

The background of the slide is decorated with numerous paint splashes in various colors, including yellow, orange, green, and blue, creating a vibrant and artistic effect.

ANALYSE ET VISUALISATION DES DONNÉES

Alexandra-Maria DOBRESCU

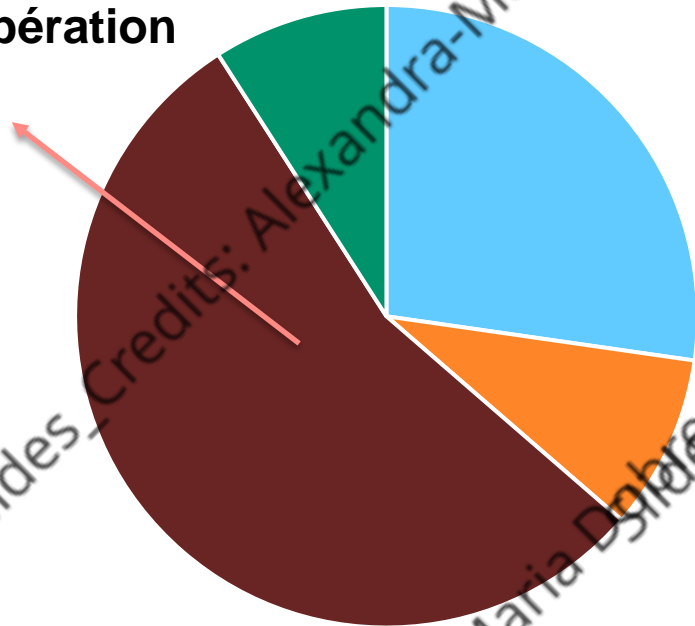
| Cours 2 |

Dis-m'en plus sur les notes

QUANTITATIVE

Travail sur place,
sans récupération

Note Finale



■ Examen ■ Tests de cours ■ Les laboratoires ■ Bonus

QUALITATIVE

Nous

RESPECT

TRANSPARENCE

TRAVAIL

APPRENTISSAGE



Introduction

- DÉFINIR L'ANALYSE DES DONNÉES
- DÉFINIR LA VISUALISATION DES DONNÉES
- APROFONDIR
- COMPARAISON
- POURQUOI L'ANALYSE DES DONNÉES EST IMPORTANTE DANS DIFFÉRENTS DOMAINES
- TYPES DE DONNÉES (STRUCTURÉES, NON STRUCTURÉES, SEMI-STRUCTURÉES)

OBJECTIF

PROCESSUS

RÉSULTATS

OUTILS

RÔLE

INTERDÉPENDANCE

Comparaison

OBJECTIF

Analyse des données

Implique l'examen, le nettoyage, la transformation et l'interprétation des données afin d'en extraire des informations significatives, d'identifier des modèles, de faire des prédictions ou de répondre à des questions de recherche spécifiques.

Vise à extraire des connaissances et à tirer des conclusions à partir des données.

Visualisation des données

Consiste principalement à représenter visuellement les données au moyen de tableaux, de graphiques, de cartes et d'autres éléments graphiques.

Son objectif principal est d'aider les utilisateurs à comprendre et à interpréter les données rapidement et efficacement.

PROCESSUS

Analyse des données

Comprend des étapes telles que le nettoyage des données, l'analyse exploratoire des données (AED), la modélisation statistique, la vérification des hypothèses et la formulation de conclusions.

Visualisation des données

Implique la sélection de techniques de visualisation appropriées, la conception de visuels et la création de tableaux ou de graphiques pour représenter visuellement les données.

Il fait souvent appel à des techniques statistiques et informatiques.

On met l'accent sur la communication des informations.

Analyse des données

Le résultat de l'analyse des données se présente généralement sous la forme de résumés statistiques, de rapports, de modèles ou de résultats numériques.

Il fournit des informations et des réponses à des questions de recherche spécifiques.

