

# Homework Lab

## Introduction to Big Data

Jie FAN  
Minh Quang HOANG

## Sommaire

<b>Introduction.....</b>	<b>1</b>
Contexte.....	1
Objectif.....	1
<b>Sources de données.....</b>	<b>2</b>
Description datasets.....	2
Origine données.....	5
<b>Architecture et gestion des données.....</b>	<b>6</b>
Schéma d'architecture.....	6
Modèle dimensionnel.....	7
<b>Traitement des données.....</b>	<b>8</b>
Infrastructures et Ressources.....	8
Étapes d'Implémentation.....	10
Étape nettoyage des données.....	12
1. Initialisation de la session Spark.....	12
2. Chargement des tables sources.....	12
3. Création des tables dimensionnels.....	12
4. Création de la table faits.....	14
Résultats du nettoyage.....	14
<b>Analyse des résultats.....</b>	<b>16</b>
Visualisation des données.....	16
1. Distribution des notes d'examen selon le genre.....	16
2. Corrélation entre l'activité physique et la note d'examen.....	17
3. Impact du niveau de bruit sur le taux de présence.....	18
4. Carte de corrélation pour plusieurs variables.....	19
<b>Conclusion.....</b>	<b>19</b>
<b>Annexes.....</b>	<b>20</b>
Code et scripts : lien vers repo Github.....	20

# Introduction

## Contexte

Dans le cadre académique, les étudiants sont confrontés à de multiples pressions, notamment des attentes élevées en matière de performance scolaire. Ces exigences, associées à des facteurs comme les troubles de santé mentale, les habitudes de sommeil, et les niveaux de stress, peuvent avoir un impact significatif sur leur bien-être et leur réussite académique.

Pour mieux comprendre ces interrelations, notre projet s'appuie sur des données réelles pour examiner comment le stress, les troubles mentaux, et d'autres facteurs liés à la vie étudiante influencent les performances scolaires. En analysant des aspects tels que le sommeil, les niveaux de stress, et les diagnostics de santé mentale, nous cherchons à proposer des stratégies pour améliorer la santé mentale et les résultats scolaires.

## Objectif

L'objectif principal de cette étude est d'analyser et de comprendre l'impact de la santé mentale et des facteurs de stress sur les performances académiques des étudiants. En exploitant les données disponibles, nous visons à identifier les schémas et corrélations clés pour proposer des solutions adaptées.

Les sous-objectifs incluent :

- Analyser les niveaux de stress et leurs relations avec les performances scolaires.
- Évaluer l'impact des troubles mentaux et des habitudes de sommeil sur la concentration et les résultats.
- Fournir des recommandations pour réduire les niveaux de stress et améliorer la santé mentale des étudiants

# Sources de données

## Description datasets

Voici un aperçu des colonnes présentes dans chaque fichier, avec une explication de leurs significations et des valeurs qu'elles contiennent :

### 1. StressLevelDataset.csv

Ce fichier contient des informations liées au stress, aux conditions de vie et à d'autres facteurs qui peuvent influencer la santé mentale. Voici les colonnes et leur signification :

- **anxiety\_level** : Niveau d'anxiété sur une échelle de 1 à 20 (valeurs plus élevées indiquent un niveau d'anxiété plus important).
- **self\_esteem** : Niveau d'estime de soi sur une échelle de 1 à 30 (valeurs plus élevées indiquent une estime de soi plus élevée).
- **mental\_health\_history** : Historique de troubles mentaux (0 = aucun historique, 1 = historique présent).
- **depression** : Niveau de dépression sur une échelle de 1 à 15.
- **headache** : Fréquence des maux de tête (1 = rare, 5 = très fréquent).
- **blood\_pressure** : Niveau de tension artérielle (1 = normal, 5 = très élevé).
- **sleep\_quality** : Qualité du sommeil sur une échelle de 1 à 5 (1 = très mauvaise, 5 = excellente).
- **breathing\_problem** : Problèmes respiratoires perçus (1 = aucun, 5 = sévère).
- **noise\_level** : Niveau de bruit dans l'environnement (1 = très calme, 5 = très bruyant).
- **living\_conditions** : Conditions de vie générales (1 = très mauvaises, 5 = excellentes).
- **basic\_needs** : Niveau de satisfaction des besoins de base (1 = non satisfait, 5 = entièrement satisfait).
- **academic\_performance** : Performance académique (1 = très faible, 5 = très élevée).
- **study\_load** : Charge d'étude perçue (1 = très légère, 5 = très lourde).
- **teacher\_student\_relationship** : Qualité de la relation enseignant-étudiant (1 = mauvaise, 5 = excellente).
- **future\_career\_concerns** : Niveau de préoccupation concernant la carrière future (1 = aucune, 5 = très préoccupé).
- **social\_support** : Niveau de soutien social perçu (1 = aucun, 5 = très élevé).
- **peer\_pressure** : Pression des pairs (1 = très faible, 5 = très forte).
- **extracurricular\_activities** : Niveau d'implication dans les activités extrascolaires (1 = aucune, 5 = très élevée).
- **bullying** : Fréquence du harcèlement (1 = jamais, 5 = très fréquent).

- **stress\_level** : Niveau global de stress sur une échelle de 1 à 10 (1 = très faible, 10 = très élevé).

## 2. Student Mental health.csv

Ce fichier documente les troubles mentaux déclarés et quelques informations démographiques. Voici les colonnes et leur signification :

- **Timestamp** : Heure et date de la réponse au questionnaire.
- **Choose your gender** : Genre de l'étudiant (Male = masculin, Female = féminin, Other = autre).
- **Age** : Âge de l'étudiant en années.
- **What is your course?** : Programme d'études suivi par l'étudiant.
- **Your current year of Study** : Année actuelle d'études (exemple : Year 1, Year 2).
- **What is your CGPA?** : Moyenne générale cumulée (exemple : 3.00 - 3.49).
- **Marital status** : Statut marital de l'étudiant (Yes = marié, No = non marié).
- **Do you have Depression?** : L'étudiant déclare-t-il souffrir de dépression (Yes = oui, No = non).
- **Do you have Anxiety?** : L'étudiant déclare-t-il souffrir d'anxiété (Yes = oui, No = non).
- **Do you have Panic attack?** : L'étudiant déclare-t-il avoir des crises de panique (Yes = oui, No = non).
- **Did you seek any specialist for a treatment?** : L'étudiant a-t-il consulté un spécialiste (Yes = oui, No = non).

