

REPUBLIQUE DU CAMEROUN
Paix-Travail-Patrie
UNIVERSITE DE YAOUNDE 1
DEPARTEMENT
D'INFORMATIQUE
BP/P.O.Box 812
Yaounde-Cameroun



REPUBLIC OF CAMEROON
Peace-Work-Fatherland
UNIVERSITY OF YAOUNDE 1
COMPUTER SCIENCES
DEPARTMENT
BP/P.O.Box 812
Yaounde-Cameroun

Methode explicable basee sur les auto-encodeurs pour la detection des anomalies dans les flux de donnees

Noms et prnoms : NAKAM YOPDUP MANUELLA KRISTEVA
Matricule : 19M2233
Niveau : Master 2
Spécialité : Sciences de Données(DS)
Encadreur : Pr Norbert TSOPZE

Superviseur: Pr Roger ETOUNDI ATSA

Contents

1 Domaine de recherche

Un domaine de recherche est un champ d'étude spécifique qui est traité par les chercheurs. Il est généralement défini par un ensemble de questions ou de problèmes qui sont étudiés, ainsi que par les méthodes et les techniques utilisées pour les étudier. Il existe plusieurs domaines de recherche mais dans la suite de notre devoir nous nous intéresserons au domaine des sciences de données.

Les sciences des données constituent un domaine interdisciplinaire qui combine des éléments des sciences naturelles, des sciences sociales, de l'ingénierie et des mathématiques. Les sciences des données se concentrent sur l'étude, l'analyse et l'interprétation des données.

1.1 Bases scientifiques des sciences de données

Afin d'exceller dans le domaine des sciences de données, des connaissances fondamentales sont nécessaires pour comprendre et pratiquer ce domaine. Entre autre on peut citer:

- **La théorie des probabilités** est utilisée pour modéliser l'incertitude et la variabilité des données.
- **La statistique** est utilisée pour analyser les données et en tirer des conclusions. Les statistiques jouent un rôle fondamental dans les sciences des données. Elles fournissent les méthodes et les outils nécessaires pour collecter, analyser et interpréter les données. Les concepts statistiques tels que la probabilité, l'échantillonnage, l'estimation et les tests d'hypothèses sont utilisés pour prendre des décisions basées sur les données.
- **L'apprentissage automatique** est utilisé pour développer des modèles qui peuvent apprendre à partir des données. L'apprentissage automatique est une branche des sciences des données qui implique la construction de modèles et d'algorithmes qui permettent aux ordinateurs d'apprendre à partir des données. Les bases statistiques et mathématiques sont utilisées pour développer des modèles d'apprentissage automatique, tandis que l'informatique fournit les outils et les infrastructures nécessaires pour entraîner et déployer ces modèles.
- **Informatique** : Les sciences des données reposent fortement sur les compétences en informatique. La programmation est essentielle pour manipuler et analyser les données à grande échelle. Des langages de programmation tels que Python, R et SQL sont couramment utilisés. Les bases de données, les systèmes de gestion de bases de données et les compétences en ingénierie logicielle sont également nécessaires pour gérer efficacement les flux de données et développer des systèmes d'analyse robustes.

- **Domaine d'application spcifique** : Les sciences des donnes reposent souvent sur des connaissances spcialises dans un domaine d'application spcifique. Par exemple, la bioinformatique, la finance, la mdecine, le marketing ou l'ingnierie peuvent necessiter des connaissances spcialises pour comprendre les donnes et dvelopper des modles appropris.

1.2 Types de recherche en sciences de donnees

La recherche en sciences des donnes est un domaine en pleine croissance, avec de nombreuses opportunit s de recherche dans une grande vari t de domaines. La recherche en sciences des donnes peut tre classe de diffrentes manires, en fonction **de l'approche, de l'objectif ou du domaine d'application**.

