Revue Systématique de la Littérature (SLR) présente plusieurs avantages par rapport aux enquêtes classiques [50] :

1. **Synthèse exhaustive et rigoureuse** : La Revue Systématique de la Littérature (SLR) permet une consolidation approfondie et méthodique des connaissances existantes sur un sujet donné, reposant sur une méthodologie transparente et reproductible.
2. **Identification objective des éléments clés** : En utilisant la SLR, il est possible d'identifier objectivement les tendances, les facteurs de succès, les risques potentiels, et les bonnes pratiques, offrant ainsi une vue d'ensemble complète du domaine étudié.
3. **Base solide pour des décisions éclairées** : La SLR autorise les praticiens et les décideurs à prendre des décisions éclairées fondées sur des preuves solides, ainsi contribuant à améliorer la qualité des projets et des processus, notamment dans des secteurs critiques tels que le secteur bancaire.

2.4.2 Importance de Surmonter les Limitations Identifiées

Il est crucial de surmonter les limitations identifiées dans une revue systématique de la littérature (SLR) afin d'améliorer la qualité et la fiabilité des résultats obtenus. En effet, en reconnaissant et en adressant ces limitations, les chercheurs peuvent renforcer

la validité interne et externe de leur étude, garantir la pertinence et la généralisabilité des conclusions, et contribuer à l'avancement des connaissances dans le domaine de recherche concerné. En outre, en prenant en compte et en corrigeant les limitations, les chercheurs peuvent accroître la crédibilité de leur travail et favoriser une meilleure acceptation de leurs résultats par la communauté scientifique.

Enfin, en surmontant les limitations identifiées, les chercheurs peuvent également ouvrir la voie à de futures recherches plus approfondies et plus robustes, permettant ainsi de progresser dans la compréhension des phénomènes étudiés [51].

2.5 Conclusion

Ce chapitre a fourni une vue d'ensemble complète de l'application de l'apprentissage automatique dans l'analyse des Big Data, en établissant le contexte et en soulignant son importance. Les fondements de la discipline, y compris les définitions et concepts clés, ainsi que les techniques d'apprentissage supervisé, non supervisé, par renforcement et profond, ont été explorés. Le chapitre a mis en lumière les défis et opportunités liés à la gestion du volume et de la variété des données, ainsi que la scalabilité des algorithmes. La revue des études classiques a révélé des limites telles que la sélection biaisée des articles et le manque de méthodologie claire, ce qui a motivé l'introduction de la Revue Systématique de la Littérature (SLR) comme approche plus rigoureuse. Les avantages de la SLR, notamment sa capacité à consolider les connaissances, identifier les lacunes de recherche, et promouvoir une méthodologie plus solide, ont été discutés. En résumé, ce chapitre jette les bases pour une exploration approfondie de l'application de l'apprentissage automatique dans l'analyse des Big Data, en insistant sur l'importance de la rigueur méthodologique.

Dans le prochain chapitre, qui constitue le cœur de notre travail, nous examinerons une revue de la littérature sur les méthodes d'apprentissage automatique pour l'analyse des Big Data.

Chapitre 3

Revue Systématique de la Littérature

3.1 Introduction

L'explosion des données numériques dans le monde moderne a entraîné une demande croissante de méthodes efficaces pour traiter et analyser ces vastes ensembles de données, communément appelés Big Data. Parmi les approches les plus prometteuses, l'apprentissage automatique se distingue par sa capacité à extraire des connaissances utiles à partir de données complexes et volumineuses.

Ce chapitre présente la méthodologie de notre étude de revue systématique de la littérature (SLR) sur les méthodes d'apprentissage automatique pour le traitement et l'analyse des Big Data. Nous commençons par exposer nos questions de recherche, qui guident notre exploration des études existantes. Ensuite, nous décrivons en détail notre méthodologie de recherche, y compris nos critères de sélection des études et notre processus de collecte et d'analyse des données. Nous pour suivons en présentant une classification des solutions d'apprentissage pour le traitement des Big Data, en mettant l'accent sur les algorithmes de classification. Nous examinerons les méthodes supervisées, non-supervisées, semi-supervisées et d'apprentissage profond, identifiant les limites de chaque approche. Enfin, nous répondons également aux questions de recherche initiales et proposons des orientations futures pour la recherche dans ce domaine en rapide évolution.

3.2 Méthodologie d'étude SLR

Le processus de réalisation d'une étude Systematic Literature Review (SLR) comprend généralement les étapes suivantes, qui sont résumées dans la Figure 3.1 :

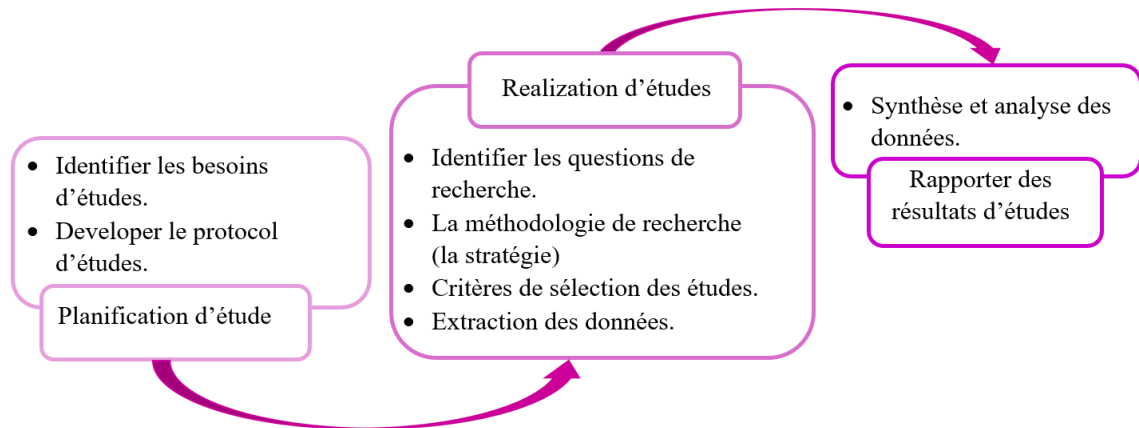


FIGURE 3.1 – Méthodologie d'étude Systématique.

- Planification de l'étude** : Identifier les besoins d'une étude systématique et développer le protocole de l'étude.
- Définition des questions de recherche** : Formuler clairement la question de recherche spécifique ou l'objectif que la SLR vise à traiter.
- Planification de la revue** : Elaborer un protocole bien défini qui décrit la stratégie de recherche, les critères de sélection, les méthodes d'extraction des données, et le plan d'analyse.
- Recherche d'études pertinentes** : Mener une recherche méthodique et approfondie dans diverses revues, actes de conférence et autres sources pertinentes pour identifier les études appropriées. La stratégie de recherche doit être transparente et clairement définie.
- Sélection et évaluation** : Appliquer des critères pour évaluer les études initialement identifiées en fonction de leur pertinence par rapport à la question de recherche, garantissant ainsi la sélection des études répondant à des critères spécifiques préalablement définis.
- Extraction et synthèse des données** : Extraire les données pertinentes des études sélectionnées, telles que les caractéristiques de l'étude, la méthodologie uti-

lisée et les résultats obtenus. Ensuite, synthétiser et analyser ces données pour identifier les tendances et les motifs émergents.

- g) **Discussions et rapport** : Interpréter les résultats synthétisés et examiner leurs implications par rapport à la question de recherche. Présenter les résultats de manière structurée pour faciliter leur compréhension et leur discussion.

3.3 Motivation et Contributions

La Revue Systématique de la Littérature (SLR) est motivée par la clarté des questions de recherche et vise à résoudre des problèmes spécifiques tout en minimisant les biais par des méthodes systématiques. Son objectif central est d'obtenir des résultats fiables pour permettre des conclusions robustes. La SLR contribue en fournissant des conclusions de haute qualité, établissant des normes méthodologiques, offrant une base solide pour la prise de décision, et réduisant le risque de biais et de publication sélective par l'enregistrement du protocole. En résumé, elle répond au besoin de produire des résultats fiables et constitue une ressource essentielle pour la prise de décision informée en recherche scientifique [52].

3.4 Méthodologie de Révision

Dans cette section, nous présentons en détail notre étude SLR, mettant en évidence les étapes clés du processus. Nous concentrons spécifiquement sur les algorithmes d'apprentissage de classification utilisés dans le traitement des Big Data.

3.4.1 Questions de Recherche

Comme illustré dans le Tableau 3.1, notre étude SLR a présenté trois questions de recherche.

TABLE 3.1 – Questions de Recherche

Questions	Objectifs
Q ₁ . Quels sont les types d'apprentissage automatique utilisés dans le traitement et l'analyse des Big Data ?	1. Présentation des méthodes utilisées pour analyser les Big Data.
Q ₂ . Quelles sont les applications de l'apprentissage automatique dans le traitement des Big Data ?	2. Utilisation de l'apprentissage automatique dans des secteurs comme la santé, la finance et l'industrie.
Q ₃ . Quelles sont les tendances émergentes dans les méthodes d'apprentissage automatique pour le traitement des Big Data ?	3. Exploration des nouvelles méthodes d'apprentissage automatique pour les Big Data.
Q ₄ . Quels sont les principaux défis et limitations des méthodes d'apprentissage automatique actuelles pour le traitement des Big Data ?	4. Analyse des obstacles dans l'utilisation de l'apprentissage automatique pour les Big Data.

3.4.2 Méthodologie de Recherche

Nous avons employé les termes clés ci-dessous pour repérer des articles pertinents liés au sujet de recherche en exploitant la base de données Google Scholar. En amalgamant des expressions alternatives et des synonymes, nous élaborons des requêtes de recherche en recourant à l'opérateur booléen "OU" et en utilisant "AND" pour fusionner les termes de recherche essentiels.

1. Mots-clés pour la recherche :

```
"Big Data" < OR > "Large-scale Data" < OR > "Massive Data"
< AND >
"Machine Learning" < OR > "ML" < OR > "Deep Learning"
< AND >
"Data Analytics" < OR > "Data Mining" < OR > "Data Science".
```

2. Recherche sur les Applications de ML dans l'Analyse de Big Data :

```
"Big Data" < OR > "Large-scale Data" < OR > "Massive Data"
< AND >
"Machine Learning" < OR > "ML" < OR > "Deep Learning"
< AND >
"Applications" < OR > "Use Cases" < OR > "Scenarios".
```

3. Recherche Générale sur l'Apprentissage Supervisé :

- "apprentissage supervisé" <OR> "supervised learning"
< AND >
- "machine learning" <OR> "ML" <OR> "apprentissage automatique".

4. Recherche sur les Algorithmes d'apprentissage Supervisé spécifiques :

- "régression linéaire" <OR> "linear regression"
< AND >
- "apprentissage supervisé" <OR> "supervised learning".

5. Recherche sur les méthodes d'apprentissage Supervisé populaires :

- "arbre de décision" <OR> "decision tree"
< AND >
- "classifieur bayésien naïf" <OR> "naive Bayes classifier"
< AND >
- "apprentissage supervisé" <OR> "supervised learning".

6. Recherche sur les modèles d'apprentissage supervisé pour la classification :

- "SVM" <OR> "Support Vector Machines" <OR> "machines à vecteurs de support"
< AND >
- "apprentissage supervisé" <OR> "supervised learning".

7. Recherche générale sur les algorithmes d'apprentissage Non Supervisé :

- "apprentissage non supervisé" <OR> "unsupervised learning"
< AND >
- "algorithmes" <OR> "méthodes"
< AND >
- "machine learning" <OR> "ML" <OR> "apprentissage automatique".

8. Recherche sur les Cartes auto-organisatrices de Kohonen :

- "cartes auto-organisatrices de Kohonen" <OR> "Kohonen self-organizing maps"
< AND >
- "apprentissage non supervisé" <OR> "unsupervised learning".

9. Recherche sur la méthode HDDC :

- "méthode HDDC" <OR> "High-Dimensional Data Clustering"
< AND >
- "algorithmes" <OR> "méthodes"
< AND >
- "apprentissage non supervisé" <OR> "unsupervised learning".

10. Recherche sur l'algorithme de Fuzzy C-Means (FCM) :

- "algorithme Fuzzy C-Means" <OR> "FCM algorithm"
< AND >
- "clustering flou" <OR> "fuzzy clustering"
< AND >
- "apprentissage non supervisé" <OR> "unsupervised learning".

11. Recherche sur l'algorithme GMM :

- "algorithme GMM" <OR> "Gaussian Mixture Model algorithm"
< AND >
- "apprentissage non supervisé" <OR> "unsupervised learning".

12. Recherche de l'algorithme de K-maense :

- "K-maens" <OR> "K-moyennes"
< AND >
- "Clustering" <OR> "Cluster Analysis"
< AND >
- "apprentissage non supervisé" <OR> "unsupervised learning".

13. Recherche Générale sur les algorithmes d'apprentissage Semi-Supervisé :

- "apprentissage semi-supervisé" <OR> "semi-supervised learning"
< AND >
- "algorithmes" <OR> "méthodes"
< AND >
- "machine learning" <OR> "ML" <OR> "apprentissage automatique".

14. Recherche sur l'algorithme de Hclustcompro :

- "algorithme de Hclustcompro" <OR> "Hclustcompro algorithm"

< AND >
- "classification ascendante hiérarchique" <OR> "CAH"
< AND >
- "apprentissage semi-supervisé" <OR> "semi-supervised learning".

15. Recherche sur l'algorithme GTC :

- "algorithme GTC" <OR> "Graph-Based Temporal Classification algorithm"
< AND >
- "reconnaissance automatique de la parole" <OR> "speech recognition"
< AND >
- "apprentissage semi-supervisé" <OR> "semi-supervised learning".

16. Recherche sur l'algorithme SimPLE :

- "algorithme SimPLE" <OR> "SimPLE algorithm"
< AND >
- "apprentissage semi-supervisé" <OR> "semi-supervised learning".

17. Recherche sur les algorithmes d'ordonnement transductif :

- "algorithme d'ordonnement transductif" <OR> "Transductive Ranking algorithm"
< AND >

- "classification bipartite" <OR> "bipartite classification"
< AND >
- "apprentissage semi-supervisé" <OR> "semi-supervised learning".

18. Recherche Générale sur les Algorithmes d'Apprentissage Profond :

- "apprentissage profond" <OR> "deep learning"
< AND >
- "algorithmes" <OR> "méthodes"
< AND >
- "machine learning" <OR> "ML" <OR> "apprentissage automatique".

19. Recherche sur l'algorithme Convolutional Neural Network (CNN) :

- "Convolutional Neural Network" <OR> "CNN"
< AND >
- "classification d'images" <OR> "traitement d'images"
< AND >
- "deep learning" <OR> "apprentissage profond".

20. Recherche sur les Réseaux Neuronaux Récursifs (RNN) :

- "Réseaux Neuronaux Récursifs" <OR> "RNN"
< AND >
- "séries temporelles" <OR> "traitement du langage naturel"
< AND >
- "deep learning" <OR> "apprentissage profond".

21. Recherche sur les Stacked Autoencoders (SAEs) :

- "Stacked Autoencoders" <OR> "SAEs"
< AND >
- "surapprentissage" <OR> "autoencoders"
< AND >

- "deep learning" <OR> "apprentissage profond".

22. Recherche sur les Réseaux Génératifs Antagonistes (GANs) :

- "Réseaux Génératifs Antagonistes" <OR> "GANs"

< AND >

- "génération d'images" <OR> "synthèse de données"

< AND >

- "deep learning" <OR> "apprentissage profond".

3.5 Classification des Solutions d'Apprentissage Automatique pour l'Analyse des Big Data

Nous examinons les méthodes d'apprentissage automatique appliquées au traitement et à l'analyse des Big Data. Cette étude vise à établir une taxonomie des algorithmes de classification des Big Data, comme présentée dans la Figure 3.2.

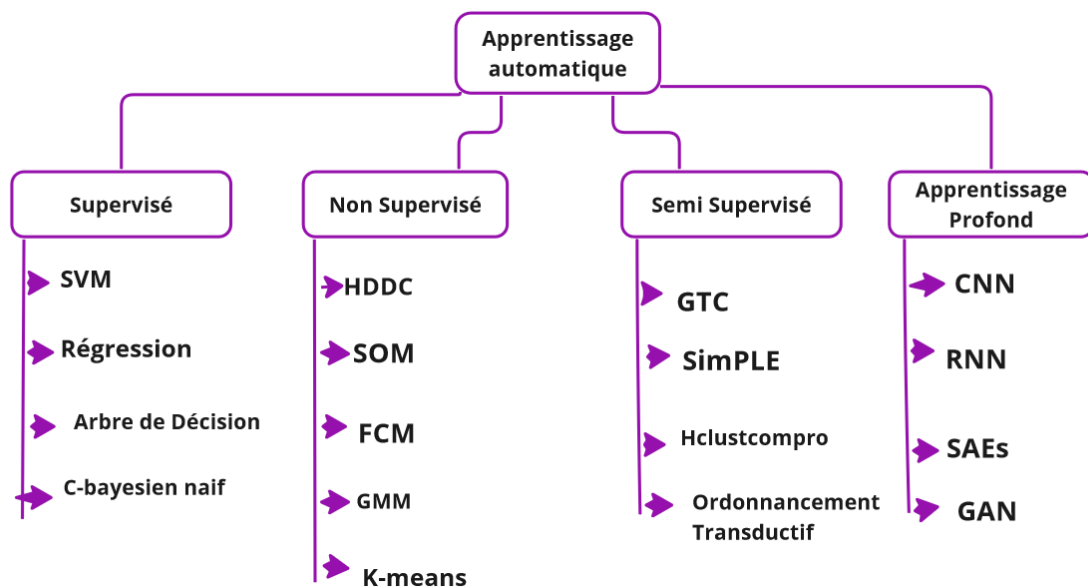


FIGURE 3.2 – Classification des Algorithmes de d'Apprentissage Automatique.

3.5.1 Les Algorithmes d'Apprentissage Automatique de Classification pour le Traitement des Big Data

Les algorithmes d'apprentissage de classification jouent un rôle crucial dans le traitement des Big Data, permettant de catégoriser et d'organiser efficacement des volumes massifs de données. Cette étude vise à examiner les différents types d'algorithmes de classification utilisés pour traiter les Big Data, en mettant en lumière leurs caractéristiques, leurs limitations.

3.5.1.1 Supervisé

L'apprentissage supervisé en machine learning utilise des données étiquetées pour entraîner un modèle à prédire des étiquettes pour de nouvelles données. Les algorithmes de classification supervisée attribuent des étiquettes de classe à des données en se basant sur des exemples d'entraînement. Dans cette section, nous explorerons les principaux algorithmes de classification supervisée appliqués aux Big Data.