## 3. Student\_sleep\_patterns.csv

Ce fichier explore les habitudes de sommeil des étudiants et leurs modes de vie. Voici les colonnes et leur signification :

- **Student\_ID** : Identifiant unique de chaque étudiant.
- **Age** : Âge de l'étudiant.
- **Gender** : Genre de l'étudiant.
- **University\_Year** : Année d'université (exemple : 1st Year, 2nd Year).
- **Sleep\_Duration** : Nombre moyen d'heures de sommeil par jour.
- **Study\_Hours** : Heures consacrées à l'étude par jour.
- **Screen\_Time** : Temps moyen passé devant un écran par jour (en heures).
- **Caffeine\_Intake** : Quantité de consommation de caféine (échelle de 1 à 5, 1 = faible, 5 = élevée).
- **Physical\_Activity** : Durée ou fréquence d'activité physique (mesurée en minutes ou sessions).

- **Sleep\_Quality** : Qualité perçue du sommeil (échelle de 1 à 10).
- **Weekday\_Sleep\_Start** : Heure moyenne de coucher en semaine (exemple : 22.30 = 22h30).
- **Weekend\_Sleep\_Start** : Heure moyenne de coucher le week-end.
- **Weekday\_Sleep\_End** : Heure moyenne de réveil en semaine.
- **Weekend\_Sleep\_End** : Heure moyenne de réveil le week-end.

#### 4. StudentPerformanceFactors.csv

Ce fichier regroupe des facteurs influençant les performances académiques. Voici les colonnes et leur signification :

- **Hours\_Studied** : Heures consacrées à l'étude chaque semaine.
- **Attendance** : Taux de présence en classe (en pourcentage).
- **Parental\_Involvement** : Niveau d'implication des parents (Low = faible, Medium = moyen, High = élevé).
- **Access\_to\_Resources** : Accès aux ressources éducatives (Low, Medium, High).
- **Extracurricular\_Activities** : Participation à des activités extrascolaires (Yes = oui, No = non).
- **Sleep\_Hours** : Nombre moyen d'heures de sommeil par nuit.
- **Previous\_Scores** : Notes obtenues précédemment (en pourcentage).
- **Motivation\_Level** : Niveau de motivation (Low, Medium, High).
- **Internet\_Access** : Accès à Internet (Yes = oui, No = non).
- **Tutoring\_Sessions** : Nombre de sessions de tutorat suivies.
- **Family\_Income** : Niveau de revenu familial (Low = faible, Medium = moyen, High = élevé).
- **Teacher\_Quality** : Qualité perçue des enseignants (Low, Medium, High).
- **School\_Type** : Type d'établissement scolaire (Public = public, Private = privé).
- **Peer\_Influence** : Influence des pairs (Positive, Neutral, Negative).
- **Physical\_Activity** : Fréquence ou durée de l'activité physique (en unités).
- **Learning\_Disabilities** : Présence de troubles d'apprentissage (Yes = oui, No = non).
- **Parental\_Education\_Level** : Niveau d'éducation des parents (High School, College, Postgraduate).
- **Distance\_from\_Home** : Distance entre le domicile et l'établissement scolaire (Near, Moderate, Far).
- **Gender** : Genre de l'étudiant.
- **Exam\_Score** : Résultat à un examen (en pourcentage).

## Origine données

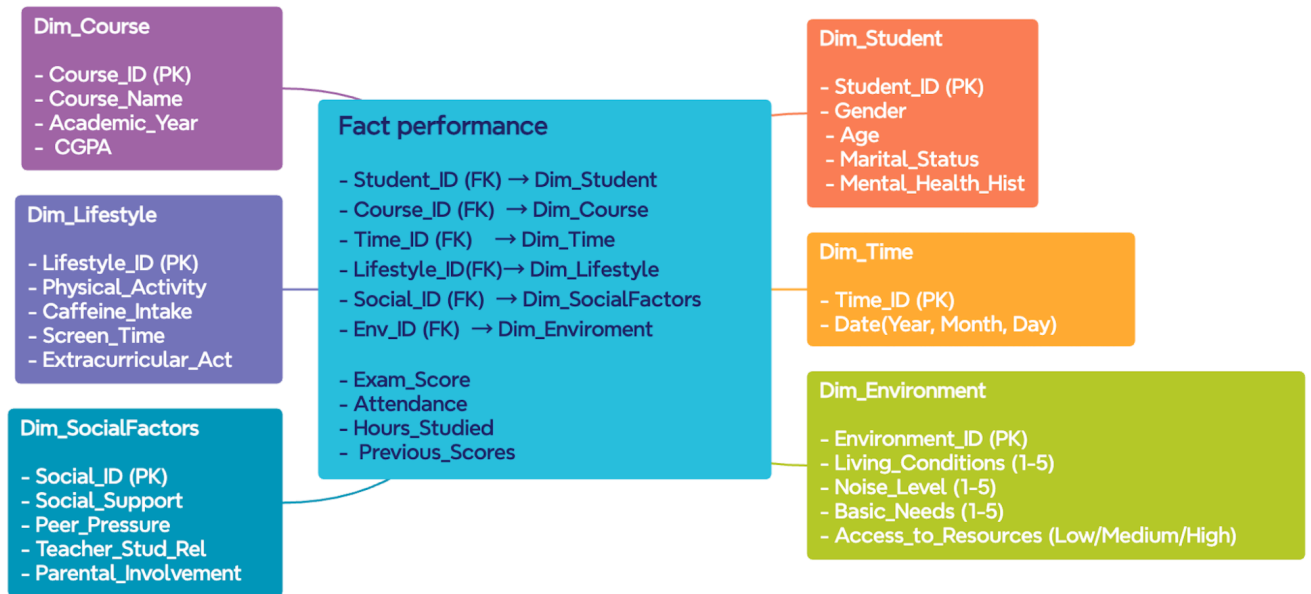
Les données utilisées dans ce projet proviennent de la plateforme Kaggle, un site web connu pour son large éventail de jeux de données publics, utilisés principalement pour des projets d'analyse de données et de visualisation. Les fichiers sélectionnés ont été choisis en raison de leur pertinence pour l'étude de la santé mentale et des performances académiques des étudiants.





