



Spécialité : DUT en Sécurité Informatique et Réseaux (SIR)

COMPTE RENDU DE

Mini Projet :

Gestion et analyse des données
académiques avec MongoDB et
PySpark

Réalisé par :

Hajar Rachid

Année universitaire 2025_2026

lien : [Mini-Projet-Big-Data.ipynb](#)

Introduction

Dans le contexte de la transformation numérique, la gestion et l'analyse des données massives occupent une place centrale, notamment dans le domaine éducatif. Les établissements d'enseignement génèrent quotidiennement un volume important de données académiques telles que les informations sur les étudiants, leurs filières et leurs notes.

Ce mini-projet s'inscrit dans le cadre du module Big Data & Bases de Données NoSQL et vise à illustrer l'utilisation d'une base de données NoSQL orientée documents (MongoDB via Mongita) ainsi que l'exploitation de l'outil PySpark pour l'analyse de données. Le projet permet de comprendre l'intérêt des technologies NoSQL et Big Data pour le traitement efficace des données académiques.

Objectifs du mini-projet

Les principaux objectifs de ce mini-projet sont :

- Comprendre le fonctionnement des bases de données NoSQL orientées documents.
- Apprendre à insérer, stocker et interroger des données JSON avec Mongita.
- Manipuler des requêtes d'exploration des données (filtrage, recherche, mise à jour, suppression).
- Exploiter PySpark pour réaliser des analyses analytiques simples sur un jeu de données.
- Utiliser SparkSQL pour interroger les données de manière déclarative.
- Produire un compte rendu structuré et professionnel du travail réalisé.

Étapes principales du mini-projet :

1. Configuration de l'environnement

L'environnement de travail a été configuré dans **Google Colab**. Les bibliothèques nécessaires ont été installées, notamment :

- PySpark pour le traitement Big Data,
- Mongita (compatible avec PyMongo) pour la base NoSQL locale,
- Pandas pour la manipulation et l'exportation des données.

La configuration de Java a également été vérifiée afin d'assurer le bon fonctionnement de PySpark.

Partie 1 – Configuration de l'environnement dans Google Colab

Nous installons ici les bibliothèques nécessaires : PySpark, PyMongo, Pandas et Mongita.

```
[1] 38s !apt-get update -q  
!apt-get install -y openjdk-11-jdk-headless -qq > /dev/null  
!pip install -q pyspark pymongo mongita pandas  
  
import os  
os.environ["JAVA_HOME"] = "/usr/lib/jvm/java-11-openjdk-amd64"  
print("Configuration terminée avec succès.")  
  
Hit:1 https://cli.github.com/packages stable InRelease  
Get:2 https://cloud.r-project.org/bin/linux/ubuntu jammy-cran40/ InRelease [3,632 B]  
Get:3 https://developer.download.nvidia.com/compute/cuda/repos/ubuntu2204/x86_64 InRelease [1,581 B]  
Get:4 http://security.ubuntu.com/ubuntu jammy-security InRelease [129 kB]  
Hit:5 http://archive.ubuntu.com/ubuntu jammy InRelease  
Get:6 https://r2u.stat.illinois.edu/ubuntu jammy InRelease [6,555 B]  
Get:7 https://cloud.r-project.org/bin/linux/ubuntu jammy-cran40/ Packages [83.6 kB]  
Get:8 http://archive.ubuntu.com/ubuntu jammy-updates InRelease [128 kB]  
Get:9 https://developer.download.nvidia.com/compute/cuda/repos/ubuntu2204/x86_64 Packages [2,159 kB]  
Get:10 https://ppa.launchpadcontent.net/deadsnakes/ppa/ubuntu jammy InRelease [18.1 kB]  
Get:11 https://r2u.stat.illinois.edu/ubuntu jammy/main all Packages [9,499 kB]  
Get:12 http://archive.ubuntu.com/ubuntu jammy-backports InRelease [127 kB]  
Hit:13 https://ppa.launchpadcontent.net/graphics-drivers/ppa/ubuntu jammy InRelease  
Get:14 http://security.ubuntu.com/ubuntu jammy-security/main amd64 Packages [3,561 kB]  
Hit:15 https://ppa.launchpadcontent.net/ubuntu/gis/ppa/ubuntu jammy InRelease  
Get:16 http://archive.ubuntu.com/ubuntu jammy-updates/universe amd64 Packages [1,598 kB]  
Get:17 http://security.ubuntu.com/ubuntu jammy-security/universe amd64 Packages [1,286 kB]  
Get:18 http://security.ubuntu.com/ubuntu jammy-security/restricted amd64 Packages [6,059 kB]  
Get:19 https://r2u.stat.illinois.edu/ubuntu jammy/main amd64 Packages [2,840 kB]  
Get:20 https://ppa.launchpadcontent.net/deadsnakes/ppa/ubuntu jammy/main amd64 Packages [38.8 kB]  
Get:21 http://archive.ubuntu.com/ubuntu jammy-updates/main amd64 Packages [3,895 kB]  
Get:22 http://archive.ubuntu.com/ubuntu jammy-updates/restricted amd64 Packages [6,273 kB]  
Fetched 37.7 MB in 8s (4.622 kB/s)
```