1. **Support Vector Machines(SVM)** : Les SVM (Support Vector Machines) sont des outils essentiels pour le traitement des Big Data, utilisés pour classer, prédire et analyser de vastes ensembles de données en recherchant l'hyperplan optimal qui maximise la marge entre les différentes classes dans l'espace de données. Ils sont particulièrement adaptés aux données linéairement séparables, mais nécessitent un choix judicieux du noyau et peuvent être coûteux en calcul pour de grandes quantités de données. Malgré ces limitations, leur capacité à minimiser l'erreur de généralisation en maximisant la marge en fait des outils puissants pour la classification binaire et multiclasse, ainsi que pour la régression, et ils sont largement utilisés dans des domaines tels que la vision par ordinateur et la bioinformatique. L'algorithme DPWSS répond aux limites des SVM traditionnels en introduisant de la randomisation dans le processus d'entraînement, ce qui garantit la confidentialité différentielle et améliore l'efficacité tout en réduisant la sensibilité au choix du noyau et en optimisant le coût computationnel pour des volumes importants de données. Cette approche permet de protéger la confidentialité des données tout en maintenant des performances de classification efficaces, faisant de l'algorithme DPWSS une contribution significative aux méthodes de SVM avec protection de

la vie privée [53].

Limitation : La limite des SVM est qu'ils peuvent être sensibles au choix du noyau, ce qui peut nécessiter une expérimentation approfondie pour obtenir de bons résultats. De plus, ils peuvent être coûteux en termes de calcul pour de grandes quantités de données, en particulier lorsque le noyau utilisé est complexe.

- Régression Logistique :** La méthode de régression logistique est une approche couramment utilisée pour prédire une variable binaire en fonction d'un ensemble de variables explicatives. Elle peut également être étendue pour traiter des problèmes de classification multiclass. Le LASSO (Least Absolute Shrinkage and Selection Operator) est une méthode de régression pénalisée qui favorise la sélection de variables significatives en réduisant les coefficients des variables moins importantes vers zéro. En appliquant une pénalité de type L1 à la fonction de coût de cette méthode, le LASSO permet d'obtenir un modèle plus simple et plus interprétable en éliminant les variables moins pertinentes. Cette approche est largement utilisée dans le traitement et l'analyse des Big Data pour la sélection de variables et la régularisation dans les modèles de régression, avec pour objectif de créer un modèle plus précis en éliminant les variables moins importantes, ce qui peut améliorer les performances prédictives du modèle [54].

Limitation : La limite de la régression logistique est sa linéarité supposée entre les variables explicatives et la variable cible. Cette contrainte peut restreindre sa capacité à modéliser des relations complexes entre les variables, ce qui peut conduire à des performances moindres dans les cas où les données présentent des interactions non linéaires.

- Arbre de Décision :** Un arbre de décision est un modèle d'apprentissage automatique qui utilise une structure en arbre pour représenter des décisions et leurs conséquences. Il fonctionne en divisant récursivement l'ensemble des données en sous-ensembles plus petits en fonction des caractéristiques des données, jusqu'à ce qu'une condition d'arrêt soit atteinte, comme une profondeur maximale de l'arbre ou un critère de pureté des nœuds. Chaque nœud de l'arbre représente une caractéristique et chaque branche une décision basée sur cette caractéristique. Les arbres de décision sont appréciés pour leur facilité d'interprétation et leur capacité à gérer des données non linéaires, ce qui en fait des outils populaires pour

la classification et la prédiction dans divers domaines. L'algorithme C4.5 est une méthode d'arbre de décision qui utilise le critère de gain d'information normalisé pour diviser les ensembles de données, évitant ainsi les erreurs de surajustement. Il utilise deux critères pour classer les tests possibles : le premier critère est le gain d'information pour minimiser l'entropie des sous-ensembles, et le deuxième critère est le ratio de gain pour diviser le gain d'information en fonction des informations sur les résultats des tests. Cela permet de sélectionner les attributs pertinents pour la classification, qu'ils soient nominaux ou numériques. L'algorithme C4.5-SHO est une technique d'optimisation du gain d'information pour les arbres de décision basés sur l'entropie, visant à améliorer la classification des données en optimisant le gain d'information, le rendant plus efficace que d'autres algorithmes tels que C4.5, ID3 et CART [55].

Limitation : La limite des arbres de décision est leur sensibilité au surapprentissage, en particulier avec des arbres profonds. De plus, ils peuvent être instables avec de petites variations dans les données d'entraînement, ce qui peut conduire à des modèles très différents pour des ensembles de données légèrement différents.

4. **Classifieur Bayésien Naïf :** Le classifieur bayésien naïf est un algorithme de classification qui repose sur le théorème de Bayes et l'hypothèse simplificatrice d'indépendance conditionnelle des caractéristiques. Il est utilisé pour attribuer une classe à un nouvel échantillon en calculant la probabilité a posteriori de chaque classe, en se basant sur les probabilités a priori des classes et les probabilités conditionnelles des caractéristiques pour chaque classe. L'échantillon est finalement assigné à la classe ayant la probabilité la plus élevée [56].

Limitation : La limite du classifieur bayésien naïf réside dans son hypothèse simpliste d'indépendance conditionnelle des caractéristiques. Cette hypothèse peut ne pas refléter la complexité des relations entre les caractéristiques dans des situations réelles, entraînant une précision de prédiction réduite.

Le Tableau 3.2 présente une comparaison des algorithmes de classification Supervisée pour le Traitement des Big Data.

TABLE 3.2 – Comparaison des algorithmes de classification Supervisé

Algorithme	Évolutivité	Efficacité	Précision	Confidentialité
SVM [53]	Moyenne	Haute	Haute	Moyenne
Régression logistique [54]	Haute	Moyenne	Moyenne	Haute
Arbres de décision [55]	Moyenne	Moyenne	Moyenne	Moyenne
Classifieur bayésien naïf [56]	Haute	Haute	Moyenne	Haute

3.5.1.2 Non Supervisé

L'apprentissage non supervisé constitue une branche fondamentale de l'intelligence artificielle, visant à extraire des informations utiles à partir de données. Parmi les différentes techniques non supervisées, plusieurs algorithmes se démarquent par leur efficacité et leur polyvalence dans la résolution de divers problèmes de traitement de données. Dans cette section, nous explorerons certains de ces algorithmes, en mettant en évidence leurs principes de fonctionnement ainsi que leurs limitations.

1. **Cartes auto-organisatrices de Kohonen (SOM) :** La Carte auto-organisatrice de Kohonen (SOM) est un réseau de neurones à apprentissage non supervisé qui permet de répartir des données en groupements similaires. Elle crée une grille bi-dimensionnelle où les neurones de la grille représentent les clusters. Les données sont associées à des neurones en fonction de leur similarité, et la topologie de la carte reflète cette similarité. Les paramètres initiaux, tels que les dimensions de la grille, influencent fortement la structure de la carte. L'algorithme de Fuzzy Self-Organizing Map (FSOM) est une extension de l'algorithme SOM, où des concepts de logique floue sont intégrés. FSOM attribue des degrés d'appartenance flous aux neurones de la grille, permettant ainsi une représentation plus nuancée des clusters. Contrairement à SOM où l'appartenance est binaire, dans FSOM, les degrés d'appartenance peuvent varier de 0 à 1, reflétant ainsi l'incertitude ou l'ambiguïté dans les données. FSOM est également un algorithme de classification non supervisée, qui organise les données en clusters sans utiliser d'étiquettes de classe préexistantes [57].

Limitation : L'algorithme des cartes auto-organisatrices de Kohonen présente une limitation notable en raison de sa sensibilité aux paramètres tels que la taille de la carte, le taux d'apprentissage et le nombre d'itérations. Des ajustements inappropriés de ces paramètres peuvent compromettre la qualité de la

représentation des données obtenue par l'algorithme.

2. **Méthode de HDDC** : La méthode HDDC (High-Dimensional Data Clustering) est un algorithme de clustering non supervisé conçu pour traiter des ensembles de données en grande dimension. Lorsque les données ne comportent pas d'étiquettes de classe pour chaque observation, l'estimation des paramètres du modèle par maximum de vraisemblance devient complexe. L'algorithme EM (Expectation-Maximization) est couramment utilisé avec la méthode HDDC pour estimer les paramètres du modèle. Il alterne entre deux étapes : l'espérance, où les valeurs attendues des variables cachées sont estimées, et la maximisation, où les paramètres du modèle sont ajustés pour maximiser la probabilité des données observées. L'algorithme EM est largement utilisé dans divers domaines, y compris le clustering, la classification et l'estimation des modèles de mélange [58].

Limitation : La méthode HDDC présente une limitation liée à l'applicabilité restreinte des méthodes statistiques traditionnelles, notamment les distributions gaussiennes multivariées, en raison de la grande dimensionnalité des données. De plus, les tailles d'échantillon souvent limitées dans le contexte des données à haute dimension posent des défis pour le contrôle des erreurs de type I et II lors de l'analyse des données omiques.

3. **La méthode des Fuzzy C-Means (FCM)** : La méthode des Fuzzy C-Means (FCM) est un algorithme de clustering non supervisé, également connu sous le nom de clustering. Contrairement aux méthodes traditionnelles de clustering, FCM permet à chaque point de données d'appartenir à plusieurs clusters simultanément, avec des degrés d'appartenance flous. L'objectif de FCM est de minimiser une fonction objective en itérant pour mettre à jour les degrés d'appartenance et les centres de clusters jusqu'à ce que les critères d'arrêt soient satisfaits. FCM est largement utilisé pour segmenter un ensemble de données en groupes homogènes, où les données à l'intérieur de chaque cluster sont similaires, mais les clusters peuvent se chevaucher [59].

Limitation : La méthode Fuzzy C-Means (FCM) est confrontée à plusieurs limitations. Elle dépend fortement de paramètres tels que le nombre de classes et le paramètre de flou, ce qui complique leur sélection et influence grandement les résultats. De plus, FCM est sensible à l'initialisation des centres de classe

et des degrés d'appartenance, pouvant entraîner des solutions sous-optimales. Sa complexité informatique élevée pour de grands ensembles de données restreint son efficacité, et l'interprétation des résultats est compliquée en raison de la nature floue des partitions générées.

4. **Algorithme GMM** : L'algorithme GMM (Modèle de Mélange Gaussien) est une méthode d'apprentissage de type clustering, utilisée pour le clustering et la modélisation probabiliste des données. Il combine plusieurs distributions gaussiennes pour représenter une distribution de probabilité. En s'appuyant sur l'algorithme de l'espérance-maximisation (EM) pour estimer les paramètres du modèle, il alterne entre l'étape d'espérance (E-step) et l'étape de maximisation (M-step) pour modéliser des distributions de données complexes. Robuste et flexible, il offre une représentation précise des données, en faisant ainsi une approche puissante pour le clustering et l'analyse de données dans divers domaines [60].

Limitation : L'algorithme GMM (modèle de Mélange Gaussien) présente plusieurs limitations. Tout d'abord, il est sensible à la sélection initiale des paramètres. Cette sensibilité peut conduire à une convergence vers des optima locaux, impactant ainsi les performances globales du modèle.

5. **Algorithme de K-means** : L'algorithme K-means est un algorithme de clustering itératif simple utilisé dans le domaine de l'apprentissage automatique pour regrouper des données non étiquetées en K clusters. Chaque observation est assignée au cluster dont le centroïde est le plus proche, ce qui en fait une technique d'apprentissage non supervisé visant à former des groupes homogènes en termes de similarité [61].

Son fonctionnement :

- (a) **Initialisation des centroïdes** : Tout d'abord, un nombre K de centroïdes est choisi au hasard dans l'espace des données. Ces centroïdes initiaux peuvent être des points de données réelles ou générées aléatoirement. Cette étape est cruciale, car les centroïdes initiaux influencent la convergence de l'algorithme et la qualité des clusters obtenus [61].
- (b) **Attribution des points au cluster le plus proche** : Une fois les centroïdes initiaux définis, chaque point de données est attribué au cluster dont le centroïde est le plus proche en termes de distance euclidienne. Cela signi-

fié que chaque point est affecté à un cluster en fonction de la similarité de ses caractéristiques avec celles du centroïde [61].

- (c) **Recalcul des centroïdes** : Après que tous les points ont été attribués à des clusters, les centroïdes sont recalculés comme la moyenne des positions de tous les points attribués à chaque cluster. En d'autres termes, chaque coordonnée du centroïde est mise à jour pour être égale à la moyenne des coordonnées des points de données attribués à ce cluster [61].
- (d) **Répétition** : Les étapes (b) et (c) sont répétées de manière itérative jusqu'à ce qu'un critère d'arrêt soit atteint. Ce critère peut être défini par un nombre maximum d'itérations, une convergence des centroïdes (c'est-à-dire que les centroïdes ne changent plus significativement entre deux itérations consécutives), ou une diminution négligeable de la variation intra-cluster [61].
- (e) **Évaluation et interprétation des clusters** : Une fois l'algorithme convergé, les clusters obtenus sont évalués en examinant la cohésion intra-cluster et la séparation inter-cluster. Les clusters peuvent être interprétés et utilisés pour segmenter les données en fonction de leurs caractéristiques communes [61].

En résumé, l'algorithme K-means progresse en ajustant itérativement les centroïdes pour minimiser la distance intra-cluster tout en maximisant la distance inter-cluster, ce qui conduit à des clusters bien définis et compacts.

Limitation : L'algorithme K-means présente une limite : c'est sa sensibilité aux valeurs initiales des centroids. Comme l'algorithme dépend des centroids initiaux, différents ensembles initiaux peuvent conduire à des partitions différentes des données. En conséquence, l'algorithme peut converger vers un optimum local plutôt que global, ce qui peut affecter la qualité de la classification. De plus, K-means suppose que les clusters sont sphériques et isotropes, ce qui peut ne pas être le cas dans des ensembles de données réels où les clusters peuvent avoir des formes irrégulières ou être de tailles différentes. Cette limitation peut entraîner des résultats sous-optimaux dans certains cas, en particulier lorsque les données ne respectent pas cette hypothèse.

Note : Nous avons sélectionné l'algorithme **K-means** comme **cas d'étude**, qui sera le sujet du chapitre 4 de notre manuscrit.

Le Tableau 3.3 présente une comparaison des algorithmes de classification Non Supervisée pour le Traitement des Big Data.

Algorithme	Évolutivité	Efficacité	Précision	Confidentialité
SOM [57]	Moyenne	Moyenne	Moyenne	Moyenne
HDCC [58]	Faible	Moyenne	Moyenne	Moyenne
FCM [59]	Moyenne	Moyenne	Moyenne	Faible
GMM [60]	Haute	Haute	Haute	Haute
K-means [61]	Moyenne	Moyenne	Moyenne	Faible

TABLE 3.3 – Comparaison des algorithmes de Classification Non Supervisé

3.5.1.3 Semi Supervisé

L'apprentissage semi-supervisé représente une approche intermédiaire entre l'apprentissage supervisé et non supervisé, où une partie des données est étiquetée tandis que l'autre partie ne l'est pas. Dans cette section, nous explorerons quelques algorithmes de classification semi-supervisée, mettant en lumière leurs mécanismes ainsi que leurs limitations.

1. **Algorithme de Hclustcompro** : L'algorithme de classification ascendante hiérarchique par compromis (hclustcompro) est une méthode d'apprentissage semi-supervisée qui prend en compte deux sources d'information pour classifier des observations. Il modifie la mesure de distance dans l'algorithme de classification ascendante hiérarchique (CAH) en combinant de manière pondérée les dissimilarités initiales des deux sources d'information. Le paramètre de mélange est crucial et détermine l'importance relative de chaque source d'information dans le processus de classification. Pour choisir le paramètre de mélange optimal, l'algorithme minimise la différence absolue des corrélations entre les dissimilarités initiales et les distances cophénétiques, avec une procédure de rééchantillonnage pour garantir la robustesse du choix. Bien qu'initialement développée pour un problème archéologique spécifique, cette méthode est flexible et peut être appliquée à d'autres domaines nécessitant une classification précise en combinant différentes sources d'information [62].

Limitation : L'algorithme Hclustcompro présente plusieurs limitations, notamment sa forte dépendance à une connaissance préalable des deux sources d'information à combiner. Cette exigence peut être difficile à satisfaire ou à modéliser

avec précision, limitant ainsi son applicabilité dans les cas où ces informations ne sont pas disponibles ou sont sujettes à des erreurs.

- 2. Algorithme de GTC :** Le GTC (Graph-Based Temporal Classification) est un algorithme de classification utilisé en reconnaissance automatique de la parole pour l'apprentissage semi-supervisé avec des informations d'étiquetage sous forme de graphe. Il présente une nouvelle fonction objectif qui généralise la populaire fonction de perte CTC (Connectionist Temporal Classification) pour accepter des automates finis pondérés, permettant ainsi de s'entraîner avec des informations d'étiquetage sous forme de graphe avec des règles de transition et des poids de transition définis par l'utilisateur. Cette approche est appliquée à des problèmes de reconnaissance automatique de la parole semi-supervisée, où un graphe de structure similaire à CTC est généré à partir d'une liste N-best de pseudo-étiquettes pour l'auto-apprentissage. L'algorithme GTC vise à trouver à la fois les meilleurs alignements temporels et les séquences d'étiquettes optimales encodées dans le graphe, offrant ainsi de meilleures performances par rapport aux méthodes d'étiquetage pseudo-standard [63].

Limitation : L'algorithme GTC présente une limitation en raison de la complexité associée à la génération et à l'utilisation de graphes pondérés dans des environnements de grande échelle. Cela peut entraîner des coûts informatiques élevés et des défis en termes d'efficacité lors de l'entraînement et de l'inférence des modèles.

- 3. Algorithme de SimPLE :** L'algorithme SimPLE (Simulated Proximal Learning with Error) est une méthode semi-supervisée de classification qui se concentre sur les relations entre les échantillons non étiquetés. Il génère des pseudo-étiquettes pour les échantillons non étiquetés en sélectionnant et en affinant les prédictions sur plusieurs variations faiblement augmentées du même échantillon. SimPLE optimise le réseau de classification en utilisant trois objectifs d'entraînement : une perte supervisée pour les données étiquetées, une perte non supervisée qui aligne les données non étiquetées avec les pseudo-étiquettes générées à partir des données faiblement augmentées, et une perte de paire qui minimise la distance statistique entre les prédictions des données fortement augmentées, basées sur la similarité et la confiance de leurs pseudo-étiquettes. Pendant les tests, SimPLE

utilise la moyenne pondérée exponentielle des poids du modèle pour faire des prédictions. En résumé, SimPLE permet au modèle de s'adapter et de tirer parti des échantillons fortement augmentés [64].

Limitation : L'algorithme SimPLE présente une limitation, car il peut être sensible au bruit présent dans les données non étiquetées. Cette sensibilité peut altérer la précision des pseudo-étiquettes générées et, par conséquent, influencer les performances de classification.