RÉSULTATS



Visualisation des données

Le résultat de la visualisation des données consiste en des diagrammes, des graphiques, des tableaux de bord et des représentations visuelles interactives qui permettent aux individus de saisir plus facilement les tendances, les modèles et les relations au sein des données.

Analyse des données

L'analyse des données utilise souvent des outils logiciels et des langages de programmation tels que R, Python, SQL et des progiciels statistiques pour effectuer des analyses et des calculs.

OUTILS



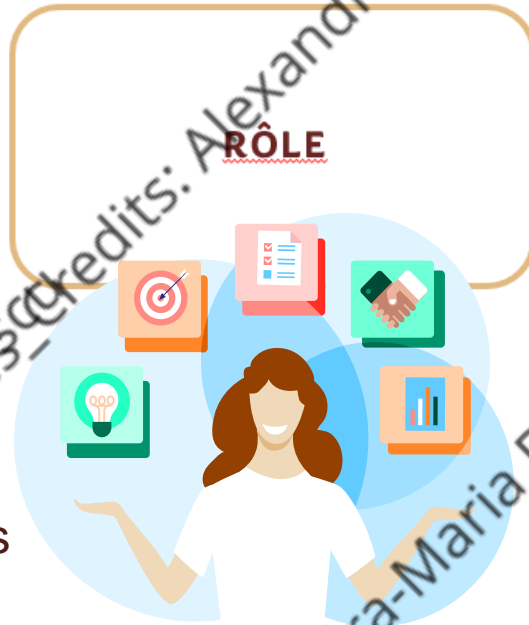
Visualisation des données

Les outils et bibliothèques de visualisation des données, tels que Tableau, D3.js, Matplotlib et ggplot2, sont utilisés pour créer des représentations visuelles des données.

Analyse des données

Les analystes et les scientifiques des données sont principalement responsables de l'analyse des données. Ils se concentrent sur l'exploration des données, la modélisation statistique et l'obtention d'informations.

RÔLE



Visualisation des données

Les spécialistes de la visualisation des données ou les analystes de données se spécialisent souvent dans la création de visualisations efficaces pour transmettre des informations sur les données à un public plus large.