2. Création du jeu de données JSON

Dans cette partie, un jeu de données au format JSON a été conçu afin de représenter des informations académiques relatives aux étudiants, telles que l'identifiant, le nom, la filière et les notes par matière. Ce format semi-structuré met en évidence la flexibilité des bases de données NoSQL orientées documents.

Partie 2 – Création du jeu de données JSON

Créons un jeu de données contenant des informations sur des étudiants et leurs notes.

```
[1] Start coding or generate with AI.

[2] os
    import json
    etudiants = [
        {
            "id": 1,
            "nom": "Amina",
            "filiere": "Informatique",
            "notes": [
                {"matiere": "Big Data", "note": 15},
                {"matiere": "Programmation", "note": 17},
                {"matiere": "Mathématiques", "note": 14}
            ],
        },
        {
            "id": 2,
            "nom": "Youssef",
            "filiere": "Réseaux",
            "notes": [
                {"matiere": "Systèmes", "note": 12},
                {"matiere": "Réseaux", "note": 14},
                {"matiere": "Sécurité", "note": 13}
            ],
        },
        {
            "id": 3,
            "nom": "Sara",
            "filiere": "Informatique",
            "notes": [
                {"matiere": "Big Data", "note": 18},
                {"matiere": "IA", "note": 17},
                {"matiere": "Programmation", "note": 16}
            ],
        }
    ]
    with open("etudiants.json", "w") as f:
        json.dump(etudiants, f, indent=2)

    print("Fichier JSON créé avec succès.")

[3] Fichier JSON créé avec succès.
```

3. Insertion des données dans la base NoSQL

L'objectif de cette partie est d'insérer le jeu de données JSON dans une base de données NoSQL locale en utilisant Mongita, une solution compatible avec MongoDB. Cette étape permet de comprendre le principe du stockage des documents dans une base orientée documents et de simuler le fonctionnement d'une base NoSQL sans configuration réseau complexe.

Partie 3 – Insertion des données dans la base NoSQL (Mongita)

Mongita fonctionne comme MongoDB mais stocke les données localement dans Colab, sans configuration réseau.

```
[1] from mongita import MongitaClientDisk
    import json

    client = MongitaClientDisk()
    db = client["enseignement_db"]
    collection = db["etudiants"]

    collection.delete_many({}) # Nettoyage
    with open("etudiants.json") as f:
        data = json.load(f)
    collection.insert_many(data)
    print("Insertion réussie dans la base Mongita locale.")

[2] Insertion réussie dans la base Mongita locale.
```

