



Alexis Balayre

Machine learning & Big Data Assignment

School of Aerospace, Transport and Manufacturing
Computational Software of Techniques Engineering

MSc
Academic Year: 2023 - 2024

Supervisor: Dr Jun Li
20th November 2023

Table of Contents

| | |
|--|-----------|
| Table of Contents | ii |
| List of Figures | iv |
| List of Tables | v |
| 1 Introduction | 1 |
| 2 Dataset Description | 2 |
| 3 Methodologies | 3 |
| 3.1 Queries Tasks | 3 |
| 3.1.1 Query 1 | 3 |
| 3.1.2 Query 2 | 4 |
| 3.1.3 Query 3 | 5 |
| 3.2 Queries Optimisation | 7 |
| 3.2.1 Common Performance Problems in Spark | 7 |
| 3.2.2 Optimisation Solutions | 7 |
| 3.3 Pipeline of the Project | 9 |
| 3.3.1 Pipeline Overview | 9 |
| 3.3.2 Pipeline Orchestration | 10 |
| 4 Results & Discussion | 11 |
| 4.1 Queries Results | 11 |
| 4.1.1 Query 1 | 11 |
| 4.1.2 Query 2 | 13 |
| 4.1.3 Query 3 | 16 |
| 4.2 Discussion of Results | 18 |
| 4.3 Ethical Considerations and Challenges | 18 |
| 5 Conclusion | 20 |
| References | 21 |
| A Documentation | 22 |
| B Query 1 Source Code | 26 |

| | | |
|----------|----------------------------|-----------|
| C | Query 2 Source Code | 29 |
| D | Query 3 Source Code | 36 |
| E | Terminal Output | 44 |

List of Figures

| | | |
|-----|--|----|
| 3.1 | Pipeline Diagram | 9 |
| 3.2 | Apache Airflow DAG Graph | 10 |
| 4.1 | Mean Daily Confirmed Cases Per Month | 12 |
| 4.2 | Top 100 Locations most affected by the pandemic | 13 |
| 4.3 | Mean Confirmed Cases By Week and Continent | 14 |
| 4.4 | Standard Deviation Confirmed Cases By Week and Continent | 14 |
| 4.5 | Max Confirmed Cases By Week and Continent | 15 |
| 4.6 | Min Confirmed Cases By Week and Continent | 15 |
| 4.7 | Top 50 Locations most affected by the pandemic | 16 |
| 4.8 | Custom KMeans Clustering on 03/2020 | 17 |
| 4.9 | Spark MLlib KMeans Clustering on 03/2020 | 17 |

List of Tables

| | | |
|-----|---|----|
| 2.1 | COVID-19 Dataset Sample | 2 |
| 2.2 | COVID-19 Dataset Statistics | 2 |
| 2.3 | COVID-19 Dataset Null Values | 2 |
| 4.1 | Query 1 Results Sample | 11 |
| 4.2 | Query 2 Results Sample | 13 |
| 4.3 | Query 3 Results Sample - Custom KMeans Clustering | 16 |
| 4.4 | Query 3 Results Sample - Spark KMeans Clustering | 16 |

Chapter 1

Introduction

During the global COVID-19 pandemic, nations took advantage of Big Data Analytics and Artificial Intelligence technologies to understand and combat the spread of the virus. This report focuses on the detailed analysis of a comprehensive global dataset of confirmed cases of COVID-19, leveraging Apache Spark, an advanced data processing framework. The main objective is to provide a clear picture of the evolution of the pandemic on a global scale. The aim is to understand not only the speed and patterns of spread of the virus. It also aims to demonstrate how Big Data tools can be used to transform, analyse and interpret massive volumes of data, offering a unique and valuable perspective in the management of global health crises. By harnessing these technologies, the aim is not only to provide insights into the current situation, but also to pave the way for faster and more effective responses to future pandemics.

Chapter 2

Dataset Description

COVID-19 Dataset

As shown in table 2.1, the “time_series_covid19_confirmed_global” dataset presents each row with key identifiers such as “province/state” and “country/region”, as well as geographical coordinates (latitude and longitude). Importantly, it traces the progression of the pandemic since its inception, with each column dated in mm/dd/yy format, revealing the total number of confirmed cases at that date. This cumulative data structure enables a dynamic and detailed analysis of the spread of the virus from 22 January 2020 to March 2023. Moreover, the dataset contains 201 distinct Countries/Regions, 92 distinct Provinces/States, and 198 null values for this last as shown in tables 2.2 and 2.3.

Table 2.1: COVID-19 Dataset Sample

| | Province/State | Country/Region | Lat | Long | 1/22/20 | ... | 3/9/23 |
|---|----------------|----------------|-----------|-----------|---------|-----|--------|
| 1 | NaN | Afghanistan | 33.93911 | 67.709953 | 0 | | 209451 |
| 2 | NaN | Albania | 41.15330 | 20.168300 | 0 | | 334457 |
| 3 | NaN | Algeria | 28.03390 | 1.659600 | 0 | | 271496 |
| 4 | NaN | Andorra | 42.50630 | 1.521800 | 0 | | 47890 |
| 5 | NaN | Angola | -11.20270 | 17.873900 | 0 | | 105288 |

Table 2.2: COVID-19 Dataset Statistics

| Column | Number of entries |
|-------------------------|-------------------|
| Distinct Province/State | 92 |
| Distinct Country/Region | 201 |

Table 2.3: COVID-19 Dataset Null Values

| Column | Number of entries |
|----------------|-------------------|
| Province/State | 198 |
| Lat | 2 |
| Long | 2 |

Chapter 3

Methodologies

3.1 Queries Tasks

3.1.1 Query 1

The first task is to calculate the average number of confirmed daily cases of people infected with COVID-19 for each country in the dataset. This task is implemented by the class **Query1** in Appendix B. The steps of the task are as follows:

1. **Pivoting the DataFrame:** The “covidDataDf” DataFrame initially contains dates in columns. This step reorganises the DataFrame so that the dates are row entries. It uses the stack function to transform each date column into a row, associated with its corresponding value. The result is a table with the following columns: “Country/Region”, “Date” and “Value”
2. **Formatting Date Column:** The “Date” column is converted from a character string to an appropriate date format, in this case “M/d/yy”.
3. **Calculating Daily Cases:** A specification window separates the data by country and sorts it by date. Daily cases are calculated by subtracting the number of cases from the previous day from the current day’s cases. If there is no data for the previous day (for example, at the start of the dataset), the number of cases is set to 0. The formula used is:

$$\text{DailyCases} = \begin{cases} 0, & \text{if PrevValue is null} \\ \text{Value} - \text{PrevValue}, & \text{otherwise} \end{cases} \quad (3.1)$$

where PrevValue is the number of cases from the previous day.

4. **Grouping and Averaging Daily Cases:** The data is grouped by “Country”, “Year” and “Month”. For each group, the average number of cases per day is calculated. This gives the average number of confirmed cases per day for each month and country.

3.1.2 Query 2

The second task is to calculate the mean, standard deviation, minimum and maximum of the number of confirmed cases daily for each week and continent. This task is implemented by the class **Query2** in Appendix C. Here are the steps performed by the task:

1. **Data Preparation** (`_prepare_data` method): The COVID data is prepared for analysis. A “Location” column is created by combining the values “Province/State” and “Country/Region” to consider the country if the state is not indicated.
2. **Assigning Continents** (`_assign_continent` method): Continents are assigned to each line of data on the basis of geographical coordinates. To do this, the method uses the GeoPandas “world_boundaries_gdf” DataFrame, loaded from a shapefile (1) containing continental boundaries, to identify the continent corresponding to a set of longitude and latitude coordinates. The method implements an internal `determine_continent` function, which searches for the matching continent by checking whether a geographic point (formed by these coordinates) lies within the boundaries of one of the continents defined in “world_boundaries_gdf”. If no match is found in this DataFrame, the method applies predefined geographic conditions (2) to approximate the continent. Finally, this logic is applied to each row of the `covidDataDf` DataFrame to add the new “Continent” column. This method was created using a GeoPandas Tutorial (3).
3. **Pivoting the DataFrame** (`_pivot_table` method): The “covidDataDf” DataFrame initially contains dates in columns. This step reorganises the DataFrame so that the dates are row entries. It uses the `stack` function to transform each date column into a row, associated with its corresponding value. The result is a table with the following columns: “Location”, “Continent”, “Date” and “Value”
4. **Computing Daily Confirmed Cases** (`_compute_daily_confirmed_cases` method): A window specification partitions the data by “Location” and orders it by “Date”. Daily confirmed cases are calculated using the formula:

$$\text{DailyCases} = \begin{cases} 0, & \text{if PrevValue is null} \\ \text{Value} - \text{PrevValue}, & \text{otherwise} \end{cases} \quad (3.2)$$

where `PrevValue` is the number of cases from the previous day.

5. **Computing Slopes** (`_compute_slopes` method): This method computes the slopes of daily confirmed cases for each location, preparing the data for linear regression. The slope is calculated using the formula:

$$\text{DailySlope} = \frac{n(\sum xy) - (\sum x)(\sum y)}{n(\sum x^2) - (\sum x)^2} \quad (3.3)$$

where:

- x is the number of days since the start of the dataset.

- y is the number of daily confirmed cases.
 - n is the number of observations.
6. **Filtering Top Affected Locations** (`--filter_top_affected` method): The top 100 affected locations are identified based on the slope of daily confirmed cases. To do this, the data is grouped by location and the maximum slope is calculated for each group. The top 100 locations are then selected.
 7. **Aggregating Data** (`--aggregate` method): For each entry, a new “WeekRange” column is added with the start and end dates of the week on which the day in the “Date” column falls. The DataFrame is then grouped by continent and week. Statistical measures like mean, standard deviation, minimum, and maximum of daily cases are calculated for each group.

3.1.3 Query 3

The last task is to calculate perform KMeans clustering on the top 50 affected locations depending on the slope of monthly confirmed cases. This task is implemented by the class **Query3** in Appendix D. Here are the steps performed by the task:

1. **Data Preparation** (`--prepare_data` method): The COVID data is prepared for analysis. A “Location” column is created by combining the values “Province/State” and “Country/Region” to consider the country if the state is not indicated.
2. **Pivoting the DataFrame** (`--pivot_table` method): The “covidDataDf” DataFrame initially contains dates in columns. This step reorganises the DataFrame so that the dates are row entries. It uses the stack function to transform each date column into a row, associated with its corresponding value. The result is a table with the following columns: “Location”, “Date” and “Value”
3. **Computing Daily Confirmed Cases** (`--compute_daily_confirmed_cases` method): Daily confirmed cases are computed using the formula:

$$\text{DailyCases} = \begin{cases} 0, & \text{if PrevValue is null} \\ \text{Value} - \text{PrevValue}, & \text{otherwise} \end{cases} \quad (3.4)$$

4. **Computing Monthly Slopes** (`--compute_monthly_slopes` method): The slopes of monthly confirmed cases for each location are computed using linear regression. The slope formula used is:

$$\text{MonthlySlope} = \frac{n(\sum xy) - (\sum x)(\sum y)}{n(\sum x^2) - (\sum x)^2} \quad (3.5)$$

where:

- x is the number of months since the start of the dataset.
- y is the number of monthly confirmed cases.
- n is the number of observations.