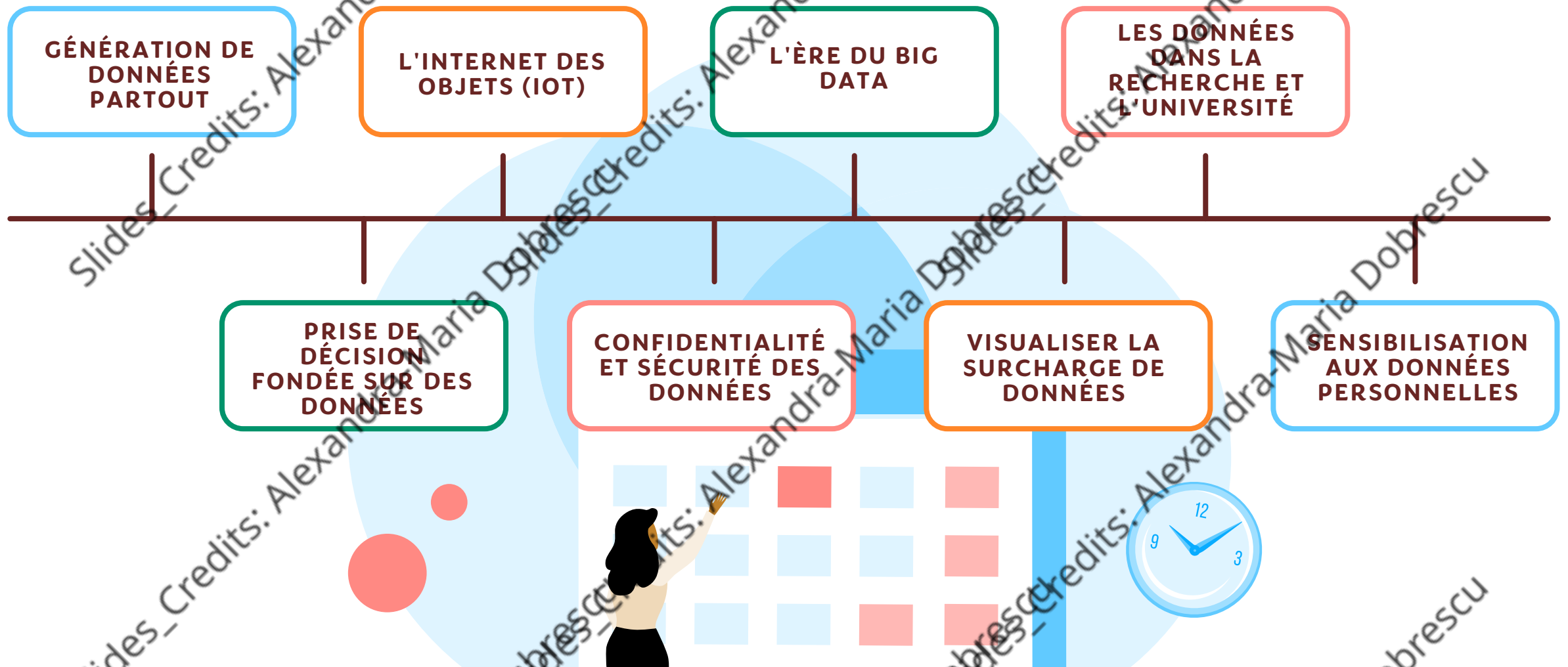
Analyse des données

INTERDÉPENDANCE

Visualisation des données

L'analyse des données peut générer des idées qui éclairent la conception des visualisations.

Les visualisations peuvent aider à présenter les résultats de l'analyse des données d'une manière claire et compréhensible.



Omniprésence des données

Génération de données partout

LES DONNÉES SONT GÉNÉRÉES EN PERMANENCE, TOUT AUTOUR DE NOUS, SOUS DIVERSES FORMES ET À PARTIR DE DIVERSES SOURCES:

Les plateformes de médias sociaux collectent des données sur les interactions des utilisateurs et leurs préférences en matière de contenu [1].



Les sites de commerce électronique enregistrent le comportement des clients, leurs achats et leur historique de navigation.



Les senseurs des appareils intelligents recueillent des données sur la température, l'emplacement et les habitudes d'utilisation [2].



L'internet des objets (IoT)

LES DISPOSITIFS DE L'IOT PRODUISENT DE GRANDES QUANTITÉS DE DONNÉES.

Flux de données continus et en temps réel : Permet de détecter immédiatement les anomalies, les irrégularités ou les événements critiques.

Exemple 1: Les thermostats intelligents surveillent en permanence la température et l'humidité [3].

Exemple 2: Les trackers de fitness suivent le rythme cardiaque et les pas.

Permettre d'améliorer considérablement l'efficacité, la réduction des coûts et la prise de décision dans différents domaines.



68 bpm

avg. resting heart rate
same as previous week



68.45

total km

▲ 20.14 km over last week

Étude de cas - L'IdO et la surveillance de l'environnement

EXERCICE : PRÉVISION DE L'INDICE DE QUALITÉ DE L'AIR (IQA)

Scénario: Imaginez que vous disposiez d'un réseau de capteurs IoT placés dans une ville pour surveiller la qualité de l'air. Ces capteurs mesurent divers polluants:

