

Revu littéraire de
Implémentation d'un service d'agrégation et d'analyse de contenu

Auteur : Sanguirè Pascal SOMDA

March 6, 2024

Chapter 1

Définitions et Revue de Littérature

1.1 Introduction

La gestion et le traitement de l'information ont toujours été des défis majeurs dans le domaine de l'informatique. Ce chapitre explore les définitions essentielles et effectue une revue de littérature pour situer notre travail dans le contexte actuel.

1.2 Définitions

1.2.1 Analyse de Contenus Textuels

L'analyse de contenu textuel consiste à extraire des informations significatives à partir de textes, utilisant des techniques et des algorithmes pour comprendre, interpréter et organiser le contenu de manière automatique ou semi-automatique.

1.2.2 Analyse de Contenus Visuels

L'analyse de contenus visuels, également connue sous le nom d'analyse d'images, utilise des techniques pour extraire des informations et des caractéristiques significatives à partir d'images, de vidéos ou d'autres contenus visuels.

1.2.3 Service d'Agrégation

Un service d'agrégation collecte, rassemble et organise automatiquement des informations provenant de différentes sources en un seul emplacement centralisé, simplifiant ainsi l'accès et la gestion des données dispersées.

1.2.4 Types de Service d'Agrégation

Les types de services incluent l'extraction d'informations, la classification de texte et l'analyse de similarité.

1.3 Revue de Littérature

1.3.1 Recherche d'Images Basée sur le Contenu

Une revue complète des techniques de recherche d'images basée sur le contenu, explorant les méthodes d'extraction de caractéristiques visuelles et les approches de comparaison.

1.3.2 Récupération d’Images basée sur la Similarité de Couleur et de Texture

Focus sur la recherche d’images basée sur la similarité de couleur et de texture, explorant les méthodes d’extraction et de mesure de similarité.

1.3.3 Enquête sur les Techniques de Recherche d’Images Basées sur le Contenu

Une enquête détaillée sur les techniques de recherche d’images basée sur le contenu, abordant les différentes méthodes d’extraction et de représentation des caractéristiques visuelles.

1.4 Synthèse

Ces études offrent une vision globale des techniques dans la recherche d’images par similarité, jetant les bases pour la mise en place d’un service de recherche d’image par similarité.

1.5 Conclusion du Chapitre

Ce chapitre a établi les définitions nécessaires et a fourni une revue de littérature pour situer notre travail dans le contexte actuel. Les chapitres suivants aborderont la méthodologie et les résultats pour répondre à l’objectif de notre recherche.

Chapter 2

Évolution des données

2.1 Introduction

Ce chapitre se concentre sur l'évolution des données et leur convergence vers l'analyse de contenu, qui permet de tirer parti de ces vastes ensembles de données pour obtenir des connaissances exploitables. Nous allons explorer comment les données ont évolué au fil du temps, en passant par les différents types de données, les formes de stockage de données à l'explosion du volume de données à l'ère numérique. Nous examinerons également les avancées technologiques qui ont permis de gérer, de stocker et de traiter ces données massives.

2.1.1 Types de données

Une donnée est une information brute, une représentation objective et factuelle d'un élément ou d'un événement. Elle peut prendre différentes formes, telles que des chiffres, du texte, des images, des vidéos, des enregistrements sonores, des codes, des symboles, etc. Les données peuvent être utilisées pour représenter des faits, des mesures, des observations, des caractéristiques ou des éléments constitutifs d'un système. Cependant, les données peuvent être organisées de manière structurée, non structurées ou de manière semi-structurées.

Données structurées

Les données structurées font référence à des informations qui sont organisées selon un format spécifique et cohérent. Elles sont généralement stockées dans des bases de données relationnelles ou dans d'autres formats structurés. Les données structurées sont organisées de manière cohérente, ce qui permet une manipulation et une analyse faciles à l'aide de requêtes et d'opérations de base de données.