# Modèle dimensionnel

schéma en étoile



# Traitement des données

## Infrastructures et Ressources

### Groupe de Ressources Azure et Compte de Stockage

Microsoft Azure

Rechercher dans les ressources, services et documents (G+)

Copilot

jie.fan@efrei.net

Accueil > projet-big-data-iceberg

Rechercher

Créer Gérer la vue Supprimer le groupe de ressources Actualiser Exporter au format CSV Ouvrir une requête Attribuer des étiquettes Déplacer

Vue d'ensemble

Journal d'activité

Contrôle d'accès (IAM)

Étiquettes

Visualiseur de ressources

Événements

Paramètres

Déploiements

Sécurité

Piles de déploiement

Stratégies

Propriétés

Verrous

Gestion des coûts

Supervision

Automatisation

Aide

Bases

Abonnement (déplacer) : Azure for Students

ID d'abonnement : 1bfa101-f6c7-4c8e-9954-d675f52540d1

Déploiements : 1 Réussite

Emplacement : West Europe

Étiquettes (modifier) : Ajouter des étiquettes

Ressources

Recommandations (3)

Filtrer un champ...

Type égal à tout

Emplacement égal à tout

Ajouter un filtre

Affichage de 1 à 1 sur 1 enregistrements.

Afficher les types masqués

Aucun regroupement

Vue liste

Nom	Type	Emplacement
iceberg2024stockage	Compte de stockage	West Europe

< Précédent Page 1 sur 1 Suivant >

Envoyer des commentaires

- **Groupe de ressources** : projet-big-data-iceberg
- **Compte de stockage** : iceberg2024stockage

### Conteneurs Blob Storage

Deux conteneurs principaux ont été créés :

Microsoft Azure

Rechercher dans les ressources, services et documents (G+)

Copilot

jie.fan@efrei.net

Accueil > projet-big-data-iceberg > iceberg2024stockage

iceberg2024stockage | Conteneurs

Rechercher

Conteneur Modifier le niveau d'accès Restaurer des conteneurs Actualiser Supprimer Envoyer des commentaires

Rechercher les conteneurs par préfixe

Afficher les conteneurs supprimés

Nom	Dernière modification	Niveau d'accès anonyme	État du bail
\$logs	26/11/2024 12:28:57	Privé	Disponible
iceberg-data	28/11/2024 17:10:46	Privé	Disponible

Vue d'ensemble

Journal d'activité

Étiquettes

Diagnostiquer et résoudre les problèmes

Contrôle d'accès (IAM)

Migration des données

Événements

Navigateur de stockage

Storage Mover

Solutions de partenaire

Stockage des données

Conteneurs

Partages de fichiers

Files d'attente

Tables

Sécurité + réseau

Mise en réseau

Front Door et CDN

Clés d'accès

1. **iceberg-data** : Contient les fichiers CSV bruts (Bronze) et les tables Iceberg (Silver).
2. **\$logs** : Contient les journaux de diagnostic.

### Contenu du conteneur **iceberg-data** :

Microsoft Azure

Rechercher dans les ressources, services et documents (G+/)

Copilot

je.fan@efrei.net

Accueil >

**iceberg-data**

Conteneur

Rechercher

Charger | Modifier le niveau d'accès | Actualiser | Supprimer | Modifier le niveau | Acquérir le bail | Résilier le bail | Afficher les instantanés

**Méthode d'authentification** : Clé d'accès (Basculer vers le compte d'utilisateur Microsoft Entra)

**Emplacement** : iceberg-data

Rechercher les objets blobs par préfixe (respect de la casse)

Ajouter un filtre

Nom	Modifié	Niveau d'accès	État de l'archive	Type d'objet blob	Taille	État du bail
<input type="checkbox"/> silver_warehouse						-
<input type="checkbox"/> silver_warehouse	12/12/2024 18:40:27	Élevé (déduit)		Objet blob de blocs	0 B	Disponible
<input type="checkbox"/> StressLevelDataset.csv	12/12/2024 17:59:00	Élevé (déduit)		Objet blob de blocs	47.58 KiB	Disponible
<input type="checkbox"/> Student Mental health.csv	12/12/2024 17:59:00	Élevé (déduit)		Objet blob de blocs	7.17 KiB	Disponible
<input type="checkbox"/> student_sleep_patterns.csv	12/12/2024 17:59:00	Élevé (déduit)		Objet blob de blocs	30.53 KiB	Disponible
<input type="checkbox"/> StudentPerformanceFactors.csv	12/12/2024 17:59:01	Élevé (déduit)		Objet blob de blocs	626.9 KiB	Disponible

- 4 Fichiers CSV originaux (bronze):
  - Student Mental health.csv,
  - StressLevelDataset.csv,
  - student\_sleep\_patterns.csv,
  - StudentPerformanceFactors.csv
- Répertoire **silver\_warehouse** : Contient les métadonnées des tables Iceberg (Silver).

Microsoft Azure

Rechercher dans les ressources, services et documents (G+/)

Copilot

je.fan@efrei.net

Accueil > projet-big-data-iceberg > iceberg2024stockage | Conteneurs >

**iceberg-data**

Conteneur

Rechercher

Charger | Modifier le niveau d'accès | Actualiser | Supprimer | Modifier le niveau | Acquérir le bail | Résilier le bail | Afficher les instantanés

**Méthode d'authentification** : Clé d'accès (Basculer vers le compte d'utilisateur Microsoft Entra)

**Emplacement** : iceberg-data / silver\_warehouse / silver

Rechercher les objets blobs par préfixe (respect de la casse)

Ajouter un filtre

Nom	Modifié	Niveau d'accès	État de l'archive	Type d'objet blob	Taille	État du bail
<input type="checkbox"/> [.]						-
<input type="checkbox"/> mental_health						-
<input type="checkbox"/> performance_factors						-
<input type="checkbox"/> sleep_patterns						-
<input type="checkbox"/> stress_levels						-
<input type="checkbox"/> mental_health	12/12/2024 18:40:27	Élevé (déduit)		Objet blob de blocs	0 B	Disponible
<input type="checkbox"/> performance_factors	12/12/2024 18:42:49	Élevé (déduit)		Objet blob de blocs	0 B	Disponible
<input type="checkbox"/> sleep_patterns	12/12/2024 18:42:01	Élevé (déduit)		Objet blob de blocs	0 B	Disponible
<input type="checkbox"/> stress_levels	12/12/2024 18:41:16	Élevé (déduit)		Objet blob de blocs	0 B	Disponible