4. **Algorithme d'Ordonnement Transductif :** Un algorithme d'ordonnement transductif est une méthode semi-supervisée utilisée dans la classification bipartite, où l'ensemble d'apprentissage comprend à la fois des instances étiquetées et non étiquetées. L'objectif est d'apprendre une fonction de score en exploitant la structure inhérente des données. Ces algorithmes construisent un graphe pondéré non orienté où les nœuds représentent les exemples et les arêtes renvoient leur similarité. Les partitions sont alors propagées à travers le graphe jusqu'à convergence, en utilisant les partitions initiales des instances étiquetées. Cela permet d'induire un ordre sur les instances non étiquetées, basé sur leur proximité dans le graphe. Toutefois, ces méthodes ne peuvent pas étiqueter les exemples absents de la phase d'apprentissage, les rendant transductives par opposition aux méthodes inductives [65].

Limitation : L'algorithme d'ordonnement transductif présente des limitations importantes. Il dépend fortement de la qualité de la similarité entre les instances, ce qui peut être difficile à déterminer dans certains cas. De plus, sa performance peut être affectée par la densité des données et la complexité de la structure du graphe construit. Ces facteurs peuvent conduire à des résultats imprécis, surtout lorsque les données sont dispersées ou présentent des structures complexes.

Le Tableau 3.4 présente une comparaison des algorithmes de classification Semi Supervisée pour le Traitement des Big Data.

TABLE 3.4 – Comparaison des algorithmes de Classification Semi Supervisé

Algorithme	Évolutivité	Efficacité	Précision	Confidentialité
Hclustcompro [62]	Moyenne	Haut	Haut	Moyenne
GTC [63]	Faible	Haut	Haut	Moyenne
SimPLE [64]	Moyenne	Haut	Moyenne	Moyenne
Ordonnancement Transductif [65]	Moyenne	Moyenne	Moyenne	Moyenne

3.5.1.4 Apprentissage Profond

L'apprentissage profond fait référence à un ensemble de techniques d'apprentissage automatique qui apprennent plusieurs niveaux de représentations dans des architectures profondes. Ces dernières années, divers algorithmes d'apprentissage profond ont été développés. Un bref aperçu des algorithmes d'apprentissage profond couramment utilisés dans l'analyse des Big Data est présenté ci-dessous.

1. **Convolutional Neural Network** : Les CNN (Convolutional Neural Networks) sont des algorithmes utilisés pour la classification d'images en exploitant des réseaux neuronaux artificiels, notamment pour des tâches complexes telles que la reconnaissance d'objets, la détection de visages et la classification d'images médicales. Ces réseaux peuvent bénéficier de modèles pré-entraînés et du transfert d'apprentissage, ce qui améliore la précision des prédictions et accélère l'apprentissage. Ils sont conçus pour imiter la structure des cerveaux humains et animaux, ce qui leur permet d'effectuer des calculs précis à travers des couches spécialisées, les rendant ainsi très efficaces dans la classification d'images.

D'autre part, l'algorithme MFFR (Multifeature Fusion Retrieval) est une méthode utilisée pour la récupération de modèles 3D qui se base sur la fusion de plusieurs caractéristiques. En utilisant des poids optimaux pour les paramètres qui contrôlent l'importance des caractéristiques sémantiques, de distribution de forme et de contexte, cet algorithme extrait des vues 2D à partir des modèles 3D, utilise un module d'attention interactif dans un CNN pour extraire des caractéristiques sémantiques, et utilise des algorithmes spécifiques pour extraire des caractéristiques globales. Les mesures de similarité entre le croquis et la vue 2D sont calculées en utilisant la distance euclidienne sur ces caractéristiques, et une somme pondérée de ces similarités est utilisée pour calculer la similarité globale. Les évaluations montrent que l'algorithme MFFR surpasse plusieurs autres

méthodes pour la récupération de modèles 3D en termes de précision et de rappel [66].

Limitation : Les limites de l’algorithme CNN incluent sa sensibilité au surapprentissage, en particulier avec des ensembles de données de petite taille. De plus, les CNN peuvent avoir du mal à traiter efficacement les données séquentielles ou temporelles, car ils ne capturent pas toujours les dépendances séquentielles de manière optimale. Enfin, les CNN peuvent être complexes en termes de calcul et de ressources nécessaires pour l’entraînement des modèles, ce qui peut limiter leur utilisation dans certaines applications.

2. **Réseaux Neuronaux Récursifs :** Les Réseaux Neuronaux Récursifs (RNN) sont des modèles d’apprentissage en profondeur efficaces pour traiter les séries temporelles et le langage naturel. Ils stockent l’information des entrées précédentes dans leur état interne pour prendre en compte les relations entre les données successives. Malgré des problèmes potentiels de disparition du gradient, des variantes comme les LSTM et les GRU ont été développées pour gérer les dépendances à long terme. Les RNN ont montré des performances remarquables dans la reconnaissance vocale, la traduction automatique et d’autres tâches de traitement du langage naturel. Ils sont adaptatifs aux données changeantes et efficaces pour extraire des caractéristiques à partir de gros volumes de données séquentielles, les rendant précieux dans divers domaines tels que la reconnaissance de la parole et la classification d’images. L’algorithme RNN-BiLSTM-CRF est une méthode d’ensemble pour la détection de la fraude à l’électricité qui combine un réseau de neurones récurrents (RNN), une mémoire à court terme bidirectionnelle longue (BiLSTM) et un champ de potentiels conditionnels (CRF). Le RNN est utilisé pour extraire automatiquement des caractéristiques importantes des données, tandis que le BiLSTM capture les dépendances séquentielles dans les données. Le CRF prend en compte les probabilités jointes de toutes les séquences d’étiquettes potentielles, permettant de modéliser les relations à longue portée dans les données séquentielles [38].

Limitation : La limite des réseaux de neurones récurrents (RNN) est leur difficulté à capturer efficacement les dépendances à long terme dans les données séquentielles en raison du problème de disparition du gradient.

3. **Stacked Autoencoders (SAEs) :** Les Stacked Autoencoders (SAEs) sont des techniques d'apprentissage profond couramment utilisées. Ils sont construits en empilant plusieurs autoencodeurs, qui sont des réseaux neuronaux feed-forward classiques. Un autoencodeur est un type de modèle d'apprentissage non supervisé composé de trois couches : une couche d'entrée, une couche cachée et une couche de sortie. L'entraînement d'un autoencodeur se fait en deux étapes : l'encodage, qui transforme les données d'entrée en une représentation cachée, et le décodage, qui reconstruit les données d'entrée à partir de la représentation cachée. Les SAEs sont généralement formés en deux étapes : la préformation, où chaque autoencodeur est formé couche par couche de manière non supervisée, et l'ajustement fin, où les poids sont mis à jour avec un ensemble d'entraînement étiqueté pour améliorer la performance du modèle [67].

Limitation : Une limite des Stacked Autoencoders (SAEs) est leur sensibilité au surapprentissage, surtout lorsque les données d'entrée sont bruitées ou incomplètes. Cette sensibilité peut compromettre leur capacité à généraliser efficacement.

4. **Réseaux Génératifs Antagonistes (GANs) :** Les Réseaux Génératifs Antagonistes (GANs) sont des modèles d'apprentissage profond composés de deux réseaux, un générateur et un discriminateur, qui s'affrontent pour générer de nouvelles données ressemblant à celles d'entraînement. Le générateur crée des données fictives à partir de bruit aléatoire, tandis que le discriminateur essaie de distinguer les données réelles des données fictives, les deux réseaux s'améliorant mutuellement jusqu'à ce que le discriminateur ne puisse plus différencier les deux types de données. Cependant, l'entraînement des GANs peut être difficile en raison de problèmes comme l'effondrement des modes, l'instabilité de l'entraînement et les problèmes de convergence. Malgré ces défis, les GANs sont largement utilisés dans divers domaines tels que la génération d'images, la traduction de style et la synthèse de données [68].

Limitation : Les GANs peuvent présenter des défis lors de leur entraînement, tels que l'effondrement des modes, l'instabilité de l'entraînement et les problèmes de convergence. Ces difficultés peuvent rendre le processus d'entraînement des GANs complexe et nécessiter des techniques avancées pour les surmonter.

Le Tableau 3.5 présente une comparaison des algorithmes de classification de type apprentissage profond pour le Traitement des Big Data.

Algorithme	Évolutivité	Efficacité	Précision	Confidentialité
CNN [66]	Haut	Haut	Haut	Moyenne
RNN [38]	Moyenne	Haut	Moyenne	Moyenne
SAEs [67]	Moyenne	Moyenne	Moyenne	Faible
GANs [68]	Haut	Moyenne	Haut	Moyenne

TABLE 3.5 – Comparaison des algorithmes de classification de type apprentissage profond .

3.6 Analyse de Résultats et Discussion

Dans cette partie, nous examinons de manière systématique et statistique les résultats de notre étude. La Figure 3.3 présente la distribution du nombre d'articles en fonction de leur année de publication, couvrant la période de 2020 à 2024. Le choix des articles inclus est déterminé selon les critères définis pour notre recherche, garantissant ainsi la pertinence et la représentativité de notre analyse. Par exemple, nous avons également calculé la moyenne et l'écart-type de l'année de publication des articles étudiés afin de mieux comprendre la tendance temporelle des recherches dans notre domaine d'étude. Ces statistiques offrent un aperçu précieux de l'évolution des travaux au fil du temps et peuvent aider à contextualiser nos résultats.

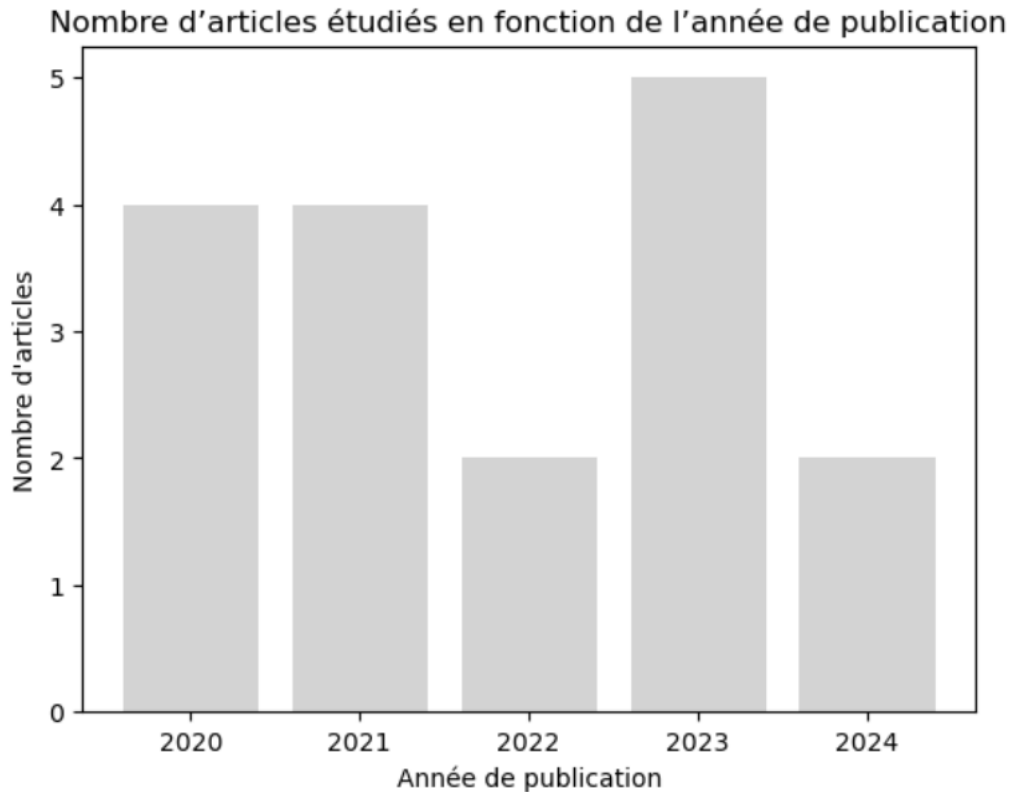


FIGURE 3.3 – Nombre d'articles étudiés en fonction de l'année de publication.

3.6.1 Réponses au Questions de Recherche :

Dans cette section, nous répondrons aux questions de recherche formulées dans le cadre de notre étude de revue systématique de la littérature (SLR).

Question 1 : Quels sont les types d'apprentissage automatique utilisés dans le traitement et l'analyse des Big Data ?

Dans le traitement et l'analyse des Big Data, les types d'apprentissage automatique utilisés.

1. **L'Apprentissage Supervisé** L'apprentissage supervisé consiste à former un algorithme à partir d'exemples de données étiquetées pour effectuer une tâche spécifique. Les algorithmes utilisent ces exemples pour apprendre à faire des prédictions sur de nouvelles données en ajustant leurs modèles en fonction des étiquettes correctes associées aux données d'entraînement. Cette méthode est largement utilisée dans des domaines tels que la classification de données, la régression linéaire, la reconnaissance d'images et la reconnaissance de la parole.
2. **L'Apprentissage non Supervisé** L'apprentissage non supervisé en machine

learning consiste à analyser et regrouper des données non étiquetées pour découvrir des modèles cachés ou des structures intrinsèques. Contrairement à l'apprentissage supervisé, il n'y a pas de guidage par des étiquettes correctes, laissant l'algorithme détecter autonome des motifs et des modèles. Les techniques telles que l'association, le clustering et la génération de modèles sont utilisées pour extraire ces structures, offrant des insights précieux sans besoin d'intervention humaine directe.

3. **L'Apprentissage Semi-Supervisé** L'apprentissage semi-supervisé est une approche intermédiaire entre l'apprentissage supervisé et non supervisé. Elle utilise à la fois des données étiquetées et non étiquetées pour entraîner un modèle, permettant ainsi de prédire les réponses pour des données non étiquetées similaires. Cette méthode combine les avantages des deux approches pour améliorer la précision des prédictions.
4. **L'Apprentissage Profond** L'apprentissage profond est une technique d'apprentissage automatique utilisant des réseaux de neurones pour analyser des données non structurées. Ces réseaux, inspirés du cerveau humain, sont composés de couches de neurones interconnectés. Plus le réseau est profond (avec de nombreuses couches), plus il peut résoudre des problèmes complexes. Les algorithmes d'apprentissage profond peuvent être supervisés (avec des données étiquetées) ou non supervisés (avec des données non étiquetées), et trouvent des applications avancées dans le traitement du langage naturel, permettant à une machine de comprendre le langage humain pour des résultats de pointe dans divers domaines.
5. **Par Renforcement** L'apprentissage par renforcement est une méthode d'apprentissage automatique où un agent IA ou un algorithme apprend de manière autonome en interagissant avec un environnement dynamique pour trouver les actions optimales à travers des essais et des erreurs successifs. Contrairement à l'apprentissage supervisé, il n'utilise pas de données étiquetées, mais vise à maximiser les récompenses. Cette méthode est utilisée dans des domaines tels que la gestion des robots, l'automatisation des usines et l'optimisation des chaînes de livraison.

Question 2 : Quelles sont les applications concrètes des méthodes d'apprentissage automatique dans le traitement des Big Data dans des domaines spécifiques comme la

santé, la finance ou l'industrie ?

Les applications concrètes des méthodes d'apprentissage automatique dans le traitement des Big Data dans des domaines spécifiques comme la santé, la finance ou l'industrie.

La Santé Les applications concrètes des méthodes d'apprentissage automatique dans le traitement des Big Data dans le domaine de la santé sont nombreuses et diverses. Voici quelques exemples :

1. **Diagnostic médical** : Les méthodes d'apprentissage automatique sont utilisées pour analyser de grandes quantités de données médicales, telles que :

Les réseaux de neurones convolutifs (CNN), les machines à vecteurs de support (SVM), les arbres de décision et les réseaux neuronaux récurrents (RNN). Les CNN sont particulièrement adaptés à l'analyse d'images médicales, comme les radiographies et les IRM, en identifiant des motifs complexes pour diagnostiquer des maladies (comme le cancer du sein et les tumeurs cérébrales). Les SVM sont utilisées pour classifier les données médicales en identifiant des frontières de décision optimales entre différentes classes de maladies (par exemple entre les patients atteints de schizophrénie et les autres). Les arbres de décision établissent des règles de décision basées sur les caractéristiques des patients et des données médicales pour aider au diagnostic (les types de cancer). Enfin, les RNN analysent les données séquentielles, comme les séries temporelles de données médicales, pour prédire l'évolution de certaines maladies.

- Ces méthodes permettent d'analyser de grandes quantités de données médicales et d'aider les médecins à diagnostiquer des maladies plus rapidement et avec une plus grande précision.

2. **Prévisions médicales** : Les méthodes d'apprentissage automatique peuvent être utilisées pour prévoir l'évolution de certaines maladies ou conditions médicales, ce qui peut aider les médecins à prendre des décisions plus éclairées sur le traitement et la gestion des patients. Parmi ces algorithmes : Les réseaux de neurones artificiels (ANN) analysent les données médicales, telles que les symptômes et les antécédents médicaux, pour prédire le risque de complications ou de rechute chez les patients. Les SVM identifient des modèles complexes dans les données

médicales, par exemple en prédisant la progression d'un cancer en fonction de caractéristiques spécifiques des patients.

3. **Recherche médicale :** Les méthodes d'apprentissage automatique peuvent être utilisées pour découvrir de nouveaux traitements potentiels et mieux comprendre les maladies. Par exemple, les réseaux de neurones artificiels (ANN) peuvent être utilisés pour analyser des données génétiques et identifier des gènes associés à des maladies spécifiques. Les machines à vecteurs de support (SVM) peuvent être utilisées pour prédire l'efficacité d'un médicament sur la base de données pharmacologiques et génétiques.
4. **la surveillance de la santé publique :** Les méthodes d'apprentissage automatique peuvent être utilisées pour détecter les tendances épidémiologiques et prévoir les épidémies de maladies. Par exemple, les modèles de réseaux de neurones récurrents (RNN) peuvent être utilisés pour analyser les données de surveillance des maladies infectieuses et détecter les tendances alarmantes, permettant aux autorités de santé de prendre des mesures préventives plus rapidement.

Finance Les applications concrètes des méthodes d'apprentissage automatique dans le traitement des Big Data dans le domaine financier sont nombreuses et diverses. Voici quelques exemples :

1. **Analyse des risques et des fraudes :** Les méthodes d'apprentissage automatique fournissent une approche puissante pour analyser les données transactionnelles et les comportements des clients, permettant ainsi d'identifier les transactions frauduleuses et de prédire les risques financiers. Parmi les algorithmes clés utilisés dans cette analyse, on trouve les SVM, les ANN, la régression logistique, les CNN et les RNN. Les SVM sont utilisés pour classer les emprunteurs en fonction de leur risque de défaut de paiement, les ANN sont efficaces pour identifier les schémas de fraude complexes, et la régression logistique est souvent utilisée pour évaluer la probabilité de défaut de paiement des emprunteurs. Les CNN peuvent être utiles pour l'analyse d'images. Ces méthodes permettent aux institutions financières d'améliorer la précision de leurs modèles d'évaluation du risque, de détecter les schémas anormaux et d'identifier les comportements suspects, contribuant ainsi de manière significative à l'analyse et à la prévention des risques et des fraudes dans le domaine financier.