Données non structurées

Les données non structurées désignent des informations qui ne sont pas organisées de manière rigide ou prédéfinie. Contrairement aux données structurées, qui sont stockées dans des formats tabulaires ou relationnels, les données non structurées ne suivent pas de schéma spécifique et peuvent varier en termes de format, de contenu et de représentation.

Données semi structurées

Les données semi-structurées se situent entre les données structurées et non structurées. Elles possèdent une certaine organisation, mais ne suivent pas un modèle de données rigide comme les données structurées.

2.1.2 Évolution des données sous l'influence des données massives

La gestion et le traitement de l'information ont toujours été des défis majeurs dans le domaine de l'informatique. Au fil du temps, les professionnels de l'informatique ont développé des compétences pour gérer des fichiers,

des répertoires, des bases de données, etc. Cependant, avec la croissance exponentielle des données et la diversification des formats de stockage, il est devenu crucial de pouvoir traiter rapidement ces informations. C'est ainsi qu'est apparu le concept du Big Data, en réponse à ce besoin croissant de stockage et de traitement efficace des données massives. Le Big Data est un concept qui fait référence à la gestion et à l'exploitation de vastes volumes de données, caractérisés par leur volume, leur variété et leur vélocité. Avec l'avènement des technologies numériques, de l'Internet et des appareils connectés, les données sont générées à une échelle sans précédent dans divers domaines tels que les réseaux sociaux, les capteurs, les transactions commerciales, les enregistrements médicaux, etc. Il s'agit d'un concept permettant de stocker un nombre indicible d'informations sur une base numérique.

2.1.3 Importance de Big Data dans le traitement des données

Le Big Data joue un rôle essentiel dans le traitement des données en permettant de gérer, d'analyser et d'exploiter des volumes massifs de données variées. Il offre de nouvelles perspectives, de meilleures prévisions et une prise de décision plus éclairée pour les organisations, leur permettant ainsi de rester compétitives et d'innover dans un environnement complexe et en constante évolution.

Données massives

Le Big Data permet de gérer des ensembles de données de grande envergure qui dépassent la capacité des systèmes traditionnels. Cela permet aux organisations de collecter, stocker et analyser des quantités massives de données provenant de diverses sources telles que les médias sociaux, les appareils connectés, les capteurs, etc.

Informations variées

Le Big Data prend en charge la diversité des données, qu'elles soient structurées, semi-structurées ou non structurées. Cela inclut des formats tels que les textes, les images, les vidéos, les fichiers audio, etc. L'analyse de ces différentes formes de données permet d'obtenir des informations approfondies et de découvrir des tendances ou des schémas cachés.

Analyse en temps réel

Le Big Data permet une analyse en temps réel des données, ce qui signifie que les organisations peuvent obtenir des informations instantanées et prendre des décisions plus rapides. Cela est particulièrement important dans des domaines tels que le commerce électronique, la surveillance de la santé, la finance, où la réactivité et la prise de décision rapide sont cruciales.

Prédiction et prévention

L'analyse des données volumineuses permet de découvrir des tendances, des corrélations et des modèles prédictifs. Cela permet aux organisations de prévoir des événements futurs, de prendre des mesures préventives et d'optimiser leurs opérations. Par exemple, dans le domaine de la santé, le Big Data peut aider à prédire les épidémies, à améliorer les soins aux patients et à prévenir les maladies.

2.1.4 Méthodes

Intégration des données

Lorsque nous évoquons l'intégration de données dans des bases de données, cela implique la possibilité de regrouper plusieurs sources de données en un seul point d'accès centralisé. Grâce à l'intégration, l'accès, l'utilisation et la description des données est possible pour la distribution et la réutilisation.

Il existe de nombreuses architectures disponibles pour l'intégration de données. Elle découle d'une combinaison de divers outils spécialisés. Cependant, il est possible de les diviser en deux catégories principales : l'intégration d'entrepôt de données et l'intégration par médiateur. Chacun a ses avantages et ses inconvénients.