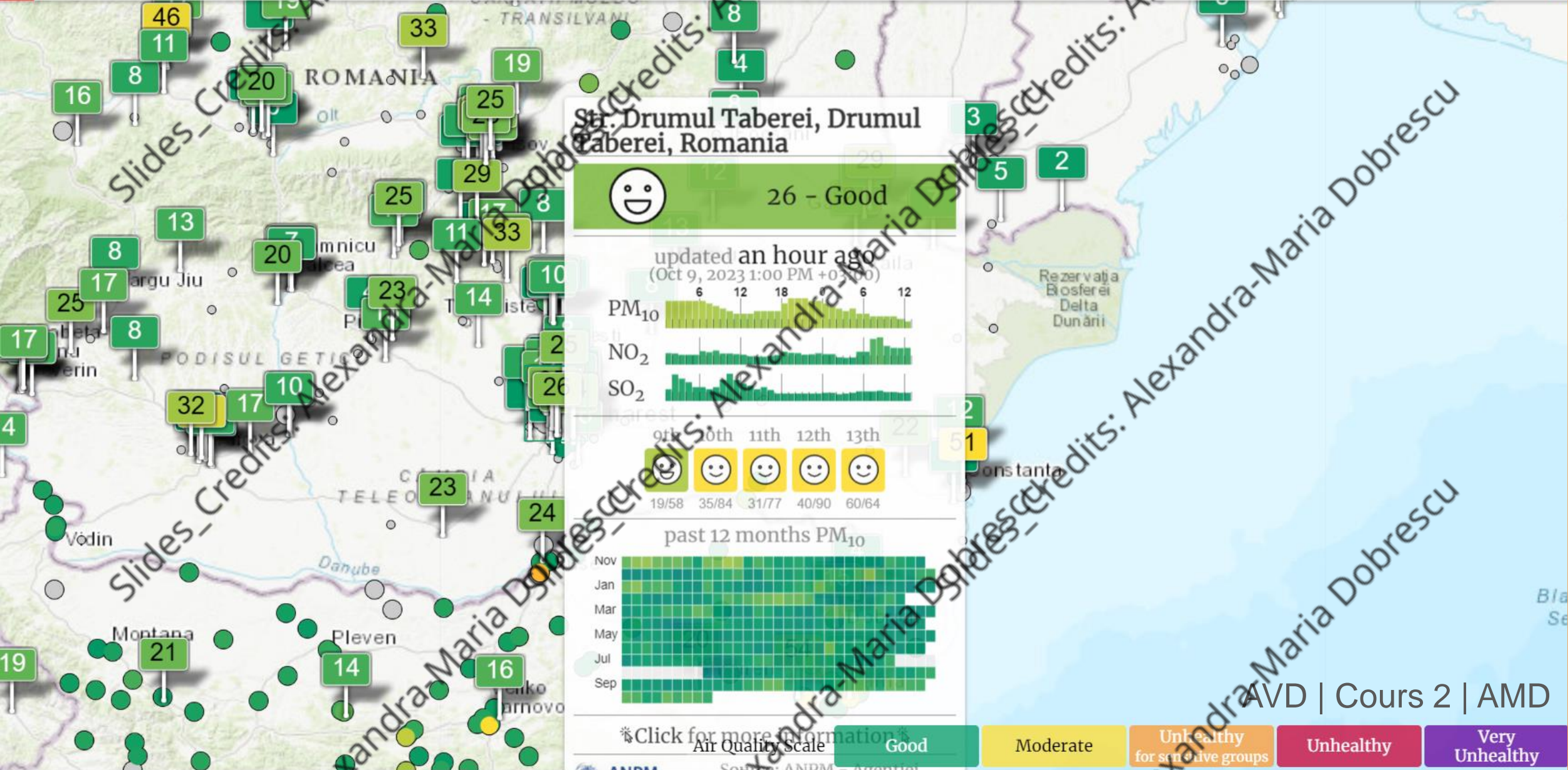
- les particules (« Particle Matter » PM2, .5),
- le monoxyde de carbone (CO)
- l'ozone (O3).

Vous souhaitez utiliser les données collectées par ces capteurs pour prédire l'indice de qualité de l'air (IQA) de la ville, qui est une mesure critique de la qualité de l'air.

Description:

1. Vous avez collecté des données au cours de l'année précédente, y compris des mesures quotidiennes de PM2.5, CO, O3, température, humidité et vitesse du vent.
2. Votre tâche consiste à créer un modèle prédictif capable d'estimer l'IQA sur la base de ces variables.