2. Prévisions de marché et trading algorithmique :

Les méthodes d'apprentissage automatique offrent un moyen efficace d'élaborer des modèles de prévision de marché et de mettre en œuvre des opérations de trading algorithmique à grande échelle. Parmi les algorithmes utilisés à cet effet, on trouve le FSOM, le FCM et le GMM. Le FSOM permet une classification non supervisée avec des degrés d'appartenance variables, tandis que le FCM autorise chaque point de données à appartenir à plusieurs clusters simultanément. De son côté, le GMM combine plusieurs distributions gaussiennes pour représenter une distribution de probabilité. Ces algorithmes sont essentiels pour segmenter les données, détecter des schémas et des tendances, et prendre des décisions en matière de prévision de marché et de trading algorithmique. Cependant, chacun présente des avantages et des limitations spécifiques, ce qui nécessite une sélection judicieuse en fonction des besoins particuliers de l'application pour une analyse efficace des marchés financiers et la gestion des risques de fraudes.

L'industrie L'apprentissage automatique industriel progresse dans la maintenance prédictive et l'optimisation des processus, réduisant les temps d'arrêt et améliorant la productivité en analysant les données des capteurs.

1. **Maintenance prédictive** : Les méthodes d'apprentissage automatique analysent les données des capteurs des équipements industriels pour prédire les pannes potentielles et recommander des actions de maintenance préventive, réduisant ainsi les temps d'arrêt imprévus et les coûts de maintenance. Les algorithmes de maintenance prédictive, tels que les cartes auto-organisatrices de Kohonen et la méthode HDDC (High-Dimensional Data Clustering), organisent les données en clusters pour détecter les tendances et les anomalies, permettant ainsi d'anticiper les besoins de maintenance. Bien que ces méthodes offrent des avantages, leur sensibilité aux paramètres et l'initialisation des centres de classe peuvent limiter leur efficacité dans la prédiction des besoins de maintenance.
2. **Optimisation des processus de fabrication** : Les méthodes d'apprentissage automatique peuvent être précieuses pour repérer les inefficacités et optimiser les processus de fabrication, ce qui conduit à une amélioration de la productivité et de la qualité des produits. Les algorithmes d'optimisation des processus de fabrication, tels que les SVM, la régression logistique, les arbres de décision et

le classifieur bayésien naïf, jouent un rôle essentiel dans le traitement des Big Data dans l'industrie. Alors que les SVM cherchent l'hyperplan optimal pour la classification, la régression logistique simplifie les modèles tout en conservant leur précision. Les arbres de décision détectent les relations non linéaires, tandis que les classifieurs Bayésiens naïfs sont efficaces pour certains types de classification. Ensemble, ces algorithmes fournissent des solutions pour améliorer la classification, la prédiction et l'analyse des données, contribuant ainsi à une meilleure efficacité des processus de fabrication.

Question 3 : Quelles sont les tendances émergentes dans les méthodes d'apprentissage automatique pour le traitement des Big Data ?

Il y a plusieurs tendances émergentes dans les méthodes d'apprentissage automatique pour le traitement des Big Data qui peuvent être identifiées :

- a) **Utilisation de l'Apprentissage Profond :** Les méthodes d'apprentissage profond, telles que les CNN, les RNN, les SAEs et les GANs, gagnent en popularité pour traiter les Big Data, en particulier dans des domaines comme la reconnaissance d'images, la traduction automatique, la reconnaissance de la parole et la génération de données.
- b) **Approches semi-supervisées et transductives :** L'apprentissage semi-supervisé et transductif devient de plus en plus important pour exploiter efficacement les données étiquetées et non étiquetées. Ces approches permettent de tirer parti de grandes quantités de données non étiquetées, ce qui est courant dans les Big Data, pour améliorer les performances des modèles.
- c) **Apprentissage Multi modale :** Les agents informatiques sont maintenant en mesure de traiter et d'interpréter des données issues de divers modes, tels que visuel, acoustique, linguistique, physiologique et tactile. Cette capacité améliore la richesse et la précision de la communication entre l'homme et la machine.
- d) **Apprentissage par Renforcement :** L'apprentissage par renforcement, qui apprend en récompensant les comportements appropriés et en punissant les mauvais comportements, devient une méthode clé. Cette approche permet aux algorithmes de s'améliorer continuellement grâce à l'interaction avec leur environnement.
- e) **Apprentissage Automatique Quantique** L'intégration de l'informatique quantique avec le machine learning permet de réaliser certaines tâches bien plus rapide-

ment que les ordinateurs traditionnels. Cette avancée promet des progrès significatifs dans le traitement des Big Data.

En résumé, les tendances émergentes dans les méthodes d'apprentissage automatique pour le traitement des Big Data évoluent au fil du temps.

Question 4 Quels sont les principaux défis et limitations des méthodes d'apprentissage automatique actuelles pour le traitement des Big Data ?

les principaux défis et limitations des méthodes d'apprentissage automatique actuelles pour le traitement des Big Data peuvent être analysés comme suit :

- a) **Évolutivité et gestion des ressources** : L'un des principaux défis est la gestion des ressources informatiques nécessaires pour traiter les Big Data. Les méthodes d'apprentissage automatique doivent être capables de s'adapter à des ensembles de données massifs tout en maintenant des temps de traitement raisonnables et en évitant les goulets d'étranglement liés aux ressources matérielles.
- b) **Complexité des données** : Les Big Data sont souvent caractérisées par leur complexité, notamment en ce qui concerne la variété, la vitesse et le volume des données. Les méthodes d'apprentissage automatique doivent être capables de gérer cette complexité et d'extraire des informations significatives malgré le bruit et les irrégularités dans les données.
- c) **Gestion de la qualité des données** : Les données massives peuvent souvent être de qualité variable, avec des valeurs manquantes, des erreurs et des incohérences. Les méthodes d'apprentissage automatique doivent être robustes face à de telles imperfections et capables de traiter des données de qualité inégale sans compromettre les performances du modèle.
- d) **Interprétabilité des modèles** : Avec l'augmentation de la complexité des modèles d'apprentissage automatique, leur interprétabilité devient un défi majeur. Il est essentiel de pouvoir comprendre et expliquer les décisions prises par les modèles, en particulier dans des domaines critiques tels que la santé et la finance.
- e) **Protection de la confidentialité et de la sécurité** : L'utilisation de données massives soulève des préoccupations concernant la confidentialité et la sécurité des informations. Les méthodes d'apprentissage automatique doivent intégrer des

mécanismes de protection des données sensibles tout en maintenant des performances élevées.

Pour améliorer les techniques actuelles d'apprentissage automatique dans le traitement des Big Data, plusieurs approches peuvent être envisagées :

- a) **Développement de techniques d'apprentissage en ligne** : Les méthodes d'apprentissage en ligne peuvent permettre de traiter efficacement les flux de données continus en adaptant les modèles de manière dynamique à mesure que de nouvelles données deviennent disponibles.
- b) **Intégration de l'apprentissage fédéré** : L'apprentissage fédéré permet de construire des modèles sur des données distribuées sans avoir besoin de les centraliser, ce qui peut contribuer à résoudre les problèmes liés à la confidentialité et à la sécurité des données.
- c) **Développement de techniques de réduction de dimensionnalité** : Les techniques de réduction de dimensionnalité peuvent aider à gérer la complexité des données en extrayant les caractéristiques les plus importantes tout en réduisant la dimensionnalité des données, ce qui peut améliorer les performances des modèles et réduire les temps de calcul.
- d) **Investissement dans la recherche sur la transparence et l'explicabilité des modèles** : Il est essentiel de développer des techniques permettant d'expliquer les décisions prises par les modèles d'apprentissage automatique, ce qui peut renforcer la confiance dans les résultats obtenus et faciliter l'adoption de ces technologies dans des domaines critiques.

En combinant ces approches et en continuant à investir dans la recherche et le développement, il est possible d'atténuer les défis et les limitations actuels de l'apprentissage automatique pour le traitement des Big Data et d'améliorer ainsi l'efficacité et la pertinence des modèles dans divers domaines d'application.

3.6.2 Considérations et Directions Futures

Les directions futures et les opportunités pour de nouvelles recherches dans le domaine des solutions d'apprentissage pour le traitement des Big Data sont vastes et pro-

metteuses. Voici quelques exemples spécifiques de domaines où des recherches peuvent être menées :

- a) **Santé et sciences de la vie** : Les Big Data dans le domaine de la santé comprennent des données provenant de dossiers médicaux électroniques, d'images médicales, de génomique, etc. Les recherches futures pourraient se concentrer sur le développement de solutions d'apprentissage qui peuvent identifier des tendances et des corrélations dans ces données pour améliorer les diagnostics, les traitements personnalisés et la découverte de médicaments.
- b) **Finance et commerce** : Dans le secteur financier, les données transactionnelles, les données de marché et les données comportementales des clients offrent des possibilités pour développer des solutions d'apprentissage qui peuvent améliorer la prévision des risques, la détection de fraudes, la personnalisation des services financiers et la gestion des investissements. De même, dans le commerce, l'analyse des Big Data peut conduire à des recommandations de produits plus précises, à une optimisation des stocks et à une compréhension plus fine du comportement des consommateurs.
- c) **Médias et divertissement** : Dans ce domaine, il existe des opportunités pour développer des solutions d'apprentissage qui peuvent analyser et comprendre les préférences des utilisateurs à partir de vastes ensembles de données telles que les historiques de visionnage, les interactions sur les réseaux sociaux et les commentaires des utilisateurs. Cela permettrait de personnaliser davantage les recommandations de contenu et d'améliorer l'engagement des utilisateurs.
- d) **Transport et logistique** : Les données provenant de capteurs embarqués, de systèmes de suivi des flottes, de données de trafic, etc., peuvent être exploitées pour améliorer l'efficacité des réseaux de transport et de logistique. Les recherches futures pourraient se concentrer sur le développement de solutions d'apprentissage pour la prévision de la demande, la planification des itinéraires, l'optimisation des opérations de logistique et la gestion du trafic.
- e) **Gouvernement et secteur public** : Les gouvernements et les organisations du secteur public collectent une grande quantité de données sur les citoyens, les infrastructures publiques, l'économie, etc. Les solutions d'apprentissage pour les Big Data dans ce domaine pourraient viser à améliorer la prestation des services publics,

la prise de décision basée sur des données probantes et la détection des fraudes.

En résumé, les opportunités de recherche dans le domaine des solutions d'apprentissage pour le traitement des Big Data sont vastes et s'étendent à de nombreux domaines d'application, offrant un potentiel significatif pour l'innovation et l'amélioration des processus dans divers secteurs.

3.7 Conclusion

Ce chapitre présente une revue systématique de la littérature offrant une analyse exhaustive des solutions d'apprentissage pour le traitement des Big Data. À travers une méthodologie rigoureuse, une variété d'approches et de techniques a été examinée, mettant particulièrement l'accent sur les algorithmes de classification supervisée, non supervisée, semi-supervisée et d'apprentissage profond. Cette étude a permis d'identifier les défis, les tendances et les lacunes du domaine, tout en mettant en évidence les contributions et les motivations essentielles à cette recherche. L'analyse des résultats a répondu aux questions de recherche initiales, ouvrant ainsi la voie à des orientations futures prometteuses pour l'évolution de ce domaine. En conclusion, cette revue systématique constitue une référence solide pour appréhender l'état actuel des solutions d'apprentissage dans le traitement des Big Data, tout en soulignant l'importance de continuer à explorer de nouvelles voies pour relever les défis à venir et saisir les opportunités émergentes.

Dans le prochain chapitre, nous examinerons une étude de cas dans le contexte de notre recherche pour mieux comprendre l'application pratique des concepts discutés. Nous utiliserons cette étude de cas pour visualiser les résultats obtenus et en discuter en détail.

Chapitre 4

Application et Évaluation d'une Méthode d'Apprentissage Automatique pour l'Analyse du Big Data : Étude de Cas

4.1 Introduction

Dans ce chapitre, nous explorons l'application et l'évaluation d'une méthode d'apprentissage automatique pour l'analyse du Big Data, mettant en avant son importance croissante. Nous commençons par définir le contexte et les objectifs de notre étude, en soulignant le rôle crucial de Google Colab et de Python dans le traitement efficace des données massives. Ensuite, nous détaillons la préparation des données, la sélection de la méthode d'apprentissage automatique, et sa pratique implémentation, mettant en évidence les techniques d'évaluation de sa performance. Enfin, nous présentons les résultats obtenus, discutons de leurs implications et conclusions sur la contribution de notre travail à l'analyse de Big Data.

4.1.1 Contexte et objectifs de l'implémentation dans l'étude de cas

Consistant en l'exploration et l'analyse du jeu de données Iris en utilisant l'algorithme de clustering K-Means, l'objectif principal est d'identifier des clusters dans les données d'Iris et d'évaluer ces clusters par rapport aux étiquettes réelles des espèces d'Iris. Cette approche vise à comprendre comment les données sont naturellement regroupées et à comparer les résultats de clustering avec les espèces d'Iris réelles pour évaluer la performance de l'algorithme. L'implémentation comprend des étapes de prétraitement des données, d'analyse exploratoire, de normalisation des données, de clustering avec K-Means, de visualisation des clusters en 2D et 3D, ainsi que d'évaluation des résultats à l'aide de différentes métriques de performance telles que l'exactitude, la pureté, l'indice de Rand ajusté, et le score de silhouette. Ces étapes offrent un cadre complet pour comprendre et évaluer l'efficacité de l'algorithme K-Means dans le contexte de l'ensemble de données Iris.

4.1.2 Importance de Colab et Python dans le traitement du Big Data

Dans le traitement et l'analyse du Big Data, l'utilisation de Google Colab et le langage de programmation Python revêt une importance particulière. Google Colab offre un environnement de développement intégré basé sur le cloud, permettant d'exécuter du code Python directement dans le navigateur sans nécessiter d'installation ni de configuration complexes. Cette plateforme offre de nombreux avantages, notamment l'accès gratuit aux ressources de calcul puissantes, la possibilité de partager et de collaborer facilement sur des notebooks, ainsi que l'intégration transparente avec d'autres services cloud de Google. Python est largement utilisé dans le domaine du Big Data en raison de sa simplicité, de sa polyvalence, et de sa richesse en bibliothèques et en outils pour le traitement et l'analyse des données. Python offre une syntaxe claire et concise, ce qui le rend accessible aux débutants tout en offrant des fonctionnalités avancées pour les utilisateurs expérimentés. De plus, Python dispose de bibliothèques telles que Pandas, NumPy, Matplotlib et Scikit-Learn, qui facilitent le traitement, la manipulation, la visualisation et l'analyse des données volumineuses caractéristiques du Big Data.

En résumé, l'utilisation de Google Colab et du langage Python joue un rôle crucial dans le traitement et l'analyse efficaces du Big Data, offrant des outils puissants et conviviaux pour les professionnels et les chercheurs travaillant dans ce domaine en constante évolution.

4.2 Configuration de l'Environnement de Développement

Dans le cadre de projets scientifiques et d'ingénierie nécessitant un environnement logiciel complexe, la configuration d'un environnement de développement approprié est essentielle pour garantir une productivité maximale et un accès efficace aux ressources informatiques. Pour répondre à ces exigences, nous avons choisi Google Colab après avoir effectué une veille technologique approfondie. Google Colab est une plateforme de notebook Jupyter hébergée en ligne, offrant un environnement de développement interactif dans le cloud.

4.2.1 Configuration de Google Colab

Pour configurer Google Colab comme environnement de développement pour des projets scientifiques et d'ingénierie nécessitant un environnement logiciel complexe, nous avons choisi Google Colab après une veille technologique en raison de sa disponibilité en ligne, de son support pour Python et de ses capacités de calcul dans le nuage. Pour commencer, nous avons accédé à Google Colab via un compte Google, puis créé un nouveau notebook Python. Ensuite, nous avons importé les bibliothèques Python nécessaires et spécifié la version de Python le cas échéant, tout en utilisant les spécificités à Colab pour charger des données depuis Google Drive, GitHub ou d'autres sources en ligne. Le montage de Google Drive a été effectué si nécessaire pour des fichiers stockés sur Google Drive. Pour une utilisation avancée, nous avons activé l'utilisation des GPU/TPU pour le calcul accéléré et assuré la sauvegarde et le partage des notebooks enregistrés sur Google Drive ou localement. Cette configuration permet une flexibilité et une reproductibilité efficace pour le développement dans le cloud, adaptée aux besoins spécifiques des projets scientifiques et d'ingénierie.

4.3 Description de l'Étude de Cas et Préparation des Données

Cette étude de cas utilise le clustering K-Means sur le jeu de données Iris pour identifier des regroupements naturels parmi les fleurs en fonction de leurs caractéristiques. Les données sont préparées, nettoyées, et analysées pour comprendre les relations entre les variables. Ensuite, K-Means est appliqué avec une visualisation en 2D et 3D, suivi de l'évaluation des performances avec des mesures telles que l'exactitude, la pureté et le score de silhouette. Une comparaison avec "AgglomerativeClustering" est également effectuée pour évaluer les résultats.

4.3.1 Présentation de l'étude de cas

Cette étude de cas explore l'application de l'algorithme de clustering K-Means au jeu de données Iris, largement utilisé en apprentissage automatique. L'objectif est de découvrir des regroupements naturels parmi les fleurs d'iris adaptés à leurs caractéristiques botaniques, telles que la longueur et la largeur des sépales et des pétales. Le processus commence par le chargement des données avec scikit-learn, suivi d'une analyse exploratoire pour comprendre les distributions et les relations entre les caractéristiques. Les données sont nettoyées, pré-traitées, normalisées, et les variables catégorielles sont encodées. En utilisant K-Means avec un nombre de clusters définis, les données sont regroupées puis visualisées en 2D et 3D pour interpréter les regroupements. Les performances sont optimisées avec diverses mesures telles que l'exactitude, la pureté, l'indice de Rand ajusté, et le score de silhouette. Une comparaison est également faite avec AgglomerativeClustering pour évaluer les différences de performance. Cette étude de cas démontre le processus complet d'application du clustering K-Means, de la préparation des données à l'évaluation des résultats, offrant ainsi un aperçu pratique de l'utilisation des techniques d'apprentissage automatique pour l'analyse de données.