Intégration par entrepôt de données

L'intégration d'un entrepôt de données est une approche dite "gourmande" ou "avide" (eager en anglais), qui consiste à créer de grandes collections de données à partir de sources disparates. De ce fait, elle mobilise d'importantes ressources informatiques. Cette approche est principalement utilisée à des fins d'analyse de données à grande échelle, y compris les métadonnées. L'intégration fait partie de cette catégorie et concentre les données en un point physique, prêt pour le traitement par des outils de traitement des documents textuels.

Cette catégorie est basée sur le processus d'intégration par copie et transformation. En effet, les outils de conversion de données doivent pouvoir traiter de grandes quantités de données en temps minimal. Parmi ces outils, nous avons le processus d'extraire, transformer, charger (Extract-Transform-Load (ETL) en anglais). Ces ETL permettent d'extraire des données en leur appliquant différentes transformations (par exemple des filtres et des modifications de la structure des documents) permettant ensuite de les charger dans un nouveau dépôt. Les nouveaux entrepôts peuvent aussi être appelés des lacs de données (data lake en anglais) si les données sont copiées sans transformation. Ces lacs sont découpés en entrepôts (data warehouse) puis en magasins de données (data mart). Les outils d'extraction, transformation et chargement comme les ETL doivent garantir un résultat stable (idempotence) afin de garantir la reproductibilité des expériences, notamment pour des utilisations intensives de données.

L'approche par entrepôt de données est la plus utilisée, notamment en entreprise, car la copie et la transformation des données rendent le dépôt de données de destination personnalisé pour les usages ciblés. Du fait de cette copie, il n'est pas garanti que ces dernières représentent la version la plus à jour des données du dépôt originel. Utiliser cette approche apporte le risque pour les chercheurs de baser leur raisonnement sur des données obsolètes.

Intégration par médiateur

L'approche médiateur, contrairement à l'approche entrepôt, ne copie pas les données, car elle permet d'aller chercher les données directement à la source. Face à la répartition des données en plusieurs points, la requête envoyée correspond généralement à un cadre de recherche plus large que celui d'un seul dépôt, et doit donc être adaptée à l'exécution par des intermédiaires : les médiateurs et les emballages de sources de données.

Un médiateur représente le processus de transformation des données du schéma global au suivant, le schéma spécifique du référentiel où les données sont stockées et sur lesquelles les requêtes sont effectuées. Par conséquent, les données ne sont pas directement accessibles. Il est ensuite transformé par des médiateurs. Les requêtes sont divisées et envoyées aux emballages de données (wrapper en anglais) pour sélectionner différentes sources correspondantes. Le rôle de l'emballage de la source de données est de faire l'interface entre le schéma global des médiateurs et la source de données : il traduit les requêtes provenant des médiateurs et renvoie le résultat via les médiateurs.

L'approche par médiateur offre l'avantage d'une vision unifiée des données, notamment par la prise en compte d'ontologies propres au domaine afin d'intégrer une terminologie précise du domaine. Cependant, selon [6], la méthode de l'approche par médiateur est plus complexe à mettre en place que celle de l'approche par entrepôt de données, mais elle est plus fiable car les requêtes accèdent directement aux dépôts de données publics d'origine, ce qui permet d'utiliser toujours les versions les plus récentes des données.

Stockage des données

Quel que soit le volume ou le type d'approche, les données seront stockées dans une base de données. Il en existe de nombreux types, à savoir les bases de données relationnelles, clés-valeurs, colonnes, documents, graphes et multimodales. Lorsqu'on parle de base de données, on voit plus les systèmes qui gèrent ces bases de données. Ces systèmes sont simplement appelés Système de Gestion de Bases de Données (SGBD) (DataBase Management System (DBMS) en anglais) et chacun présente des avantages et des inconvénients d'après les dires de [7].