World's Air Pollution: Real-time Air Quality Index



Étude de cas - L'IdO et la surveillance de l'environnement

EXERCICE : PRÉVISION DE L'INDICE DE QUALITÉ DE L'AIR (IQA)

Étapes:

1. Divisez votre ensemble de données en un ensemble de formation et un ensemble de test (par exemple, 80 % pour le « training » et 20 % pour le test).
2. Utilisez la régression linéaire ou une autre technique de régression appropriée pour construire un modèle prédictif. Vous pouvez utiliser des outils tels que la bibliothèque scikit-learn de Python [4].
3. Entraînez votre modèle à l'aide de l'ensemble de données d'entraînement. Utilisez des variables telles que les PM_{2,5}, le CO, l'O₃, la température, l'humidité et la vitesse du vent comme caractéristiques d'entrée et l'IQA comme variable cible.
4. Évaluez les performances du modèle à l'aide de l'ensemble de données de test. Calculez des mesures telles que l'erreur absolue moyenne (MAE) et le R au carré (R²) pour évaluer la capacité de votre modèle à prédire l'IQA.
5. Interprétez les résultats. Discutez de l'importance de chaque variable d'entrée dans la prédiction de l'IQA. Par exemple, une concentration plus élevée de PM_{2.5} entraîne-t-elle un IQA plus élevé ?
6. Examinez les possibilités d'optimisation du modèle. Pourriez-vous améliorer sa précision en incluant des variables supplémentaires ou en utilisant une technique de régression différente ?

L'ère du Big Data

LES BIG DATA DÉSIGNENT DES ENSEMBLES DE DONNÉES EXTRÊMEMENT VOLUMINEUX ET COMPLEXES QUI NE PEUVENT ÊTRE GÉRÉS OU ANALYSÉS EFFICACEMENT À L'AIDE D'OUTILS ET DE MÉTHODES DE TRAITEMENT DE DONNÉES TRADITIONNELS.

Dans le traitement traditionnel des données, les systèmes RDBMS (Relational Database Management System) tels qu'Oracle, MySQL ou Microsoft SQL Server ont été les solutions privilégiées pour le stockage et l'analyse des données structurées.

Ces systèmes sont excellents pour:

- traiter des données structurées avec des schémas bien définis
- des volumes de données relativement modérés.

Big Data: comprennent généralement de grandes quantités de **données non structurées** ou **semi-structurées**, en plus des données structurées.



L'ère du Big Data

LES LIMITATIONS DU RDBMS

1. **Scalabilité:** Les bases de données traditionnelles peuvent avoir du mal à évoluer horizontalement pour gérer les volumes massifs de données générés dans les scénarios de big data.
2. **Performance:** Au fur et à mesure que les données augmentent, les requêtes sur un RDBMS traditionnel peuvent devenir lentes et inefficaces.
3. **Complexité:** La structuration et la gestion des données dans une base de données relationnelle traditionnelle peuvent devenir complexes, en particulier lorsqu'il s'agit de divers types de données.
4. **Coût:** La mise à l'échelle des bases de données traditionnelles peut être coûteuse, tant en termes de matériel que de licences logicielles.

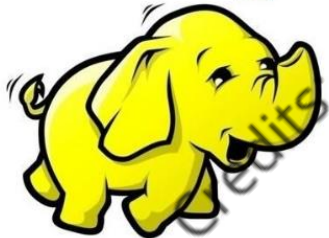


L'ère du Big Data

POUR GÉRER L'ÉCHELLE ET LA VARIÉTÉ DES DONNÉES VOLUMINEUSES DE MANIÈRE PLUS EFFICACE, ON A:

Des technologies Big Data telles que:

hadoop



Quelle technologie choisir ?

HADOOP

- Traitement de grands ensembles de données dans des environnements où la taille des données dépasse la mémoire disponible.
- Mise en place d'une infrastructure d'analyse de données avec un budget limité.
- Exécution de tâches ne nécessitant pas de résultats immédiats et pour lesquelles le temps n'est pas un facteur limitant.
- Traitement par batch avec des tâches exploitant les opérations de lecture et d'écriture sur disque.
- Analyse de données historiques et d'archives.