4.3.2 Dataset

Le jeu de données Iris est un ensemble de données classique utilisé en apprentissage automatique et en statistiques. Il contient des mesures de quatre caractéristiques de

fleurs d'iris (longueur et largeur du sépale et du pétale) ainsi que l'espèce de fleur (iris setosa, iris versicolor ou iris virginica). Cet ensemble de données est souvent utilisé pour illustrer des concepts tels que la classification supervisée. L'objectif est de prédire l'espèce de fleur d'iris en fonction de ses caractéristiques optimisées.

4.3.3 Nettoyage et pré-traitement des données

Le nettoyage et le pré-traitement des données sont des étapes cruciales en analyse de données et en apprentissage automatique. Ces processus visent à préparer les données brutes afin de les rendre compatibles avec l'utilisation de modèles d'apprentissage automatique. Les étapes clés incluent la gestion des données manquantes par suppression ou imputation, la détection et la suppression des doublons, ainsi que le traitement des valeurs aberrantes qui pourraient altérer les résultats. Les données numériques sont souvent normalisées ou standardisées pour les mettre à la même échelle, et les variables catégorielles sont converties en variables numériques par encodage. La réduction de la dimensionnalité peut également être utilisée pour simplifier les données tout en préservant les informations essentielles. Enfin, les données sont généralement divisées en ensembles d'entraînement et de test pour évaluer les modèles. Ces étapes sont essentielles pour garantir la qualité des données, améliorer les performances des modèles, et assurer des analyses fiables en apprentissage automatique.

4.3.4 Analyse exploratoire des données

L'analyse exploratoire des données (AED) est une étape cruciale dans tout projet d'analyse de données. Elle vise à comprendre la nature, la structure, et les relations présentes dans un ensemble de données. L'AED comprend plusieurs tâches importantes, notamment l'examen des statistiques descriptives, la visualisation des données à l'aide de graphiques, l'analyse des corrélations entre les variables, et éventuellement l'utilisation de techniques de réduction de dimensionnalité. L'AED permet de détecter les tendances, les motifs, les valeurs aberrantes, et de formuler des hypothèses préliminaires qui guideront la modélisation ultérieure et l'analyse approfondie des données. En résumé, l'AED est une étape essentielle pour comprendre les caractéristiques et les potentiels cachés des données avant d'entreprendre des analyses plus avancées.

4.4 Sélection et Adaptation de la Méthode d'Apprentissage Automatique

Dans cette section, nous mettons en lumière la justification du choix de l'algorithme ainsi que l'adaptation de celui-ci aux spécificités de l'étude de cas.

4.4.1 Justification du choix de l'algorithme

Le choix de l'algorithme K-Means pour le clustering des données repose sur sa simplicité conceptuelle, son efficacité computationnelle et sa capacité à gérer efficacement des ensembles de données de taille modérée à grande. K-Means est largement utilisé en raison de sa facilité d'implémentation et de son temps d'exécution linéaire, ce qui le rend scalable pour des applications pratiques. Il produit des clusters bien définis et interprétables, adaptés à diverses tâches d'analyse exploratoire ou de segmentation.

4.4.2 Adaptation de l'algorithme aux spécificités de l'étude de cas

Nous avons adapté le nombre de clusters dans l'algorithme K-Means en utilisant la méthode du coude et l'analyse de la silhouette. Voici comment nous avons procédé :

1. **Choix du nombre de clusters (K) :** Nous avons décidé d'utiliser 3 clusters dans notre étude de cas, correspondant aux trois espèces d'iris dans le jeu de données. Ce choix a été fait en connaissance des trois espèces d'iris dans le jeu de données, mais dans d'autres cas, il aurait été possible d'explorer plusieurs valeurs de K pour trouver le nombre optimal de clusters.
2. **Méthode du coude :** Nous avons utilisé la méthode du coude pour déterminer le nombre optimal de clusters. Cette méthode consiste à tracer la somme des carrés des distances entre les points et leur centre de cluster le plus proche en fonction du nombre de clusters. Le coude dans le graphique représente le point où l'ajout d'un autre cluster n'explique pas beaucoup plus de variance. Nous avons observé visuellement que le coude était à $K=3$, ce qui a confirmé notre choix initial de 3 clusters.

3. **Analyse de la silhouette** : Nous avons également utilisé l'analyse de la silhouette pour évaluer la qualité de la segmentation des clusters. La silhouette mesure à quel point un point de données est similaire à son propre cluster par rapport aux autres clusters. Nous avons calculé le score de silhouette pour différentes valeurs de K et avons choisi celle qui maximisait ce score, ce qui a également abouti à $K = 3$ comme le meilleur nombre de clusters.

4.5 Implémentation de la Méthode

Dans cette section, nous détaillons l'implémentation de l'algorithme K-Means sur le jeu de données Iris, en passant en revue chaque étape du processus. De plus, nous examinons diverses méthodes d'évaluation visant à mesurer la qualité des clusters par rapport aux étiquettes de classe réelle du jeu de données.

4.5.1 Détails de l'implémentation de l'algorithme

Pour détailler l'implémentation de l'algorithme K-Means sur le jeu de données Iris, nous allons examiner chaque étape du processus, y compris la normalisation des données, l'initialisation et l'entraînement de l'algorithme K-Means, la visualisation des clusters obtenus, et l'évaluation des résultats.

Étape 1 : Charger les bibliothèques nécessaires

Pour commencer, chargeons les bibliothèques nécessaires pour effectuer les opérations de clustering, normaliser les données, et visualiser les résultats.

```
1 import pandas as pd
2 import numpy as np
3 import matplotlib.pyplot as plt
4 from sklearn.cluster import KMeans, AgglomerativeClustering
5 from sklearn.preprocessing import StandardScaler
6 from sklearn import datasets
7 from sklearn.metrics import accuracy_score, confusion_matrix,
8 classification_report, adjusted_rand_score, silhouette_score,
   homogeneity_score
9 from sklearn.metrics import confusion_matrix
10 from sklearn.decomposition import PCA
```



```
11 from mpl_toolkits.mplot3d import Axes3D
```

Étape 2 : Charger le jeu de données Iris

Nous importons le jeu de données Iris à partir de scikit-learn, puis nous stockons les caractéristiques dans un DataFrame appelé data. Les véritables étiquettes de classes sont également stockées dans un tableau appelé true_labels.

```
1 iris = datasets.load_iris()
2 data = pd.DataFrame(iris.data, columns=iris.feature_names)
3 true_labels = iris.target # Vraies étiquettes des classes
4 # Aperçu des données
5 print(data.head())
```

Étape 3 : Analyse exploratoire des données (EDA)

L'exploration initiale des données implique d'abord d'examiner la distribution des caractéristiques pour comprendre leur répartition. Ensuite, nous visualisons les relations entre les caractéristiques à l'aide de graphiques comme les histogrammes et les nuages de points. Enfin, nous analyserons les corrélations entre les caractéristiques pour détecter les éventuelles relations entre les variables.

```
1 # Visualisation des distributions des caractéristiques
2 data.hist(bins=20, figsize=(12, 8))
3 plt.suptitle('Distribution des caractéristiques des données Iris')
4 plt.show()
5
6 # Visualisation des relations entre les caractéristiques
7 pd.plotting.scatter_matrix(data, c=true_labels,
8 figsize=(12, 12), marker='o', hist_kwds={'bins': 20}, alpha=0.5)
9 plt.suptitle('Relations entre les caractéristiques')
10 plt.show()
11
12 # Examen des corrélations entre les caractéristiques
13 correlation_matrix = data.corr()
14 print("Matrice de corrélation des caractéristiques :")
15 print(correlation_matrix)
```

Étape 4 : Normalisation des données

Il est important de normaliser les données pour que toutes les caractéristiques soient à la même échelle. Utilisons StandardScaler pour normaliser les caractéristiques.

```

1 scaler = StandardScaler()
2 scaled_features = scaler.fit_transform(data)

```

Étape 5 : Application de l'algorithme K-Means pour le clustering

Procédons à l'initialisation et à l'entraînement de l'algorithme K-Means afin d'obtenir des clusters à partir des données normalisées.

```

1 num_clusters = 3 # Choix du nombre de clusters (correspondant aux
   trois esp ces)
2
3 # Initialiser et entra ner l'algorithme K-Means
4 kmeans = KMeans(n_clusters=num_clusters, random_state=42)
5 kmeans.fit(scaled_features)
6
7 # Ajouter les labels de cluster aux donn es
8 data['Cluster'] = kmeans.labels_

```

Étape 6 : Visualisation des clusters

Visualisons les clusters formés en utilisant les deux premières caractéristiques, telles que 'sepal length' et 'sepal width', pour une représentation en 2D et 3D afin d'illustrer la répartition des données par cluster.

```

1 # Visualiser les clusters avec K-Means
2 plt.figure(figsize=(10, 8))
3 plt.scatter(data['sepal length (cm)'], data['petal length (cm)'], c=
   data['Cluster'], cmap='viridis')
4 plt.xlabel('Sepal Length (cm)')
5 plt.ylabel('Petal Length (cm)')
6 plt.title('Clusters avec K-Means dans le jeu de donn es Iris')
7 plt.show()
8
9 # Visualisation des clusters en 3D avec K-Means
10 fig = plt.figure()
11 ax = fig.add_subplot(111, projection='3d')
12 ax.scatter(data['sepal length (cm)'], data['sepal width (cm)'],
13 data['petal length (cm)'], c=data['Cluster'], cmap='viridis')
14 ax.set_xlabel('Sepal Length (cm)')
15 ax.set_ylabel('Sepal Width (cm)')
16 ax.set_zlabel('Petal Length (cm)')
17 plt.title('Visualisation 3D des clusters avec K-Means

```

```

18 dans le jeu de données Iris')
19 plt.show()

```

Étape 7 : Analyser les caractéristiques moyennes des clusters

Le code affiche les caractéristiques moyennes des clusters.

```

1 print(data.groupby('Cluster').mean())

```

Étape 8 : Évaluation des résultats de K-Means

Les performances de K-Means sont réalisées à l'aide de différentes métriques.

```

1 # Définition de la fonction purity_score
2 def purity_score(y_true, y_pred):
3     cm = confusion_matrix(y_true, y_pred)
4     return np.sum(np.amax(cm, axis=0)) / np.sum(cm)
5
6 # Calcul de l'exactitude (accuracy)
7 accuracy = accuracy_score(true_labels, data['Cluster'])
8 print(f"Accuracy avec K-Means: {accuracy:.2f}")
9
10 # Calcul de la pureté des clusters
11 purity = purity_score(true_labels, data['Cluster'])
12 print(f"Purity avec K-Means: {purity:.2f}")
13
14 # Calcul de l'Indice de Rand Ajusté (Adjusted Rand Index)
15 ari = adjusted_rand_score(true_labels, data['Cluster'])
16 print(f"Indice de Rand Ajusté avec K-Means: {ari:.2f}")
17
18 # Calcul du score de silhouette
19 silhouette = silhouette_score(scaled_features, data['Cluster'])
20 print(f"Score de silhouette avec K-Means: {silhouette:.2f}")
21
22 # Classification report
23 report = classification_report(true_labels, data['Cluster'])
24 print("Rapport de classification avec K-Means:")
25 print(report)

```

Étape 9 : Évaluation par courbe d'inertie avec K-Means

Le code calcule la courbe d'inertie pour évaluer le bon nombre de clusters.

```

1 inertia_values = []
2 for k in range(1, 11):

```

```

3     kmeans_model = KMeans(n_clusters=k, random_state=42)
4     kmeans_model.fit(scaled_features)
5     inertia_values.append(kmeans_model.inertia_)
6 plt.figure()
7 plt.plot(range(1, 11), inertia_values, marker='o')
8 plt.xlabel('Nombre de clusters')
9 plt.ylabel('Inertie')
10 plt.title('Courbe d\'inertie pour K-Means')
11 plt.show()

```

Étape 10 : Visualisation de la distribution des clusters dans chaque classe

Le code visualise la distribution des classes réelles pour chaque cluster.

```

1 for cluster in np.unique(data['Cluster']):
2     plt.figure(figsize=(8, 6))
3     mask = data['Cluster'] == cluster
4     plt.hist(true_labels[mask], bins=3, alpha=0.7, label=f'Cluster {
5         cluster}')
6     plt.xlabel('Classe r elle')
7     plt.ylabel('Nombre')
8     plt.title(f'Distribution des classes r elles pour le Cluster {
9         cluster}')
10    plt.legend()
11    plt.show()

```

Étape 11 : Comparaison avec AgglomerativeClustering

Le code compare les résultats avec K-Means à l'aide de l'algorithme Agglomerative-Clustering

```

1 # Comparaison avec un autre algorithme de clustering :
2   AgglomerativeClustering
3 agglo = AgglomerativeClustering(n_clusters=num_clusters)
4 agglo_labels = agglo.fit_predict(scaled_features)
5
6 # Calcul de la pureté pour AgglomerativeClustering
7 agglo_purity = purity_score(true_labels, agglo_labels)
8 print(f"Purity avec AgglomerativeClustering: {agglo_purity:.2f}")
9
10 # Calcul de l'Indice de Rand Ajusté (Adjusted Rand Index) pour
11   AgglomerativeClustering
12 agglo_ari = adjusted_rand_score(true_labels, agglo_labels)

```

```

11 print(f"Indice de Rand Ajust avec AgglomerativeClustering: {
    agglo_ari:.2f}")
12
13 # Calcul du score de silhouette pour AgglomerativeClustering
14 agglo_silhouette = silhouette_score(scaled_features, agglo_labels)
15 print(f"Score de silhouette avec AgglomerativeClustering: {
    agglo_silhouette:.2f}")

```

Étape 12 : Comparaison des résultats entre K-Means et Agglomerative-Clustering

Le code compare les résultats des deux algorithmes.

```

1 # Comparaison des r sultats entre K-Means et AgglomerativeClustering
2 print(f"Purity de K-Means: {purity:.2f}, Purity d'
    AgglomerativeClustering: {agglo_purity:.2f}")
3 print(f"ARI de K-Means: {ari:.2f}, ARI d'AgglomerativeClustering: {
    agglo_ari:.2f}")
4 print(f"Score de silhouette de K-Means: {silhouette:.2f}, Score de
    silhouette d'AgglomerativeClustering: {agglo_silhouette:.2f}")

```

4.5.2 Méthodes d'évaluation utilisées

Dans le code présenté pour le clustering K-Means sur le jeu de données Iris, plusieurs méthodes d'évaluation sont utilisées pour évaluer la qualité des clusters obtenus. Ces méthodes permettent de mesurer à quel point les clusters correspondent aux étiquettes des classes dans le véritable jeu de données. Voici les principales méthodes d'évaluation utilisées :

1. **Accuracy (Exactitude)** : L'accuracy (exactitude) est une mesure de la qualité globale du clustering qui compare les étiquettes prédites des clusters avec les véritables étiquettes des classes. Elle est calculée à l'aide de `accuracy_score(true_labels, data['Cluster'])`, où `true_labels` sont les étiquettes réelles des classes dans le jeu de données et `data['Cluster']` sont les étiquettes prédites des clusters obtenues à partir du clustering K-Means. L'exactitude est comprise entre 0 (aucune correspondance) et 1 (correspondance parfaite).
2. **Purity (Pureté)** : La pureté mesure l'homogénéité des clusters en évaluant dans quelle mesure chaque cluster contient uniquement des points appartenant à

une seule classe. Elle est calculée avec la fonction `purity_score(true_labels, data['Cluster'])`. Une pureté de 1 indique une correspondance parfaite entre les clusters et les classes réelles.

3. **Adjusted Rand Index (Indice de Rand Ajusté) :** L'Indice de Rand Ajusté (ARI) est une mesure de similarité entre les clusters obtenus et les vraies étiquettes de classes, modifiées pour le hasard. Il corrige le fait que l'ARI peut donner des valeurs proches de zéro pour des regroupements aléatoires en utilisant une correction pour le hasard. Il est calculé à l'aide de `adjusted_rand_score(true_labels, data['Cluster'])`. Un ARI de 0 indique un regroupement aléatoire, tandis qu'un ARI proche de 1 indique un regroupement parfaitement corrélé avec les étiquettes réelles.
4. **Silhouette Score (Score de Silhouette) :** Le score de silhouette mesure la similarité des objets au sein d'un même cluster par rapport à ceux des autres clusters. Un score de silhouette proche de 1 indique des clusters bien séparés, tandis qu'un score proche de -1 indique des clusters qui se chevauchent. Il est calculé avec `silhouette_score(scaled_features, data['Cluster'])`.
5. **Cohésion Interne (Intra-cluster cohesion) :** La cohésion interne mesure à quel point les points au sein d'un cluster sont similaires entre eux. Il est généralement calculé comme la moyenne des distances entre tous les points d'un cluster par rapport à leur centre (par exemple, le centre de gravité ou la moyenne). Une cohésion interne élevée indique que les points au sein d'un cluster sont étroitement regroupés.
6. **Séparation Inter-clusters (Inter-cluster séparation) :** La séparation inter-clusters mesure à quel point les clusters sont distincts les uns des autres. Elle est généralement calculée comme la distance entre les centres de deux clusters différents. Une bonne séparation inter-clusters indique que les clusters sont bien séparés et distincts les uns des autres..

Ces métriques d'évaluation permettent d'analyser la qualité du clustering K-Means en comparant les clusters obtenus avec les véritables étiquettes des classes. Elles permettent de déterminer la robustesse du modèle de clustering et de choisir les paramètres optimaux pour obtenir les meilleurs résultats.

4.6 Résultats

Dans cette section, nous allons effectuer une analyse des résultats, des étapes d'exécution du code K-Means avec des captures d'écran.

1. Charger le jeu de données Iris :

La Figure 4.1 montre les cinq premières lignes du DataFrame pandas, qui contient les mesures de longueur et de largeur des sépales et des pétales de différentes espèces d'iris. Ce code charge le jeu de données Iris, avec les noms de colonnes des caractéristiques (longueur des sépales (sepal length), largeur des sépales (sepal width), longueur des pétales (petal length), largeur des pétales (petal width)). Les données sont organisées dans un DataFrame, affichant les mesures de chaque caractéristique pour chaque iris.

	sepal length (cm)	sepal width (cm)	petal length (cm)	petal width (cm)
0	5.1	3.5	1.4	0.2
1	4.9	3.0	1.4	0.2
2	4.7	3.2	1.3	0.2
3	4.6	3.1	1.5	0.2
4	5.0	3.6	1.4	0.2

FIGURE 4.1 – Mesures des sépales et des pétales de différentes espèces d'iris

2. Analyse exploratoire des données :

Dans cette section, nous allons voir les résultats du code pour la distribution des caractéristiques des données Iris ainsi que la matrice de corrélation des caractéristiques.