Systèmes de gestion de base de données relationnelles

Les SGBD relationnelles tels que PostgreSQL 21 ou MySQL 22 offrent la possibilité d'effectuer des jointures entre les tables de données. Leur principal avantage réside dans leur stabilité et leur utilisation du langage

de requête structuré (SQL). Toutefois, cette rigidité peut rendre l'insertion de nouveaux éléments complexe, nécessitant la conception d'une base de données, la mise en place d'un schéma de tables et de relations, ainsi que la normalisation des données avant insertion, qui requièrent des compétences en informatique, notamment en conception et développement d'après [7]. Bien que ce modèle soit stable, les requêtes impliquant plusieurs jointures peuvent être lentes, et l'intégration de données non structurées reste un défi. Pour gérer efficacement de grandes quantités de données, il est recommandé d'utiliser des outils de parallélisation, également appelés "horizontal scalability", qui permettent de répartir la charge de travail sur plusieurs machines physiques afin de paralléliser les traitements.

Systèmes de gestion de base de données NoSQL

Les systèmes de gestion de bases de données qui ne reposent pas sur le modèle relationnel ont été regroupés sous le nom de "NoSQL" ou "non-SQL" par abus de langage, car ils ne sont pas basés sur le langage SQL. Cependant, de plus en plus de systèmes de gestion de bases de données peuvent maintenant utiliser le langage SQL pour interroger leurs données, même s'ils ne sont pas basés sur un modèle relationnel. Par conséquent, l'acronyme "NoSQL" ou "non-SQL" est maintenant préféré pour signifier "pas seulement SQL" ou "Not Only SQL" en anglais. Les systèmes de gestion de bases de données NoSQL ont été développés pour offrir des alternatives aux limites des systèmes de gestion de bases de données relationnelles, notamment en ce qui concerne le traitement de grandes quantités de données. En effet, les SGBD NoSQL ont généralement de meilleures performances que les SGBD relationnels pour la gestion de données massives. De plus, les SGBD NoSQL possèdent plusieurs fonctionnalités intéressantes, telles que la possibilité d'être distribués sur plusieurs machines (scalabilité horizontale en anglais) et la prise en charge de la réplication. Ces types de systèmes sont donc mieux adaptés à la gestion de données massives. Il existe plusieurs types de systèmes de gestion de bases de données NoSQL, tels que les systèmes clé-valeur, document et graphe. Différents critères doivent être pris en compte pour choisir le type de SGBD NoSQL approprié, tels que les performances ou la combinaison avec d'autres outils pour créer un système cohérent. Par exemple, MongoDB et Elasticsearch sont deux systèmes de gestion de bases de données NoSQL orientés vers les documents, capables de gérer des données semi-structurées sous des formats tels que XML et JSON. En combinant deux SGBD orientés documents (l'un pour le stockage de données et l'autre pour la recherche de données).

Exploitation des données

L'exploitation des données, également appelée "data mining" ou "data analytics", fait référence au processus d'exploration, d'analyse et d'utilisation des données pour en extraire des informations précieuses, des modèles, des tendances ou des relations cachées. L'objectif de l'exploitation des données est de découvrir des insights significatifs et exploitables à partir des données brutes, afin de prendre des décisions éclairées et d'améliorer les performances ou les résultats d'une entreprise ou d'une organisation. Le processus d'exploitation des données comprend plusieurs étapes.

Collecte des données

La collecte des données implique la collecte et l'agrégation de données provenant de différentes sources telles que des bases de données, des fichiers, des systèmes, des capteurs, des médias sociaux, des sites web, etc. Les données peuvent être structurées, semi-structurées ou non structurées.

Préparation des données

Cette étape consiste à nettoyer, filtrer et transformer les données pour les rendre cohérentes, de haute qualité et adaptées à l'analyse ultérieure. Cela peut impliquer l'élimination des valeurs manquantes, la suppression des duplications, la normalisation des formats, la conversion des données dans un format commun, etc.

Exploration des données

C'est l'étape où les données sont explorées en utilisant des techniques d'analyse statistique et des méthodes visuelles pour découvrir des tendances, des motifs ou des anomalies. Cela peut inclure des analyses descriptives telles que des statistiques récapitulatives, des visualisations graphiques, des tests de corrélation.