# Étapes d'Implémentation

## Étape 1 : Montage dans Databricks

- Le conteneur **iceberg-data** a été monté sur le chemin **dbfs:/mnt/iceberg-data** dans Databricks. Cela permet d'accéder directement aux fichiers dans Databricks sans spécifier de chemin Azure.

```
▶ 06:34 PM (12s) 1 Python [ ] ⋮

# 1. info storage
storage_name = "iceberg2024stockage"
container_name = "iceberg-data"
access_key = "6u5KBSn4tf5V5ak1R0CGA2H6dKEojvQIy74085wkvmmKuyMA5UILRDw3f3kwRHL+nIV0ehuInAX9+ASt3LgZFw=="
mount_point_key = f"/mnt/iceberg-data"

# 2. Azure Blob
try:
    dbutils.fs.mount(
        source=f"wasbs://{container_name}@{storage_name}.blob.core.windows.net",
        mount_point=mount_point_key,
        extra_configs={f"fs.azure.account.key.{storage_name}.blob.core.windows.net": access_key}
    )
    print(f"Successfully mounted {container_name} to {mount_point_key}")
except Exception as e:
    print(f"Already mounted or failed to mount: {e}")

# 3. verify
try:
    files = dbutils.fs.ls(mount_point_key)
    print(f"Files in {mount_point_key}:")
    display(files)
except Exception as e:
    print(f"Failed to list files: {e}")
```

## Étape 2 : Transformation en Tables Iceberg (Silver)

- Les fichiers CSV bruts (Bronze) sont lus via Spark, nettoyés, puis stockés sous forme de tables Iceberg dans le répertoire `silver_warehouse`.

```

# set up Spark
spark.conf.set("spark.sql.catalog.silver_catalog", "org.apache.iceberg.spark.SparkCatalog")
spark.conf.set("spark.sql.catalog.silver_catalog.type", "hadoop")
spark.conf.set("spark.sql.catalog.silver_catalog.warehouse", "dbfs:/mnt/iceberg-data/silver_warehouse")

# Read the CSV files from the already mounted DBFS path
df_mental_health = (
    spark.read
    .format("csv")
    .option("header", "true")
    .option("inferSchema", "true")
    .load("dbfs:/mnt/iceberg-data/Student Mental health.csv")
)

# Read the CSV files from the already mounted DBFS path
df_mental_health.writeTo("silver_catalog.silver.mental_health").createOrReplace()

# The above steps will create the corresponding Iceberg table metadata and data files in the specified warehouse path (dbfs:/mnt/iceberg-data/silver_warehouse)

df_stress_levels = (
    spark.read
    .format("csv")
    .option("header", "true")
    .option("inferSchema", "true")
    .load("dbfs:/mnt/iceberg-data/StressLevelDataset.csv")
)
df_stress_levels.writeTo("silver_catalog.silver.stress_levels").createOrReplace()

df_sleep_patterns = (
    spark.read
    .format("csv")
    .option("header", "true")
    .option("inferSchema", "true")
    .load("dbfs:/mnt/iceberg-data/student_sleep_patterns.csv")
)
df_sleep_patterns.writeTo("silver_catalog.silver.sleep_patterns").createOrReplace()

df_performance = (
    spark.read
    .format("csv")
    .option("header", "true")
    .option("inferSchema", "true")
    .load("dbfs:/mnt/iceberg-data/StudentPerformanceFactors.csv")
)
df_performance.writeTo("silver_catalog.silver.performance_factors").createOrReplace()

```

### Étape 3 : Vérification des Tables Iceberg

```
1-iceberg-format-open-table Python %  
File Edit View Run Help Last edit was 24 days ago
```

```
> ✓ 12/30/2024 (Sat) 2  
  
%sql  
SHOW TABLES IN silver_catalog.silver; /* verify tables */  
  
+ [ ] pyspark.sql.dataframe.DataFrame = (database: string, tableName: string ... 1 more field)  
  
Table +  
┌─ database ─┐ ┌─ tableName ─┐ ┌─ IsTemporary ─┐  
1 | silver      | sleep_patterns | false          |  
2 | silver      | mental_health  | false          |  
3 | silver      | stress_levels  | false          |  
4 | silver      | performance_facto... | false          |  
  
↓ 4 rows | 4.87 seconds runtime  
Refreshed 23 days ago
```

```
> ✓ 12/30/2024 (Thu) 3  
  
%sql  
DESCRIBE TABLE silver_catalog.silver.mental_health; /* verify tables */  
  
SELECT * FROM silver_catalog.silver.mental_health LIMIT 10;  
  
+ [ ] Spark Jobs  
+ [ ] pyspark.sql.dataframe.DataFrame = [Timestamp: string, Choose your gender: string ... 9 more fields]  
  
Table +  
┌─ Timestamp ─┐ ┌─ Choose your gender ─┐ ┌─ Age ─┐ ┌─ What is your course? ─┐ ┌─ Your current year of Study ─┐ ┌─ What is your CGPA? ─┐ ┌─ Marital status ─┐ ┌─ Do you have Depression? ─┐ ┌─ Do you have Anxiety? ─┐  
1 | 8/7/2020 12:02 | Female              | 18    | Engineering             | year 1                | 3.00 - 3.49        | No                     | Yes                    | No                    |  
2 | 8/7/2020 12:04 | Male                | 21    | Islamic education       | year 2                | 3.00 - 3.49        | No                     | No                     | No                    |  
3 | 8/7/2020 12:05 | Male               | 19    | BIT                    | Year 1                | 3.00 - 3.49        | No                     | Yes                    | Yes                   |  
4 | 8/7/2020 12:08 | Female              | 22    | Laws                   | year 3                | 3.00 - 3.49        | Yes                    | Yes                    | No                    |  
5 | 8/7/2020 12:13 | Male                | 23    | Mathematics            | year 4                | 3.00 - 3.49        | No                     | No                     | No                    |  
6 | 8/7/2020 12:31 | Male                | 19    | Engineering             | Year 2                | 3.50 - 4.00        | No                     | No                     | No                    |  
7 | 8/7/2020 12:32 | Female              | 23    | Pendidikan Islam       | year 2                | 3.50 - 4.00        | Yes                    | Yes                    | No                    |  
8 | 8/7/2020 12:33 | Female              | 18    | BCS                    | year 1                | 3.50 - 4.00        | No                     | No                     | Yes                   |  
9 | 8/7/2020 12:35 | Female              | 19    | Human Resources         | Year 2                | 2.50 - 2.99        | No                     | No                     | No                    |  
  
↓ 10 rows | 12.16 seconds runtime  
Refreshed 23 days ago
```

# Étape nettoyage des données

Dans cette section, nous détaillons les étapes de nettoyage des données effectuées dans le cadre de notre projet. Ce processus est crucial pour garantir la qualité des données en vue de leur exploitation dans la construction des tables de dimensions et de faits.