5. **Filtering Top Affected Locations** (`--filter_top_affected` method): The top 50 affected locations are identified based on the slope of monthly confirmed cases. To do this, the data is grouped by location and the mean slope is calculated for each group. The top 50 locations are then selected.
6. **Applying Custom Clustering** (`--apply_custom_clustering` method): A custom KMeans algorithm is applied to cluster the locations based on their slopes. This source code was adapted from the following Medium article (4).
7. **Applying Clustering** (`--apply_clustering` method): The standard KMeans algorithm from Spark MLlib is used for clustering the locations.

3.2 Queries Optimisation

Once the 3 queries had been built, an optimisation step was carried out to improve task execution time. This section will look at the various difficulties encountered and how to overcome them.

3.2.1 Common Performance Problems in Spark

The five most common performance issues encountered in Apache Spark, known as the 5 Ss (5) are:

1. **Spill:** When there is not enough RAM memory to process the current data, Spark is forced to move some data to the hard disk, a process known as “spilling”.
2. **Skew:** The “skew” problem arises when data is not distributed evenly across partitions in Spark. Some nodes end up with a much heavier workload than others, creating bottlenecks and slowing down overall processing.
3. **Shuffle:** During complex operations such as joins or clusters, Spark redistributes the data so that the corresponding elements are on the same node. Shuffling is costly in terms of performance because it involves intensive data transfer over the network.
4. **Storage:** Inefficient storage management, such as processing many small files or using non-optimised file formats, can lead to a high number of I/O operations, which slows down performance. In addition, a poor storage strategy can also lead to shuffling and asymmetry problems.
5. **Serialization:** “Serialization” in Spark refers to the conversion of objects into a format that can be easily transmitted over the network or stored on disk. Inefficient serialization and deserialization processes can significantly slow down data transfer between cluster nodes and increase overall processing time.

3.2.2 Optimisation Solutions

Here is a list of techniques for optimising Apache Spark processes and overcoming performance problems (6):

1. **Efficient partitioning:** Efficient partitioning in Apache Spark is essential for the balanced distribution of data across cluster nodes. This distribution plays a crucial role in optimising performance, particularly for shuffle operations. By adjusting the number of partitions with `repartition()` or `coalesce()`, it is possible to balance the workload between nodes, which is effective for join and aggregation operations, for example.
2. **Persist/Cache and Early Filtering:** The judicious use of caching and early filtering goes hand in hand in optimising Spark processes. By using `persist()` or `cache()`, frequently accessed data is retained in memory, reducing the time required

for repeated operations. At the same time, early data filtering, before resource-intensive operations such as joins or shuffles, can significantly reduce the amount of data processed, speeding up the whole process.

3. Data Format Choice: The choice of data format is an often underestimated but crucial aspect of Spark optimisation. Columnar formats such as Parquet or ORC are preferable for efficient read and write operations. By choosing the right data format, you can significantly improve read performance and reduce the amount of storage space required.

4. Minimising Shuffle and Broadcast Operations: Minimising shuffle operations is crucial to improving performance in Spark. Shuffles, which are necessary for operations such as joins and groupings, can be costly in terms of performance because they involve moving large amounts of data around the network. In addition, broadcasting is a powerful technique for joins involving a large table and a smaller table.

3.3 Pipeline of the Project

3.3.1 Pipeline Overview

The pipeline of this project is composed of four main components: **data ingestion**, **query 1**, **query 2** and **query 3**.

The **ingestion** task retrieves the last version of the data set from the source and stores it in the data lake (CSV file stored in local storage). Then, the **queries 1, 2 & 3** tasks retrieves the data from the data lake and performs the queries on the data.

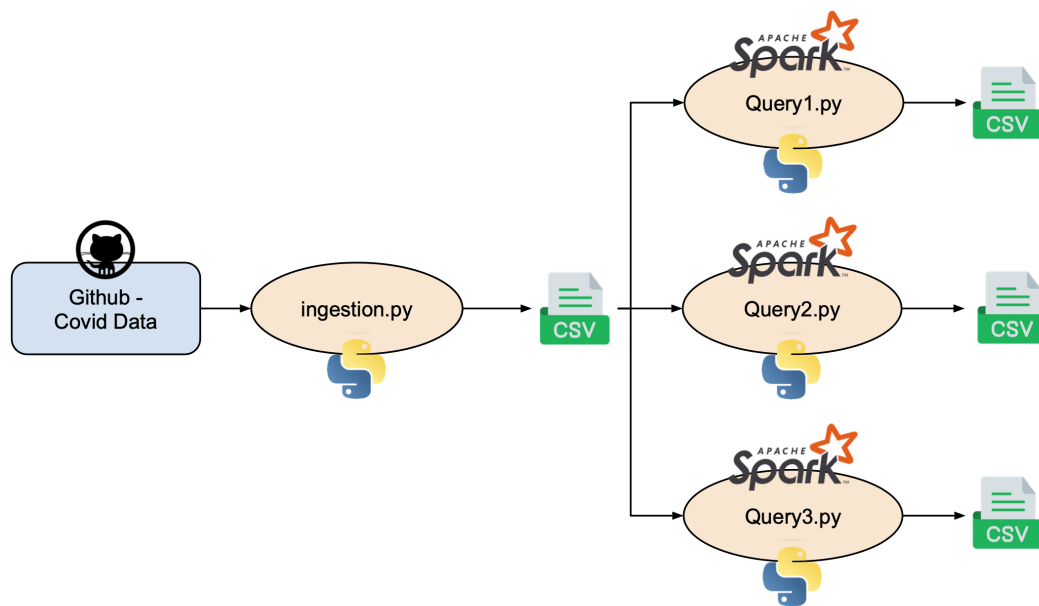


Figure 3.1: Pipeline Diagram

3.3.2 Pipeline Orchestration

In order to orchestrate and automate the pipeline, a scheduled task must be run every day to retrieve the latest version of the dataset and run the tasks when a new daily row is added at 23:59 UTC to the dataset.

To perform this task, a DAG (Directed Acyclic Graph) was created using Apache Airflow. The DAG is scheduled to run every day at 00:00 UTC and is composed of four tasks: **ingestion**, **query 1**, **query 2** and **query 3**. The screenshot below 3.2 shows the DAG graph of the pipeline in the web interface of Apache Airflow.

The benefits of using a workflow platform such as Apache Airflow are its ability to schedule and automate the pipeline, as well as its ability to monitor the pipeline and send alerts if a task fails.

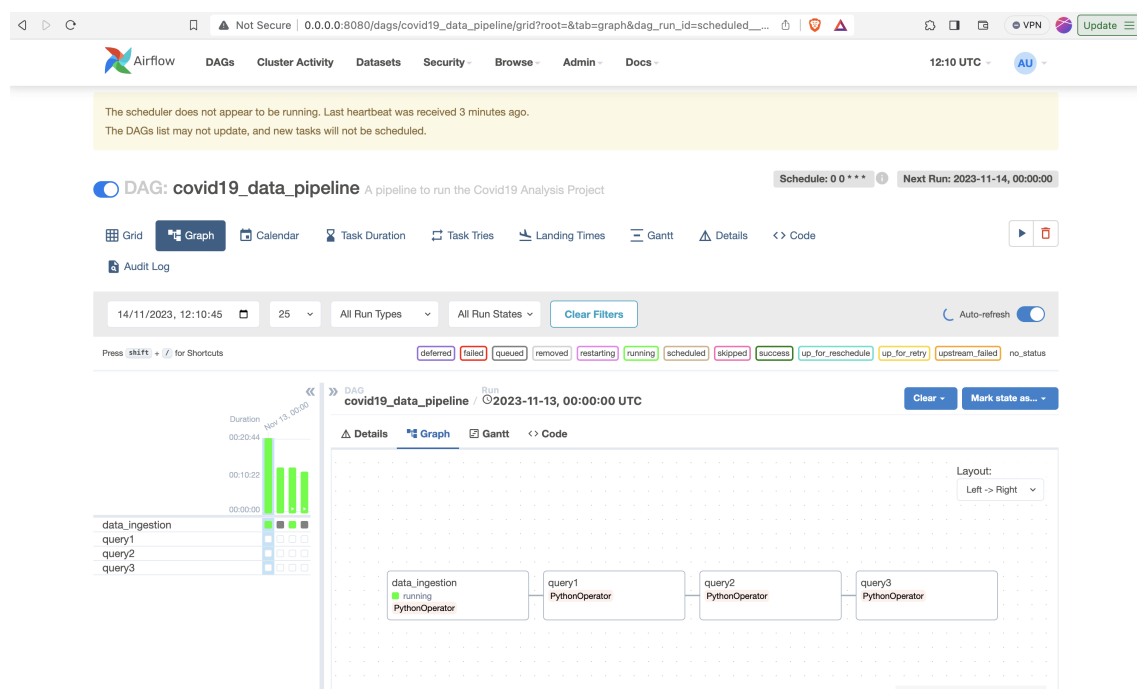


Figure 3.2: Apache Airflow DAG Graph

Chapter 4

Results & Discussion

4.1 Queries Results

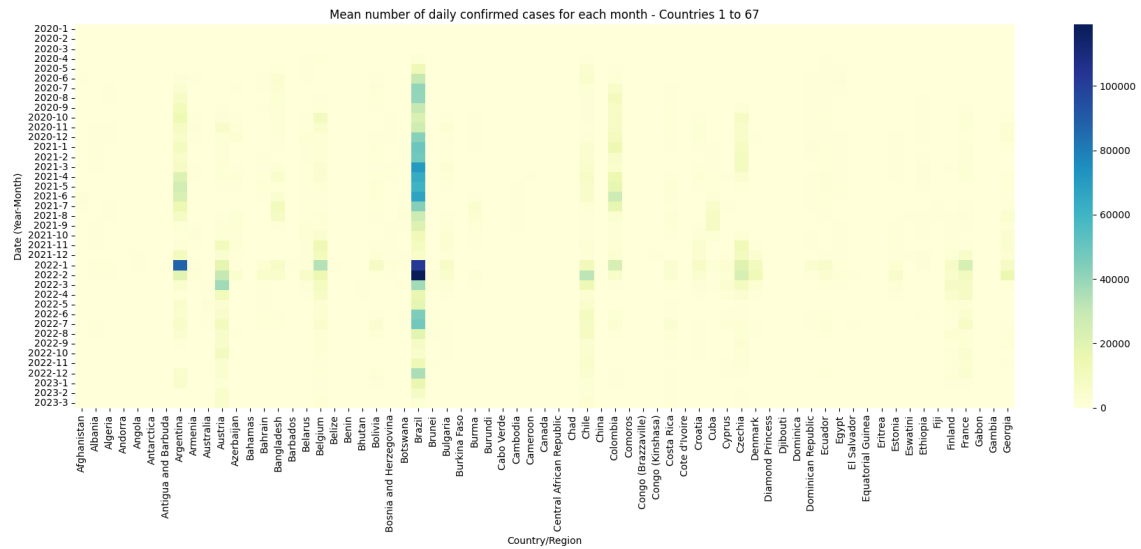
The programme was last run on 18 November 2023. The appendix E shows the output of the program on the terminal.