Métadonnées

Définition et importance

Les métadonnées sont des données qui fournissent des informations sur d'autres données selon [8]. Elles décrivent les caractéristiques, les propriétés ou les attributs des données, ce qui permet de les organiser, de les identifier et de les gérer de manière plus efficace. Les métadonnées peuvent inclure des informations telles que la date de création, l'auteur, la taille du fichier, le format, la résolution, les mots-clés, les droits d'accès, les informations de localisation, etc.

Dans le contexte des données, les métadonnées jouent un rôle crucial dans leur gestion et leur utilisation. Elles facilitent la recherche, la récupération et l'organisation des données en permettant aux utilisateurs de comprendre leur contenu et leur contexte. Les métadonnées peuvent également être utilisées pour garantir l'intégrité et la qualité des données, ainsi que pour faciliter leur échange et leur partage entre différents systèmes.

Gestion des métadonnées

Collecte et création de métadonnées

Les métadonnées peuvent être collectées automatiquement lors de la création ou de la capture des données, ou bien elles peuvent être ajoutées manuellement par les utilisateurs. Il est essentiel de déterminer quelles informations sont pertinentes et nécessaires pour décrire les données de manière adéquate.

Stockage des métadonnées

Les métadonnées peuvent être stockées de différentes manières, en fonction des besoins et des exigences du système. Elles peuvent être incluses dans les fichiers de données eux-mêmes, stockées dans des bases de données dédiées, ou gérées par des systèmes de gestion de métadonnées. Le stockage des métadonnées doit garantir leur accessibilité, leur cohérence et leur sécurité.

Utilisation des métadonnées

Les métadonnées sont utilisées pour améliorer la découverte, l'accès et l'utilisation des données. Elles permettent aux utilisateurs de rechercher des données spécifiques, de comprendre leur provenance, leur qualité et leur pertinence, et de prendre des décisions informées sur leur utilisation. Les métadonnées sont également cruciales dans le contexte de l'intégration des données, de la gestion de la confidentialité et de la conformité réglementaire.

Modèles de métadonnées

Les modèles de métadonnées définissent la structure et les normes pour représenter les métadonnées dans un système. Ils déterminent les types d'informations pouvant être incluses, les relations entre les différentes informations, et les règles de gestion des métadonnées. Les modèles de métadonnées peuvent être spécifiques à un domaine ou génériques, et ils sont souvent basés sur des normes industrielles pour assurer l'interopérabilité.

Interopérabilité des métadonnées

L'interopérabilité des métadonnées est la capacité des systèmes et des applications à échanger et à utiliser des métadonnées de manière cohérente et efficace. Elle est essentielle pour assurer la collaboration, l'intégration des données et la gestion des informations à l'échelle de l'entreprise. L'adoption de normes de métadonnées communes facilite l'interopérabilité en fournissant un langage commun et des conventions de représentation des métadonnées.

Conclusion

En conclusion, les métadonnées jouent un rôle fondamental dans la gestion, la découverte et l'utilisation efficace des données. Elles offrent un contexte essentiel pour comprendre les informations, garantissent la

qualité des données et facilitent l'intégration des données à l'échelle de l'entreprise. La collecte, le stockage, l'utilisation et l'interopérabilité des métadonnées sont des aspects clés pour maximiser leur valeur et soutenir les processus décisionnels informés.

Analyse prédictive

L'analyse prédictive est une branche de l'analyse des données qui vise à prédire les résultats futurs en utilisant des modèles statistiques, des algorithmes d'apprentissage automatique et d'autres méthodes d'analyse des données. Elle repose sur l'idée d'extraire des informations à partir de données existantes pour anticiper les tendances, les comportements ou les résultats futurs. L'analyse prédictive a des applications dans divers domaines, notamment le marketing, la finance, la santé, les opérations commerciales, la gestion des ressources humaines, etc.