En fonction de **l'approche**, la recherche en sciences des donnes peut tre divise en deux catgories principales :

- **La recherche fondamentale** vise dvelopper de nouvelles connaissances sur les donnes et les mthodes d'analyse des donnes. Elle est gnralement mene par des universitaires dans des laboratoires de recherche.
- **La recherche applique** vise rsoudre des problmes concrets l'aide des donnes. Elle est gnralement mene par des entreprises, des organisations publiques ou des organismes de recherche applique.

En fonction de **l'objectif**, la recherche en sciences des donnes peut tre divise en deux catgories principales :

- **La recherche exploratoire** vise comprendre les donnes et identifier des tendances ou des schmas. Elle est gnralement utilise pour orienter la recherche future.
- **La recherche confirmatoire** vise tester des hypothses ou valider des modles. Elle est gnralement utilise pour prendre des dcisions ou pour dvelopper des applications.

En fonction du **domaine d'application**, la recherche en sciences des donnes peut tre divise en de nombreuses sous-catgories, telles que :

- **La recherche en sant** utilise les donnes pour amliorer la prvention, le diagnostic et le traitement des maladies.

- **La recherche en finance** utilise les données pour prendre des décisions d'investissement et de gestion des risques.
- **La recherche en marketing** utilise les données pour comprendre les comportements des consommateurs et pour développer des stratégies marketing efficaces.
- **La recherche en sécurité** utilise les données pour détecter les menaces et les incidents de sécurité.

1.3 Méthodologie de recherche en sciences de données

La méthodologie de recherche en sciences des données implique une approche systématique pour mener des études et des projets de recherche basés sur les données. Bien que la méthodologie puisse varier en fonction du contexte spécifique et des objectifs de recherche. De manière générale la méthodologie de recherche en sciences des données comprend généralement les étapes suivantes :

- **Définition du problème de recherche:** La première étape consiste à définir le problème de recherche. Cela implique de clarifier la question de recherche, de préciser les objectifs de la recherche et d'identifier les données et les méthodes d'analyse nécessaires.
- **Collecte des données:** La deuxième étape consiste à collecter les données. Les données peuvent être collectées à partir de diverses sources, telles que des bases de données, des enquêtes, des observations ou des expériences.
- **Nettoyage des données :** Une fois que les données ont été collectées, elles doivent être nettoyées. Cela implique de corriger les erreurs, de supprimer les données aberrantes et de formater les données de manière à ce qu'elles soient prêtes pour l'analyse.
- **Sélection des méthodes et des modèles :** Sélectionnez les méthodes et les modèles d'analyse appropriés en fonction de la nature de votre problème de recherche. Cela peut inclure des techniques de statistiques, d'apprentissage automatique, de fouille de données, de modélisation prédictive ou d'autres approches analytiques.
- **Analyse des données :** La quatrième étape consiste à analyser les données. L'analyse des données peut être réalisée à l'aide de diverses méthodes, telles que l'analyse statistique, l'apprentissage automatique ou l'intelligence artificielle.
- **Interprétation des résultats :** La cinquième étape consiste à interpréter les résultats de l'analyse des données. Cela implique de tirer des conclusions sur les données et de les mettre en relation avec le problème de recherche.

- **Communication des résultats :** La sixième étape consiste à communiquer les résultats de la recherche. Cela peut être fait à travers des publications scientifiques, des conférences ou des rapports.

La méthodologie de recherche en sciences des données est un processus flexible qui peut être adapté aux besoins spécifiques de chaque projet de recherche.