4.1.1 Query 1

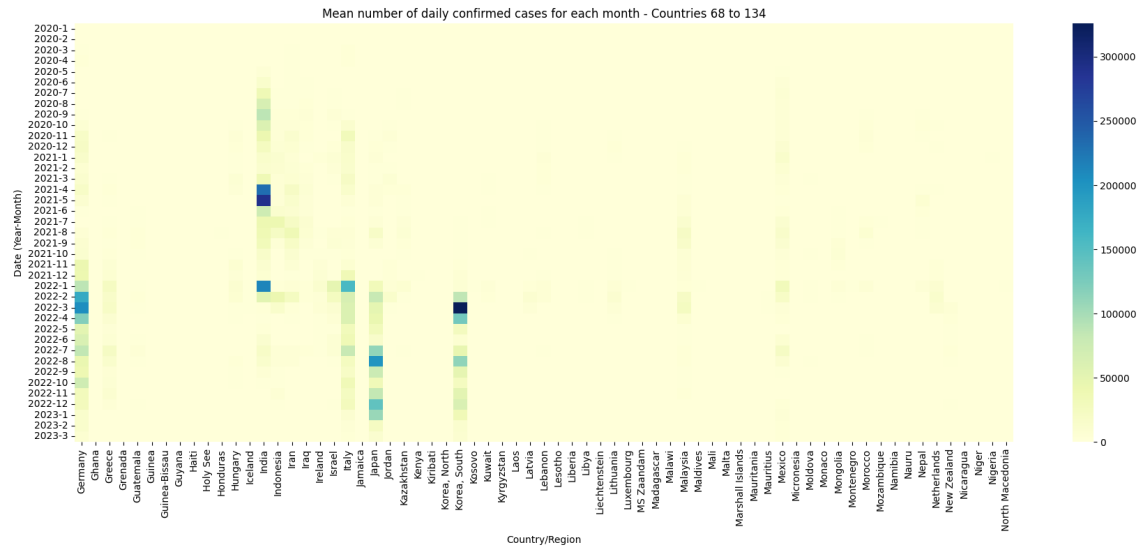
The first query takes around 3 seconds, and the table 4.1 shows a sample of the data calculated during the task. The results are consistent with what was expected. For example, the figure 4.1a shows that Brazil was heavily impacted by the pandemic, reaching an average peak in February 2022. The same is true for Korea in figure 4.1b and the United States in figure 4.1c, which were heavily impacted.

Table 4.1: Query 1 Results Sample

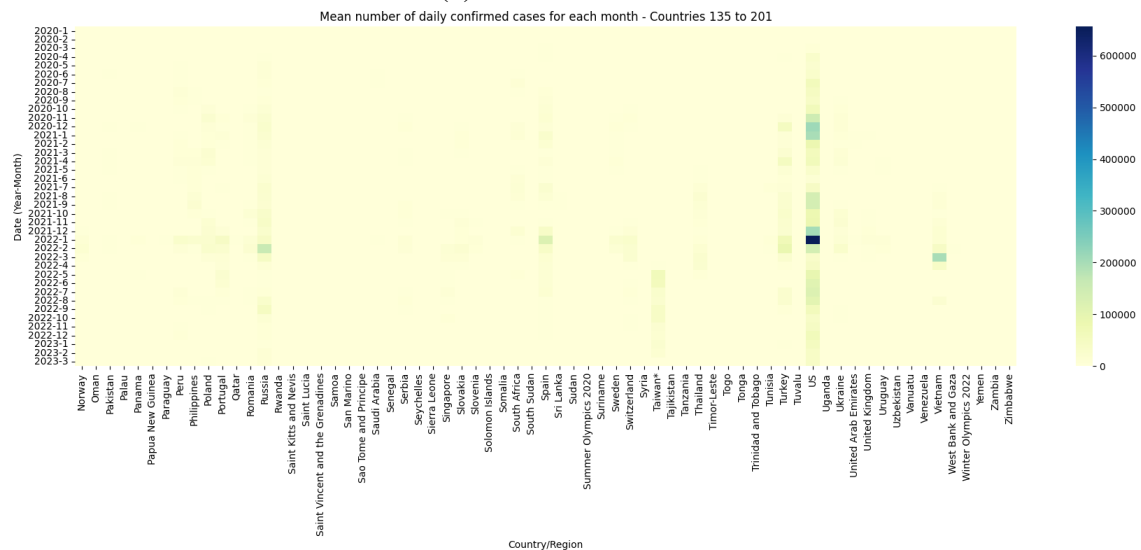
| | Country/Region | Year | Month | Average |
|---|----------------|------|-------|--------------------|
| 1 | Afghanistan | 2020 | 1 | 0.0 |
| 2 | Afghanistan | 2020 | 2 | 0.1724137931034483 |
| 3 | Afghanistan | 2020 | 3 | 5.193548387096774 |
| 4 | Afghanistan | 2020 | 4 | 55.36666666666667 |
| 5 | Afghanistan | 2020 | 5 | 430.741935483871 |



(a) Countries 1 to 67



(b) Countries 68 to 134



(c) Countries 135 to 201

Figure 4.1: Mean Daily Confirmed Cases Per Month

4.1.2 Query 2

The second query takes around 15 seconds, and the table 4.2 shows a sample of the data calculated during the task. Locations used to compute the statistics are shown on the map of figure 4.2. The area of the circles is proportional to how the location has been affected by the pandemic.

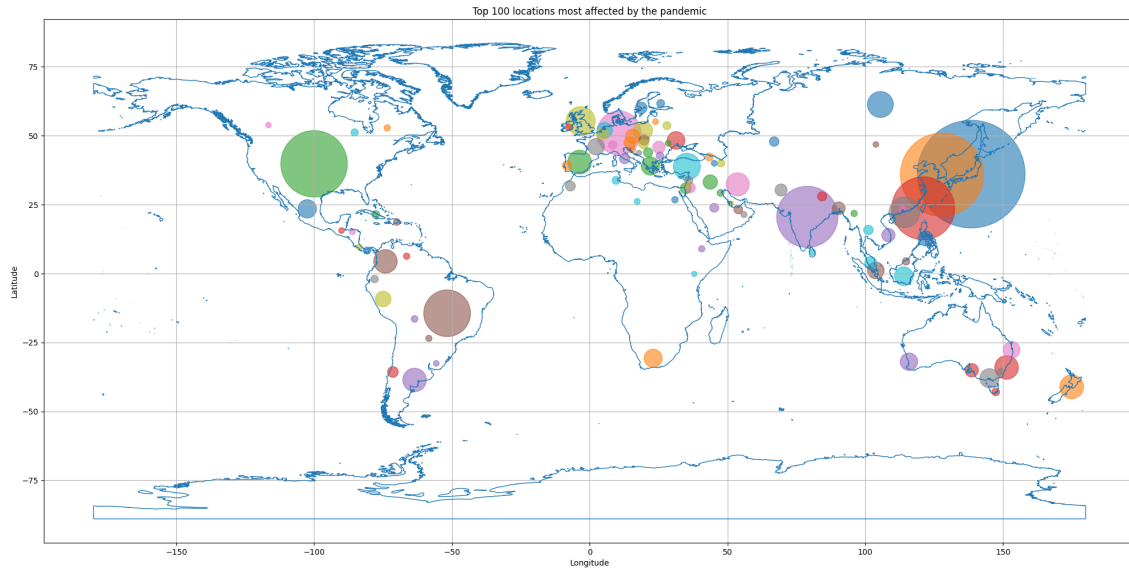


Figure 4.2: Top 100 Locations most affected by the pandemic

Table 4.2: Query 2 Results Sample

| | Continent | WeekRange | Mean | Std | Min | Max |
|----|-----------|-----------------------|----------|----------|-----|-----|
| 1 | Africa | 19/01/2020-25/01/2020 | 0.0 | 0.0 | 0 | 0 |
| 2 | Africa | 26/01/2020-01/02/2020 | 0.0 | 0.0 | 0 | 0 |
| 3 | Africa | 02/02/2020-08/02/2020 | 0.0 | 0.0 | 0 | 0 |
| 4 | Africa | 09/02/2020-15/02/2020 | 0.0 | 0.0 | 0 | 0 |
| 5 | Africa | 16/02/2020-22/02/2020 | 0.0 | 0.0 | 0 | 0 |
| 6 | Africa | 23/02/2020-29/02/2020 | 0.0 | 0.0 | 0 | 0 |
| 7 | Africa | 01/03/2020-07/03/2020 | 0.02857 | 0.16903 | 0 | 1 |
| 8 | Africa | 08/03/2020-14/03/2020 | 0.65714 | 1.73108 | 0 | 9 |
| 9 | Africa | 15/03/2020-21/03/2020 | 2.74285 | 4.53964 | 0 | 17 |
| 10 | Africa | 22/03/2020-28/03/2020 | 12.62857 | 15.33754 | 0 | 59 |
| 11 | Africa | 29/03/2020-04/04/2020 | 14.97142 | 19.31242 | 0 | 82 |