Modèles de données

Définition

Un modèle de données est une représentation structurée des concepts, des relations, des contraintes et des règles qui définissent la manière dont les données sont stockées, traitées et organisées dans un système d'information. Il sert de cadre pour comprendre et communiquer la structure logique des données, facilitant ainsi la conception, la mise en œuvre et la maintenance des bases de données.

Types de modèles de données

Il existe plusieurs types de modèles de données, chacun adapté à des besoins spécifiques de modélisation et de traitement des données. Les principaux types de modèles de données sont les suivants :

Modèle conceptuel de données (MCD) Le modèle conceptuel de données représente les concepts et les relations métier sans se soucier des détails techniques de mise en œuvre. Il offre une vue abstraite des entités, de leurs attributs et des relations entre elles. Le modèle conceptuel de données est souvent utilisé lors des phases initiales de conception des bases de données.

Modèle logique de données (MLD) Le modèle logique de données se concentre sur la conversion du modèle conceptuel en une représentation plus détaillée et technique. Il spécifie les tables, les colonnes, les clés primaires, les clés étrangères et d'autres éléments nécessaires à la mise en œuvre d'une base de données relationnelle. Le modèle logique de données est généralement utilisé dans le processus de conception détaillée de la base de données.

Modèle physique de données (MPD) Le modèle physique de données décrit la manière dont les données sont réellement stockées et organisées au niveau physique, y compris les types de données, les index, les partitions, etc. Il est étroitement lié à la mise en œuvre technique de la base de données, tenant compte des performances et de l'efficacité du stockage.

Utilisation des modèles de données

Les modèles de données sont largement utilisés dans le domaine de la gestion de bases de données et du développement logiciel. Ils jouent un rôle crucial dans les processus de conception, de mise en œuvre et de maintenance des systèmes d'information. Voici quelques-unes des utilisations principales des modèles de données :

Conception de bases de données Les modèles conceptuels aident les concepteurs de bases de données à comprendre et à représenter les besoins métier de manière abstraite. Les modèles logiques guident la conception détaillée des bases de données relationnelles, tandis que les modèles physiques détaillent l'implémentation technique.

Communication Les modèles de données servent de moyen de communication entre les différentes parties prenantes impliquées dans le développement logiciel et la gestion de bases de données. Ils fournissent une représentation visuelle et structurée des structures de données.

Documentation Les modèles de données servent de documentation essentielle pour comprendre la structure, les contraintes et les relations des bases de données. Ils facilitent la maintenance et les évolutions des systèmes d'information au fil du temps.

Analyse et optimisation Les modèles de données aident à analyser la structure des bases de données et à identifier des opportunités d'optimisation en termes de performances, d'efficacité de stockage et de gestion des accès.

Migration et intégration Lors de la migration ou de l'intégration de systèmes, les modèles de données facilitent la compréhension des différences structurelles entre les bases de données, contribuant ainsi à une transition en douceur.

Conclusion

En conclusion, les modèles de données sont des outils essentiels dans le domaine de la gestion de bases de données et du développement logiciel. Ils fournissent une représentation structurée des concepts, des relations et des règles qui guident la conception, la communication, la documentation, l'analyse et l'optimisation des bases de données. Les modèles de données sont un élément central pour garantir la cohérence, la qualité et la pérennité des systèmes d'information.

Chapter 3

Outils et langage de programmation utilisés

3.1 Introduction

Dans le cadre de ce chapitre, nous explorerons les différentes technologies et ressources qui jouent un rôle essentiel dans le domaine de la recherche d'images. Ces outils sont conçus pour faciliter diverses tâches et processus liés à notre sujet d'étude, offrant des fonctionnalités avancées pour améliorer notre compréhension, notre analyse et notre prise de décision.