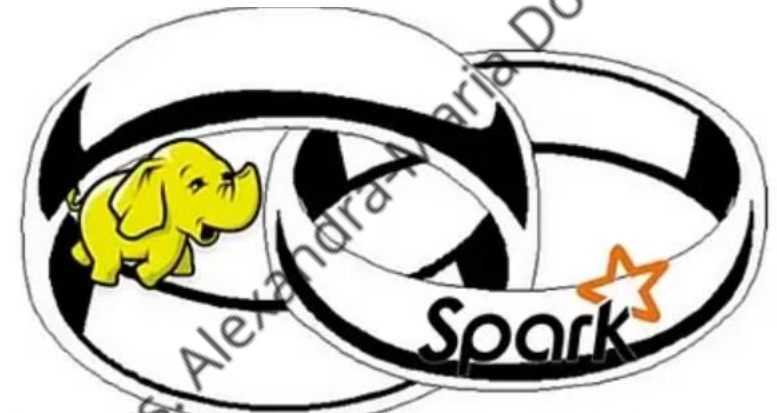
SPARK

- L'analyse des données de flux en temps réel.
- Lorsque le temps est compté, Spark fournit des résultats rapides grâce à des calculs en mémoire.
- Traitement des chaînes d'opérations parallèles à l'aide d'algorithmes itératifs.
- Le traitement parallèle des graphes pour modéliser les données.
- Toutes les applications d'apprentissage automatique.

Quelle technologie choisir ?

Hadoop et Spark fonctionnent mieux ensemble. En tant que successeur, Spark n'est pas là pour remplacer Hadoop, mais pour utiliser ses fonctionnalités afin de créer un nouvel écosystème amélioré.

En combinant les deux, Spark peut tirer parti des fonctionnalités qui lui manquent, telles qu'un système de fichiers.



**Not Mutually Exclusive But
Better Together**

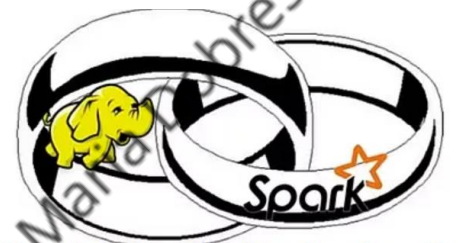


Quelle technologie choisir ?

En combinant les deux, Spark peut tirer parti des fonctionnalités qui lui manquent, telles qu'un système de fichiers.

Sans Hadoop, les applications professionnelles risquent de passer à côté de données historiques cruciales que Spark ne peut pas traiter.

Hadoop stocke une grande quantité de données à l'aide d'un matériel abordable et effectue ensuite des analyses, tandis que Spark apporte un traitement en temps réel pour traiter les données entrantes.



Not Mutually Exclusive But
Better Together

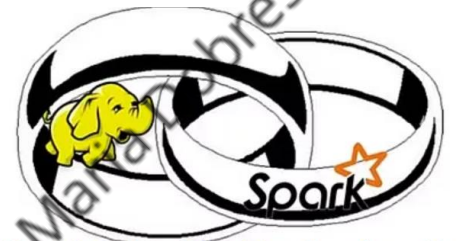


Quelle technologie choisir ?

Spark tire également parti des avantages de Hadoop en matière de sécurité et de gestion des ressources.

Vous pouvez exécuter automatiquement les charges de travail Spark en utilisant toutes les ressources disponibles.

Avec YARN, la mise en cluster de Spark et la gestion des données sont beaucoup plus faciles.



Not Mutually Exclusive But
Better Together

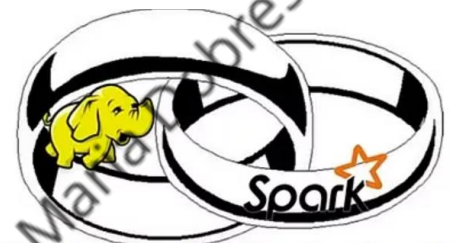


Quelle technologie choisir ?