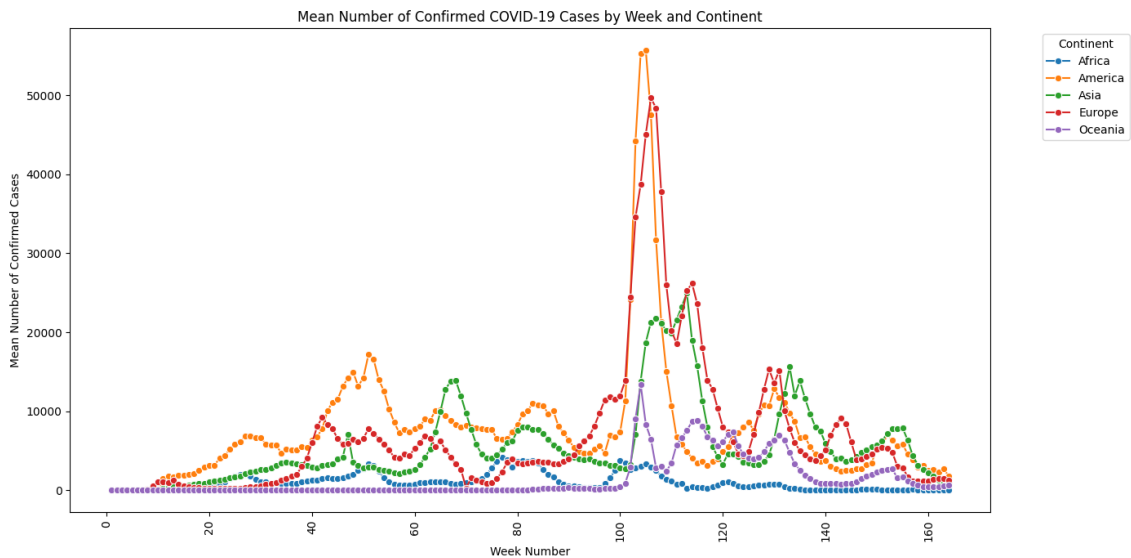


Figure 4.3: Mean Confimed Cases By Week and Continent

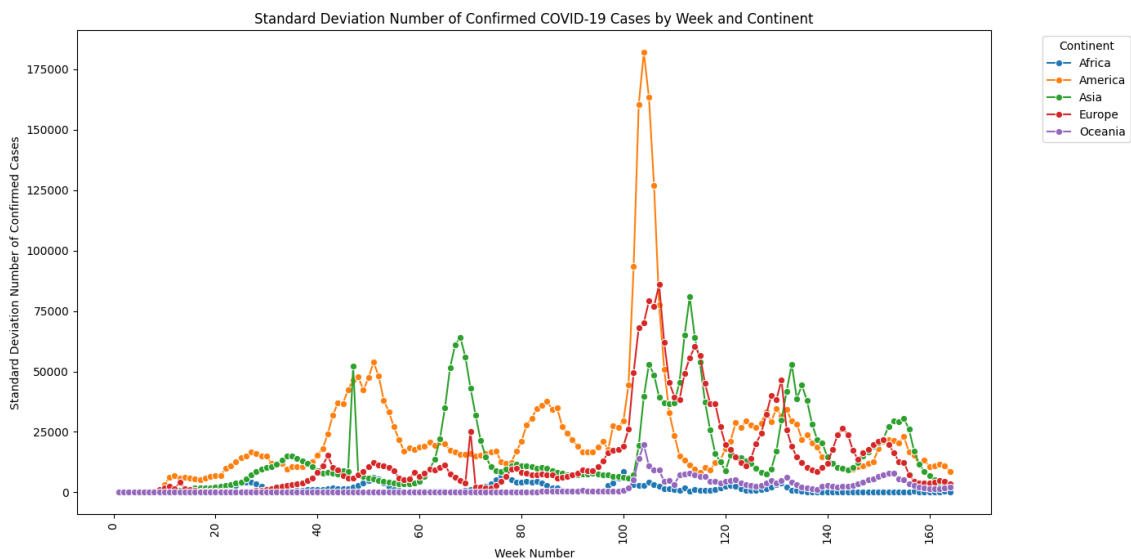


Figure 4.4: Standard Deviation Confimed Cases By Week and Continent

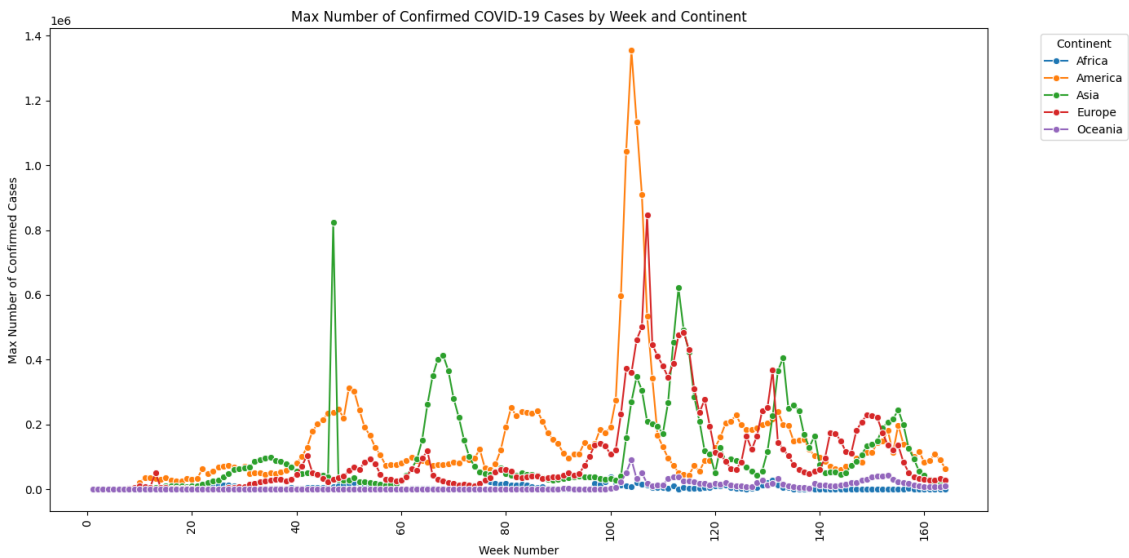


Figure 4.5: Max Confimed Cases By Week and Continent

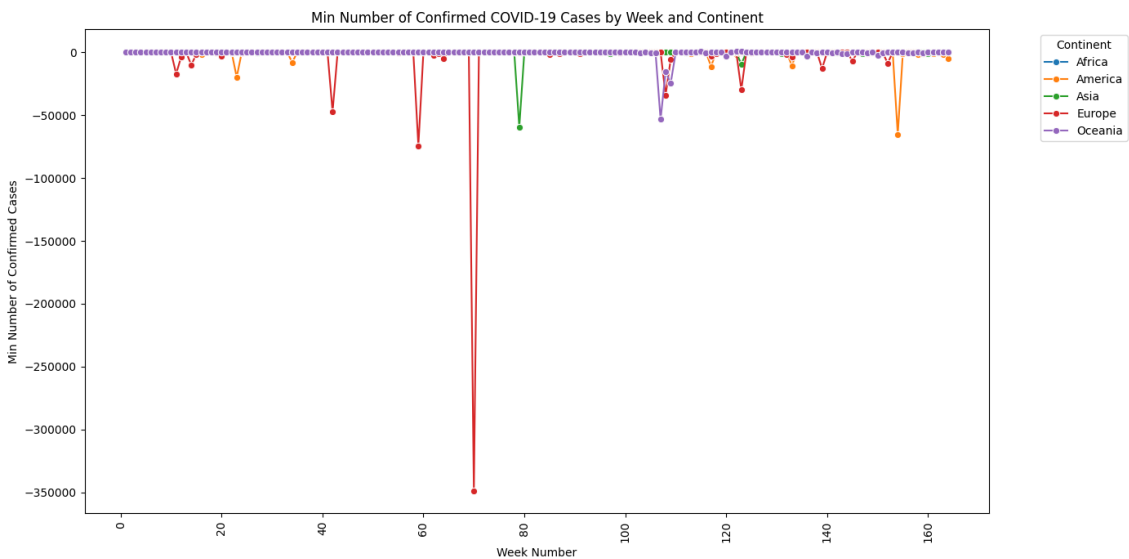


Figure 4.6: Min Confimed Cases By Week and Continent

4.1.3 Query 3

The third query takes around 3 minutes as clustering with the custom implementation takes 60 seconds and clustering with the Spark MLlib implementation takes 110 seconds. The table 4.3 shows a sample of the data calculated during the task with the custom implementation and the table 4.4 with the Spark MLlib implementation. Locations used to compute the statistics are shown on the map of figure 4.7. The area of the circles is proportional to how the location has been affected by the pandemic.

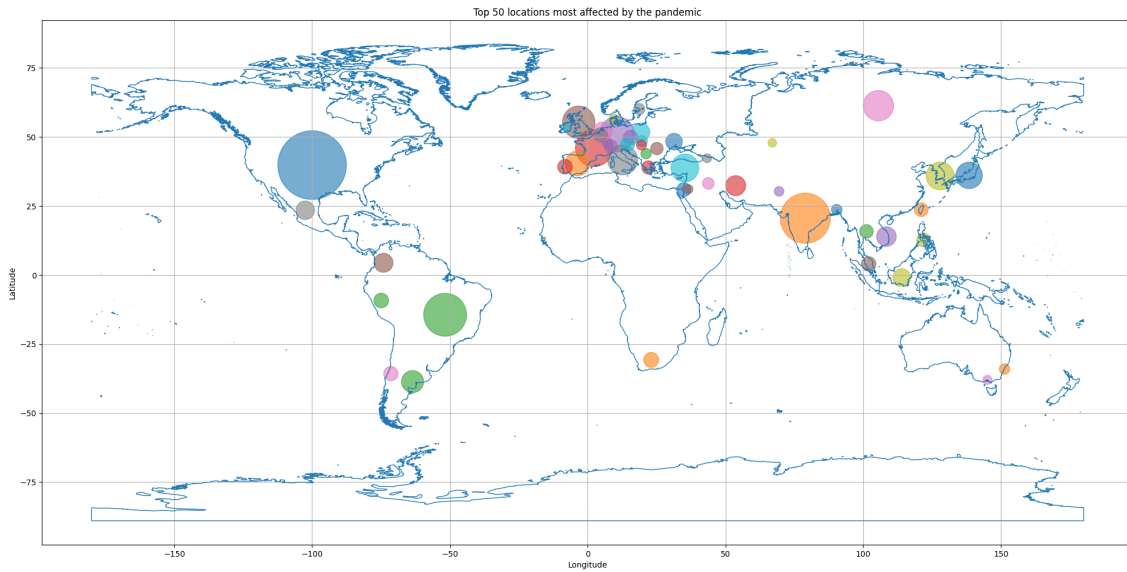


Figure 4.7: Top 50 Locations most affected by the pandemic

Table 4.3: Query 3 Results Sample - Custom KMeans Clustering

| | Location | Month | Cluster |
|---|-----------|---------|---------|
| 1 | Argentina | 2020-01 | 2 |
| 2 | Austria | 2020-01 | 2 |
| 3 | Brazil | 2020-01 | 2 |
| 4 | Czechia | 2020-01 | 2 |
| 5 | France | 2020-01 | 1 |

Table 4.4: Query 3 Results Sample - Spark KMeans Clustering

| | Location | Month | Cluster |
|---|-----------|---------|---------|
| 1 | Argentina | 2020-01 | 2 |
| 2 | Austria | 2020-01 | 0 |
| 3 | Brazil | 2020-01 | 2 |
| 4 | Czechia | 2020-01 | 2 |
| 5 | France | 2020-01 | 1 |

The figures below 4.8 and 4.9 show the clusters of the top 50 locations most affected by the pandemic in March 2020. The clusters are represented by different colours.

