

RÉPUBLIQUE DU SÉNÉGAL
UNIVERSITÉ GASTON BERGER DE SAINT-LOUIS

INSTITUT  POLYTECHNIQUE
DE SAINT-LOUIS



Projet de Big Data

Présenté par :

Bara Sow

Professeur : M. Mboup

Année universitaire : 2023-2024

Plan

- I. Introduction
- II. Etude de comparaison des différents formats de fichiers
 1. Parquet
 2. ORC
 3. Avro
 4. Apache Arrow
- III. Choisir le format de fichier approprié
- IV. Conclusion

I. Introduction

Dans le domaine du Big Data, où les volumes de données atteignent des tailles gigantesques, la gestion et le stockage efficaces de ces données deviennent des défis majeurs. C'est là qu'entrent en jeu les formats de fichiers, qui jouent un rôle crucial en permettant de structurer, d'organiser et de compresser ces données pour une utilisation optimale. Parmi les formats de fichiers les plus répandus pour le Big Data, on trouve Parquet, ORC, Avro et Apache Arrow. Chacun d'entre eux présente des caractéristiques et des atouts distincts, et le choix du format le plus approprié dépend des besoins spécifiques de chaque cas d'utilisation.

II. Etude de comparaison des différents formats de fichiers

1. Parquet : Un format de fichier performant pour l'analyse de données volumineuses

Parquet s'impose comme un format de fichier incontournable dans le domaine de l'analyse de données volumineuses. Son architecture en colonnes le rend particulièrement adapté aux requêtes complexes et aux charges de travail en lecture seule, offrant de nombreux avantages aux utilisateurs :

Avantages :

- **Stockage optimisé** : En organisant les données par colonnes, Parquet facilite un accès rapide aux informations spécifiques requises pour les analyses, réduisant ainsi considérablement les temps de lecture.
- **Compression efficace** : Les fichiers Parquet bénéficient d'une compression performante, permettant de minimiser l'espace de stockage nécessaire et de réduire les coûts associés.
- **Gestion de données complexes** : La structure flexible de Parquet s'adapte aisément aux types de données complexes, y compris les données imbriquées et les schémas évolutifs.

- **Adoption répandue** : Intégré nativement à l'écosystème Hadoop, Parquet profite d'une large adoption au sein de la communauté Big Data, facilitant l'interopérabilité et la collaboration.
- **Parfaitement adapté aux lectures fréquentes** : Pour les charges de travail en lecture seule (WORM - Write Once, Read Many), Parquet excelle, offrant des performances optimales pour les requêtes fréquentes sur des ensembles de données volumineux.

Inconvénients :

- **Performances pour requêtes simples** : Sur les requêtes simples, Parquet peut s'avérer moins performant que d'autres formats de fichiers.
- **Complexité du schéma** : La gestion de schémas complexes peut avoir un impact sur les performances, nécessitant une attention particulière lors de la conception des structures de données.
- **Mises à jour de données** : Les mises à jour de données dans les fichiers Parquet impliquent un traitement supplémentaire, ce qui peut s'avérer moins efficace que d'autres approches pour des workflows d'écriture fréquents.

2. ORC : Un format de fichier performant pour les environnements Hadoop

ORC (Optimized Row Columnar) s'impose comme un format de fichier performant et flexible pour le stockage et le traitement de données volumineuses, particulièrement dans les environnements Hadoop. Son architecture en colonnes et sa prise en charge de types de données complexes en font un choix judicieux pour une large gamme d'analyses de données.

Avantages :

- **Format flexible et extensible** : ORC offre une grande adaptabilité aux différents types de données et structures de schémas, permettant une évolution aisée des ensembles de données.
- **Compression efficace** : Les fichiers ORC bénéficient d'une compression performante, réduisant considérablement l'espace de stockage nécessaire et les coûts associés.
- **Prise en charge des types de données complexes** : La structure flexible d'ORC s'adapte aisément aux types de données complexes, y compris les données imbriquées et les schémas évolutifs.
- **Optimisé pour Hive** : ORC est nativement intégré à Apache Hive et fonctionne de manière optimale avec d'autres outils Hadoop, facilitant l'intégration et l'analyse des données.
- **Bon choix pour les charges de travail mixtes** : ORC offre de bonnes performances pour les requêtes fréquentes et les opérations de lecture/écriture, ce qui le rend adapté à une variété de cas d'utilisation.

Inconvénients :

- **Performances pour les analyses complexes** : Sur certaines analyses complexes, Parquet peut s'avérer plus performant qu'ORC.

- **Structure colonnaire** : L'organisation en colonnes peut ne pas être idéale pour toutes les requêtes, notamment celles nécessitant un accès aléatoire aux données.
- **Adoption** : ORC est moins largement adopté que Parquet, ce qui peut limiter son interopérabilité avec certains outils et workflows.