## 1. Initialisation de la session Spark

```
spark = SparkSession.builder \
    .appName("Data Cleaning and Dimension/Fact Table Creation") \
    .config("spark.sql.catalog.silver_catalog", "org.apache.iceberg.spark.SparkCatalog") \
    .config("spark.sql.catalog.silver_catalog.type", "hadoop") \
    .config("spark.sql.catalog.silver_catalog.warehouse",
"dbfs:/mnt/iceberg-data/silver_warehouse") \
    .getOrCreate()
```

- **Objectif** : Configurer et initialiser une session Spark.
- **Explication** : Cette étape initialise une session Spark avec des configurations spécifiques pour utiliser le catalogue "Iceberg". Cela nous permet de lire et d'écrire des tables directement dans le *data lake* situé dans le répertoire `dbfs:/mnt/iceberg-data/silver_warehouse`.

## 2. Chargement des tables sources

```
df_mental_health = spark.table("silver_catalog.silver.mental_health")
df_sleep_patterns = spark.table("silver_catalog.silver.sleep_patterns")
df_stress_levels = spark.table("silver_catalog.silver.stress_levels")
df_performance_factors = spark.table("silver_catalog.silver.performance_factors")
```

- **Objectif** : Charger les données brutes depuis le catalogue "Iceberg".
- **Explication** : Quatre tables (santé mentale, habitudes de sommeil, niveaux de stress, et facteurs de performance) sont chargées dans des *DataFrames* Spark. Ces données constituent les sources pour les étapes ultérieures.

## 3. Création des tables dimensionnels

### *Ajout d'un identifiant synthétique*

Dans chaque table de dimension, on conserve deux références au même identifiant :

- **Synthetic\_ID** (le même que dans les données source/fact) pour pouvoir effectuer des jointures (JOIN) précises avec la table de faits.
- `col("Synthetic_ID").alias("XXX_ID")` (par ex. `Student_ID`, `Course_ID`) pour définir la clé primaire de la dimension et respecter la logique du modèle en étoile (Star Schema).

Ainsi, `Synthetic_ID` est la référence de correspondance avec la table de faits, tandis que `XXX_ID` est la clé propre à la dimension.

### *Nettoyage et standardisation des colonnes*

Les transformations suivantes ont été appliquées au niveau des tables pour harmoniser les colonnes :

- **Renommage des colonnes** : Par exemple, `Choose your gender` devient `Gender`, `What is your course?` devient `Course_Name`, etc. Cela permet de simplifier la lecture des données et d'uniformiser les noms des colonnes.
- **Conversion des types de données** : Pour certaines colonnes comme les dates (`Timestamp`), une conversion explicite en format `yyyy-MM-dd` a été effectuée pour les rendre compatibles avec les outils d'analyse.

```

# Create Dimension Tables

# Dim_Student
# Perform the join using the synthetic key
dim_student = (
    df_mental_health_with_id
    .join(df_stress_levels_with_id, on="Synthetic_ID", how="inner")
    .select(
        # ID for join in Fact table
        col("Synthetic_ID"),

        # primary key
        col("Synthetic_ID").alias("Student_ID"),

        col("Choose your gender").alias("Gender"),
        col("Age"),
        col("Marital status").alias("Marital_Status"),
        col("mental_health_history").alias("Mental_Health_Hist")
    )
)

dim_student.writeTo("silver_catalog.dimensions.dim_student").createOrReplace()

```

## 4. Création de la table faits

```
from pyspark.sql.functions import col

fact_performance = []
df_performance_factors_with_id.alias("f")

# 1) Join Dim_Student
.join(
    dim_student.select("Synthetic_ID", "Student_ID".alias("s"),
                      one=col("f.Synthetic_ID") == col("s.Synthetic_ID")),
    how="inner"
)

# 2) Join Dim_Course
.join(
    dim_course.select("Synthetic_ID", "Course_ID".alias("c"),
                     one=col("f.Synthetic_ID") == col("c.Synthetic_ID")),
    how="inner"
)

# 3) Join Dim_Time
.join(
    dim_time.select("Synthetic_ID", "Time_ID".alias("t"),
                   one=col("f.Synthetic_ID") == col("t.Synthetic_ID")),
    how="inner"
)

# 4) Join Dim_Lifestyle
.join(
    dim_lifestyle.select("Synthetic_ID", "Lifestyle_ID".alias("l"),
                       one=col("f.Synthetic_ID") == col("l.Synthetic_ID")),
    how="inner"
)

# 5) Join Dim_SocialFactors
.join(
    dim_social_factors.select("Synthetic_ID", "Social_ID".alias("so"),
                             one=col("f.Synthetic_ID") == col("so.Synthetic_ID")),
    how="inner"
)

# 6) Join Dim_Environment
.join(
    dim_environment.select("Synthetic_ID", "Environment_ID".alias("env"),
                           one=col("f.Synthetic_ID") == col("env.Synthetic_ID")),
    how="inner"
)

# 7)
.select(
    col("f.Exam_Score"),
    col("f.Attendance"),
    col("f.Hours_Studied"),
    col("f.Previous_Scores"),

    # foreign key
    col("s.Student_ID"),
    col("c.Course_ID"),
    col("t.Time_ID"),
    col("l.Lifestyle_ID"),
    col("so.Social_ID"),
    col("env.Environment_ID")
)
```

## Résultats du nettoyage

Après cette étape, les données sont :

1. **Prêtes pour les jointures** : Grâce à l'ajout des identifiants synthétiques.
2. **Normalisées** : Les noms de colonnes sont cohérents et les formats de données sont standardisés.
3. **Sans doublons** : Les colonnes clés sont dédoublées afin de réduire les incohérences dans les analyses futures.



Accueil >

iceberg-data

Conteneur

Rechercher

Charger

Modifier le niveau d'accès

Actualiser

Supprimer

Modifier le niveau

Acquérir le bail

Résilier le bail

Afficher les instantanés

Vue d'ensemble

Diagnostiquer et résoudre les problèmes

Contrôle d'accès (IAM)

Paramètres

Méthode d'authentification : Clé d'accès (Basculer vers le compte d'utilisateur Microsoft Entra)

Emplacement : iceberg-data / silver\_warehouse

Rechercher les objets blobs par préfixe (respect de la casse)

Afficher les objets blob supprimés

Ajouter un filtre

Nom	Modifié	Niveau d'accès	État de l'archive	Type d'objet blob	Taille	État du bail
<input type="checkbox"/> [.]						...
<input type="checkbox"/> dimensions						-
<input type="checkbox"/> fact						-
<input type="checkbox"/> silver						-
<input type="checkbox"/> dimensions	27/12/2024 20:24:27	Élevé (déduit)		Objet blob de blocs	0 B	Disponible
<input type="checkbox"/> fact	27/12/2024 23:12:11	Élevé (déduit)		Objet blob de blocs	0 B	Disponible
<input type="checkbox"/> silver	12/12/2024 18:40:27	Élevé (déduit)		Objet blob de blocs	0 B	Disponible