3.2 Pile ELK

ELK est un ensemble de logiciels comprenant Elasticsearch, Logstash et Kibana. Chacun de ces outils joue un rôle spécifique dans le processus qui nous permet de rechercher, analyser et visualiser des données provenant de diverses sources et formats, en garantissant la fiabilité et la sécurité en temps réel. Chaque outil a sa propre fonctionnalité distincte, mais ils fonctionnent ensemble de manière harmonieuse pour offrir une expérience intégrée.

3.2.1 Elasticsearch

ElasticSearch, lancé en 2010 par Shay Banon et commercialisé par Elastic, a gagné en popularité en tant que moteur de stockage, de recherche et d'analyse de contenu distribué. Bien qu'il soit open source malgré son aspect commercial, ElasticSearch se démarque par l'utilisation du format JSON pour stocker les données, éliminant ainsi le besoin de support de stockage externe. Il repose sur Apache Lucene pour l'indexation et la recherche de contenu, tandis que Logstash et Kibana sont utilisés pour l'analyse des données.

L'interrogation d'ElasticSearch se fait via une API REST, accessible par le protocole HTTP, avec le support des normes HTML, REST et JSON, facilitant son intégration avec d'autres applications. Les fonctionnalités distribuées permettent des recherches rapides, généralement en moins d'une seconde. ElasticSearch offre un support NoSQL, combinant stockage, indexation, recherche et analyse des données. Son orientation documentaire en fait un choix attrayant, et son API REST offre une flexibilité d'interaction. La distribution et l'évolutivité sont des différences notables par rapport à d'autres outils similaires, justifiant ainsi son choix.

3.2.2 Kibana

Kibana est une plateforme gratuite qui aide à visualiser et analyser les données stockées dans Elasticsearch. C'est un outil convivial largement utilisé pour créer des tableaux de bord interactifs, des graphiques, des cartes géographiques, et d'autres représentations visuelles des données. Il est souvent utilisé pour la surveillance des systèmes, l'analyse des journaux et la visualisation des données en temps réel.

3.2.3 Python

Python, créé en 1991 par Guido van Rossum, est un langage de programmation polyvalent, lisible et facile à apprendre. Sa simplicité en fait un choix populaire pour les débutants en programmation, et il est largement utilisé dans divers domaines tels que le développement web, la science des données, l'intelligence artificielle et l'automatisation des tâches. La polyvalence de Python lui permet de créer des scripts simples ou des programmes complexes, notamment pour le traitement d'images et le calcul de similarités entre celles-ci. Dans notre cas, Python a été utilisé pour extraire les métadonnées et les caractéristiques visuelles des images, en utilisant des formules mathématiques. De plus, le framework Flask a été employé pour développer une application web permettant aux utilisateurs de visualiser des images similaires à leurs requêtes.

3.2.4 Flask

Flask, un framework web léger en Python, simplifie la création rapide et efficace d'applications web. Adoptant le modèle de développement "micro-framework", Flask offre uniquement les fonctionnalités de base nécessaires, permettant aux développeurs de choisir les bibliothèques et outils complémentaires selon leurs besoins. Il propose des fonctionnalités telles que le routage des URL, la gestion des sessions, la création de templates HTML, l'accès aux bases de données, la gestion des formulaires, et bien d'autres. Célèbre pour sa simplicité et sa flexibilité, Flask permet aux développeurs d'écrire un code clair et concis tout en offrant une grande liberté dans la structuration et le développement de leurs applications.

3.2.5 Bibliothèque

cv2 (OpenCV)

OpenCV, une bibliothèque renommée pour le traitement d'images et la vision par ordinateur en Python. Le module cv2 d'OpenCV permet de manipuler et de traiter des images, réalisant des opérations de vision par ordinateur telles que la détection d'objets, le suivi, la reconnaissance faciale, etc.

os

Le module os offre des fonctionnalités pour interagir avec le système d'exploitation. Il facilite la manipulation des fichiers et des répertoires, l'accès aux variables d'environnement, la création de processus, etc.

Numpy

NumPy, une bibliothèque essentielle pour le calcul scientifique en Python, fournit des structures de données performantes, notamment des tableaux multidimensionnels, ainsi que des fonctions mathématiques pour la manipulation de ces tableaux.