3. Avro : Un format de fichier binaire flexible pour les échanges de données

Avro s'impose comme un format de fichier binaire flexible et indépendant du langage, particulièrement adapté aux flux de données et aux échanges de données entre différents systèmes. Sa structure flexible et ses capacités de compression en font un choix judicieux pour de nombreux cas d'utilisation.

Avantages :

- **Format binaire flexible** : Avro est indépendant du langage de programmation, permettant une utilisation transparente entre différents systèmes et plateformes. De plus, sa structure évolutive facilite l'adaptation aux changements de schémas de données.
- **Stockage des schémas en JSON** : Les schémas Avro sont définis en JSON, offrant une lisibilité et une compatibilité accrues, simplifiant la gestion et l'interprétation des données.
- **Compression efficace** : Les fichiers Avro bénéficient d'une compression performante, réduisant l'espace de stockage nécessaire et les coûts associés, particulièrement important pour les flux de données volumineux.
- **Bien adapté aux flux de données** : La conception d'Avro le rend idéal pour les échanges de données entre différents systèmes, garantissant une transmission efficace et fiable des informations.
- **Utilisation avec Kafka** : Avro est largement utilisé avec Apache Kafka, une plateforme de streaming de données populaire, offrant une intégration transparente et des performances optimisées.

Inconvénients :

- **Performances pour les analyses complexes** : Pour les analyses complexes impliquant des agrégations et des jointures, Parquet ou ORC peuvent offrir de meilleures performances.
- **Format binaire** : Le format binaire d'Avro peut nécessiter un traitement supplémentaire pour la lecture et l'interprétation des données par rapport aux formats textuels.
- **Adoption dans l'écosystème Hadoop** : Avro est moins couramment utilisé au sein de l'écosystème Hadoop que Parquet et ORC, ce qui peut limiter son intégration avec certains outils et workflows.

4. Apache Arrow : Un format de fichier en mémoire performant pour l'analyse de données

Apache Arrow s'impose comme un format de fichier en mémoire performant, conçu pour optimiser les opérations analytiques sur des ensembles de données volumineux. Sa structure colonnaire et sa prise en charge de types de données complexes en font un outil puissant pour les data scientists et les analystes de données.

Avantages :

- **Format colonnaire en mémoire** : L'organisation des données en colonnes permet un accès rapide et efficace aux informations spécifiques requises pour les analyses, réduisant considérablement les temps de traitement.
- **Prise en charge des types de données complexes et imbriqués** : Arrow gère efficacement les structures de données complexes, y compris les tableaux imbriqués, les données JSON et les types de données personnalisés.
- **Couche d'abstraction** : Arrow offre une couche d'abstraction qui simplifie l'interaction avec différents formats de fichiers, permettant une utilisation transparente de diverses sources de données.
- **Réduction de la surcharge CPU** : La structure en mémoire d'Arrow minimise les transferts de données entre la mémoire et le disque, réduisant la surcharge du processeur et améliorant les performances globales des analyses.
- **Adapté aux bibliothèques de calcul analytique** : Arrow est particulièrement performant avec Apache Spark, une plateforme de traitement de données distribuées, offrant une intégration transparente et une accélération significative des analyses.

Inconvénients :

- **Adoption** : En tant que format relativement récent, Arrow n'a pas encore une adoption aussi large que des formats plus établis comme Parquet ou ORC.
- **Intégration** : L'utilisation d'Arrow peut nécessiter une intégration avec les outils et bibliothèques existants, ce qui peut impliquer un certain effort de développement.
- **Compatibilité** : Arrow peut ne pas être compatible avec tous les systèmes hérités et les workflows existants, nécessitant une évaluation au cas par cas.

III. Choisir le format de fichier approprié

Le choix du format de fichier le plus approprié pour vos données est une décision cruciale qui impacte directement la gestion efficace, les analyses performantes et l'interopérabilité optimale de vos informations. Divers formats présentent des caractéristiques et des avantages distincts, rendant la sélection tributaire de vos besoins spécifiques.

Cas d'utilisation : Les formats colonnaires comme Parquet et ORC sont plus performants pour les analyses complexes, tandis qu'Avro est mieux adapté aux flux de données et aux échanges de données.

Performance : Parquet et Arrow offrent généralement des performances plus élevées pour les analyses complexes, tandis que ORC est un bon choix pour les charges de travail mixtes.

Adoption : Parquet et ORC sont largement adoptés dans l'écosystème Hadoop, tandis qu'Avro est plus populaire pour les applications Kafka et Arrow gagne en popularité pour les analyses in-memory.

Complexité du schéma : Des schémas complexes peuvent affecter les performances de Parquet, tandis qu'ORC est plus flexible.

Intégration : Assurez-vous que le format choisi est compatible avec vos outils et bibliothèques existants.

IV. Conclusion

Chaque format de fichier présente ses propres avantages et inconvénients. En analysant les besoins spécifiques de votre cas d'utilisation, vous pouvez choisir le format le plus approprié pour optimiser le stockage, les performances et l'interopérabilité de vos données.