1.4 Les grands noms du domaine

Le domaine des sciences des données est vaste et en constante évolution, et de nombreux experts ont contribué de manière significative au développement. Voici quelques-uns des pionniers du domaine des sciences des données :

- **Geoffrey Hinton** : Geoffrey Hinton est considéré comme l'un des pionniers de l'apprentissage profond (deep learning). Ses travaux sur les réseaux de neurones artificiels ont grandement contribué à l'avancement de l'intelligence artificielle et de l'apprentissage automatique.
- **Yann LeCun** : Yann LeCun est un chercheur en intelligence artificielle et directeur de l'AI Research chez Facebook. Il est connu pour ses contributions dans le domaine de la vision par ordinateur et de l'apprentissage profond, notamment pour avoir développé le concept des réseaux de neurones convolutifs (CNN).
- **Andrew Ng** : Andrew Ng est un chercheur et entrepreneur spécialisé dans le domaine de l'apprentissage automatique. Il a cofondé Google Brain et a été l'un des cofondateurs de Coursera. Ses travaux ont contribué à populariser l'apprentissage automatique et rendre les connaissances en sciences des données plus accessibles grâce à des cours en ligne.
- **Judea Pearl** : Judea Pearl est un informaticien et philosophe reconnu pour ses travaux fondamentaux sur le raisonnement causal et les réseaux bayésiens. Ses contributions ont permis d'élargir les capacités de la modélisation causale dans les sciences des données.
- **DJ Patil** : DJ Patil est un data scientist et entrepreneur qui a joué un rôle clé dans la popularisation du terme "scientifique des données". Il a occupé le poste de Chief Data Scientist des États-Unis sous l'administration Obama et a contribué à promouvoir l'utilisation des données pour la prise de décisions dans divers domaines.

- **Fei-Fei Li** : Fei-Fei Li est une chercheuse en intelligence artificielle et vision par ordinateur. Elle a jou un rôle important dans le développement de méthodes d'apprentissage profond pour la reconnaissance d'images et a contribué à la création de la base de données d'images ImageNet, qui a été largement utilisée dans le domaine de la vision par ordinateur.
- **Hadley Wickham** : Hadley Wickham est un statisticien et programmeur renommé, connu pour ses contributions majeures dans le domaine de l'analyse de données et de la visualisation. Il a développé des packages populaires tels que ggplot2 et dplyr en langage de programmation R, qui sont largement utilisés par les scientifiques des données.

Nous avons aussi des grands noms plus récents tels que:

- **Andrew Ng** est un informaticien américain qui est considéré comme l'un des pionniers de l'apprentissage profond. Il a cofondé Coursera, une plateforme d'apprentissage en ligne, et a été nommé directeur de l'intelligence artificielle chez Baidu.
- **Ilya Sutskever** est un informaticien russe qui est également un pionnier de l'apprentissage profond. Il a cofondé OpenAI, une organisation à but non lucratif axée sur la recherche et au développement de l'intelligence artificielle, et a été directeur de l'intelligence artificielle chez Tesla.
- **Yann LeCun** est un informaticien français qui est également un pionnier de l'apprentissage profond. Il est professeur d'informatique à l'Université de New York et directeur du laboratoire de vision et de reconnaissance de formes. Quoc Le est un informaticien américain qui est un expert en apprentissage automatique et en traitement du langage naturel. Il est directeur de l'ingénierie chez Google AI.
- **Timnit Gebru** est une informaticienne éthiopienne-américaine qui est une experte en apprentissage automatique et en justice sociale. Elle est cofondatrice de Black in AI, une organisation qui vise à promouvoir la diversité et l'inclusion dans le domaine de l'intelligence artificielle.
- **Kaggle Community** : Bien que ce ne soit pas une personne spécifique, la communauté Kaggle a joué un rôle majeur dans le domaine des sciences des données depuis 2010. Kaggle est une plateforme en ligne populaire qui héberge des compétitions de science des données et permet aux chercheurs de partager des ensembles de données, de collaborer et de repousser les limites de l'apprentissage automatique.

