

Lecture : Exploiter l'IA générative dans le processus d'ingénierie des données

Introduction

Dans le paysage actuel axé sur les données, l'intégration de l'intelligence artificielle (IA) et des algorithmes génétiques (AG) est devenue une force transformatrice, remodelant les approches traditionnelles de l'ingénierie des données. De la collecte de données à l'analyse et au-delà, GenAI propose une suite de solutions innovantes qui optimisent les processus, améliorent la prise de décision et ouvrent des perspectives sans précédent.

Cette lecture se penche sur le rôle central de GenAI à travers les étapes du cycle de vie de l'ingénierie des données, illustrant comment ses capacités adaptatives et évolutives révolutionnent l'efficacité, l'évolutivité et l'efficience dans la gestion et l'exploitation des données. À travers des exemples concrets et des informations sur le secteur, nous explorons comment GenAI permet aux organisations de relever des défis complexes en matière de données, de stimuler l'innovation et de garder une longueur d'avance dans le paysage actuel de l'ingénierie des données en évolution rapide.

Exploiter l'IA générative à différentes étapes de l'ingénierie des données

1. Collecte de données :

- *Découverte automatisée des données :*
formez les modèles GenAI sur les données et la documentation existantes pour identifier et catégoriser automatiquement les sources de données potentielles (par exemple, les API, les bases de données, les capteurs) dans toute l'organisation, économisant ainsi un temps et des efforts précieux par rapport à la découverte manuelle. Cette automatisation peut être particulièrement bénéfique dans les environnements complexes avec de nombreuses sources de données.

2. Ingestion de données :

- *Génération de code pour les pipelines de données :*
en fonction des sources et des formats identifiés, vous pouvez utiliser les modèles GenAI pour générer des extraits de code pour les scripts d'extraction et de transformation des données, réduisant ainsi considérablement le temps de développement et minimisant les erreurs par rapport au codage manuel. Les modèles GenAI permettent aux ingénieurs de données de se concentrer sur des tâches plus stratégiques.
- *Détection et correction d'anomalies :*
Entraînez les modèles GenAI sur des échantillons de données propres pour identifier et traiter les incohérences et les erreurs (par exemple, les valeurs manquantes, les valeurs aberrantes) lors de l'ingestion des données, garantissant ainsi la qualité des données dès le départ et rationalisant les processus en aval.

3. Stockage des données :

- *Prédiction du schéma de données :*
utilisez GenAI pour analyser les modèles d'utilisation des données et prévoir les besoins d'accès futurs. Cela permet de recommander des formats et des structures de stockage optimaux, d'optimiser l'efficacité du stockage et de faciliter une récupération plus rapide des données lorsque cela est nécessaire.

4. Informatique:

- *Nettoyage automatisé des données :*
entraînez les modèles GenAI sur des échantillons de données propres pour identifier et corriger automatiquement les incohérences et les anomalies dans le flux de données. Cela peut impliquer des tâches telles que l'imputation de valeurs manquantes, la correction des fautes de frappe et l'identification et la gestion des valeurs aberrantes, réduisant ainsi considérablement l'effort manuel requis pour le nettoyage des données.

5. Intégration des données :

- *Alignement des schémas :*
exploitez GenAI pour analyser et suggérer des mappages entre différents formats de données provenant de diverses sources, facilitant ainsi une intégration transparente. Cela peut être particulièrement utile lorsque vous traitez des structures et des formats de données disparates.
- *Génération de données synthétiques :*
Générez des données synthétiques qui reflètent la structure et les relations des données réelles, facilitant ainsi l'intégration tout en protégeant les informations sensibles. Cette génération peut s'avérer cruciale pour permettre le partage et la collaboration des données tout en respectant les réglementations relatives à la confidentialité des données.

6. Modélisation des données :

- *Suggestion d'ingénierie des fonctionnalités :*
utilisez GenAI pour analyser les données et suggérer des fonctionnalités potentielles à inclure dans le modèle de données. Cette analyse peut impliquer l'identification des relations entre les fonctionnalités existantes, la recommandation de transformations de

fonctionnalités et la suggestion de fonctionnalités entièrement nouvelles basées sur les données, améliorant ainsi potentiellement les performances et la précision du modèle.

