

Lecture : Étude de cas sur les mises en œuvre réussies de l'IA générative dans l'ingénierie des données

Défi

The Headline est un média de premier plan avec une forte présence en ligne et des publications imprimées. The Headline est confronté à des défis de gestion des données en raison du stockage cloisonné des données et des formats de données incohérents entre les différents services (par exemple, les plateformes d'actualités en ligne, les publications imprimées, le marketing sur les réseaux sociaux). Cela entrave leur capacité à :

- **Obtenez une vue complète de l'engagement de l'audience** : ils manquent d'une vue unifiée du comportement des utilisateurs sur différentes plateformes, ce qui rend difficile la compréhension des préférences de l'audience et l'optimisation de la diffusion de contenu.
- **Personnaliser le contenu et la publicité** : Le manque de données utilisateur intégrées rend difficile la personnalisation des recommandations de contenu et le ciblage efficace de la publicité.
- **Générer des informations basées sur les données** : la difficulté d'accéder aux données provenant de sources disparates et de les analyser entrave leur capacité à générer des informations basées sur les données pour la prise de décision stratégique.

Objectif du projet

Affectez une équipe d'ingénieurs de données pour concevoir et mettre en œuvre une plateforme de données unifiée (UDP) qui intègre les données provenant de diverses sources, rationalise la gestion des données et permet des capacités d'analyse avancées.

Responsabilités des ingénieurs de données

- **Identification et évaluation des sources de données**
Identifiez toutes les sources de données de l'organisation (par exemple, le trafic du site Web, les données des médias sociaux, la CRM et les données de vente) et évaluez leur qualité et leur cohérence.
- **Développement de pipelines de données**
Concevez et créez des pipelines de données pour extraire, transformer et charger (ETL) des données provenant de diverses sources dans l'UDP dans un format standardisé.
- **Assurance qualité des données**
Mettre en œuvre des techniques de nettoyage et de validation des données pour garantir l'exactitude et la cohérence des données au sein de l'UDP.
- **Entreposage et stockage de données**
Concevoir et mettre en œuvre une architecture d'entrepôt de données au sein de l'UDP pour un stockage et une récupération efficaces des données.
- **Accès et sécurité des données**
Développer des mécanismes pour un accès sécurisé aux données sur l'UDP tout en respectant les réglementations sur la confidentialité des données.

Résultats attendus

- **Stockage de données intégré et centralisé**
Une plateforme unique hébergeant toutes les données pertinentes pour une gestion et une analyse simplifiées.
- **Amélioration de la qualité et de la cohérence des données**
Les formats de données standardisés et les processus d'assurance qualité garantissent des données fiables pour la prise de décision.
- **Informations utilisateur améliorées**
Les données utilisateur unifiées permettent une analyse complète de l'audience et facilitent le contenu personnalisé et les stratégies publicitaires.
- **Prise de décision basée sur les données**
Un accès facile à des données propres et intégrées permet une prise de décision éclairée dans toute l'organisation.

En mettant en œuvre avec succès le projet UDP, The Headline peut libérer le potentiel de ses données, acquérir une compréhension plus approfondie de son public et stimuler la croissance de son entreprise grâce à des stratégies axées sur les données.

Solution

L'équipe d'ingénierie des données de The Headline a mené à bien le projet UDP en intégrant l'IA générative (GenAI) à différentes étapes, démontrant ainsi son potentiel pour rationaliser la gestion des données et révéler des informations précieuses. Voici comment ils ont exploité GenAI :

1. Identification et évaluation des sources de données

- *Découverte automatisée des données* :
les modèles GenAI ont été formés sur des données et une documentation existantes pour identifier et catégoriser automatiquement les sources de données potentielles dans l'ensemble de l'organisation, économisant ainsi du temps et des efforts par rapport à la découverte manuelle.

2. Développement du pipeline de données

- *Génération de code pour les pipelines de données* :
sur la base des sources et des formats de données identifiés, les modèles GenAI ont été utilisés pour générer des extraits de code pour les pipelines ETL, réduisant ainsi le temps de développement et minimisant les erreurs par rapport au codage manuel.

3. Assurance qualité des données

- *Détection et correction des anomalies* :
les modèles GenAI ont été formés sur des échantillons de données propres pour identifier et signaler les incohérences et les anomalies dans le flux de données entrant, permettant ainsi un nettoyage et une correction automatisés des données.

4. Entreposage et stockage de données

- *Optimisation du schéma* :
GenAI a analysé les modèles d'utilisation des données et prédit les futurs besoins d'accès aux données pour recommander et optimiser automatiquement le schéma de l'entrepôt de données pour un stockage et une récupération efficaces.

5. Accès et sécurité des données

- *Génération de données synthétiques* :
pour fournir un accès sécurisé aux données sensibles à des fins d'analyse, GenAI a généré des données synthétiques réalistes qui ont préservé les distributions et les relations de données sans révéler d'informations réelles sur les utilisateurs.
- *Automatisation du contrôle d'accès aux données* :
les modèles GenAI ont aidé à définir et à mettre en œuvre des contrôles d'accès des utilisateurs en fonction des rôles et de la sensibilité des données, garantissant la sécurité des données et la conformité aux réglementations.

Avantages de l'utilisation de GenAI

- **Efficacité accrue**
L'automatisation de tâches telles que la découverte de données, la génération de code et la détection d'anomalies a considérablement réduit le temps de développement et les besoins en ressources.
- **Amélioration de la qualité des données**
Le nettoyage des données et la génération de données synthétiques optimisés par GenAI ont assuré l'exactitude des données et facilité l'accès sécurisé pour l'analyse.
- **Des délais d'accès aux informations plus rapides**
Des pipelines de données rationalisés et des contrôles de qualité des données automatisés ont permis un accès plus rapide à des données propres et fiables pour l'obtention d'informations et la prise de décision.

Conclusion

Cette utilisation innovante de l'IA générative par l'équipe d'ingénierie des données de The Headline démontre le potentiel de cette technologie pour révolutionner la gestion des données et permettre aux organisations de libérer tout le potentiel de leurs données pour une prise de décision éclairée et la croissance de leur entreprise.

Auteur(s)

[Abhishek Gagneja](#)

© IBM Corporation. Tous droits réservés.