3.2.6 Conclusion

En conclusion, ce chapitre a présenté divers outils et un langage de programmation couramment utilisés dans le domaine de traitement images et recherche d'images. Nous avons exploré des outils tels que Elasticsearch, Logstash et kibana et un langage tels que Python, en mettant en évidence leurs caractéristiques, leurs utilisations et leurs avantages respectifs.

Chapter 4

Réalisation

4.1 Introduction

Au cours de cette réalisation, nous avons relevé plusieurs défis techniques, tels que la collecte et le traitement des données d'images, la conception, ainsi que l'optimisation des performances et de la précision de la recherche. En utilisant des méthodes telles que la segmentation d'images, l'extraction de caractéristiques visuelles et la comparaison des similarités.

4.1.1 Architecture du modèle

Image Indexation -> Calcul des métadonnées -> Signature La requête -> Calcul des métadonnées + Signature
-> Le comparateur entre en action => Montre les images similaires

4.1.2 Bases d'images

PlantVillage est une plateforme en ligne qui vise à aider les agriculteurs à diagnostiquer et à traiter les maladies des plantes

4.1.3 Extraction et calcul de similarité

On peut segmenter une image en la convertissant de BGR vers HSV (En lui ajoutant une couleur) pour ressortir les différences d'images. Puis après traitement, on utilise Kibana pour la visualisation d'image

4.1.4 Conclusion

En conclusion, l'implémentation d'un service d'agrégation et d'analyse de contenus est une étape clé pour tirer des informations précieuses à partir de grandes quantités de données. Ce service permet de collecter, d'organiser et d'analyser des contenus provenant de diverses sources, ce qui facilite la recherche d'informations pertinentes et la prise de décisions éclairées.

4.2 Conclusion générale

En conclusion générale, l'intégration des données est un aspect fondamental de la gestion de l'information dans divers domaines tels que la recherche, l'entreprise, la santé, etc. Les organisations s'efforcent de tirer parti de leurs données pour prendre des décisions éclairées et obtenir un avantage concurrentiel. Cependant, l'intégration des données présente des défis importants liés à la diversité des sources, aux formats variés, aux changements constants et aux volumes massifs de données.

Deux approches principales pour l'intégration des données sont l'approche par entrepôt de données et l'approche par médiateur. Chacune a ses avantages et ses inconvénients, et le choix entre les deux dépend des besoins spécifiques de l'organisation.

Le stockage des données repose sur l'utilisation de systèmes de gestion de base de données (SGBD). Les SGBD relationnelles, tels que PostgreSQL et MySQL, offrent stabilité et utilisation du langage SQL, tandis que les SGBD NoSQL sont plus adaptés à la gestion de grandes quantités de données et offrent une meilleure évolutivité.

L'exploitation des données, y compris l'analyse prédictive et l'exploration des données, vise à découvrir des insights significatifs à partir des données brutes. Cela implique la collecte, la préparation et l'exploration des données pour en extraire des informations précieuses.

Les métadonnées jouent un rôle crucial dans la gestion des données, fournissant des informations sur les caractéristiques, la qualité et le contexte des données. La collecte, le stockage, l'utilisation et l'interopérabilité des métadonnées sont des aspects clés pour maximiser leur valeur.

Enfin, les modèles de données sont des outils essentiels pour la conception, la communication, la documentation, l'analyse et l'optimisation des bases de données. Ils offrent une représentation structurée des concepts et des relations qui guident les processus liés aux bases de données.

En résumé, l'intégration des données est un domaine complexe mais essentiel, et les organisations doivent adopter des approches, des outils et des pratiques adaptés pour relever les défis et exploiter pleinement le potentiel de leurs données. L'évolution continue de la technologie, des normes et des méthodologies dans le domaine de l'intégration des données nécessite une adaptation constante pour rester à la pointe de ce domaine en évolution constante.