7. Transformation des données :

- *Génération de code pour des transformations complexes :*
générez des extraits de code pour des transformations de données complexes en fonction de règles définies par l'utilisateur ou de modèles appris à partir de données existantes. Les extraits de code peuvent automatiser des tâches telles que la normalisation des données, l'agrégation et la création de fonctionnalités, ce qui permet aux ingénieurs de données de se concentrer sur des tâches de manipulation de données plus complexes.

8. Analyse des données :

- *Exploration des données et découverte de modèles :*
entraînez les modèles GenAI à identifier les modèles et les relations cachés dans les données, en suggérant des pistes potentielles pour une analyse plus approfondie. Cette exploration des données peut impliquer des tâches telles que la détection d'anomalies, l'identification de corrélations et la découverte de clusters, fournissant des informations précieuses qui pourraient être manquées par les méthodes d'analyse traditionnelles.
- *Génération de rapports automatisée :*
générez des rapports préliminaires contenant des informations clés basées sur des modèles prédéfinis et des résultats d'analyse de données. Les modèles peuvent automatiser l'étape initiale de création de rapports, permettant aux ingénieurs de données de se concentrer sur l'affinement de l'analyse et de fournir une interprétation plus approfondie des résultats.

9. Visualisation des données :

- *Suggestion automatique de graphiques :*
en fonction des données et de l'analyse, proposez des formats de visualisation de données appropriés (par exemple, des graphiques à barres, des diagrammes de dispersion, des cartes thermiques) pour communiquer efficacement des informations. Les graphiques peuvent aider les parties prenantes non techniques à comprendre des données complexes et à prendre des décisions éclairées.

10. Gouvernance et sécurité des données :

- *Génération de données synthétiques :*
Générez des données synthétiques pour l'accès et l'analyse des utilisateurs, protégeant ainsi les informations sensibles et respectant les réglementations en matière de confidentialité des données. Les données synthétiques permettent un accès plus large aux données à des fins d'analyse et de prise de décision tout en atténuant les risques pour la confidentialité.
- *Recommandation de contrôle d'accès automatisé aux données :*
utilisez GenAI pour analyser les rôles des utilisateurs et la sensibilité des données, en suggérant des politiques de contrôle d'accès appropriées. Cette analyse rationalise les processus de gouvernance des données et garantit que les données ne sont accessibles qu'aux utilisateurs autorisés en fonction de leurs besoins spécifiques.

11. Suivi et optimisation :

- *Détection d'anomalies dans les pipelines de données :*
Formez les modèles GenAI pour surveiller les pipelines de données et identifier les problèmes potentiels tels que les erreurs ou les retards, facilitant ainsi la maintenance proactive. Cette maintenance garantit la fluidité du flux de données et évite les perturbations dans les processus en aval.
- *Suggestions d'optimisation des performances :*
analysez les flux de travail de traitement et de stockage des données avec GenAI et recommandez des optimisations pour une gestion des données plus rapide et plus efficace. Cette gestion des données peut impliquer l'identification des goulots d'étranglement, la suggestion d'algorithmes alternatifs et l'optimisation de l'allocation des ressources, améliorant ainsi l'efficacité globale du processus d'ingénierie des données.

Conclusion

Il est important de noter que l'efficacité des outils GenAI dépend de la qualité et de la pertinence des données d'entraînement. Cependant, en intégrant stratégiquement GenAI tout au long du cycle de vie des données, les ingénieurs de données peuvent améliorer considérablement l'efficacité, améliorer la qualité des données et extraire des informations précieuses de leurs données, ce qui permet en fin de compte une prise de décision basée sur les données et stimule la croissance de l'entreprise. À mesure que le domaine de l'IA générative continue d'évoluer, son potentiel à révolutionner tous les aspects de l'ingénierie des données continuera de s'étendre, permettant aux ingénieurs de données de devenir des partenaires encore plus stratégiques pour favoriser la réussite des organisations.

Auteur(s)

[Abhishek Gagneja](#)

© IBM Corporation. Tous droits réservés.