Accueil >

iceberg-data

Conteneur

Rechercher

Charger

Modifier le niveau d'accès

Actualiser

Supprimer

Modifier le niveau

Acquérir le bail

Résilier le bail

Afficher les instantanés

Vue d'ensemble

Diagnostiquer et résoudre les problèmes

Contrôle d'accès (IAM)

Paramètres

Méthode d'authentification : Clé d'accès (Basculer vers le compte d'utilisateur Microsoft Entra)

Emplacement : iceberg-data / silver\_warehouse / fact

Rechercher les objets blobs par préfixe (respect de la casse)

Afficher les objets blob supprimés

Ajouter un filtre

Nom	Modifié	Niveau d'accès	État de l'archive	Type d'objet blob	Taille	État du bail
<input type="checkbox"/> [.]						...
<input type="checkbox"/> fact_performance						-
<input type="checkbox"/> fact_performance	27/12/2024 23:12:11	Élevé (déduit)		Objet blob de blocs	0 B	Disponible

Accueil >

iceberg-data

Conteneur

Rechercher

Charger

Modifier le niveau d'accès

Actualiser

Supprimer

Modifier le niveau

Acquérir le bail

Résilier le bail

Afficher les instantanés

Vue d'ensemble

Diagnostiquer et résoudre les problèmes

Contrôle d'accès (IAM)

Paramètres

Méthode d'authentification : Clé d'accès (Basculer vers le compte d'utilisateur Microsoft Entra)

Emplacement : iceberg-data / silver\_warehouse / dimensions

Rechercher les objets blobs par préfixe (respect de la casse)

Afficher les objets blob supprimés

Ajouter un filtre

Nom	Modifié	Niveau d'accès	État de l'archive	Type d'objet blob	Taille	État du bail
<input type="checkbox"/> [.]						...
<input type="checkbox"/> dim_course						-
<input type="checkbox"/> dim_environment						-
<input type="checkbox"/> dim_lifestyle						-
<input type="checkbox"/> dim_social_factors						-
<input type="checkbox"/> dim_student						-
<input type="checkbox"/> dim_time						-
<input type="checkbox"/> dim_course	27/12/2024 20:25:11	Élevé (déduit)		Objet blob de blocs	0 B	Disponible
<input type="checkbox"/> dim_environment	27/12/2024 21:26:49	Élevé (déduit)		Objet blob de blocs	0 B	Disponible
<input type="checkbox"/> dim_lifestyle	27/12/2024 20:44:42	Élevé (déduit)		Objet blob de blocs	0 B	Disponible
<input type="checkbox"/> dim_social_factors	27/12/2024 21:20:54	Élevé (déduit)		Objet blob de blocs	0 B	Disponible
<input type="checkbox"/> dim_student	27/12/2024 20:24:27	Élevé (déduit)		Objet blob de blocs	0 B	Disponible
<input type="checkbox"/> dim_time	27/12/2024 20:25:56	Élevé (déduit)		Objet blob de blocs	0 B	Disponible

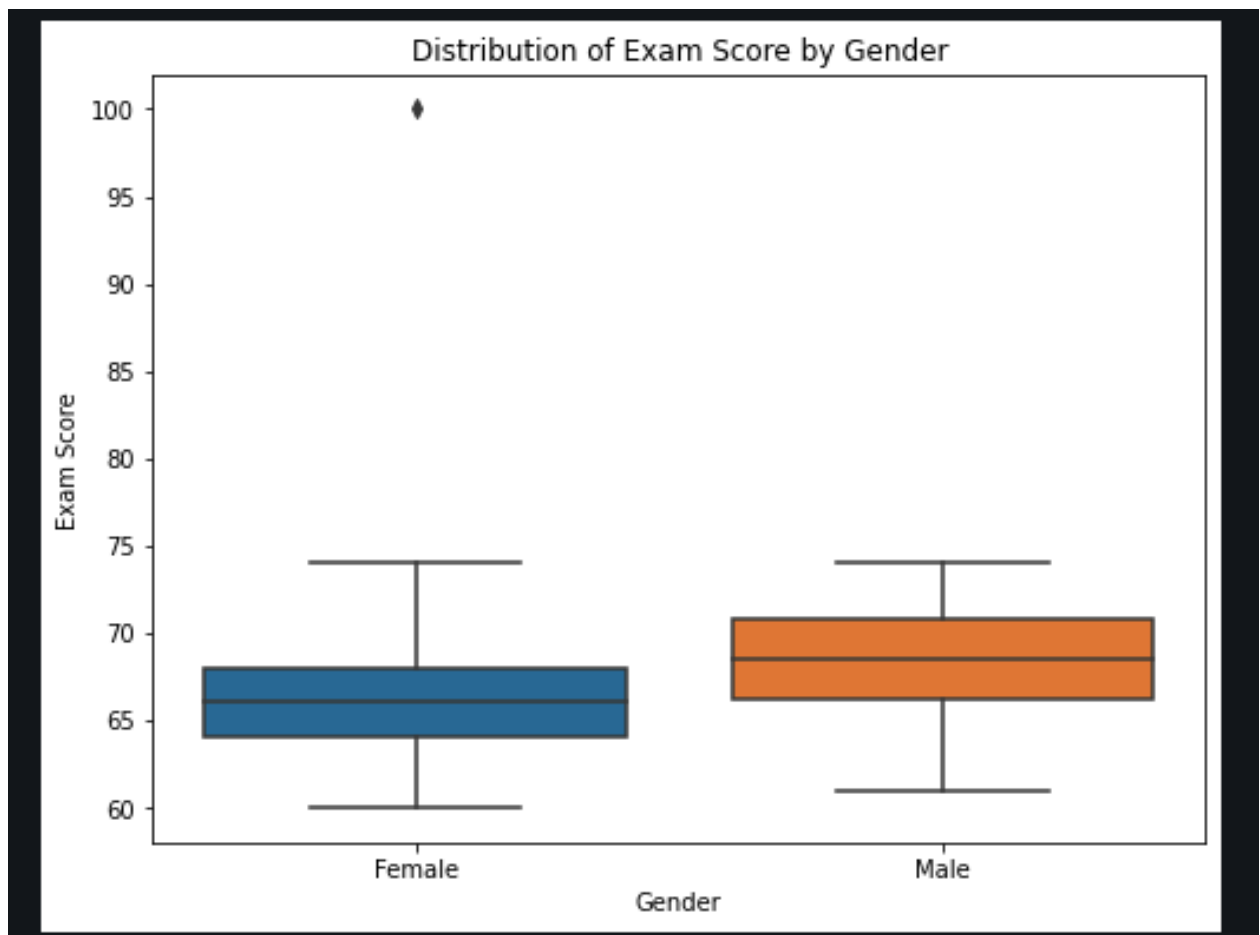
# Analyse des résultats

## Visualisation des données

Nous avons effectué une jointure entre la table de faits et les différentes tables de dimensions afin de rassembler toutes les informations nécessaires sur un seul DataFrame. Cela nous permet d'analyser conjointement les facteurs de style de vie, l'environnement, le soutien social et la performance académique des étudiants.