La Figure 4.2 montre des histogrammes illustrant la distribution des caractéristiques des fleurs. Chaque histogramme représente une caractéristique différente : la longueur (Sepal length) et la largeur (Sepal width) des sépales, ainsi que la longueur (petal length) et la largeur (petal width) des pétales. Chaque barre dans les histogrammes correspond à un intervalle de valeurs (bin) pour la caractéristique respective. La hauteur de chaque barre indique le nombre d'échantillons qui tombent dans cet intervalle. Ces visualisations permettent de visualiser rapidement la répartition des valeurs pour chaque caractéristique, ce qui peut être utile pour identifier des tendances ou des motifs dans les données .

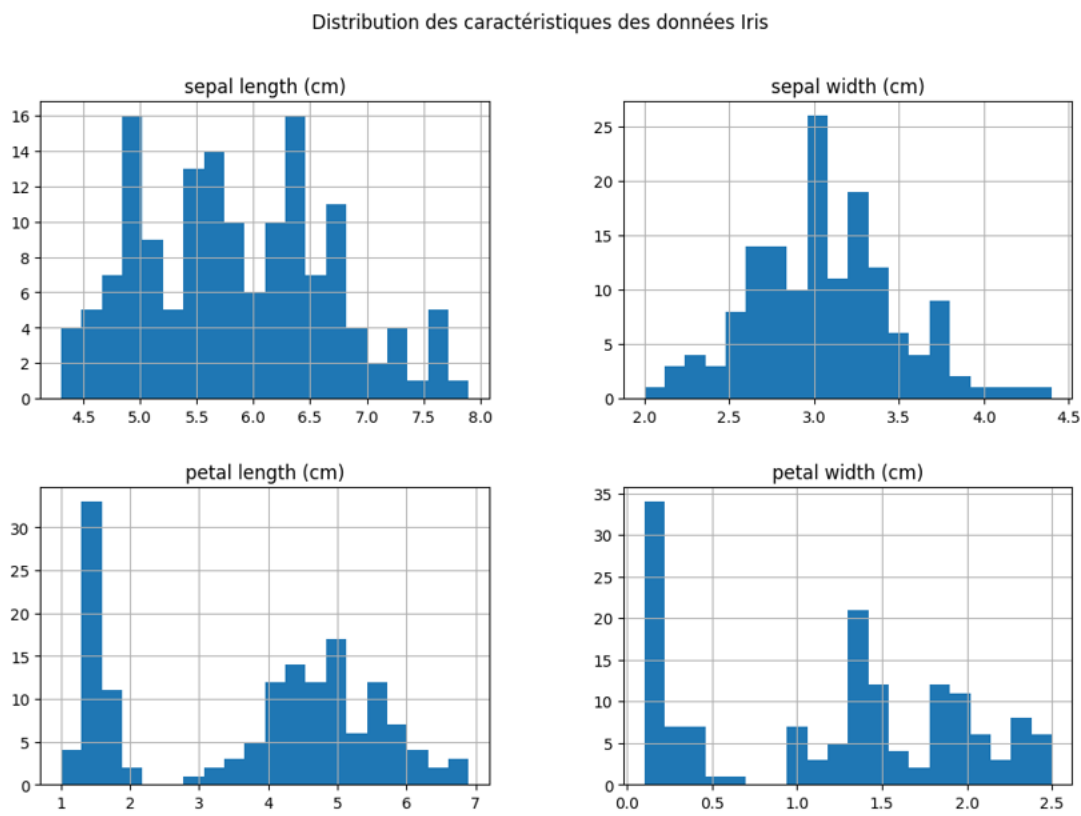


FIGURE 4.2 – Distribution des caractéristiques des données Iris

3. Relation entre les caractéristiques :

La Figure 4.3 représente les relations entre les différentes caractéristiques des fleurs Iris, telles que la longueur et la largeur des sépales et des pétales. Chaque point dans la matrice de dispersion représente la relation entre deux de ces caractéristiques. Par exemple, on peut observer comment la longueur des pétales est liée à leur largeur ou comment la longueur des sépales est liée à la longueur des pétales, et ainsi de suite. Cette visualisation nous permet d'identifier des schémas qui pourraient nous aider dans la classification des différentes espèces de fleurs Iris.

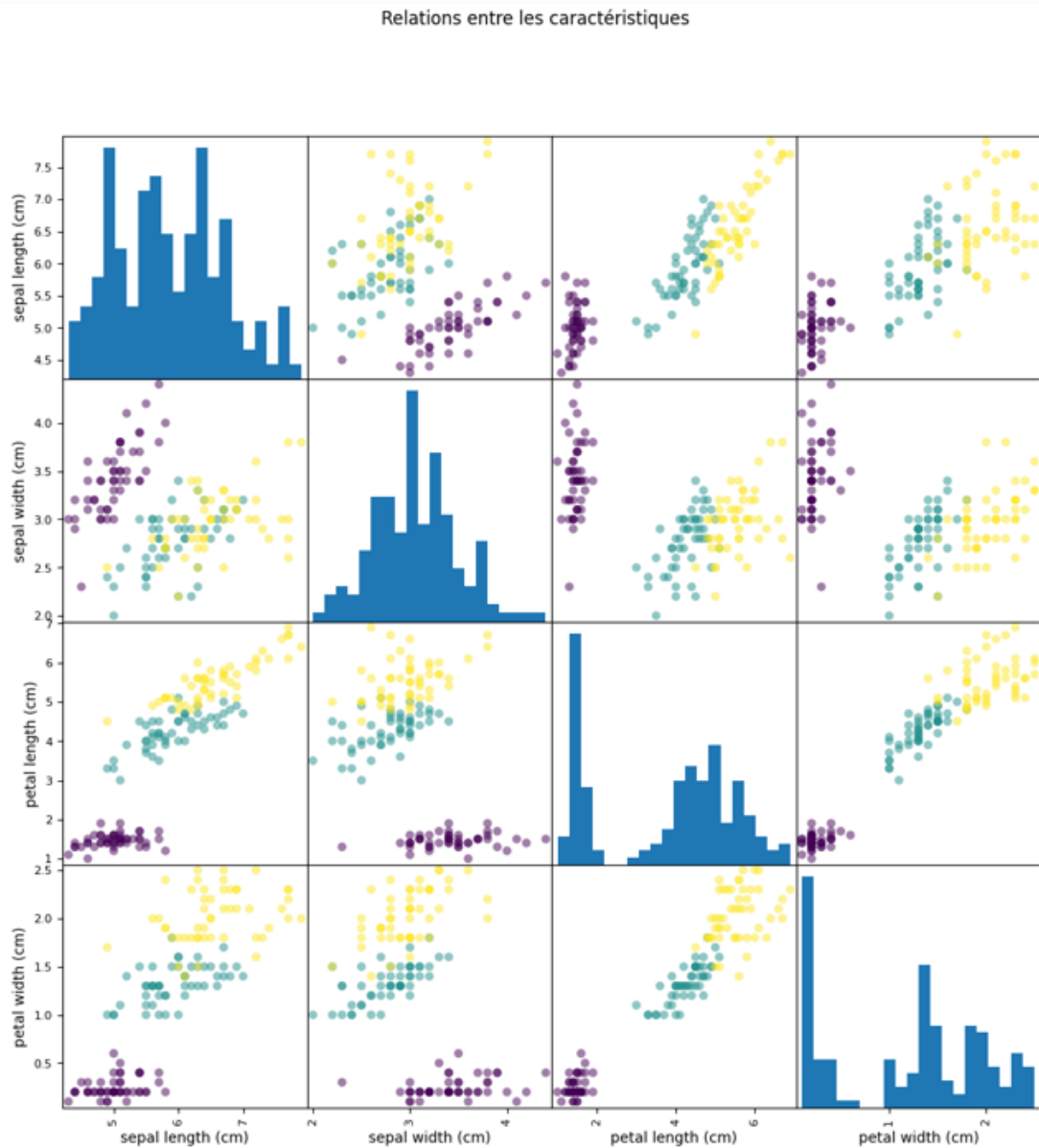


FIGURE 4.3 – Relation entre les caractéristiques

La Figure 4.4 montre la matrice de corrélation des caractéristiques du jeu de données Iris révèle des relations significatives entre les différentes mesures de l'Iris. Les résultats montrent une forte corrélation positive entre la longueur et la largeur des pétales, ainsi qu'entre la longueur des sépales et les dimensions des pétales. En revanche, la largeur des sépales présente une corrélation plus faible avec les autres caractéristiques. Ces observations permettent des schémas cohérents dans les mesures de l'Iris, ce qui peut être utile pour identifier les caractéristiques les plus informatives pour la classification des espèces d'Iris.

	sepal length (cm)	sepal width (cm)	petal length (cm)	\
sepal length (cm)	1.000000	-0.117570	0.871754	
sepal width (cm)	-0.117570	1.000000	-0.428440	
petal length (cm)	0.871754	-0.428440	1.000000	
petal width (cm)	0.817941	-0.366126	0.962865	
	petal width (cm)			
sepal length (cm)	0.817941			
sepal width (cm)	-0.366126			
petal length (cm)	0.962865			
petal width (cm)	1.000000			

FIGURE 4.4 – Matrice de corrélation des caractéristiques

4. Analyser les caractéristiques moyennes des clusters :

Dans cette partie, nous analyserons les résultats du code concernant les moyennes des caractéristiques des iris dans chaque cluster. Les résultats affichent les valeurs moyennes de chaque caractéristique pour chaque cluster, comme illustré dans la Figure4.5.

- **Cluster 0** : Les iris de ce cluster ont en moyenne une longueur de sépale de 6,78 cm, une largeur de sépale de 3,10 cm, une longueur de pétale de 5,51 cm, et une largeur de pétale de 1,97 cm.
- **Cluster 1** : Les iris de ce cluster ont en moyenne une longueur de sépale de 5,01 cm, une largeur de sépale de 3,43 cm, une longueur de pétale de 1.46 cm, et une largeur de pétale de 0.25 cm.
- **Cluster 2** : Les iris de ce cluster ont en moyenne une longueur de sépale de 5.80 cm, une largeur de sépale de 2.67 cm, une longueur de pétale de 4,37 cm, et une largeur du pétale de 1,41 cm.

D'après ces résultats, les moyennes des caractéristiques de l'iris dans chaque groupe peuvent être comparées. Ces informations sont utiles pour caractériser les différents groupes d'iris identifiés dans l'ensemble de données.

	sepal length (cm)	sepal width (cm)	petal length (cm)	\
Cluster				
0	6.780851	3.095745	5.510638	
1	5.006000	3.428000	1.462000	
2	5.801887	2.673585	4.369811	

	petal width (cm)
Cluster	
0	1.972340
1	0.246000
2	1.413208

FIGURE 4.5 – Caractéristiques moyennes des iris dans chaque cluster

5. Évaluation des résultats de K-Means :

La Figure 4.6 montre L'évaluation des performances du modèle K-Means sur l'ensemble de données Iris avec lumière des résultats contrastés. Alors que la pureté des clusters est élevée, atteignant 83%, ce qui suggère une bonne homogénéité des regroupements, les autres métriques révèlent des défis importants. La précision, mesurant la proportion d'observations correctement prédites, est extrêmement faible à 9%, indiquant une capacité très limitée du modèle à prédire correctement les classes de l'iris. De même, l'Indice de Rand Ajusté de 0,62 montre une similarité seulement raisonnable entre les regroupements trouvés par le modèle et les vrais regroupements des données. Le score de silhouette de 0.46 suggère une certaine structure dans les clusters mais laisse place à une significative.

- **Classification Report (Rapport de classification)** : Le rapport de classification fournit des mesures de précision, de rappel (rappel) et de f1-score pour chaque classe, ainsi que leur moyenne. Il montre également le support, c'est-à-dire le nombre d'occurrences de chaque classe dans l'ensemble de données. Dans ce cas, les scores de précision, de rappel, et de f1-score sont tous très faibles, ce qui suggère que le modèle ne performe pas bien dans la classification des classes d'iris.

```

Accuracy avec K-Means: 0.09
Purity avec K-Means: 0.83
Indice de Rand Ajusté avec K-Means: 0.62
Score de silhouette avec K-Means: 0.46
Rapport de classification avec K-Means:
      precision    recall  f1-score   support

     0         0.00     0.00     0.00        50
     1         0.00     0.00     0.00        50
     2         0.26     0.28     0.27        50

 accuracy                   0.09        150
 macro avg         0.09     0.09     0.09        150
 weighted avg     0.09     0.09     0.09        150

```

FIGURE 4.6 – Évaluation des performances du clustering K-Means

6. Évaluation par courbe d'inertie avec K-Means

La Figure 4.7 représente les résultats du code pour l'analyse de la courbe d'inertie pour le clustering K-Means. Le résultat affiche la courbe d'inertie en traçant le nombre de clusters sur l'axe des x et l'inertie correspondante sur l'axe des y. Cette courbe permet de visualiser la décroissance de l'inertie en fonction du nombre de clusters et peut aider à déterminer le nombre optimal de clusters à utiliser pour le clustering K-Means.

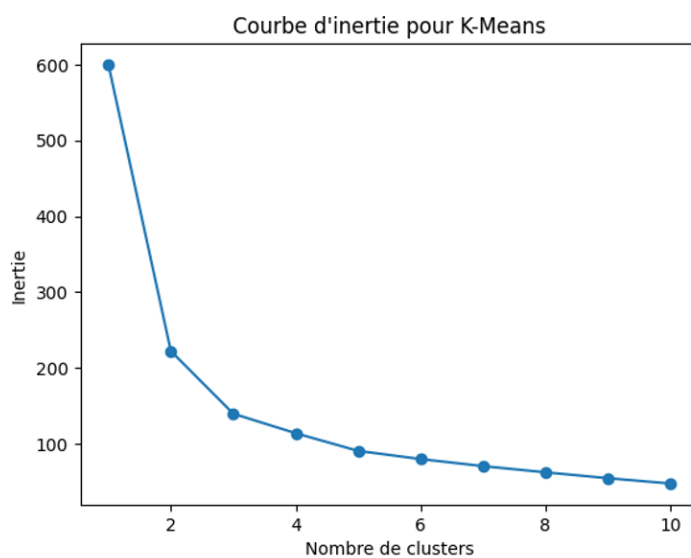


FIGURE 4.7 – Analyse de la courbe d'inertie pour le clustering K-Means.

7. Visualisation de la distribution des clusters dans chaque classe

La Figure 4.8 montre Les histogrammes des résultats,illustrant la répartition des classes réelles des iris pour chaque cluster identifié par l’algorithme de clustering. Ces histogrammes permettent de visualiser comment les différentes classes d’iris sont réparties au sein de chaque cluster identifié. :

- **Cluster 0** : Cet histogramme montre la distribution des classes réelles des iris dans le cluster 0. La hauteur de chaque barre indique le nombre d’iris appartenant à chaque classe réelle (0, 1 ou 2) dans ce cluster.
- **Cluster 1** : Cet histogramme représente la répartition des classes réelles des iris dans le cluster 1.
- **Cluster 2** : cet histogramme montre la distribution des classes réelles des iris dans le cluster 2.

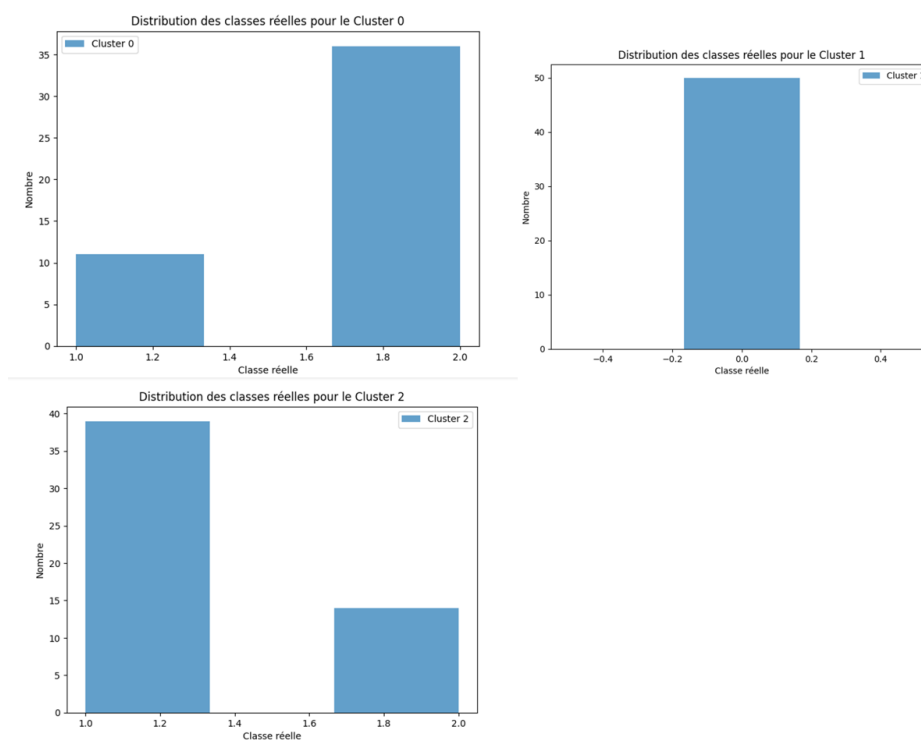


FIGURE 4.8 – la distribution des clusters dans chaque classe

4.6.1 Visualisation des Clusters

La visualisation des clusters obtenus à partir de l’algorithme K-Means sur le jeu de données Iris peut être réalisée à l’aide de graphiques. Voici les différentes visualisations réalisées dans le code :

a) Visualisation des clusters dans l'espace 2D

```

1 # Visualiser les clusters avec K-Means (2D)
2 plt.figure(figsize=(10, 8))
3 plt.scatter(data['sepal length (cm)'], data['petal length (cm)'], c=
    data['Cluster'], cmap='viridis')
4 plt.xlabel('Longueur du S pale (cm)')
5 plt.ylabel('Longueur du P tale (cm)')
6 plt.title('Clusters avec K-Means dans le jeu de donn es Iris')
7 plt.show()