Spark tire également parti des avantages de Hadoop en matière de sécurité et de gestion des ressources.

Vous pouvez exécuter automatiquement les charges de travail Spark en utilisant toutes les ressources disponibles.

Avec **YARN**, la mise en cluster de Spark et la gestion des données sont beaucoup plus faciles.



Not Mutually Exclusive But
Better Together

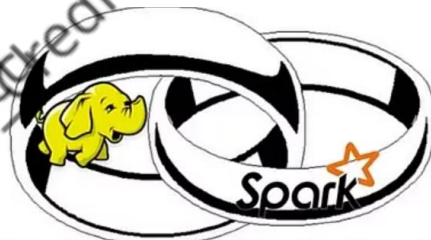


Quelle technologie choisir ?

CETTE COLLABORATION PERMET D'OBTENIR LES MEILLEURS RÉSULTATS EN MATIÈRE:

- D'ANALYSE DE DONNÉES TRANSACTIONNELLES RÉTROACTIVES,
- D'ANALYSE AVANCÉE
- DE TRAITEMENT DES DONNÉES IOT/IDO.

TOUS CES CAS D'UTILISATION SONT POSSIBLES DANS UN SEUL ENVIRONNEMENT.



Not Mutually Exclusive But
Better Together

ZeYre

Quelle BD choisir [5] ?



- Système de gestion de base de données **NoSQL** gratuit et open-source, distribué, à large stockage de colonnes.
- Prévu pour traiter de grandes quantités de données sur de nombreux serveurs de base, offrant une haute disponibilité sans point de défaillance unique.



- MongoDB est un programme de base de données non relationnelle (c'est-à-dire **NoSQL**) orienté vers les documents (document-oriented) et multiplateforme (cross-platform).
- Offre une grande vitesse, une grande disponibilité et une grande évolutivité.

Caractéristiques [5]	Cassandra	MongoDB
Modèle de données	<u>Base de données NoSQL</u> à grand nombre de colonnes qui utilise un format tabulaire pour le stockage des données. Elle est conçue pour traiter des quantités massives <u>de données structurées.</u>	<u>Base de données NoSQL</u> basée sur des documents qui stocke les données dans des documents BSON flexibles de type JSON.
Langage de requête (Query Language)	CQL (Cassandra Query Language) pour les requêtes, qui est similaire à SQL mais avec quelques différences dues à sa nature NoSQL.	Langage de requête puissant pour des requêtes flexibles et complexes qui peuvent être exprimées dans une syntaxe de type JSON.

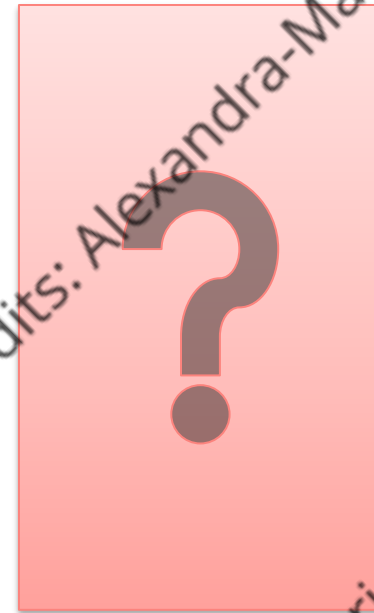
Caractéristiques [5]	Cassandra	MongoDB
Évolutivité (Scalability)	<p>Connu pour son excellente <u>évolutivité horizontale</u>, ce qui le rend adapté aux déploiements distribués à grande échelle.</p>	<p>Prend également en charge l'<u>évolutivité horizontale</u>, mais il se peut qu'il n'évolue pas de manière aussi transparente que Cassandra dans des scénarios de données extrêmement volumineuses et à grande vitesse.</p>
Modèle de Consistance (Consistency Model)	<p>Offre des niveaux de cohérence réglables, ce qui vous permet de choisir entre une <u>consistance forte</u> et une <u>consistance éventuelle</u> en fonction des besoins de votre application.</p>	<p>Offre généralement une cohérence forte au sein d'un document unique, mais permet une <u>consistance éventuelle</u> lorsqu'il est distribué.</p>