4. Requêtes d'exploration des données

Cette partie est consacrée à l'exploration et à la manipulation des données stockées dans la base NoSQL. Différentes requêtes ont été réalisées afin d'afficher, filtrer et rechercher les informations académiques. Ces opérations permettent de démontrer la souplesse des bases NoSQL dans l'accès et la manipulation des données.

```
[4] os
# Afficher tous les étudiants
print("\nListe complète des étudiants :")
for doc in collection.find():
    print(doc)

# Afficher seulement le nom et la filière (projection manuelle en Python)
print("\nNoms et filières des étudiants :")
for doc in collection.find():
    print(("nom": doc.get("nom"), "filiere": doc.get("filiere")))

# Rechercher les étudiants d'une filière donnée
print("\nEtudiants en Informatique :")
for doc in collection.find({"filiere": "Informatique"}):
    print(doc)

# Afficher les étudiants ayant une note > 15 en Programmation
print("\nEtudiants ayant une note > 15 en Programmation :")
for doc in collection.find():
    for note in doc.get("notes", []):
        if note["matiere"] == "Programmation" and note["note"] > 15:
            print(f"\t{doc['nom']} : {note['note']}")

# Compter le nombre total d'étudiants
print("\nNombre total d'étudiants :", collection.count_documents({}))
```



```
[5] os
Liste complète des étudiants :
{'id': 1, 'nom': 'Amina', 'filiere': 'Informatique', 'notes': [{'matiere': 'Big Data', 'note': 15}, {'matiere': 'Programmation', 'note': 17}, {'matiere': 'Mathématiques', 'note': 14}], '_id': ObjectId('693033dee09b4b3f2ecbcc61')}
{'id': 2, 'nom': 'Youssef', 'filiere': 'Réseaux', 'notes': [{'matiere': 'Systèmes', 'note': 12}, {'matiere': 'Réseaux', 'note': 14}, {'matiere': 'Sécurité', 'note': 13}], '_id': ObjectId('693033dee09b4b3f2ecbcc61')}
{'id': 3, 'nom': 'Sara', 'filiere': 'Informatique', 'notes': [{'matiere': 'Big Data', 'note': 18}, {'matiere': 'IA', 'note': 17}, {'matiere': 'Programmation', 'note': 16}], '_id': ObjectId('693033dee09b4b3f2ecbcc62')}

Noms et filières des étudiants :
{'nom': 'Amina', 'filiere': 'Informatique'}
{'nom': 'Youssef', 'filiere': 'Réseaux'}
{'nom': 'Sara', 'filiere': 'Informatique'}

Etudiants en Informatique :
{'id': 1, 'nom': 'Amina', 'filiere': 'Informatique', 'notes': [{'matiere': 'Big Data', 'note': 15}, {'matiere': 'Programmation', 'note': 17}, {'matiere': 'Mathématiques', 'note': 14}], '_id': ObjectId('693033dee09b4b3f2ecbcc61')}
{'id': 3, 'nom': 'Sara', 'filiere': 'Informatique', 'notes': [{'matiere': 'Big Data', 'note': 18}, {'matiere': 'IA', 'note': 17}, {'matiere': 'Programmation', 'note': 16}], '_id': ObjectId('693033dee09b4b3f2ecbcc62')}

Etudiants ayant une note > 15 en Programmation :
Amina : 17
Sara : 16

Nombre total d'étudiants : 3
```

5 . Exportation des données vers PySpark

Dans cette partie, les données stockées dans la base NoSQL ont été exportées vers un fichier JSON afin de pouvoir être exploitées par PySpark. Cette étape assure la liaison entre la base de données NoSQL et l'outil Big Data, permettant ainsi de passer d'un stockage des données à une phase d'analyse.

```
[5] os
Partie 5 – Exportation du dataset vers PySpark

import pandas as pd

# Extraire tous les documents (sans projection)
docs = list(collection.find())

# Supprimer le champ _id qui pose problème (ajout automatique par MongoDB)
for doc in docs:
    if "_id" in doc:
        del doc["_id"]

# Exporter vers JSON
pd.DataFrame(docs).to_json("etudiants_export.json", orient="records")
print("Export vers etudiants_export.json terminé.")

Export vers etudiants_export.json terminé.
```