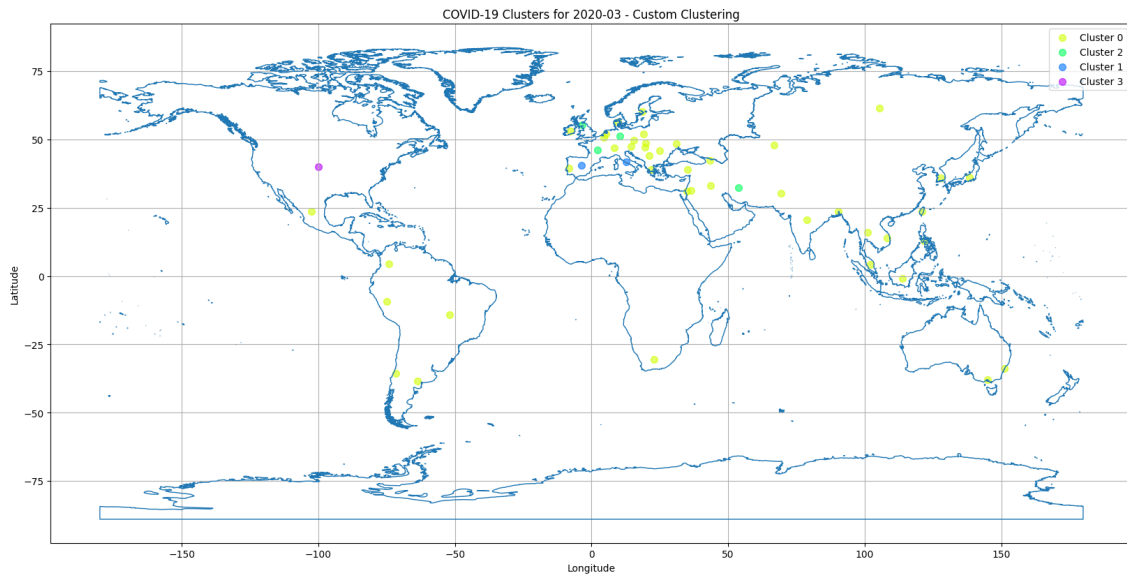


Figure 4.8: Custom KMeans Clustering on 03/2020

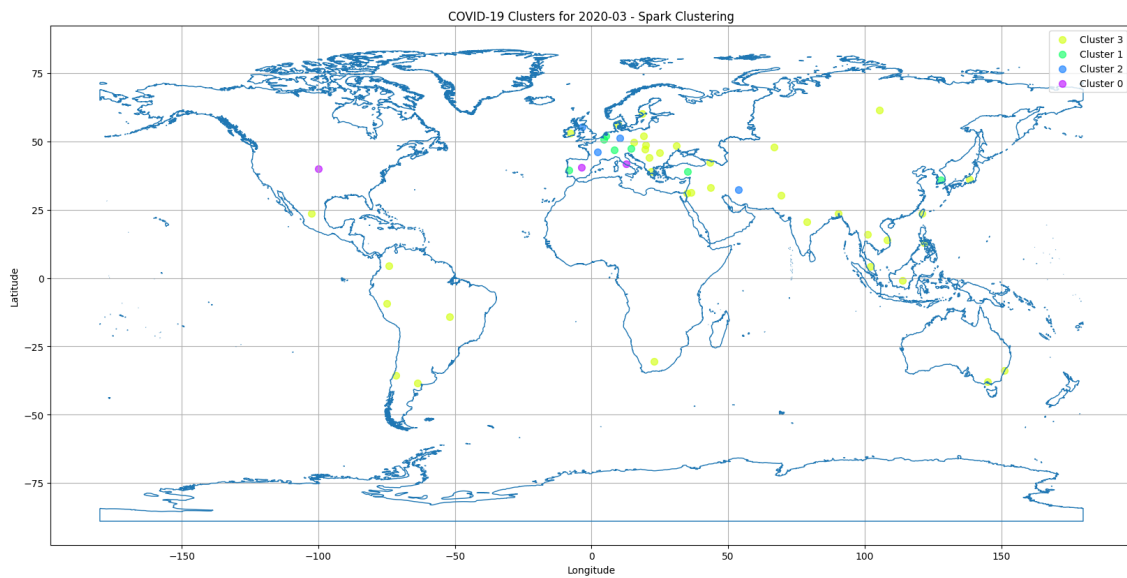


Figure 4.9: Spark MLlib KMeans Clustering on 03/2020

4.2 Discussion of Results

As the previous results show, it seems that scripts not using Spark are faster. While these differences in execution time may seem surprising at first glance, they can in fact be attributed to several factors:

1. **Data size and structure:** The DataFrame, with its 289 entries and 1147 columns occupying 2.5 MB of memory, is relatively modest in size. Pandas is particularly efficient at processing such large amounts of data in memory on a single node. Spark, on the other hand, is designed for distributed processing of large datasets. In this case, the overheads of distributing the data and managing the Spark environment may outweigh the benefits of using it for smaller datasets.
2. **Operational efficiency:** Pandas performs vectorised operations that are optimised for speed, especially with datasets that fit easily into memory on a single computer. Spark, while powerful for processing large volumes of data, introduces an initial overhead for distributing data and configuring the distributed environment, which can slow down processing for smaller datasets.
3. **Complexity of the environment:** Complexity of the environment: Running operations in a Spark environment involves initializing a cluster (even in local mode), distributing tasks, and managing distributed memory, which adds extra processing time compared to running in-memory Pandas directly.

Thus in this project run locally, Pandas is faster than Spark, due to its efficient management of in-memory operations on a single node. Spark's distributed processing overhead makes it less efficient for such tasks.

4.3 Ethical Considerations and Challenges

In the context of the COVID-19 pandemic, the use of Machine Learning and Big Data, while beneficial for crisis management, raises complex ethical issues and challenges.

Firstly, the management of health data, which is essential for tracking and preventing the spread of the virus, creates major risks for privacy. preventing the spread of the virus, creates major privacy risks. Massive data collection and analysis can lead to intrusive surveillance, where the boundaries between public security and individual privacy become blurred. This can constitute an infringement of the right to privacy. In addition, some countries have taken the opportunity to monitor their populations. For example, China has used facial recognition technology to track people's movements and apply quarantine measures.

Secondly, algorithmic biases represent a major challenge. Machine learning systems, although powerful, can incorporate and amplify existing biases in the data. In the context of the pandemic, this could result in inequitable distribution of medical resources, biased diagnoses or discriminatory public health policies, disproportionately affecting certain

social or ethnic groups.

Thirdly, transparency and accountability are crucial issues. Complex and often opaque algorithms make it difficult to understand and challenge AI-based decisions. This raises questions of governance and regulation: who is responsible for errors or harm caused by automated decisions? How can these technologies be effectively controlled and regulated?

Finally, there is the overall challenge of striking a balance between the benefits of using these technologies and respect for fundamental rights. Striking the right balance between effective management of the pandemic and protection of individual freedoms is a delicate exercise, requiring careful ethical reflection and judicious regulation.

Chapter 5

Conclusion

This study demonstrated the effectiveness of Big Data and Machine Learning technologies in the analysis and management of the COVID-19 pandemic. By exploiting the potential of Apache Spark to process large datasets, the analysis revealed significant trends in the spread of the virus, providing essential information for public health decisions. The methodologies adopted provided an in-depth understanding of the dynamics of the pandemic, highlighting geographical and temporal variations in confirmed cases. At the same time, the ethical challenges raised, such as privacy protection and algorithmic fairness, underline the need for a balanced and regulated approach to the use of information technologies. In short, this project illustrates how the judicious integration of AI and Big Data can transform our response to health emergencies, while reminding us of the importance of ethical responsibility in the application of technological advances.

References

1. ESRI. World Continents. ArcGIS Hub; 2013. Available at: <https://hub.arcgis.com/datasets/esri::world-continents/about>. Last Updated: February 9, 2023, (Accessed: October 29, 2023).
2. GISGeography. World Map with Latitudes and Longitudes; 2023. Available at: <https://gisgeography.com/world-map-with-latitudes-and-longitudes/>. Last Updated: October 29, 2023, (Accessed: October 29, 2023).
3. Luna JC. GeoPandas Tutorial: An Introduction to Geospatial Analysis. DataCamp; 2023. Available at: <https://www.datacamp.com/tutorial/geopandas-tutorial-geospatial-analysis>. (Accessed: October 29, 2023).
4. Dabbura I. K-means Clustering: Algorithm, Applications, Evaluation Methods, and Drawbacks. Towards Data Science. 2018 Sep. Available at: <https://towardsdatascience.com/k-means-clustering-algorithm-applications-evaluation-methods-and-drawbacks-aa03e644b48a>. (Accessed: November 8, 2023).
5. Ippolito PP. Apache Spark Optimization Techniques. Towards Data Science. 2023 Jan. Available at: <https://towardsdatascience.com/apache-spark-optimization-techniques-fa7f20a9a2cf>. (Accessed: November 8, 2023).
6. (NNK) N. Spark Performance Tuning & Best Practices; 2023. Available at: <https://sparkbyexamples.com/spark/spark-performance-tuning/>. Last modified: February 7, 2023, (Accessed: November 8, 2023).

Appendix A

Documentation

1. Project tree

```
Covid19AnalysisProjectLib /
  NonOptimised /
    CustomKMeans.py
    CustomQuery1.py
    CustomQuery2.py
    utils.py
  Queries /
    Query1.py
    Query2.py
    Query3.py
  ingestion.py
  Spark.py
dags /
  covid19AnalysisProject.py
data /
  covid_data /
    <last-ingestion-date>.csv
  World_Continents.zip
results /
  figures /
    mean-daily-confirmed-cases-per-month-1.png
    mean-daily-confirmed-cases-per-month-2.png
    mean-daily-confirmed-cases-per-month-3.png
    max-confirmed-cases-by-week-and-continent.png
    mean-confirmed-cases-by-week-and-continent.png
    min-confirmed-cases-by-week-and-continent.png
    std-confirmed-cases-by-week-and-continent.png
    covid-19-clusters-custom.png
    covid-19-clusters-spark.png
  query1 /
    mean_confirmed_cases_per_month_pd.csv
    mean_daily_cases_per_month_spark /
      part-00000-<hash>.csv
```

```
query2 /
    weekly_stats_pd.csv
    weekly_stats_spark /
        part-00000-<hash>.csv
query3 /
    clusters_all_months /
        part-00000-<hash>.csv
    clusters_custom_all_months /
        part-00000-<hash>.csv
main.py
README.md
requirements.txt
setup.py
visualisation.ipynb
```

2. Getting Started

To run the program, follow these steps:

1. Create a virtual environment using `python3 -m venv venv`.
2. Activate the virtual environment using `source venv/bin/activate`.
3. Install the required dependencies using `pip3 install -r requirements.txt`.
4. Run the program using `python3 main.py`.
5. Visualise the results using `visualisation.ipynb` (Jupyter Notebook).