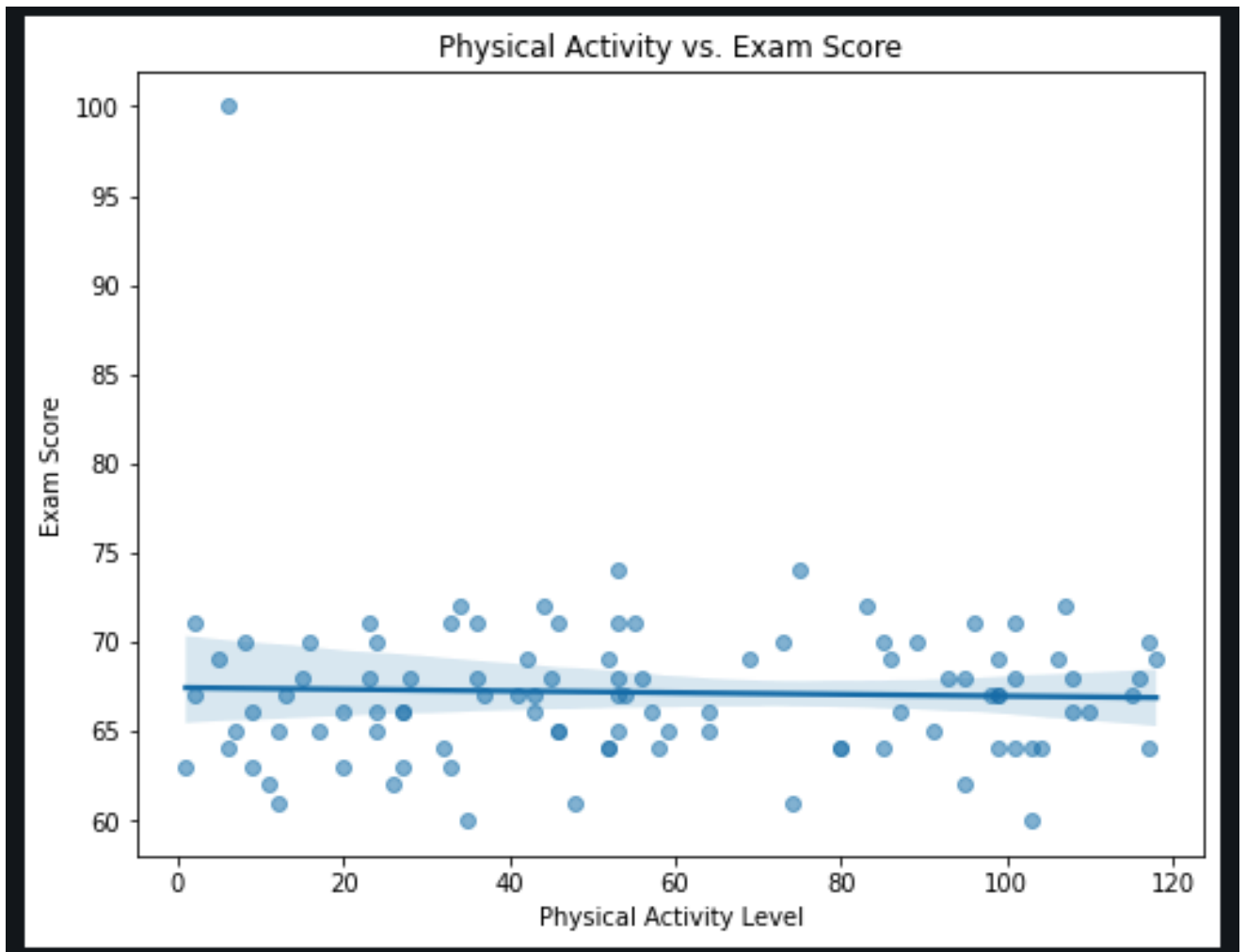
### 1. Distribution des notes d'examen selon le genre

**Interprétation possible :** Déterminer si l'un des genres obtient systématiquement des notes plus élevées ou si les distributions sont semblables.