6 . Analyse des données avec PySpark

L'objectif de cette partie est d'analyser les données académiques à l'aide de PySpark. Après le chargement du fichier JSON dans un DataFrame, plusieurs opérations analytiques ont été réalisées, telles que l'explosion des tableaux de notes, le calcul des moyennes par filière et l'évaluation des performances globales des étudiants.

```
[6] 20s    from pyspark.sql import SparkSession
          from pyspark.sql.functions import explode, col, avg

        spark = SparkSession.builder.appName("MiniProjet_Etudiants").getOrCreate()
        df = spark.read.json("etudiants_export.json")
        df.printSchema()
        df.show(truncate=False)

    root
    |-- filiere: string (nullable = true)
    |-- id: long (nullable = true)
    |-- nom: string (nullable = true)
    |-- notes: array (nullable = true)
    |   |-- element: struct (containsNull = true)
    |   |   |-- matiere: string (nullable = true)
    |   |   |-- note: long (nullable = true)

    +-----+-----+-----+
    |filiere|id |nom  |notes
    +-----+-----+-----+
    |Informatique|1 |Amina |[[{Big Data, 15}, {Programmation, 17}, {Mathématiques, 14}]]
    |Réseaux|2 |Youssef|[{{Systèmes, 12}, {Réseaux, 14}, {Sécurité, 13}}]
    |Informatique|3 |Sara  |[{Big Data, 18}, {IA, 17}, {Programmation, 16}]
    +-----+-----+-----+
```



```
[7] 3s    df_explode = df.withColumn("note", explode(col("notes")))
          result = df_explode.groupBy("filiere").agg(avg(col("note.note")).alias("moyenne_filiere"))
          result.show()

    +-----+-----+
    |filiere| moyenne_filiere|
    +-----+-----+
    |Informatique|16.166666666666668|
    | Réseaux| 13.0|
    +-----+-----+
```

7 .Manipulation, analyse et requêtes SQL

Cette partie constitue l'étape finale du mini-projet. Elle regroupe des opérations avancées de manipulation des données avec Mongita, des analyses détaillées avec PySpark ainsi que l'utilisation de SparkSQL pour interroger les données de manière déclarative. Cette partie permet de consolider l'ensemble des compétences acquises tout au long du projet et de démontrer l'intérêt combiné du NoSQL et du Big Data dans l'analyse des données académiques.

Partie 7 – Travail à faire

--- A.1 Ajouter 2 nouveaux étudiants et réinsérer ---

```
[10]  Os
    import json
    from pymongo import MongoClientDisk

    # connexion Mongita (même DB / collection que tol)
    client = MongoClientDisk()
    db = client['enseignement_db']
    collection = db['etudiants']

    # Nouveaux étudiants à ajouter
    nouveaux = [
        {
            "id": 4,
            "nom": "Khalid",
            "filiere": "Informatique",
            "notes": [
                {"matiere": "Big Data", "note": 14},
                {"matiere": "Programmation", "note": 15},
                {"matiere": "Mathématiques", "note": 13}
            ]
        },
        {
            "id": 5,
            "nom": "Nadia",
            "filiere": "Réseaux",
            "notes": [
                {"matiere": "Systèmes", "note": 16},
                {"matiere": "Réseaux", "note": 17},
                {"matiere": "Sécurité", "note": 15}
            ]
        }
    ]

    # Charger le fichier existant et ajouter les nouveaux (ou reconstruire le fichier)
    with open("etudiants.json", "r") as f:
        data = json.load(f)

    # Ajouter les nouveaux (évite les doublons d'id)
    ids_existantes = {d["id"] for d in data}
    for nd in nouveaux:
        if nd["id"] not in ids_existantes:
            data.append(nd)

    # Écrire le fichier mis à jour
    with open("etudiants.json", "w") as f:
        json.dump(data, f, indent=2, ensure_ascii=False)

    # Réinsertion dans Mongita : on nettoie et on insert tout
    collection.delete_many({})
    collection.insert_many(data)
    print("Insertion des documents (avec 2 nouveaux) effectuée. Total inséré :", collection.count_documents({}))
```