```

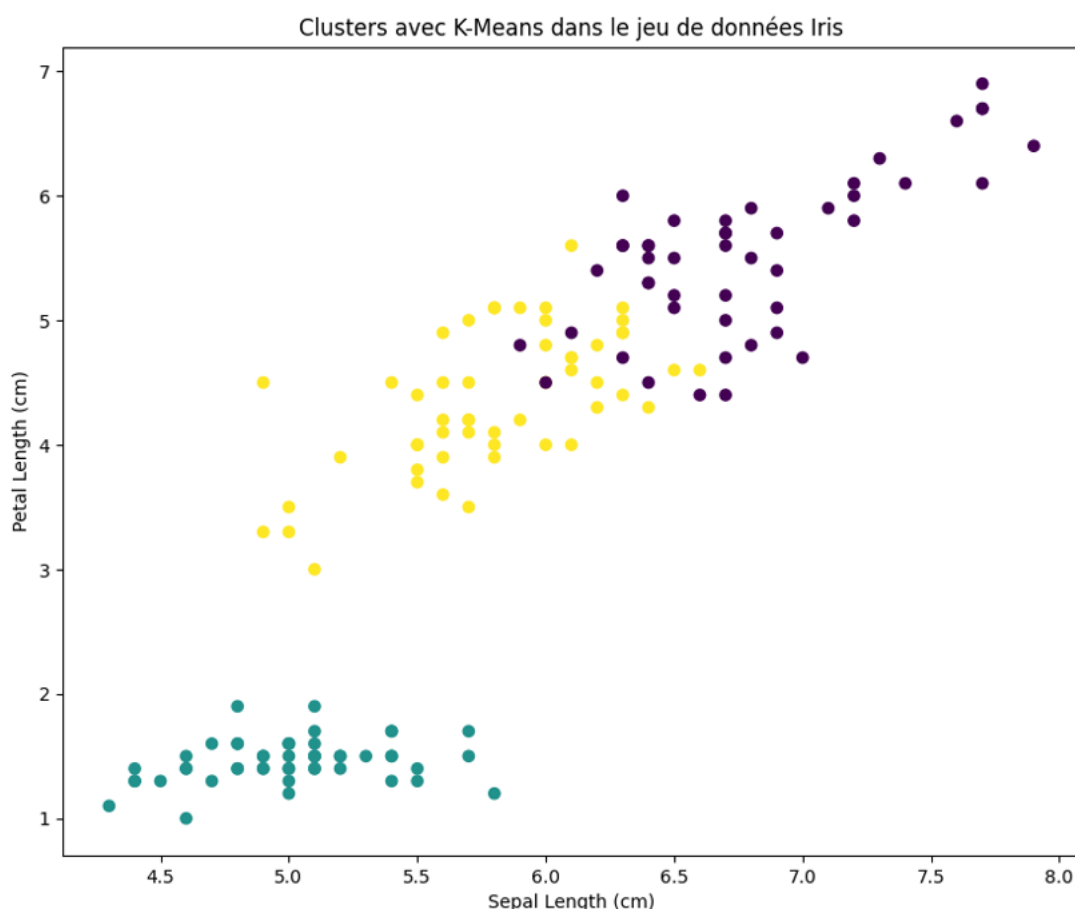


FIGURE 4.9 – Visualiser les clusters avec K-Means (2D)

La Figure 4.9 illustre cette visualisation en utilisant les colonnes 'sepal length (cm)' (longueur du sépale) et 'petal length (cm)' (longueur du pétale) pour représenter les clusters colorés par la colonne 'Cluster' attribuée par K-Means. Chaque point dans le graphique représente une observation, et sa couleur indique à quel cluster elle appartient.

b) Visualisation des clusters dans l'espace 3D

```

1 # Visualisation des clusters en 3D avec K-Means
2 fig = plt.figure()
3 ax = fig.add_subplot(111, projection='3d')
4 ax.scatter(data['sepal length (cm)'], data['sepal width (cm)'],
5 data['petal length (cm)'], c=data['Cluster'], cmap='viridis')
6 ax.set_xlabel('Sepal Length (cm)')
7 ax.set_ylabel('Sepal Width (cm)')
8 ax.set_zlabel('Petal Length (cm)')
9 plt.title('Visualisation 3D des clusters avec K-Means dans le jeu
10 de donn es Iris')
11 plt.show()

```

Visualisation 3D des clusters avec K-Means dans le jeu de données Iris

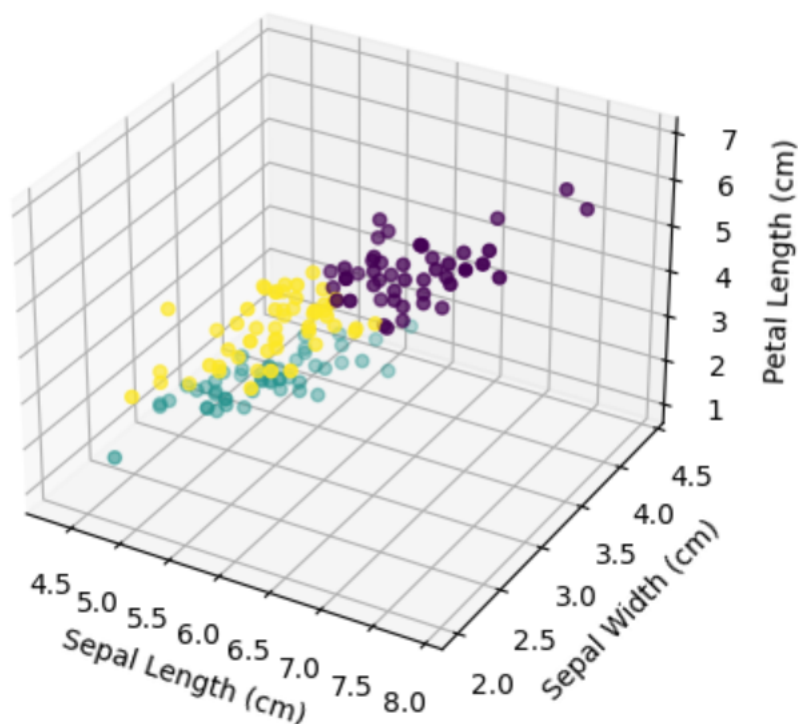


FIGURE 4.10 – Visualisation 3D des clusters avec K-Means dans le jeu de données Iris

La Figure 4.10 illustre cette visualisation en utilisant les colonnes 'sepal length (cm)' (longueur du sépale), 'sepal width (cm)' (largeur du sépale), et 'petal length (cm)' (longueur du pétale) pour représenter les grappes en 3D, avec chaque axe correspondant à une caractéristique différente. Les points sont colorés en fonction du cluster auquel ils appartiennent.

Ces graphiques permettent de visualiser la répartition des clusters dans l'espace en

fonction des caractéristiques sélectionnées. Ils offrent un aperçu visuel des performances de l’algorithme K-Means dans la segmentation des données en groupes distincts.

4.6.2 Évaluation des Clusters

Le Tableau 4.1 présente les performances du code implémentant l’algorithme K-Means sur le jeu de données Iris, évaluant les clusters obtenus. :

Métrique	Valeur
Exactitude	0.09
Pureté	0.83
Indice de Rand Ajusté	0.62
Score de Silhouette	0.46

TABLE 4.1 – Évaluation des performances de l’algorithme K-Means.

4.7 Comparaison avec ”Agglomerative” Clustering

Pour élargir notre compréhension des méthodes de clustering, nous comparons les résultats de l’algorithme K-Means avec ceux de l’algorithme ”Agglomerative Clustering” sur le jeu de données Iris. L’algorithme ”Agglomerative Clustering” est une méthode hiérarchique qui fusionne progressivement les points de données en clusters. Dans les étapes suivantes, nous importons l’algorithme ”Agglomerative Clustering” à l’aide de la bibliothèque Scikit-learn, nous effectuons le clustering sur les mêmes caractéristiques que celles utilisées pour K-Means, puis nous visualisons les clusters obtenus pour une comparaison directe avec les résultats de K-Means.

La comparaison avec l’algorithme ”Agglomerative Clustering” est réalisée dans les étapes 11 et 12 du code.

1. **Comparaison avec ”Agglomerative” Clustering :** Dans cette étape, l’algorithme `AgglomerativeClustering` est utilisé pour effectuer le clustering, et les performances sont évaluées de la même manière que pour K-Means.

```

1 # Comparaison avec un autre algorithme de clustering
  AgglomerativeClustering
2 aggro = AgglomerativeClustering(n_clusters=num_clusters)
3 aggro_labels = aggro.fit_predict(scaled_features)
4

```



```

5 # Calcul de la pureté pour AgglomerativeClustering
6 aggro_purity = purity_score(true_labels, aggro_labels)
7 print(f"Purity avec AgglomerativeClustering: {aggro_purity:.2f}")
8
9 # Calcul de l'Indice de Rand Ajusté (Adjusted Rand Index)
10 #pour AgglomerativeClustering
11 aggro_ari = adjusted_rand_score(true_labels, aggro_labels)
12 print(f"Indice de Rand Ajusté avec AgglomerativeClustering: {
    aggro_ari:.2f}")
13
14 # Calcul du score de silhouette pour AgglomerativeClustering
15 aggro_silhouette = silhouette_score(scaled_features, aggro_labels
    )
16 print(f"Score de silhouette avec AgglomerativeClustering:
17 {aggro_silhouette:.2f}")

```

Nous utilisons l'algorithme `AgglomerativeClustering` pour effectuer le clustering sur les données normalisées `scaled_features`. Nous évaluons ensuite les performances de cet algorithme en calculant plusieurs métriques telles que la pureté, l'indice de Rand ajusté, et le score de silhouette. Enfin, nous comparons ces performances avec celles de l'algorithme *K*-Means pour déterminer lequel fonctionne le mieux pour le jeu de données Iris en termes de qualité des clusters et de similarité avec les regroupements réels.

La Figure 4.11 compare les résultats de *K*-Means à ceux d'`AgglomerativeClustering` en termes de pureté, indice de Rand ajusté, et score de silhouette.

```

Purity avec AgglomerativeClustering: 0.83
Indice de Rand Ajusté avec AgglomerativeClustering: 0.62
Score de silhouette avec AgglomerativeClustering: 0.45

```

FIGURE 4.11 – Comparaison entre *K*-Means et `AgglomerativeClustering`

2. **Comparaison des performances entre *K*-Means et "Agglomerative" Clustering** : Dans cette étape, les performances des deux algorithmes sont comparées en termes de pureté, indice de Rand amélioré, et score de silhouette.

```

1 # Comparaison des résultats entre K-Means et
    AgglomerativeClustering
2 print(f"Purity de K-Means: {purity:.2f},

```

```

3 Purity d'AgglomerativeClustering: {agglo_purity:.2f}")
4 print(f"ARI de K-Means: {ari:.2f},
5 ARI d'AgglomerativeClustering: {agglo_ari:.2f}")
6 print(f"Score de silhouette de K-Means: {silhouette:.2f},
7 Score de silhouette d'AgglomerativeClustering: {agglo_silhouette
  :.2f}")

```

Nous comparons les performances de l'algorithme K-Means avec celles de l'algorithme AgglomerativeClustering. Nous examinerons les métriques de performance telles que la pureté, l'indice de Rand ajusté, et le score de silhouette pour les deux algorithmes. Cette comparaison nous permet de déterminer lequel des deux algorithmes fournit les clusters les plus cohérents et les plus similaires aux étiquettes réelles du jeu de données Iris.

La Figure 4.12 montre une comparaison finale des performances des deux algorithmes.

```

Purity de K-Means: 0.83, Purity d'AgglomerativeClustering: 0.83
ARI de K-Means: 0.62, ARI d'AgglomerativeClustering: 0.62
Score de silhouette de K-Means: 0.46, Score de silhouette d'AgglomerativeClustering: 0.45

```

FIGURE 4.12 – Comparaison des résultats entre K-Means et AgglomerativeClustering

Ces étapes permettent de comparer les performances de K-Means et d'Agglomerative Clustering sur le jeu de données Iris en termes de pureté, d'indice de Rand ajusté, et de score de silhouette.

4.8 Discussion

La mise en œuvre de l'algorithme de clustering K-Means sur le jeu de données Iris a produit des résultats significatifs et prometteurs. En examinant les caractéristiques moyennes de chaque cluster, nous pouvons interpréter les regroupements identifiés et comprendre les similarités et les différences entre eux. Les métriques de performance telles que l'exactitude, la pureté, l'Indice de Rand Ajusté et le score de silhouette fournissent une évaluation quantitative de la qualité des clusters, et la courbe d'inertie nous aide à choisir judicieusement le nombre de clusters. En comparant les résultats avec l'algorithme Agglomerative Clustering, nous validons la robustesse de notre approche. Enfin, cette méthode offre des perspectives d'application dans divers domaines, soulignant

son utilité pour la segmentation de la clientèle, la détection d'anomalies et d'autres analyses exploratoires de données, tout en permettant une meilleure compréhension des structures sous-jacentes des ensembles de données non étiquetées.

4.8.1 Analyse critique des performances du modèle

Bien que le modèle réussisse à segmenter efficacement les données du jeu de données Iris en clusters significatifs, certaines limitations et défis subsistent. Parmi ceux-ci, on peut citer la sensibilité de l'algorithme K-Means aux valeurs aberrantes et la nécessité de déterminer le nombre optimal de clusters à l'avance.

La sensibilité des K-Means aux valeurs aberrantes peut entraîner une affectation incorrecte des observations aux clusters, ce qui peut nuire à la qualité globale de la segmentation. Pour atténuer ce problème, il est souvent nécessaire de prétraiter les données en éliminant ou en traitant les valeurs aberrantes de manière appropriée.

De plus, le choix du nombre de clusters est un aspect crucial de l'utilisation de l'algorithme K-Means. Bien que dans ce cas, le nombre de clusters corresponde aux trois espèces d'iris présentes dans le jeu de données, dans d'autres scénarios, déterminer le nombre optimal de clusters peut être plus complexe et nécessiter des techniques telles que la méthode du coude ou la méthode de silhouette.

Ces défis soulignent l'importance de la préparation des données et du réglage des paramètres lors de l'utilisation de l'algorithme K-Means. Il est essentiel de comprendre les caractéristiques des données et de choisir judicieusement les paramètres de l'algorithme pour obtenir des résultats fiables et significatifs. De plus, il est recommandé d'explorer d'autres algorithmes de clustering et de comparer leurs performances pour s'assurer que le modèle choisi répond adéquatement aux besoins de l'analyse.

4.8.2 Implications des résultats pour l'étude de cas

Les résultats de l'étude de cas sur le jeu de données Iris ont plusieurs implications importantes pour l'analyse et la compréhension des données botaniques. Voici quelques-unes des implications clés :

1. **Identification des espèces d'iris :** Les clusters identifiés par l'algorithme de clustering peuvent correspondre aux différentes espèces d'iris présentes dans le

jeu de données. Cela peut aider les botanistes et les chercheurs à mieux comprendre les caractéristiques distinctives de chaque espèce et à les utiliser pour l'identification et la classification des plantes dans la nature.

2. **Détection de similitudes et de différences** : L'analyse de clustering peut révéler des similitudes et des différences entre les différentes espèces d'iris en termes de caractéristiques morphologiques, telles que la longueur et la largeur des pétales et des sépales. Cela peut être utile pour étudier l'évolution des plantes et comprendre les relations phylogénétiques entre les espèces.
3. **Validation des données existantes** : En comparant les clusters identifiés par l'algorithme avec les étiquettes réelles des espèces d'iris, on peut valider la qualité et la précision des données existantes. Cela peut aider à identifier d'éventuelles erreurs ou incohérences dans les données et à améliorer la fiabilité des ensembles de données botaniques.
4. **Prédiction et classification** : Une fois que le modèle de clustering est entraîné sur un ensemble de données étiquetées, il peut être utilisé pour prédire la classe ou l'espèce d'iris à laquelle appartient une nouvelle observation en fonction de ses caractéristiques morphologiques. Cela peut être utile pour automatiser le processus d'identification des plantes dans les applications de classification botanique.
5. **Guidance pour la recherche future** : Les insights tirés de l'analyse de clustering peuvent orienter la recherche future sur les iris et d'autres plantes en mettant en évidence les caractéristiques les plus importantes pour la classification et en identifiant les zones de recherche potentielles pour approfondir la compréhension de la diversité botanique.

En résumé, les résultats de l'étude de cas sur le jeu de données Iris ont des implications importantes pour la botanique en fournissant des informations précieuses sur la classification et la compréhension des espèces d'iris, ainsi que des orientations pour la recherche future dans ce domaine.

4.9 Conclusion

En conclusion, ce chapitre a fourni une vue d'ensemble détaillée des différentes étapes impliquées dans l'application et l'évaluation d'une méthode d'apprentissage au-

tomatique pour l'analyse de Big Data, à travers une étude de cas spécifique. Il a débuté par l'introduction du contexte et des objectifs de l'implémentation, mettant en avant l'importance cruciale de Colab et Python dans ce domaine. Par la suite, la configuration de l'environnement de développement et la préparation des données ont été abordées, notamment le nettoyage et l'analyse exploratoire. Le chapitre a ensuite justifié la sélection et l'adaptation de l'algorithme, suivi de la mise en œuvre détaillée et des résultats obtenus, incluant la visualisation et l'évaluation des clusters. Enfin, une discussion critique des performances du modèle a été menée, soulignant les implications pour l'étude de cas et identifiant des pistes d'amélioration pour de futures recherches.

4.9.1 Résumé des contributions principales et de l'impact sur les objectifs de l'étude de cas

L'étude de cas sur l'algorithme de clustering K-Means appliqué au jeu de données Iris offre une contribution significative en introduisant les principes fondamentaux de l'algorithme et en démontrant son application pratique. En utilisant le jeu de données Iris, l'étude détaille le processus de clustering, interprète les résultats obtenus en termes de séparation des espèces d'iris en clusters distincts, et évalue la performance de l'algorithme. Cette analyse contribue à approfondir la compréhension de K-Means et de son efficacité dans le regroupement de données, tout en soulignant son impact sur la segmentation efficace des données Iris.

4.9.2 Réflexion sur l'importance de l'approche choisie pour le traitement et l'analyse de Big Data

L'algorithme de clustering K-Means, choisi pour le traitement et l'analyse de Big data, revêt une importance capitale, se révélant crucial à plusieurs niveaux. Tout d'abord, le choix de cet algorithme impacte directement la capacité à découvrir des structures et des patterns au sein des données. Une sélection judicieuse peut conduire à des regroupements pertinents, fournissant ainsi des informations précieuses sur les relations entre les données. De plus, le prétraitement des données, tel que la normalisation, garantit que les caractéristiques sont sur la même échelle, ce qui est essentiel pour des résultats de clustering précis. L'analyse exploratoire des données (EDA) permet de

comprendre la nature des données, d'identifier des tendances et des corrélations, et de guider le choix des algorithmes et des paramètres appropriés. Enfin, l'évaluation des performances de l'algorithme, avec des métriques telles que l'exactitude, la pureté, et le score de silhouette, est essentielle pour mesurer la qualité des clusters et la pertinence des résultats obtenus. En somme, l'approche choisie pour le traitement et l'analyse de Big Data joue un rôle déterminant dans la génération d'informations exploitables à partir de grandes quantités de données, permettant ainsi de prendre des décisions éclairées et de découvrir des connaissances cachées.