Définir de nouveaux mots-clés

CONSISTANCE FORTE



CONSISTANCE ÉVENTUELLE



Définir de nouveaux mots-clés

CONSISTANCE FORTE

Imaginez que vous disposez d'un compte bancaire dont le solde est de 100 dollars. Vous effectuez deux transactions simultanées: l'une pour retirer 30 dollars et l'autre pour vérifier le solde de votre compte. Dans un système à forte cohérence, le système garantit que lorsque vous vérifiez votre solde immédiatement après avoir retiré 30 \$, vous verrez le solde mis à jour de 70 \$.

En effet, la cohérence forte garantit que toutes les opérations semblent se produire instantanément, comme s'il n'existait qu'une seule copie de votre compte.

CONSISTANCE ÉVENTUELLE

Dans un autre scénario, disons que vous avez publié une photo sur une plateforme de médias sociaux. Vous la consultez ensuite immédiatement. Il est possible que vous ne voyiez pas immédiatement votre propre publication lorsque vous la consultez juste après l'avoir publiée. Toutefois, au bout d'un certain temps (quelques secondes ou quelques minutes), toutes les répliques des données du système distribué se synchroniseront et vous finirez par voir votre propre message.

En effet, la cohérence éventuelle autorise les incohérences temporaires car elle privilégie la disponibilité et la faible latence par rapport à la cohérence immédiate.

Définir de nouveaux mots-clés

CONSISTANCE FORTE

CONSISTANCE ÉVENTUELLE

Dans les deux exemples, la cohérence forte garantit que vous voyez l'état immédiat et cohérent des données, tandis que la cohérence à terme autorise des incohérences temporaires mais garantit que les données convergeront vers un état cohérent au fil du temps. Le choix entre les deux dépend des exigences spécifiques de l'application et des compromis entre la cohérence, la disponibilité et la latence.

Caractéristiques [5]	Cassandra	MongoDB
Flexibilité du modèle de données (Data Model Flexibility)	<p>Bien adapté aux données de séries temporelles et à l'enregistrement d'événements, où la structure des données reste relativement constante.</p>	<p>La conception sans schéma (schema-less design) est plus flexible et polyvalente, ce qui la rend adaptée à un large éventail d'applications et de données dynamiques.</p>
Cas d'utilisation (Use Cases)	<p>Souvent utilisé dans des scénarios où la haute disponibilité, la tolérance aux pannes et les charges de travail à forte intensité d'écriture sont cruciales, telles que les données chronologiques, les données de capteurs et les analyses en temps réel.</p>	<p>Populaire pour les applications nécessitant de la flexibilité, des requêtes complexes et un développement rapide, telles que les systèmes de gestion de contenu, les catalogues et les profils d'utilisateurs.</p>

Who Uses Cassandra



Jelvix

Who Uses MongoDB



jelvix.com

Bibliographie

- [1] https://www.freepik.com/premium-vector/round-social-media-icons-network-platforms-logos_12628744.htm
- [2] <https://www.taggdigital.com/blog/what-are-smartwatch-sensors-and-how-do-they-function>
- [3] <https://www.goodhousekeeping.com/uk/product-reviews/tech/g685723/best-smart-heating-thermostats/#product-e8b0b9d6-3938-4403-8422-24f8da27ff43>
- [4] Kramer, O., & Kramer, O. (2016). Scikit-learn. Machine learning for evolution strategies, 45-53.
- [5] <https://www.geeksforgeeks.org/difference-between-cassandra-and-mongodb/>
- [6] <https://www.image-net.org/about.php>
- [7] G. Miller, R. Beckwith, C. Felbaum, D. Gross, K. Miller, Introduction to WordNet: An On-line Lexical Database (1993)
- [8] Knafllic, C. N. (2015). Storytelling with data: A data visualization guide for business professionals. John Wiley & Sons.