Insertion des documents (avec 2 nouveaux) effectuée. Total inséré : 5

--- A.2 Afficher tous les étudiants insérés ---

```
[11]  Os
    print("Tous les étudiants :")
    for doc in collection.find():
        print(doc)
```

Tous les étudiants :

```
{'id': 1, 'nom': 'Amina', 'filiere': 'Informatique', 'notes': [{'matiere': 'Big Data', 'note': 15}, {'matiere': 'Programmation', 'note': 17}, {'matiere': 'Mathématiques', 'note': 14}], '_id': ObjectId('693035a7e09b4b3f2ecbc64')}, {'id': 2, 'nom': 'Youssef', 'filiere': 'Réseaux', 'notes': [{'matiere': 'Systèmes', 'note': 12}, {'matiere': 'Réseaux', 'note': 14}, {'matiere': 'Sécurité', 'note': 13}], '_id': ObjectId('693035a7e09b4b3f2ecbc64')}, {'id': 3, 'nom': 'Sara', 'filiere': 'Informatique', 'notes': [{'matiere': 'Big Data', 'note': 18}, {'matiere': 'IA', 'note': 17}, {'matiere': 'Programmation', 'note': 16}], '_id': ObjectId('693035a7e09b4b3f2ecbc65')}, {'id': 4, 'nom': 'Khalid', 'filiere': 'Informatique', 'notes': [{'matiere': 'Big Data', 'note': 14}, {'matiere': 'Programmation', 'note': 15}, {'matiere': 'Mathématiques', 'note': 13}], '_id': ObjectId('693035a7e09b4b3f2ecbc65')}, {'id': 5, 'nom': 'Nadia', 'filiere': 'Réseau', 'notes': [{'matiere': 'Systèmes', 'note': 16}, {'matiere': 'Réseau', 'note': 17}, {'matiere': 'Sécurité', 'note': 15}], '_id': ObjectId('693035a7e09b4b3f2ecbc67')}
```

--- A.3 Afficher seulement noms et filières (filtrage en Python) ---

```
[11]  Os
    print("\nNoms et filières :")
    for doc in collection.find():
        print("{'nom': doc.get('nom'), 'filiere': doc.get('filiere')}")
```

Noms et filières :

```
{'nom': 'Amina', 'filiere': 'Informatique'}
{'nom': 'Youssef', 'filiere': 'Réseaux'}
{'nom': 'Sara', 'filiere': 'Informatique'}
{'nom': 'Khalid', 'filiere': 'Informatique'}
{'nom': 'Nadia', 'filiere': 'Réseaux'}
```

--- A.4 Rechercher étudiants de la filière "Réseaux" ---

```
[12]  Os
    print("\nétudiants en Réseaux :")
    for doc in collection.find({"filiere": "Réseaux"}):
        print(doc)
```

Étudiants en Réseaux :

```
{'id': 2, 'nom': 'Youssef', 'filiere': 'Réseaux', 'notes': [{"matiere': 'Systèmes', 'note': 12}, {"matiere': 'Réseaux', 'note': 14}, {"matiere': 'Sécurité', 'note': 13}], '_id': ObjectId('693035a7e09b4b3f2ecbc64')}, {'id': 5, 'nom': 'Nadia', 'filiere': 'Réseaux', 'notes': [{"matiere': 'Systèmes', 'note': 16}, {"matiere': 'Réseaux', 'note': 17}, {"matiere': 'Sécurité', 'note': 15}], '_id': ObjectId('693035a7e09b4b3f2ecbc67')}
```

✓ --- A.5 Étudiants dont la note en "Programmation" > 15 ---

```
[13] 0s print("\nétudiants avec note > 15 en Programmation :")
for doc in collection.find():
    for note in doc.get("notes", []):
        if note["matiere"] == "Programmation" and note["note"] > 15:
            print(f"({doc['nom']} : {note['note']})")
```