Conclusion Générale

Après avoir mené une revue systématique de la littérature sur l'application des méthodes d'apprentissage automatique au traitement et à l'analyse des Big Data, ainsi qu'une étude de cas spécifique utilisant l'algorithme K-means sur le jeu de données Iris, plusieurs conclusions importantes peuvent être tirées.

Tout d'abord, cette revue systématique a permis de mettre en lumière la diversité des techniques et des algorithmes disponibles pour aborder les défis posés par les Big Data. Chaque méthode d'apprentissage automatique présente ses propres avantages et défis, et il est crucial de sélectionner la technique la plus appropriée en fonction des caractéristiques spécifiques des données et des objectifs de l'analyse.

L'étude de cas sur l'algorithme K-means a illustré à la fois l'efficacité et les limites des méthodes d'apprentissage non supervisé pour la classification des données. Bien que K-means soit largement utilisé pour sa simplicité et sa rapidité de convergence, il présente des défis tels que la sensibilité aux valeurs aberrantes et la dépendance à l'initialisation des centroids. Cette étude a souligné l'importance de la normalisation des données et de la visualisation des clusters pour une interprétation efficace des résultats.

Enfin, cette revue et cette étude de cas soulignent l'importance continue de la recherche dans le domaine de l'apprentissage automatique appliqué aux Big Data. Les avancées technologiques rapides et l'explosion des données exigent des méthodes d'analyse sophistiquées et robustes pour extraire des informations significatives à partir de vastes ensembles de données. L'exploration de nouvelles approches, telles que les réseaux de neurones convolutifs, et l'accent mis sur l'optimisation des techniques existantes sont essentielles pour relever ces défis et exploiter pleinement le potentiel des Big Data.

En conclusion, les méthodes d'apprentissage automatique offrent des outils puissants pour le traitement et l'analyse des Big Data. Cette étude souligne l'importance

d'une approche systématique et rigoureuse dans la sélection et l'application de ces techniques, ainsi que la nécessité d'une recherche continue pour répondre aux défis en constante évolution posés par les Big Data. L'avenir de ce domaine promet des avancées significatives qui pourraient transformer notre capacité à exploiter les vastes quantités de données disponibles et à générer des insights précieux pour la science, l'industrie et la société dans son ensemble.

Bibliographie

- [1] A. K. N. Derouaz, “Proposition d’une technique efficace pour l’analyse du big data basée sur l’intelligence artificielle,” 2023. Soutenu publiquement le 22/06/2023.
- [2] O. Polo, “Big data : Une révision,” *École National d’Ingenieurs de Metz, France*, October 2015. Master’s thesis.
- [3] Z. Sayah, *Intégration de la Sémantique dans le Big Data*. Doctorat lmd en informatique, Université de Ouargla, 2021.
- [4] H. Medfouni, “Validation de clustering des données dans un contexte big data,” 2018.
- [5] A. Anis, “Les enjeux de l’usage de la big data analytics sur l’efficacité et l’efficience de l’audit interne,” mémoire de fin d’études, Université de Mouloud Mammeri de Tizi Ouzou, 2022.
- [6] B. Boudjehem, “Le traitement des données manquantes dans le « big data » médical,” Master’s thesis, Université de 8 Mai 1945 – Guelma, Juin 2022. Mémoire de Master, Faculté des Mathématiques, d’Informatique et des Sciences de la matière.
- [7] M. I. A. Bourahdoun, “Impact des méthodes analytiques dans le contexte des données massives,” Master’s thesis, Université 8 Mai 1945 – Guelma, 2020.
- [8] R. Feldman and J. Sanger, *The Text Mining Handbook : Advanced Approaches in Analyzing Unstructured Data*. New York, NY : Cambridge University Press, 2007.
- [9] A. C. Eberendu, “Unstructured data : an overview of the data of big data,” *International Journal of Emerging Trends & Technology in Computer Science*, vol. 38, no. 1, pp. 46–49, 2016.
- [10] J. I. Maletic and A. Marcus, “Data cleansing : Beyond integrity analysis,” in *Proceedings of the 2000 Conference on Information Quality*, (Campus Box 526429

- Memphis, TN 38152), Division of Computer Science, The Department of Mathematical Sciences, The University of Memphis, 2000.
- [11] K. Nagorny, P. Lima-Monteiro, J. Barata, and A. W. Colombo, “Big data analysis in smart manufacturing : A review,” *International Journal of Communications, Network and System Sciences*, vol. 10, pp. 31–58, 2017.
- [12] A. Menasria, “Vers une nouvelle méthode de stockage de données « big data » dans un environnement smart city,” Master’s thesis, Université [Nom de votre université], Juin 2022. Mémoire de Fin d’études - Master, Filière : Informatique, Option : Science et technologie de l’information et de la communication.
- [13] G. Saint-Cirgue, *Apprendre le Machine Learning en une semaine*. machinelearning.com, 2019.
- [14] M. Prakash, G. Padmapriy, and M. V. Kumar, “A review on machine learning big data using r,” in *2018 Second International Conference on Inventive Communication and Computational Technologies (ICICCT)*, pp. 1873–1877, 2018.
- [15] Y. Zhao, Y. Li, X. Zhang, G. Geng, W. Zhang, and Y. Sun, “A survey of networking applications applying the software defined networking concept based on machine learning,” *IEEE Access*, vol. 7, pp. 95397–95417, 2019.
- [16] A. Allmang, “Qu’est-ce que l’apprentissage non supervisé?,” 2023.
- [17] E. Delsol, “Avec le deep learning, jcdcaux modélise l’eye tracking publicitaire,” *Le Monde Informatique*, Octobre 2023.
- [18] A. Krajnc, “Tout savoir sur le renforcement learning,” 2023. CEO & Fondateur. Dernière mise à jour le 9 novembre 2023.
- [19] L. Schapira, “Deep learning ou apprentissage profond : qu’est-ce que c’est?,” 2022. Data Science.
- [20] A. non spécifié, “Big data transformer : l’immensité de l’apprentissage automatique et l’analyse du big data,” 2023. Mis à jour le 17 Nov 2023, réductions de FasterCapital.
- [21] P. Ongsulee, V. Chotchaung, E. Bamrungsi, and T. Rodcheewit, “Big data, predictive analytics and machine learning,” in *2018 16th International Conference on ICT and Knowledge Engineering (ICTKE)*, pp. 1–6, 2018.

- [22] G. E. Melo-Acosta, F. Duitama-Muñoz, and J. D. Arias-Londoño, “Fraud detection in big data using supervised and semi-supervised learning techniques,” in *2017 IEEE Colombian Conference on Communications and Computing (COLCOM)*, pp. 1–6, 2017.
- [23] S. Dwivedi and V. S. K. Roshni, “Recommender system for big data in education,” in *2017 5th National Conference on E-Learning E-Learning Technologies (ELELTECH)*, pp. 1–4, 2017.
- [24] A. P. Jain and P. Dandannavar, “Application of machine learning techniques to sentiment analysis,” in *2016 2nd International Conference on Applied and Theoretical Computing and Communication Technology (iCATccT)*, pp. 628–632, 2016.
- [25] K. Barznji, *BIG DATA SENTIMENT ANALYSIS USING MACHINE LEARNING ALGORITHMS*. PhD thesis, University of Chemical Technology and Metallurgy, Sofia, Bulgaria, May 2018. PhD Research Scholar at UCTM.
- [26] Z. Zong and C. Hong, “On application of natural language processing in machine translation,” in *2018 3rd International Conference on Mechanical, Control and Computer Engineering (ICMCCE)*, pp. 506–510, 2018.
- [27] A. Ali, J. Qadir, R. Rasool, *et al.*, “Big data for development : applications and techniques,” *Big Data Analytics*, vol. 1, no. 2, 2016.
- [28] I. K. Nti, J. A. Quarcoo, J. Aning, and G. K. Fosu, “A mini-review of machine learning in big data analytics : Applications, challenges, and prospects,” *Big Data Mining and Analytics*, vol. 5, no. 2, pp. 81–97, 2022.
- [29] S. Shetty, S. Shetty, C. Singh, and A. Rao, *Supervised Machine Learning : Algorithms and Applications*, pp. 1–16. 02 2022.
- [30] C. Zhuang, S. Yan, A. Nayebi, M. Schrimpf, M. C. Frank, J. J. DiCarlo, and D. L. K. Yamins, “Unsupervised neural network models of the ventral visual stream,” *Edited by Marlene Behrmann, Carnegie Mellon University, Pittsburgh, PA*, 2020. Approved on December 9, 2020 (Received for review on July 7, 2020).
- [31] K. Nguyen, C. Fookes, and S. Sridharan, “Improving deep convolutional neural networks with unsupervised feature learning,” in *2015 IEEE International Conference on Image Processing (ICIP)*, pp. 2270–2274, 2015.

- [32] S. Sakib, N. Ahmed, A. J. Kabir, and H. Ahmed, “An overview of convolutional neural network : Its architecture and applications,” *To be filled*, 2024. Corresponding author : 15305026@iubat.edu.
- [33] S. Aladdin, S. El-Tantawy, M. M. Fouda, and A. S. Tag Eldien, “Marla-sg : Multi-agent reinforcement learning algorithm for efficient demand response in smart grid,” *IEEE Access*, vol. 8, pp. 210626–210639, 2020.
- [34] M. Sughasiny and J. Rajeshwari, “Application of machine learning techniques, big data analytics in health care sector – a literature survey,” in *2018 2nd International Conference on I-SMAC (IoT in Social, Mobile, Analytics and Cloud) (I-SMAC)I-SMAC (IoT in Social, Mobile, Analytics and Cloud) (I-SMAC)*, 2018 2nd International Conference on, pp. 741–749, 2018.
- [35] A. A. Mohamed, “La détection des webshells pour assurer la sécurité des serveurs web,” Master’s thesis, Faculté des Mathématiques, d’Informatique et des Sciences de la matière, Département d’Informatique, Septembre 2021. Mémoire de Fin d’études Master en Informatique, Filière : Informatique, Option : Système Informatique.
- [36] L. DE MATTEIS, S. JANNY, S. NATHAN, and W. SHU-QUARTIER, “Introduction à l’apprentissage automatique,” *Culture Sciences de l’Ingénieur*, mai 2022.
- [37] H. Mezili, “Vers une amélioration de la détection d’intrusion par les méthodes de sélection des fonctionnalités à l’aide des arbres de décision,” Master’s thesis, Université Ibn Khaldoun - Tiaret, Octobre 2021. Mémoire de Master en Réseaux et Télécommunications.
- [38] B.-A. Hammou, A. Al Lahcen, and S. Mouline, “Towards a realtime processing framework based on improved distributed recurrent neural network variants with fasttext for social big data analytics,” *Information Processing Management*, 2020.
- [39] C. Talon, E. Dautrême, E. Remy, Y. Dirat, and C. D. L. Strat, “Analyse de différents algorithmes de classification par apprentissage automatique sur un cas d’usage du domaine nucléaire,” in *Congrès LambdaMu21 ”Maîtrise des risques et transformation numérique : opportunités et menaces”*, (Reims, France), October 2018.

- [40] R. Mazouzi, C. D. Runz, and H. Akdag, “Un système collectif d’utilisation d’un grand ensemble de classifieurs sur le cloud pour la classification de big data,” January 2016. Published on ResearchGate.
- [41] M. Severo and R. Lamarche-Perrin, “L’analyse des opinions politiques sur twitter : Défis et opportunités d’une approche multi-échelle,” *Revue française de sociologie*, vol. 59, no. 3, pp. 507–532, 2018.
- [42] A. Bouzouane, “Séminaire – apprentissage automatique pour le bigdata (8inf912).” Plan de cours, 2020.
- [43] M. de l’Enseignement Supérieur et de la Recherche Scientifique, “Optimisation automatique dans le compilateur tiramisu : Amélioration de la scalabilité par la sélection des paramètres,” Juin 2022. Soutenue le 30 juin 2022.
- [44] J. Xie, Z. Song, Y. Li, Y. Zhang, H. Yu, J. Zhan, Z. Ma, Y. Qiao, J. Zhang, and J. Guo, “A survey on machine learning-based mobile big data analysis : Challenges and applications,” *Wireless Communications and Mobile Computing*, vol. 2018, p. 8738613, 2018.
- [45] H. A. Selmy, H. K. Mohamed, and W. Medhat, “Big data analytics deep learning techniques and applications : A survey,” *Information Systems*, vol. 120, p. 102318, 2024. Contents lists available at ScienceDirect.
- [46] A. H. Gandomi, F. Chen, and L. Abualigah, “Machine learning technologies for big data analytics,” *Electronics*, vol. 11, no. 3, p. 421, 2022. Received : 24 January 2022, Accepted : 28 January 2022, Published : 30 January 2022.
- [47] W. Li, Y. Chai, F. Khan, S. R. U. Jan, S. Verma, V. G. Menon, Kavita, and X. Li, “A comprehensive survey on machine learning-based big data analytics for iot-enabled smart healthcare system,” *Journal of Medical Systems*, January 6 2021. Published online : 6 January 2021, Accepted : 22 November 2020.
- [48] A. M. Rahmani, E. Azhir, S. Ali, M. Mohammadi, O. H. Ahmed, M. Y. Ghafour, S. H. Ahmed, and M. Hosseinzadeh, “Artificial intelligence approaches and mechanisms for big data analytics : a systematic study,” *PeerJ Computer Science*, vol. 7, p. e488, 2021. Submitted 17 November 2020, Accepted 20 March 2021, Published 14 April 2021.

- [49] K. W. AMARA and A. OUZANDJA, “Étude systématique de la littérature pour les algorithmes de clusterisation pour les réseaux véhiculaires,” Master’s thesis, Université de Mohamed El Bachir El Ibrahimi de Bordj Bou Arréridj, Faculté des Mathématiques et d’Informatique, Département d’informatique, République Algérienne Démocratique et Populaire, Juin 2023. Mémoire de Master en informatique, spécialité Réseaux & Multimédias.
- [50] I. E. Khammal, A. Maizate, and E. Ziyati, “Projets informatiques d’intégration dans le secteur bancaire : Planning agile vs planning projet en cascade,” in *Colloque sur les Objets et systèmes Connectés*, (CASABLANCA, Maroc), Ecole Supérieure de Technologie de Casablanca (Maroc), Institut Universitaire de Technologie d’Aix-Marseille (France), Jun 2019. Submitted on 25 Sep 2019.
- [51] F. Bordignon, “Rédaction, soumission, peer-review : un retour d’expérience multi-perspectives de la publication d’un data paper,” *GTSO*, 2022.
- [52] A. Nambiemaa, J. Fouquetb, J. Guilloteaub, and b. A. Descathaa, “La revue systématique et autres types de revue de la littérature : qu’est-ce que c’est, quand, comment, pourquoi?,” *Archives des Maladies Professionnelles et de l’Environnement*, 2021. ARTICLE IN PRESS.
- [53] D. M. Abdullah and A. M. Abdulazeez, “Machine learning applications based on svm classification a review,” *Qubahan Academic Journal*, vol. 1, no. 2, pp. 81–90, 2021.
- [54] C. Starbuck, “Logistic regression,” in *The Fundamentals of People Analytics : With Applications in R*, pp. 223–238, Springer, 2023.
- [55] S. Abbas, Z. Jalil, A. R. Javed, I. Batool, M. Z. Khan, A. Noorwali, T. R. Gadekallu, and A. Akbar, “Bcd-wert : a novel approach for breast cancer detection using whale optimization based efficient features and extremely randomized tree algorithm,” *PeerJ Computer Science*, vol. 7, p. e390, Mar. 2021.
- [56] M. L. D. De Lara, “Persistent homology classification algorithm,” *PeerJ Computer Science*, vol. 9, p. e1195, jan 2023.
- [57] J. Villaneau, F. Saïd, and A. Achour, “Modélisation et interprétation des catégories taxonomiques des animaux et aliments chez des enfants d’âge

- préscolaire,” in *Extraction et Gestion des connaissances (EGC)*, (Bruxelles, Belgique), pp. 245–252, hal-02469861, 2020.
- [58] J. Rahnenführer, R. D. Bin, A. Benner, F. Ambrogi, L. Lusa, A.-L. Boulesteix, E. Migliavacca, H. Binder, S. Michiels, W. Sauerbrei, L. McShane, and for topic group “High-dimensional data” (TG9) of the STRATOS initiative, “Statistical analysis of high-dimensional biomedical data : a gentle introduction to analytical goals, common approaches and challenges,” *BMC Medicine*, vol. 21, no. 182, 2023. Open Access article licensed under a Creative Commons Attribution 4.0 International License.
- [59] B. Albert, V. Antoine, and J. Koko, “Optimisation de fuzzy c-means (fcm) clustering par la méthode des directions alternées (admm),” in *Extraction et Gestion des Connaissances : Actes de la conférence EGC’2023*, vol. 39, BoD-Books on Demand, 2023.
- [60] W. Lu, D. Ding, X. Wu, and G. Yuan, “An efficient 1 iteration learning algorithm for gaussian mixture model and gaussian mixture embedding for neural network,” *arXiv preprint arXiv :2308.09444v2*, 2023.
- [61] J. Prezelj, J. Murovec, S. Huemer-Kals, K. Häsler, and P. Fischer, “Identification of different manifestations of nonlinear stick–slip phenomena during creep groan braking noise by using the unsupervised learning algorithms k-means and self-organizing map,” *Mechanical Systems and Signal Processing*, vol. 166, p. 108349, 2022.
- [62] L. Bellanger, A. Coulon, and P. Husi, “Une méthode de classification ascendante hiérarchique par compromis : hclustcompro,” in *9e Conférence Internationale Francophone sur la Science des Données (CIFSD)* (M. Quafafou, ed.), (Marseille, France), CIFSD, June 2021.
- [63] N. Moritz, T. Hori, and J. Le Roux, “Semi-supervised speech recognition via graph-based temporal classification,” *arXiv preprint arXiv :2010.15653v2*, 2021. [cs.LG] 16 Feb 2021.
- [64] Z. Hu, Z. Yang, X. Hu, and R. Nevatia, “Simple : Similar pseudo label exploitation for semi-supervised classification,” in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2020.

- [65] V. Truong and M.-R. Amini, “Apprentissage semi-supervisé de fonctions d’ordonnement,” Master’s thesis, Laboratoire d’Informatique de Paris 6, 104, Avenue du Président Kennedy, 75016 Paris, France, 2024.
- [66] M. Guzel, M. Kalkan, E. Bostanci, K. Acici, and T. Asuroglu, “Cloud type classification using deep learning with cloud images,” *PeerJ Computer Science*, vol. 10, p. e1779, jan 2024.
- [67] M. Hajirahimova and A. Aliyeva, “A survey on deep learning in big data analytics,” 2020.
- [68] S. Shurrab and R. Duwairi, “Self-supervised learning methods and applications in medical imaging analysis : a survey,” *PeerJ Computer Science*, vol. 8, p. e1045, jul 2022.