3. Detailed Features of Classes and Functions

Classes

Spark: Class to create a Spark session and load the data set into a DataFrame. This class has the following methods:

- `__init__(self, appName, master)`: Constructor to initialise the class with an application name and a master URL.
- `getSpark(self)`: Method to return the Spark session.
- `getSparkDf(self, path)`: Method to load the data set into a DataFrame and return it, given a path.
- `stopSpark(self)`: Method to stop the Spark session.

Query1: Class to perform Query 1, which is to calculate the average number of confirmed daily cases of people infected with COVID-19 for each country in the dataset, and save the results to a CSV file, using Spark. This class has the following methods:

- `__init__(self, sparkSession, covidDataDf)`: Constructor to initialise the class with a Spark session and a DataFrame containing COVID-19 data.
- `run(self)`: Method to run the query.

Query2: Class to perform Query 2, which is to calculate the mean, standard deviation, minimum and maximum of the number of confirmed cases daily for each week and continent, and save the results to a CSV file, using Spark. This class has the following methods:

- `__init__(self, sparkSession, covidDataDf)`: Constructor to initialise the class with a Spark session and a DataFrame containing COVID-19 data.
- `run(self)`: Method to run the query.

Query3: Class to perform Query 3, which is to perform clustering on the top 50 affected locations based on the maximum slope of monthly confirmed cases, and save the results to a CSV file, using Spark. This class has the following methods:

- `__init__(self, sparkSession, covidDataDf)`: Constructor to initialise the class with a Spark session and a DataFrame containing COVID-19 data.
- `run(self)`: Method to run the query.

CustomKMeans: Class to perform a custom KMeans clustering algorithm. This class has the following methods:

- `__init__(self, k, max_iter)`: Constructor to initialise the class with the number of clusters and the maximum number of iterations.
- `fit(self, X)`: Method to fit the model to the data.
- `predict(self, X)`: Method to predict the cluster of each data point.

CustomQuery1: Class to perform Query1, using Pandas without Spark. This class has the following methods:

- `__init__(self, covidDataDf)`: Constructor to initialise the class with a DataFrame containing COVID-19 data.
- `run(self)`: Method to run the query.

CustomQuery2: Class to perform Query2, using Pandas without Spark. This class has the following methods:

- `__init__(self, covidDataDf)`: Constructor to initialise the class with a DataFrame containing COVID-19 data.
- `run(self)`: Method to run the query.

Functions

`ingestion.py`

- `fetch_covid_data()`: Function to ingest the data set from the source and store it in the data lake.

`utils.py`

- `getContinentByCoordinates(longitude, latitude)`: Function used by CustomQuery2 class to assign a continent to a location based on its coordinates.
- `compute_slope(row)`: Function used by CustomQuery2 class to compute the slope of a location.
- `format_date_range(date)`: Function used by CustomQuery2 class to format a date range.

Appendix B

Query 1 Source Code

```
1 from pyspark.sql import functions as F
2 from pyspark.sql.window import Window
3 import time
4
5
6 class Query1:
7     """
8     Class to run Query 1.
9
10    Parameters:
11        sparkSession (SparkSession): Spark Session
12        covidDataDf (DataFrame): Covid DataFrame
13
14    Attributes:
15        sparkSession (SparkSession): Spark Session
16        covidDataDf (DataFrame): Covid DataFrame
17    """
18
19    def __init__(self, sparkSession, covidDataDf):
20        self.sparkSession = sparkSession # Spark Session
21        self.covidDataDf = covidDataDf # Covid DataFrame
22
23    def run(self):
24        """
25        Runs Query 1 and writes the result to a CSV file.
26        """
27        # Start the timer
28        startTime = time.time()
29
30        # Get the date columns (Remove State/Province, Country/
31        # Region, Lat, Long columns)
32        date_cols = self.covidDataDf.columns[4:]
33
34        # Create the stack expression to pivot the table to make
35        # the date columns into rows
36        stack_expr = (
37            "stack("
38            + str(len(date_cols))
39            + ", "
```

```
38         + ", ".join(["'" + x + "'", '"' + x + '"'] for x in
date_cols]))
39         + ") as (Date, Value)"
40     )
41
42     # Pivot the table to make the date columns into rows
43     covidDataDf_pivot = self.covidDataDf.select(
44         "Country/Region", F.expr(stack_expr)
45     ).withColumn(
46         "Date", F.to_date("Date", "M/d/yy")
47     ) # Convert the Date column to a date type
48
49     # Window specification to get the previous day's cases for
each country
50     windowSpec = Window.partitionBy("Country/Region").orderBy(
("Date")
51
52     # Calculate daily confirmed cases
53     covidDataDf_daily = (
54         covidDataDf_pivot.withColumn(
55             "DailyCases", F.col("Value") - F.lag("Value").over(
windowSpec)
56         )
57         .fillna({"DailyCases": 0}) # Fill the null values with
0
58         .drop("Value") # Drop the Value column
59     )
60
61     # Cache the dataframe
62     covidDataDf_daily.cache()
63
64     # Group by country and month to get average daily cases
65     df_grouped = (
66         covidDataDf_daily.groupBy(
67             "Country/Region",
68             F.year("Date").alias("Year"),
69             F.month("Date").alias("Month"),
70         ) # Group by country, year, and month
71         .agg(
72             F.avg("DailyCases").alias("Average")
73         ) # Compute the average daily cases
74         .orderBy(
75             "Country/Region", "Year", "Month"
76         ) # Order by country, year, and month
77     )
78
79     # Unpersist the dataframe
80     covidDataDf_daily.unpersist()
81
82     # Write the result to a CSV file
83     try:
84         df_grouped.write.csv(
85             "results/query1/mean_daily_cases_per_month_spark",
86             header=True,
87             mode="overwrite",
88         )
```


Appendix B. Query 1 Source Code

```
89     except Exception as e:
90         print(f"Failed to write the file: {e}")
91
92     # Stop the timer
93     endTime = time.time()
94     print(f"Query 1 took {endTime - startTime} seconds.")
```

Appendix C

Query 2 Source Code

```
1 import geopandas as gpd
2 from shapely.geometry import Point
3 from pyspark.sql import functions as F
4 from pyspark.sql.types import StringType
5 from pyspark.sql.window import Window
6 import time
7
8
9 class Query2:
10     """
11     Class to run Query 2.
12
13     Parameters:
14         sparkSession (SparkSession): Spark Session
15         covidDataDf (DataFrame): Covid DataFrame
16
17     Attributes:
18         sparkSession (SparkSession): Spark Session
19         covidDataDf (DataFrame): Covid DataFrame
20         world_boundaries_gdf (GeoDataFrame): Shapefile of world
21         boundaries (continents)
22     """
23
24     def __init__(self, sparkSession, covidDataDf):
25         self.sparkSession = sparkSession # Spark Session
26         self.covidDataDf = covidDataDf # Covid DataFrame
27         self.world_boundaries_gdf = gpd.read_file(
28             "data/World_Continents.zip"
29         ) # Shapefile of world boundaries (continents)
30
31     def __prepare_data(self):
32         """
33         Prepares the COVID data for processing.
34         """
35         # Location column : Province/State if not null, else
36         # Country/Region
37         self.covidDataDf = (
38             self.covidDataDf.withColumn(
39                 "Location", F.coalesce(F.col("Province/State"), F.
40                     col("Country/Region"))
```

```

38         ) # Coalesce the Province/State and Country/Region
columns
39         .filter(
40             F.col("Lat").isNotNull() & F.col("Long").isNotNull
41         )
42         ) # Filter out the rows with null values in Lat and
Long columns
43         .drop(
44             "Province/State", "Country/Region"
45         ) # Drop the Province/State and Country/Region columns
46     )
47     def __assign_continent(self):
48         """
49         Assign continents to each data row based on coordinates.
50         """
51         # Broadcast the world boundaries dataframe to all the nodes
in the cluster for faster processing
52         broadcasted_gdf = self.sparkSession.sparkContext.broadcast(
53             self.world_boundaries_gdf
54         )
55
56         def determine_continent(longitude, latitude):
57             """
58             Determine the continent of a given set of coordinates.
59
60             Parameters:
61                 longitude (double): Longitude of the location (
in degrees)
62                 latitude (double): Latitude of the location (in
degrees)
63
64             Returns:
65                 (str): Continent of the location
66             """
67             point = Point(
68                 longitude, latitude
69             ) # Create a Point object from the given coordinates
70             matching_continents = broadcasted_gdf.value[
71                 broadcasted_gdf.value.geometry.contains(point)
72             ][
73                 "CONTINENT"
74             ] # Get the matching continents from the world
boundaries geodataframe
75             if len(matching_continents) > 0:
76                 # If the continent is Australia, return Oceania
77                 if matching_continents.values[0] == "Australia":
78                     return "Oceania"
79                 # If the continent is North America or South
America, return America
80                 if (
81                     matching_continents.values[0] == "North America
"
82                     or matching_continents.values[0] == "South
America"
83                 ):

```

```

84         return "America"
85         return matching_continents.values[0]
86         # If the continent is not found, determine the
continent based on the coordinates (Approximate)
87         elif -34 <= latitude <= 37 and -17 <= longitude <= 51:
88             return "Africa"
89         elif 34 <= latitude <= 82 and -25 <= longitude <= 60:
90             return "Europe"
91         elif -56 <= latitude <= 71 and -168 <= longitude <=
-34:
92             return "America"
93             elif -1 <= latitude <= 81 and (25 <= longitude <= 171):
94                 return "Asia"
95             elif -50 <= latitude <= 10 and 110 <= longitude <= 180:
96                 return "Oceania"
97             elif -90 <= latitude <= -60:
98                 return "Antarctica"
99             return None
100
101         # Create a UDF to determine the continent of a given set of
coordinates
102         udf_determine_continent = F.udf(
103             lambda long, lat: determine_continent(long, lat),
StringType()
104         )
105
106         # Assign continents to each data row based on coordinates
107         self.covidDataDf = self.covidDataDf.withColumn(
108             "Continent", udf_determine_continent(F.col("Long"), F.
col("Lat")))
109         ).filter(
110             F.col("Continent").isNotNull()
111         ) # Filter out the rows with null values in the Continent
column
112
113         def __pivot_table(self):
114             """
115             Pivot the table to make the date columns into rows.
116             """
117             # Get the date columns (Remove Location, Continent, Lat,
Long columns)
118             date_cols = self.covidDataDf.columns[2:-2]
119
120             # Create the stack expression to pivot the table to make
the date columns into rows
121             stack_expr = (
122                 "stack("
123                 + str(len(date_cols))
124                 + ", "
125                 + ", ".join(["'" + x + "'", '" + x + "'" for x in
date_cols])
126                 + ") as (Date, Value)"
127             )
128
129             # Pivot the table to make the date columns into rows
130             self.covidDataDf = self.covidDataDf.select(