1.5 Conferences et journaux dans le domaine de la data sciences

Il existe plusieurs confrences et journaux renomms dans le domaine des sciences des donnes. Voici quelques exemples parmi les plus influents :

Conferences :

- **Conference on Neural Information Processing Systems (NeurIPS)** : NeurIPS est l'une des principales confrences en apprentissage automatique et en intelligence artificielle. Elle rassemble des chercheurs, des praticiens et des experts du monde entier pour prsenter et discuter des dernires avances dans le domaine.
- **International Conference on Machine Learning (ICML)** : L'ICML est une confrence de premier plan en apprentissage automatique qui se concentre sur la prsentation de travaux de recherche novateurs et sur l'change d'ides entre les scientifiques des donnes.
- **International Conference on Data Mining (ICDM)** : L'ICDM est une confrence de premier plan en fouille de donnes et en extraction de connaissances partir de grandes bases de donnes. Elle rassemble des chercheurs et des praticiens pour partager leurs travaux sur les techniques, les mthodologies et les applications de l'exploration de donnes.
- **International Conference on Very Large Data Bases (VLDB)** : VLDB est une confrence majeure dans le domaine des bases de donnes et de la gestion de donnes grande chelle. Elle offre une plateforme pour prsenter des recherches sur les systmes de gestion de bases de donnes, l'analyse de donnes et les applications lies aux mgadonnes.
- **ACM SIGKDD Conference on Knowledge Discovery and Data Mining (KDD)** : KDD est une confrence de premier plan en science des donnes et en fouille de donnes. Elle explore les dernires avances dans les domaines de la dcouverte de connaissances, de l'apprentissage automatique et de l'exploration de donnes.
- **ACM SIGAI Conference on Artificial Intelligence (AAAI)** : La confrence AAAI sur l'intelligence artificielle (AAAI) est une confrence internationale de premier plan consacr l'avancement de la recherche sur l'intelligence artificielle (IA). Organise chaque anne par l'Association for the Advancement of Artificial Intelligence (AAAI), elle est considre comme l'une des confrences les plus prestigieuses dans le domaine de l'IA. Elle aborde les sujets suivants : Apprentissage automatique (Machine Learning), Apprentissage profond (Deep Learning), Vision par ordinateur

(Computer Vision), Traitement automatique du langage naturel (Natural Language Processing), Raisonnement et planification (Reasoning and Planning), Robotique et agents autonomes (Robotics and Autonomous Agents), Apprentissage par renforcement (Reinforcement Learning) et éthique de l'IA (AI Ethics).

Journaux :

- **Journal of Machine Learning Research (JMLR)** : JMLR est un journal en accès libre qui publie des articles de recherche originaux dans le domaine de l'apprentissage automatique et de l'intelligence artificielle. Il est considéré comme l'un des journaux les plus prestigieux dans le domaine de l'apprentissage automatique.
- **IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI)** : TPAMI est un journal de l'IEEE qui se concentre sur la publication d'articles de recherche sur la vision par ordinateur, la reconnaissance de formes et l'analyse d'images.
- **Data Mining and Knowledge Discovery (DMKD)** : DMKD est un journal qui couvre les aspects théoriques et pratiques de la fouille de données, de l'extraction de connaissances et de l'exploration de données.
- **Journal of Big Data** : Ce journal se concentre sur la publication d'articles de recherche sur les défis et les opportunités liés aux données massives (big data), y compris les méthodologies, les outils et les applications.
- **ACM Transactions on Knowledge Discovery from Data (TKDD)** : TKDD est un journal qui couvre les aspects de la découverte de connaissances à partir des données, y compris les algorithmes, les modèles, les méthodologies et les applications.
- **Knowledge and Information Systems (KAIS)** : est une revue scientifique internationale comitée de lecture publiée par Springer Nature. Elle est consacrée à la recherche sur les systèmes de connaissance et d'information.
- **L'ACM Transactions on Intelligent Systems and Technology (TIST)** : est une revue scientifique internationale comitée de lecture publiée par l'Association for Computing Machinery (ACM). Elle est consacrée à la recherche et à la pratique des systèmes intelligents et de la technologie.
- **IEEE Transactions on Knowledge and Data Engineering (TKDE)** : est une revue scientifique internationale comitée de lecture publiée par l'Institute of Electrical and Electronics Engineers (IEEE). Elle est considérée comme l'une des principales revues dans le domaine des sciences des données et de l'ingénierie des connaissances.