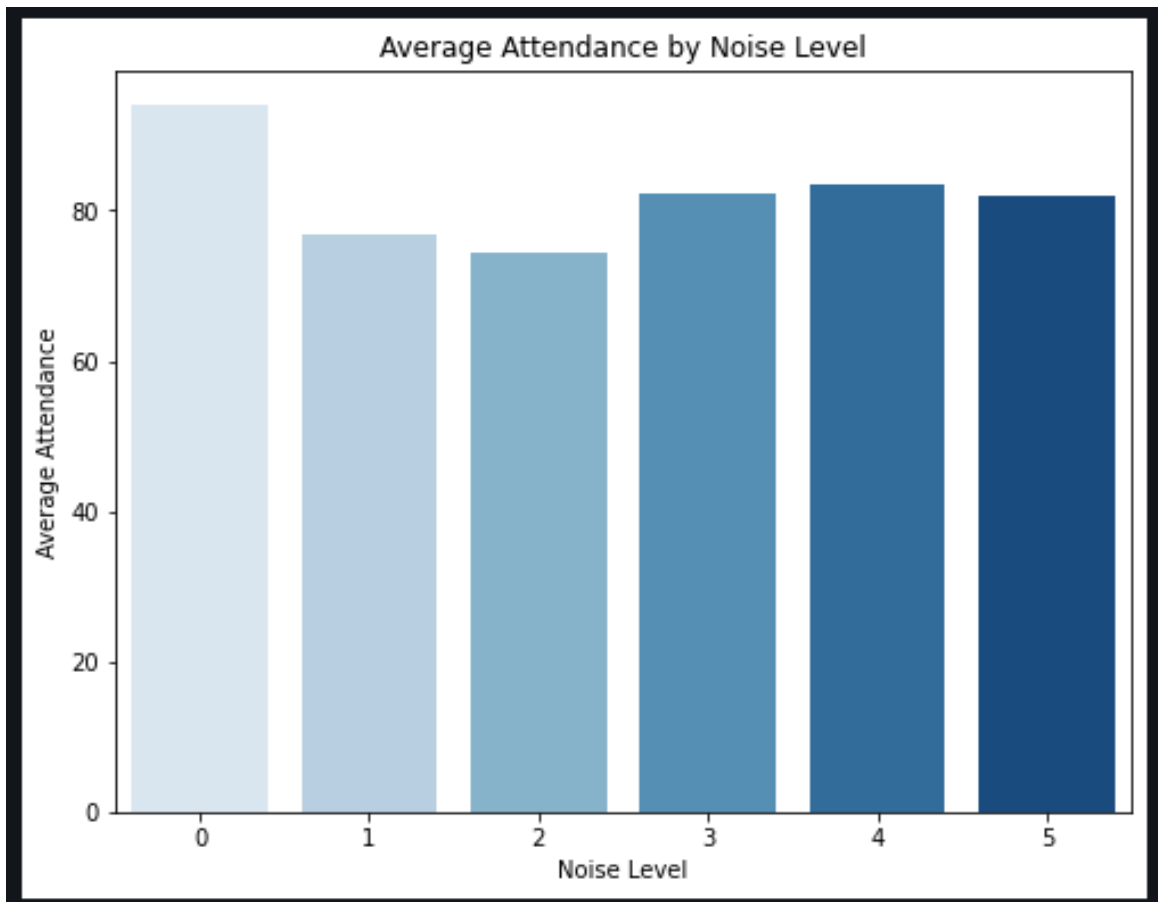
## 2. Corrélation entre l'activité physique et la note d'examen

**Interprétation possible :** Une pente positive suggère qu'augmenter l'activité physique coïncide avec de meilleures notes.



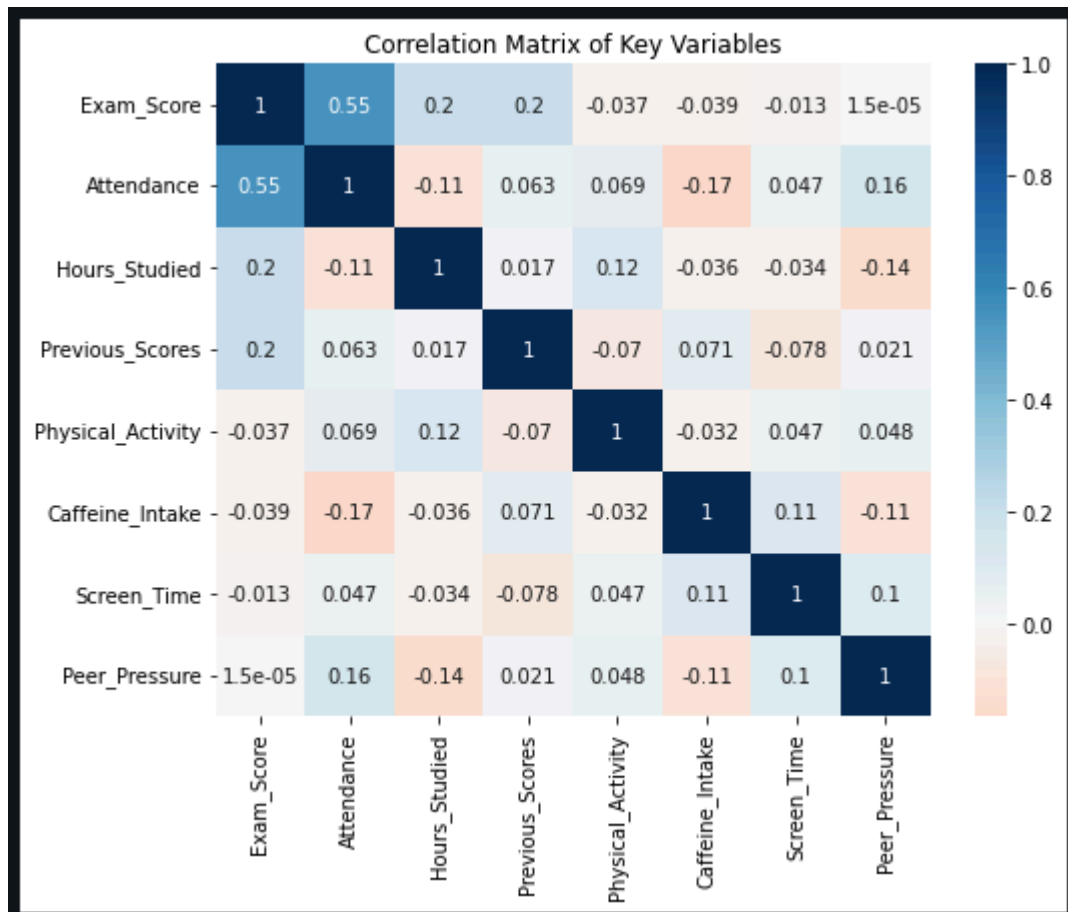
### 3. Impact du niveau de bruit sur le taux de présence

**Interprétation possible :** Un environnement plus calme contribuerait-il à un meilleur engagement ?



#### 4. Carte de corrélation pour plusieurs variables

**Interprétation possible :** Identifier rapidement quelles variables sont corrélées positivement ou négativement.



## Conclusion

Ce projet Big Data axé sur la santé mentale et les performances académiques des étudiants nous a démontré l'importance d'une architecture de données pour une meilleure compréhension et un suivi continu. Grâce aux différentes couches (Bronze, Silver et Gold), nous avons pu nettoyer, transformer puis organiser nos données de façon fiable, garantissant leur cohérence et leur intégrité tout au long du pipeline.

L'utilisation d'Iceberg et de modèles en étoile (tables de faits et dimensions) nous a permis de structurer les informations clés (sommeil, stress, environnement, etc.) et de faciliter l'exploration analytique. Les visualisations finales ont mis en évidence certaines corrélations

significatives, par exemple entre le niveau de stress et la réussite académique, ou encore l'impact positif de bonnes conditions de vie sur les notes et la santé mentale.

En définitive, l'architecture mise en place et les analyses menées offrent une vision de l'état de la santé mentale et de la performance des étudiants.

## Annexes

Code et scripts : <https://github.com/Jie01236/Projet-Big-Data.git>