```

```

131         "Location", "Continent", F.expr(stack_expr)
132     ).withColumn(
133         "Date", F.to_date("Date", "M/d/yy")
134     ) # Convert the Date column to a date type
135
136     def __compute_daily_confirmed_cases(self):
137         # Window specification to get the previous day's cases for
138         each country
139         windowSpec = Window.partitionBy("Location").orderBy("Date")
140
141         # Calculate daily confirmed cases
142         self.covidDataDf = (
143             self.covidDataDf.withColumn(
144                 "DailyCases", F.col("Value") - F.lag("Value").over(
145                     windowSpec)
146             ) # Calculate the daily cases
147             .fillna({"DailyCases": 0}) # Fill the null values with
148             0
149             .drop("Value") # Drop the Value column
150         )
151
152     def __compute_slopes(self):
153         """
154         Compute the slopes of the daily confirmed cases for each
155         location.
156         """
157         # Determine the start date for computing DaysSinceStart
158         start_date = self.covidDataDf.agg(F.min("Date")).first()[0]
159         # Start date
160         days_since_start = F.datediff(
161             F.col("Date"), F.lit(start_date)
162         ) # Days since start
163
164         # Prepare the data for linear regression
165         windowSpec = Window.partitionBy("Location")
166         regression_df = (
167             self.covidDataDf.withColumn(
168                 "DaysSinceStart", days_since_start
169             ) # Days since start (x)
170             .withColumn(
171                 "DaysSquared", F.pow(F.col("DaysSinceStart"), 2)
172             ) # Days squared
173             .withColumn("n", F.col("DaysSinceStart") + 1) # Number
174             of observations (n)
175             .withColumn(
176                 "CasesTimesDays", F.col("DaysSinceStart") * F.col("
177                 DailyCases")
178             ) # Cases * Days
179             .withColumn(
180                 "sumDays", F.sum("DaysSinceStart").over(windowSpec)
181             ) # Sum of DaysSinceStart
182             .withColumn(
183                 "sumCases", F.sum("DailyCases").over(windowSpec)
184             ) # Sum of DailyCases
185             .withColumn(

```

```

179         "sumDaysSquared", F.sum("DaysSquared").over(
windowSpec)
180     ) # Sum of DaysSquared
181     .withColumn(
182         "sumCasesTimesDays", F.sum("CasesTimesDays").over(
windowSpec)
183     ) # Sum of CasesTimesDays
184 )
185
186 # Compute the slope of the linear regression
187 numerator = F.col("n") * F.col("sumCasesTimesDays") - F.col(
("sumDays") * F.col(
188     "sumCases"
189 ) # Numerator of the slope
190 denominator = F.col("n") * F.col("sumDaysSquared") - F.pow(
F.col("sumDays"), 2
191 ) # Denominator of the slope
192 self.covidDataDf = regression_df.withColumn(
193     "DailySlope",
194     F.when(denominator == 0, 0).otherwise(
195         numerator / denominator
196     ), # Daily slope
197 ).select(
198     "Location", "Date", "DailyCases", "DailySlope", "
Continent"
199 ) # Select the required columns
200
201
202 def __filter_top_affected(self):
203     """
204     Select the top 100 affected locations.
205     """
206     # Aggregate the slopes by location, calculating the maximum
slope
207     top100LocationsDf = (
208         self.covidDataDf.groupBy("Location") # Group by
location
209         .agg(F.max("DailySlope").alias("MaxSlope")) # Maximum
slope
210         .orderBy(F.desc("MaxSlope")) # Order by the maximum
slope
211         .limit(100) # Select the top 100 locations
212     )
213
214     # Join the top 100 countries back to the original dataframe
215     self.covidDataDf = self.covidDataDf.join(
216         F.broadcast(
217             top100LocationsDf
218         ), # Broadcast the top 100 locations dataframe to all
the nodes in the cluster for faster processing
219         ["Location"],
220     )
221
222 def __aggregate(self):
223     """
224     Perform statistics calculations on the top 100 affected
locations.

```

```

225     """
226     # Group by continent and week
227     self.covidDataDf = (
228         self.covidDataDf.withColumn(
229             "WeekStart",
230             F.date_sub(
231                 F.col("Date"), F.dayofweek(F.col("Date")) - 1
232             ), # Starting date of the week
233         )
234         .withColumn(
235             "WeekEnd", F.date_add(F.col("WeekStart"), 6) #
Ending date of the week
236         )
237         .withColumn(
238             "WeekRange",
239             F.concat(
240                 F.date_format(F.col("WeekStart"), "dd/MM/yyyy")
241                 ,
242                 F.lit(" - "),
243                 F.date_format(F.col("WeekEnd"), "dd/MM/yyyy"),
244             ),
245         )
246
247     # Aggregate the data by continent and week
248     weekly_stats = (
249         self.covidDataDf.groupBy("Continent", "WeekRange")
250         .agg(
251             F.mean("DailyCases").alias("Mean"), # Mean of
DailyCases
252             F.stddev("DailyCases").alias("Std"), # Standard
deviation of DailyCases
253             F.min("DailyCases").alias("Min"), # Minimum of
DailyCases
254             F.max("DailyCases").alias("Max"), # Maximum of
DailyCases
255         )
256         .orderBy(
257             "Continent", F.min("WeekStart")
258         ) # Order by continent and week start date
259     )
260
261     # Write the aggregated data to a CSV file
262     try:
263         weekly_stats.write.csv(
264             "results/query2/weekly_stats_spark", header=True,
mode="overwrite"
265         )
266     except Exception as e:
267         print(f"Failed to write the file: {e}")
268
269     def run(self):
270         """
271         Runs the query 2 and writes the results to a CSV file.
272         """
273         startTime = time.time()

```

```
274         self.__prepare_data()
275         self.__assign_continent()
276         self.__pivot_table()
277         self.covidDataDf.cache() # Cache the dataframe in memory
    for faster processing
278         self.__compute_daily_confirmed_cases()
279         self.__compute_slopes()
280         self.__filter_top_affected()
281         self.__aggregate()
282         self.covidDataDf.unpersist() # Unpersist the dataframe
    from memory
283         endTime = time.time()
284         print(f"Query 2 took {endTime - startTime} seconds to
    complete.")
```


Appendix D

Query 3 Source Code

```
1 from pyspark.sql import functions as F
2 from pyspark.sql.window import Window
3 from pyspark.ml.feature import VectorAssembler
4 from pyspark.ml.clustering import KMeans
5
6 import numpy as np
7 import time
8
9 from Covid19AnalysisProjectLib.NonOptimised.CustomKMeans import
   CustomKMeans
10
11
12 class Query3:
13     """
14     Class to run Query 3.
15
16     Parameters:
17         sparkSession (SparkSession): Spark Session
18         covidDataDf (DataFrame): Covid DataFrame
19
20     Attributes:
21         sparkSession (SparkSession): Spark Session
22         covidDataDf (DataFrame): Covid DataFrame
23     """
24
25     def __init__(self, sparkSession, covidDataDf):
26         self.sparkSession = sparkSession # Spark Session
27         self.covidDataDf = covidDataDf # Covid DataFrame
28
29     def __prepare_data(self):
30         """
31         Prepares the COVID data for processing.
32         """
33         # Location column : Province/State if not null, else
   Country/Region
34         self.covidDataDf = (
35             self.covidDataDf.withColumn(
36                 "Location",
37                 F.coalesce(
```

```

38         F.col("Province/State"), F.col("Country/Region"
39     )
40     ), # If Province/State is null, use Country/Region
41     )
42     .filter(
43         F.col("Lat").isNotNull() & F.col("Long").isNotNull
44     )
45     ) # Filter out rows with null Lat and Long values
46     .drop(
47         "Province/State", "Country/Region", "Lat", "Long"
48     ) # Drop unnecessary columns
49
50 def __pivot_table(self):
51     """
52     Pivots the table to make the date columns into rows.
53     """
54     # Get the date columns (Remove Location column)
55     date_cols = self.covidDataDf.columns[0:-1]
56
57     # Create the stack expression to pivot the table to make
58     the date columns into rows
59     stack_expr = (
60         "stack("
61         + str(len(date_cols))
62         + ", "
63         + ", ".join(["'" + x + "'", '" + x + "' for x in
64 date_cols])
65         + ") as (Date, Value)"
66     )
67
68     # Pivot the table to make the date columns into rows
69     self.covidDataDf = self.covidDataDf.select(
70         "Location", F.expr(stack_expr)
71     ).withColumn(
72         "Date", F.to_date("Date", "M/d/yy")
73     ) # Convert the Date column to a date type
74
75 def __compute_daily_confirmed_cases(self):
76     """
77     Compute the daily confirmed cases for each location.
78     """
79     # Window specification to get the previous day's cases for
80     each location
81     windowSpec = Window.partitionBy("Location").orderBy("Date")
82     # Calculate daily confirmed cases
83     self.covidDataDf = (
84         self.covidDataDf.withColumn(
85             "DailyCases", F.col("Value") - F.lag("Value").over(
86 windowSpec)
87         ) # Calculate the difference between the current and
88         previous day's cases
89         .fillna({"DailyCases": 0}) # Fill the null values with
90         0
91         .drop("Value") # Drop the Value column
92     )

```

```

86
87     def __compute_monthly_slopes(self):
88         """
89         Compute the slopes of the monthly confirmed cases for each
location.
90         """
91         # Extract the month from the date
92         self.covidDataDf = self.covidDataDf.withColumn(
93             "Month", F.date_format("Date", "yyyy-MM")
94         )
95
96         # Determine the start date for computing MonthsSinceStart
97         start_date = self.covidDataDf.agg(F.min("Date")).first()[0]
98         months_since_start_expr = (F.year("Date") - F.year(F.lit(
start_date))) * 12 + (
99             F.month("Date") - F.month(F.lit(start_date))
100         ) # Calculate the number of months since the start date
101         self.covidDataDf = self.covidDataDf.withColumn(
102             "MonthsSinceStart", months_since_start_expr
103         ) # Add the MonthsSinceStart column
104
105         # Calculate monthly total cases for each location and
include 'MonthsSinceStart'
106         monthly_totals = self.covidDataDf.groupBy(
107             "Location", "Month", "MonthsSinceStart"
108         ).agg(F.sum("DailyCases").alias("MonthlyCases"))
109
110         # Window specification to get the cumulative sum of
MonthsSinceStart and MonthlyCases for each location
111         window_spec = Window.partitionBy("Location").orderBy("Month
")
112
113         # Prepare the dataframe for linear regression
114         regression_df = (
115             monthly_totals.withColumn(
116                 "x", F.col("MonthsSinceStart")
117             ) # MonthsSinceStart (x)
118             .withColumn(
119                 "MonthsSquared", F.pow(F.col("MonthsSinceStart"),
2)
120             ) # MonthsSinceStart squared
121             .withColumn(
122                 "n", F.col("MonthsSinceStart") + 1
123             ) # Number of observations (n)
124             .withColumn(
125                 "CasesTimesMonths",
126                 F.col("MonthsSinceStart")
127                 * F.col("MonthlyCases"), # MonthlyCases *
MonthsSinceStart
128             )
129             .withColumn(
130                 "sumMonths", F.sum("MonthsSinceStart").over(
window_spec)
131             ) # Cumulative sum of MonthsSinceStart
132             .withColumn(
133                 "sumCases", F.sum("MonthlyCases").over(window_spec)
134             ) # Cumulative sum of MonthlyCases