2 Les axes de recherche

2.1 Les différents axes de recherche

En data science, il existe plusieurs axes de recherche qui englobent un large ventail de domaines et de sujets. Voici quelques-uns des principaux axes de recherche :

- **Apprentissage automatique et apprentissage profond** : Cette branche explore les algorithmes et les modèles pour permettre aux ordinateurs d'apprendre à partir de données et de prendre des décisions sans être explicitement programmés.
- **Analyse de données massives (Big Data)** : Ce domaine se concentre sur les méthodes, les outils et les infrastructures pour gérer, traiter et analyser des ensembles de données massives et complexes.
- **Recherche sur les réseaux neuronaux** : Il s'agit d'une approfondie des architectures, des méthodes d'entraînement et des applications des réseaux neuronaux, y compris les réseaux convolutifs, récurrents et d'autres architectures avancées.
- **Visualisation de données** : Cette discipline cherche à développer des méthodes pour représenter visuellement des données complexes de manière à en extraire des informations significatives et à faciliter la prise de décision.
- **Exploration de données (Data Mining)** : L'exploration de données vise à découvrir des schémas, des tendances ou des relations significatives au sein de grands ensembles de données pour en tirer des informations utiles.
- **Analyste prédictive** : Utilisation de techniques statistiques et d'apprentissage automatique pour prédire des événements futurs ou des tendances en se basant sur des données historiques.
- **Sécurité et confidentialité des données** : Cette recherche se concentre sur la protection des données sensibles contre les accès non autorisés, les atteintes à la vie privée et les cybermenaces.
- **Traitement du langage naturel (NLP)** : Cette branche se focalise sur le développement de modèles et d'algorithmes pour permettre aux ordinateurs de comprendre, interpréter et générer un langage humain.
- **Optimisation et apprentissage par renforcement** : Cette recherche porte sur le développement d'algorithmes permettant de prendre des décisions séquentielles dans des environnements complexes pour atteindre des objectifs spécifiques.

- **Data Science thique et responsable** : Un domaine mergeant qui se concentre sur l'ethique de la collecte, de l'utilisation et de la diffusion des données, ainsi que sur les implications sociales de l'analyse des données.
- **Vision par ordinateur (Computer Vision)** : Il s'agit d'un domaine de recherche qui se concentre sur l'analyse, l'interprétation et la compréhension des images et des vidéos par les machines. Cela comprend la reconnaissance d'objets, la détection de motifs, la segmentation d'images, etc.

Ces axes de recherche en data science ne sont pas exhaustifs et peuvent souvent se chevaucher. De plus, de nouvelles branches émergent constamment mesure que la discipline évolue et que de nouveaux défis se présentent.

2.2 Axe de recherche choisi

Après avoir fait les cours fouille de données en Master 1 je me suis intéressée aux flux de données et en particulier à la détection d'anomalies dans les flux de données. J'ai donc lu de nombreux articles sur le sujet tel que:

- **Review of Anomaly Detection Algorithms for Data Streams, Tianyuan Lu, Lei Wang * and Xiaoyong Zhao[?]**: explore les principaux défis liés à la détection d'anomalies dans les flux de données, tels que la nature dynamique des données, la limitation des ressources et la nécessité de détecter rapidement les anomalies.
- **DeepStream: Autoencoder-Based Stream Temporal Clustering and Anomaly Detection, Shimon Harush, Yair Meidan, Asaf Shabtai[?]** : présente une méthode novatrice appelée DeepStream, qui utilise des autoencodeurs pour la détection d'anomalies et le regroupement temporel dans les flux de données.
- **Méthodes parallèles pour le traitement des flux de données continus, Mme Ge Song[?]** : travaille sur les méthodes parallèles pour le traitement des flux de données continus.
- **Détection d'anomalies dans les flux de données par structure d'indexation et approximation : Application à l'analyse en continu des flux de messages du système d'information de la SNCF, Lucas Foulon[?]** : semble proposer une approche qui combine la structure d'indexation et l'approximation pour détecter les anomalies dans les flux de messages du système d'information de la SNCF.