Étudiants avec note > 15 en Programmation :
Amina : 17
Sara : 16

✓ --- A.6 Modifier la note de "Big Data" d'un étudiant (+1 point) et afficher ---

```
[14] 0s # Exemple : augmenter Big Data de Sara (id 3)
doc = collection.find_one({"id": 3})
if doc:
    # modifier dans la structure Python
    for note in doc["notes"]:
        if note["matiere"] == "Big Data":
            note["note"] = note["note"] + 1 # +1 point
    # Mettre à jour le document complet (remplacer)
collection.replace_one({"id": doc["id"]}, doc)
print(f"Nouvelle note Big Data pour {doc['nom']} :",
      next(n["note"] for n in doc["notes"] if n["matiere"]=="Big Data"))
else:
    print("Étudiant id 3 introuvable.")
```

✓ --- A.7 Supprimer un étudiant de test ajouté par erreur ---

```
[15] 0s # Exemple : supprimer l'étudiant de test ayant id = 999 si existe
res = collection.delete_one({"id": 999})
print("Suppression (id 999) - documents supprimés :", res.deleted_count)
```

Suppression (id 999) - documents supprimés : 0

✓ --- A.8 Compter le nombre total d'étudiants ---

```
[16] 0s total = collection.count_documents({})
print("Nombre total d'étudiants :", total)
```

Nombre total d'étudiants : 5

✓ --- A.9 Compter le nombre d'étudiants par filière (dictionnaire Python) ---

```
[17] 0s from collections import defaultdict
compte = defaultdict(int)
for doc in collection.find():
    compte[doc["filiere"]] += 1
print("Nombre d'étudiants par filière :", dict(compte))
```

Nombre d'étudiants par filière : {'Informatique': 3, 'Réseaux': 2}

--- A.10 Trouver l'étudiant ayant la meilleure moyenne générale (calcul en Python) --

```
[18] 0s
def moyenne(doc):
    notes = [n["note"] for n in doc.get("notes", [])]
    return sum(notes)/len(notes) if notes else 0

meilleure = None
best_avg = -1
for doc in collection.find():
    avg = moyenne(doc)
    if avg > best_avg:
        best_avg = avg
        meilleure = doc

if meilleure:
    print("Étudiant avec la meilleure moyenne :", meilleure["nom"], "->", round(best_avg, 2))
else:
    print("Aucun étudiant trouvé.")

Étudiant avec la meilleure moyenne : Sara -> 17.33
```

✓ B. Analyse et transformation (PySpark)

```
[19] 0s
# Réexporter vers etudiants_export.json (supprimer _id si présent)
docs = list(collection.find())
for d in docs:
    if "_id" in d:
        del d["_id"]

import pandas as pd
pd.DataFrame(docs).to_json("etudiants_export.json", orient="records", force_ascii=False)
print("Export vers etudiants_export.json terminé. Documents exportés : ", len(docs))
```

Export vers etudiants_export.json terminé. Documents exportés : 5

--- B.1 Charger et afficher schéma + types ---

```
[20] 16
from pyspark.sql import SparkSession
from pyspark.sql.functions import explode, col, avg, round as spark_round

spark = SparkSession.builder.appName("MiniProjet_Etudiants_Partie7").getOrCreate()
df = spark.read.json("etudiants_export.json")

print("Schéma du DataFrame :")
df.printSchema()
df.show(truncate=False)

Schéma du DataFrame :
root
 |-- filiere: string (nullable = true)
 |-- id: long (nullable = true)
 |-- nom: string (nullable = true)
 |-- notes: array (nullable = true)
 |   |-- element: struct (containsNull = true)
 |   |   |-- matiere: string (nullable = true)
 |   |   |-- note: long (nullable = true)

+-----+-----+
|filiere|id    |nom   |notes
+-----+-----+
|Informatique|1|Amina |[{Big Data, 15}, {Programmation, 17}, {Mathématiques, 14}]
|Réseaux     |2|Youssef|[Systèmes, 12], {Réseaux, 14}, {Sécurité, 13}]
|Informatique|3|Sara   |[{Big Data, 19}, {IA, 17}, {Programmation, 16}]
|Informatique|4|Khald  |[{Big Data, 14}, {Programmation, 15}, {Mathématiques, 13}]
|Réseaux     |5|Nadia  |[{Systèmes, 16}, {Réseaux, 17}, {Sécurité, 15}]
+-----+-----+
```