```

```

134         .withColumn(
135             "sumMonthsSquared", F.sum("MonthsSquared").over(
window_spec)
136         ) # Cumulative sum of MonthsSquared
137         .withColumn(
138             "sumCasesTimesMonths",
139             F.sum("CasesTimesMonths").over(
140                 window_spec
141             ), # Cumulative sum of MonthlyCases *
MonthsSinceStart
142         )
143     )
144
145     # Calculate the slope for each location
146     numerator = F.col("n") * F.col("sumCasesTimesMonths") - F.
col(
147         "sumMonths"
148     ) * F.col("sumCases")
149     denominator = F.col("n") * F.col("sumMonthsSquared") - F.
pow(
150         F.col("sumMonths"),
151         2,
152     )
153     self.covidDataDf = regression_df.withColumn(
154         "MonthlySlope",
155         F.when(denominator == 0, 0).otherwise(numerator /
denominator),
156     ).select("Location", "Month", "MonthlyCases", "MonthlySlope
")
157
158     def __filter_top_affected(self):
159         """
160         Select the top 50 affected locations.
161         """
162         # Aggregate the slopes by location, calculating the maximum
or average slope
163         top50LocationsDf = (
164             self.covidDataDf.groupBy("Location")
165             .agg(
166                 F.mean("MonthlySlope").alias("MeanSlope")
167             ) # Get the maximum slope for each location
168             .orderBy(F.desc("MeanSlope")) # Order by descending
slope
169             .limit(50) # Select the top 50 locations
170         )
171
172     # Join the top 50 locations back to the original dataframe
173     self.covidDataDf = (
174         self.covidDataDf.join(
175             F.broadcast(
176                 top50LocationsDf
177             ), # Broadcast the top 50 locations dataframe to
all nodes for efficient join
178             ["Location"],
179         )
180         .drop(

```

```

181         "Date", "MonthsSinceStart", "DailyCases"
182     ) # Drop the unnecessary columns
183     .dropDuplicates(["Location", "Month"]) # Drop
duplicate rows
184     .orderBy("Month", "Location") # Order by month and
location
185     )
186
187     def __apply_custom_clustering(self):
188         """
189         Apply custom clustering to the data for each month using
the CustomKMeans model.
190         """
191         # Start the timer
192         start_time_custom = time.time()
193
194         # Initialize CustomKMeans model
195         custom_kmeans = CustomKMeans(n_clusters=4, max_iter=20,
seed=1, tol=1e-4)
196
197         # Processing each month in parallel
198         all_clusters = []
199         unique_months = self.covidDataDf.select("Month").distinct()
.collect()
200         for month_row in unique_months:
201             # Filter the data for the current month
202             month = month_row["Month"]
203             monthly_data = self.covidDataDf.filter(F.col("Month")
== month)
204
205             # Convert the MonthlySlope column to a numpy array
206             slope_data = np.array(
207                 monthly_data.select("MonthlySlope").collect()
208             ).reshape(-1, 1)
209
210             # Train and predict using the custom KMeans model
211             custom_kmeans.fit(slope_data)
212             cluster_assignments = custom_kmeans.predict(slope_data)
213
214             # Create a dataframe with the cluster assignments
215             monthly_data = monthly_data.withColumn(
216                 "row_index", F.monotonically_increasing_id()
217             )
218             clusters_df = self.sparkSession.createDataFrame(
219                 [
220                     (int(i), int(cluster_assignments[i]))
221                     for i in range(len(cluster_assignments))
222                 ],
223                 ["row_index", "Cluster"],
224             )
225
226             # Join the cluster assignments to the original
dataframe
227             clusters = monthly_data.join(clusters_df, "row_index").
select(
228                 "Location", "Month", "Cluster"

```

```
229         )
230
231         # Append to all_clusters list
232         all_clusters.append(clusters)
233
234         # Union all clusters and repartition for efficient writing
235         final_clusters_df = self.sparkSession.createDataFrame(
236             self.sparkSession.sparkContext.emptyRDD(), all_clusters
237             [0].schema
238         )
239         for clusters in all_clusters:
240             final_clusters_df = final_clusters_df.union(clusters)
241 # Union all clusters
242
243         # Repartition by month for efficient writing to a single
244 CSV file
245         final_clusters_df = final_clusters_df.repartition("Month")
246
247         # Write the aggregated results to a single CSV file
248 try:
249     final_clusters_df.write.csv(
250         "results/query3/clusters_custom_all_months",
251         header=True,
252         mode="overwrite",
253     )
254 except Exception as e:
255     print(f"Failed to write the file: {e}")
256
257     print(
258         f"Clustering completed by custom implementation in {
259 time.time() - start_time_custom} seconds"
260     )
261
262 def __apply_clustering(self):
263     """
264     Apply clustering to the data for each month using Spark
265 MLlib.
266     """
267     # Start the timer
268     start_time_spark = time.time()
269
270     # Initialize VectorAssembler and KMeans model outside the
271 loop
272     vec_assembler = VectorAssembler(
273         inputCols=["MonthlySlope"], outputCol="features"
274     )
275     kmeans = (
276         KMeans()
277         .setK(4) # Number of clusters
278         .setSeed(1) # Random seed
279         .setFeaturesCol("features") # Input features
280         .setPredictionCol("Cluster") # Output cluster
281     )
282
283     # Vectorise the MonthlySlope column
```

```

278         self.covidDataDf = vec_assembler.transform(self.covidDataDf
279     )
280     # Apply clustering to the data for each month using Spark
281     MLlib
282     all_clusters = []
283     unique_months = self.covidDataDf.select("Month").distinct()
284     .collect()
285     for month_row in unique_months:
286         # Filter the data for the current month
287         month = month_row["Month"]
288         monthly_data = self.covidDataDf.filter(F.col("Month")
289     == month)
290
291         # Train and predict using the Spark MLlib KMeans model
292         model = kmeans.fit(monthly_data)
293         clusters = model.transform(monthly_data).select(
294             "Location", "Month", "Cluster"
295         )
296
297         # Append to all_clusters list
298         all_clusters.append(clusters)
299
300     # Union all clusters and repartition for efficient writing
301     final_clusters_df = self.sparkSession.createDataFrame(
302         self.sparkSession.sparkContext.emptyRDD(), all_clusters
303     [0].schema
304     )
305     for clusters in all_clusters:
306         final_clusters_df = final_clusters_df.union(clusters)
307     # Union all clusters
308
309     # Repartition by month for efficient writing to a single
310     CSV file
311     final_clusters_df = final_clusters_df.repartition("Month")
312
313     # Write the aggregated results to a single CSV file
314     try:
315         final_clusters_df.write.csv(
316             "results/query3/clusters_all_months", header=True,
317             mode="overwrite"
318         )
319     except Exception as e:
320         print(f"Failed to write the file: {e}")
321
322     # Stop the timer
323     stop_time_spark = time.time()
324     print(f"Clustering completed by Spark MLlib in {
325     stop_time_spark - start_time_spark} seconds")
326
327     def run(self):
328         """
329         Runs the query 3 and writes the results to a CSV file.
330         """
331         # Execute methods in sequence
332         self.__prepare_data()

```

```
325         self.__pivot_table()
326         self.__compute_daily_confirmed_cases()
327         self.__compute_monthly_slopes()
328         self.covidDataDf.cache() # Cache the dataframe to speed up
processing
329         self.__filter_top_affected()
330         self.__apply_custom_clustering()
331         self.__apply_clustering()
332         self.covidDataDf.unpersist() # Unpersist the dataframe
```


Appendix E

Terminal Output

```
1 Setting default log level to "WARN".
2 To adjust logging level use sc.setLogLevel(newLevel). For SparkR,
  use setLogLevel(newLevel).
3 23/11/18 18:04:06 WARN NativeCodeLoader: Unable to load native-
  hadoop library for your platform... using builtin-java classes
  where applicable
4
5 ----- STEP 1: FETCHING THE LATEST DATA -----
6
7 1. Deleting the existing data...
8   -> Deleting file: 2023-11-18.csv
9 Done deleting the existing data.
10
11 2. Fetching the latest data...
12 Done fetching the latest data.
13
14 3. Writing the data to a CSV file: 2023-11-18.csv
15 Done writing the data to a CSV file.
16
17
18 ----- STEP 2: LOADING THE DATA INTO A SPARK DATAFRAME
  -----
19
20 Done loading the data into a Spark dataframe.
21
22
23 ----- STEP 3: Query 1 -----
24
25 23/11/18 18:04:11 WARN SparkStringUtils: Truncated the string
  representation of a plan since it was too large. This behavior
  can be adjusted by setting 'spark.sql.debug.maxToStringFields'.
26 Query 1 took 3.1029131412506104 seconds.
27
28 ----- STEP 4: Query 2 -----
29
30 23/11/18 18:04:20 WARN GarbageCollectionMetrics: To enable non-
  built-in garbage collector(s) List(G1 Concurrent GC), users
  should configure it(them) to spark.eventLog.gcMetrics.
  youngGenerationGarbageCollectors or spark.eventLog.gcMetrics.
  oldGenerationGarbageCollectors
```

Appendix E. Terminal Output

```
31 Query 2 took 15.06498098373413 seconds to complete.
32
33 ----- STEP 5: Query 3 -----
34
35 Clustering completed by custom implementation in 58.84758687019348
   seconds
36 23/11/18 18:05:32 WARN InstanceBuilder: Failed to load
   implementation from:dev.ludovic.netlib.blas.JNIBLAS
37 23/11/18 18:05:32 WARN InstanceBuilder: Failed to load
   implementation from:dev.ludovic.netlib.blas.VectorBLAS
38 23/11/18 18:07:17 WARN DAGScheduler: Broadcasting large task binary
   with size 2020.2 KiB
39 Clustering completed by Spark MLlib in 109.29923582077026 seconds
```