- **Anomaly Detection of Time Series with Smoothness-Inducing Sequential Variational Auto-Encoder[?]** : prsente une approche novatrice pour la dtection d'anomalies dans les sries temporelles. Les auteurs proposent l'utilisation d'un auto-encodeur variationnel sequentiel (Sequential Variational Auto-Encoder, SVAE) qui intgre une rgularisation favorisant la rgularit et la continuit des sries temporelles.

C'est donc sur l'article **Anomaly Detection of Time Series with Smoothness-Inducing Sequential Variational Auto-Encoder**, Longyuan Li, Junchi Yan, Member, IEEE,, Haiyang Wang, and Yaohui Jin Member, IEEE, que j'ai decide de m'appuyer avec les conseils de Dr Jiechieu Florentin. Cet artvcle traitant sur les auto-encodeurs; l'axe de recherche choisis est donc : **l'apprentissage automatique**

3 Justification du theme

Le theme qui m'a ete propose est: **Methode explicable basee sur les auto-encodeurs pour la detection des anomalies dans les flux de donnees**

3.1 Definitions de concepts

- **Auto-encodeurs:** sont une forme de rseau de neurones utiliss pour apprendre des reprsentations efficaces des donnes en tentant de reconstruire les entres la sortie. Ils sont compos de deux parties principales : un encodeur et un dcodeur. Ils sont souvent utiliss pour **la compression de donnes, la dtection d'anomalies, la rduction de dimension, la gnration de donnes**, etc.
- **Flux de donnees:** est une squence continue de donnes qui est gnre et traite de manire continue et progressive. Les flux de donnes peuvent provenir de diverses sources, telles que **des capteurs, des transactions, des vnements ou des mdias sociaux**.
- **Anomalies:** se rfrent des observations, des comportements ou des vnements qui diffrent significativement du modle ou du schma habituel des donnes. Elles peuvent tre **le rsultat d'erreurs, de dfauts, de manipulations malveillantes ou simplement de variations inhabituelles dans les donnes**.

3.2 Objectifs du theme

il sera donc question tout au long de mon travail de :

- Trouver un jeu de donnees representant un flux de donnees tels que Apache SPARK :
- En fonction du contexte du jeu de donnees choisis identifier ce qui est considere comme anomalie
- On va appliquer l'algorithme de l'article choisi sur notre jeu de donnees
- On va observer les limites et essayer d'apporter un exemple d'amelioration tout en portant l'interet sur l'explicabilite du modele
- On va appliquer la methode proposee sur notre modele et faire des tests
- Puis on va comparer les resultats

3.3 Les apports du theme

Tout theme de recherche se veut d'un apport soit pour le domaine de la recherche soit pour la vie courante:

- Concernant la recherche, les auto-encodeurs sont generalement des boites noires donc obtenir une methode explicable est un veritable atout pour faire valider les resultats par un expert et justifier son utilisation.
- D'autres parts les applications de la detection des anomalies dans les flux de donnees peuvent jouer un role important dans de nombreux domaines(sante, securite,)

4 Conclusion

Le choix d'un theme de recherche demande de suivre un ensemble d'étapes qui sont autant importantes les unes comme les autres. Le but de cette approche est de se rassurer que le projet respecte les normes de la recherche tels que l'axe de recherche. De plus l'importance et l'utilité du theme doivent être un indispensable dans le choix de celui-ci.