--- B.2 / B.3 Créer df_explode correctement (matière, note) ---

```
[22] 1s
# --- B.2 / B.3 Créer df_explode correctement (matière, note) ---
from pyspark.sql.functions import explode

# explode renvoie une struct; on l'extract ensuite en colonnes matière/note
df_explode_struct = df.select("filiere", explode(col("notes")).alias("note_struct"))
df_explode = df_explode_struct.select(
    "filiere",
    col("note_struct.matiere").alias("matiere"),
    col("note_struct.note").alias("note")
)

print("Aperçu df_explode :")
df_explode.show(truncate=False)

Aperçu df_explode :
+-----+-----+
|filiere|matiere|note|
+-----+-----+
|Informatique|Big Data|15 |
|Informatique|Programmation|17 |
|Informatique|Mathématiques|14 |
|Réseaux|Systèmes|12 |
|Réseaux|Réseaux|14 |
|Réseaux|Sécurité|13 |
|Informatique|Big Data|19 |
|Informatique|IA|17 |
|Informatique|Programmation|16 |
|Informatique|Big Data|14 |
|Informatique|Programmation|15 |
|Informatique|Mathématiques|13 |
|Réseaux|Systèmes|16 |
|Réseaux|Réseaux|17 |
|Réseaux|Sécurité|15 |
+-----+-----+
```

--- B.4 Calculer la moyenne des notes par filière ---

```
[23] 2s
moy_par_filiere = df_explode.groupby("filiere").agg(spark_round(avg("note"), 2).alias("moyenne_filiere"))
moy_par_filiere.show()

+-----+
| filiere|moyenne_filiere|
+-----+
|Informatique|      15.56|
|Réseaux|        14.5|
+-----+
```

--- B.5 Calculer la moyenne des notes par matière ---

```
[24] 1s
moy_par_matiere = df_explode.groupby("matiere").agg(spark_round(avg("note"), 2).alias("moyenne_matiere"))
moy_par_matiere.orderBy(col("moyenne_matiere").desc()).show(truncate=False)

+-----+-----+
|matiere|moyenne_matiere|
+-----+
|IA|17.0
|Big Data|16.0
|Programmation|16.0
|Réseaux|15.5
|Systèmes|14.0
|Sécurité|14.0
|Mathématiques|13.5
+-----+
```

--- B.6 Identifier la filière ayant la meilleure moyenne ---

```
[25] 0s
    best_filiere = moy_par_filiere.orderBy(col("moyenne_filiere").desc()).limit(1)
    print("Filière avec la meilleure moyenne :")
    best_filiere.show()
```

Filière	moyenne_filiere
Informatique	15.56

--- B.7 Identifier la matière la mieux réussie ---

```
[26] 1s
    best_matiere = moy_par_matiere.orderBy(col("moyenne_matiere").desc()).limit(1)
    print("Matière la mieux réussie :")
    best_matiere.show()
```

Matière	moyenne_matiere
IA	17.0

--- B.8 Moyenne générale (toutes filières) ---

```
[27] 0s
    moyenne_generale = df_explode.agg(spark_round(avg("note"),2).alias("moyenne_generale"))
    moyenne_generale.show()
```

moyenne_generale
15.13

--- B.9 Ajouter une colonne "mention" selon la moyenne individuelle ---

```
[28] 3s
    # ===== B.9 Ajouter une colonne "mention" (CODE UNIQUE PRÊT À COLLER) =====
    from pyspark.sql.functions import explode, col, avg as spark_avg, round as spark_round, when

    # 1) Explode
    df_explode = df.select(
        "id",
        "nom",
        "filiere",
        explode(col("notes")).alias("note_struct")
    ).select(
        "id",
        "nom",
        "filiere",
        col("note_struct.matiere").alias("matiere"),
        ...
    )

    # 2) Calcul de la moyenne par étudiant
    moy_par_etudiant = df_explode.groupby(
        "id", "nom"
    ).agg(
        spark_round(spark_avg("note"), 2).alias("moyenne")
    )

    # 3) Ajout de la colonne "mention"
    resultat_final = moy_par_etudiant.withColumn(
        "mention",
        when(col("moyenne") > 16, "Excellent")
            .when((col("moyenne") >= 14) & (col("moyenne") <= 16), "Bien")
            .when((col("moyenne") >= 12) & (col("moyenne") < 14), "Assez-bien")
            .otherwise("Passable")
    )

    # 4) Affichage final
    resultat_final.orderBy(col("moyenne").desc()).show(truncate=False)
```

--- C.1 Créer table temporaire ---

```
[31] 0s # Utiliser df (ou df_explode selon la requête)
df.createOrReplaceTempView("etudiants_sql")
df_explode.createOrReplaceTempView("etudiants_notes_sql")
```

--- C.21 Requête SQL : afficher les étudiants de la filière "Informatique" ---

```
[32] 0s q1 = spark.sql("""
SELECT id, nom, filiere
FROM etudiants_sql
WHERE filiere = 'Informatique'
ORDER BY id
""")
q1.show()
```

id	nom	filiere
1	Amina	Informatique
3	Sara	Informatique
4	Khalid	Informatique

--- C.22 Requête SQL : moyenne des notes par matière ---

```
[33] 1s q2 = spark.sql("""
SELECT matiere, ROUND(AVG(note),2) as moyenne_matiere
FROM etudiants_notes_sql
GROUP BY matiere
ORDER BY moyenne_matiere DESC
""")
q2.show(truncate=False)
```

matiere	moyenne_matiere
IA	17.0
Big Data	16.0
Programmation	16.0
Réseaux	15.5
Systèmes	14.0
Sécurité	14.0
Mathématiques	13.5

--- C.23 Requête SQL : filière ayant la meilleure moyenne ---

```
[34] 0s q3 = spark.sql("""
SELECT filiere, ROUND(AVG(note),2) as moyenne_filiere
FROM etudiants_notes_sql
GROUP BY filiere
ORDER BY moyenne_filiere DESC
LIMIT 1
""")
q3.show()
```

filiere	moyenne_filiere
Informatique	15.56

--- C.24 Requête SQL : compter le nombre total d'étudiants dans chaque filière ---

```
[35] 1s q4 = spark.sql("""
SELECT filiere, COUNT(*) as nombre_etudiants
FROM etudiants_sql
GROUP BY filiere
ORDER BY nombre_etudiants DESC
""")
q4.show()
```

filiere	nombre_etudiants
Informatique	3
Réseaux	2

Conclusion

Ce mini-projet a permis de mettre en pratique les notions fondamentales des bases de données NoSQL et du Big Data dans un contexte concret lié au domaine éducatif. L'utilisation de MongoDB (via Mongita) a montré la flexibilité et la simplicité des bases orientées documents pour le stockage des données semi-structurées.

Par ailleurs, l'exploitation de PySpark et SparkSQL a mis en évidence la puissance des outils Big Data pour effectuer des analyses rapides et efficaces sur des volumes importants de données. Ce projet a également renforcé notre compréhension des opérations de transformation, d'agrégation et d'analyse des données.

En conclusion, les technologies NoSQL et Big Data représentent aujourd'hui des solutions indispensables pour la gestion et l'analyse des données académiques à grande échelle, offrant performance, flexibilité et scalabilité.