

**Big Data Platforms**  
**Level: M**  
**Module Code:MMI227050**  
**Coursework 2 (Resit)**

**Cloud data pipeline Design Report**  
**MAHEK KATHIRIYA**

## Contents

<b>Part B. Cloud Data Pipeline Design Report</b> .....	2
<b>1. Overall Concept Design</b> .....	2
<b>1.1 Data Format and Schema</b> .....	2
<b>1.2 Transformations</b> .....	2
<b>1.3 Analyses and Visualizations</b> .....	2
<b>Step-by-Step Cloud Data Pipeline Design</b> .....	3
<b>2. Platforms</b> .....	3
<b>2.1 Data Store: Google Cloud Storage (GCS)</b> .....	3
<b>2.2 File System: Google Cloud Storage (GCS)</b> .....	4
<b>2.3 Analytic Engine: Google Dataproc</b> .....	4
Pyspark code: .....	4
<b>2.4 Visualization: Google Data Studio</b> .....	10
<b>3. Design Illustration</b> .....	10
<b>4. Component Descriptions</b> .....	11
<b>5. Visualization Examples</b> .....	11
<b>6. Justification</b> .....	12
<b>7. References</b> .....	12

## Part B. Cloud Data Pipeline Design Report

The objective of this report was to outline the high-level design that was employed to upscale the initial Databricks data pipeline implementation into a cloud-based architecture. This endeavor predominantly revolved around Google Cloud Platform (GCP) as the designated provider. The overarching ambition was to augment data processing and analysis capacities, while simultaneously upholding principles of scalability, dependability, and cost-efficiency.

### 1. Overall Concept Design

The purpose of this report was to delineate the high-level design that had been employed to elevate the primary Databricks data pipeline implementation into a cloud-based framework. The emphasis in this design was on Google Cloud Platform (GCP) as the chosen provider. The intention was to enhance data processing and analytical capabilities, while also ensuring scalability, dependability, and cost-effectiveness.

#### 1.1 Data Format and Schema

The data in its original incarnation was presented in CSV format. This dataset encompassed specifications of mobile devices, encompassing attributes such as battery power, clock speed, RAM, and more. The schema of the data was composed of features such as "battery\_power," "clock\_speed," "ram," in conjunction with categorical attributes like "blue," "dual\_sim," and more.

#### 1.2 Transformations

In terms of the ETL process, a series of transformations were applied:

- Numeric columns were converted into the double data type.

- Numeric features such as battery\_power, clock\_speed, ram, and mobile\_wt underwent standard scaling, executed through the utilization of Google Dataproc's Spark clusters.
- Categorical attributes, including "blue," "dual\_sim," and others, were subjected to encoding using methodologies like StringIndexer and OneHotEncoder.

#### 1.3 Analyses and Visualizations

The projected analyses and visualizations were tailored towards uncomplicated operations such as filtering, projection, and aggregation. For instance, focal points

encompassed the identification of devices possessing specific attributes (such as high RAM or 4G support), alongside the aggregation of average battery power across diverse categories.

## Step-by-Step Cloud Data Pipeline Design

### 2. Platforms

Each fundamental component within the expanded pipeline necessitated an apt cloud service provider, with careful consideration given to the deployability of these platforms on the provider's infrastructure. The chosen platform for this design endeavor was Google Cloud Platform (GCP). The subsequent components and services were harnessed to drive the pipeline:

#### 2.1 Data Store: Google Cloud Storage (GCS)

Google Cloud Storage emerged as the prime choice for the data store, housing both the original dataset and the refined data. GCS provided a scalable and economical solution for object storage, ensuring both data durability and ease of access.

Google Cloud

BDP-CW2

data studio

Search

4

?

:

M

LEARN Tutorial

?

×

Cloud Storage

Buckets

Monitoring

Settings

Bucket details

REFRESH

HELP ASSISTANT

LEARN

bdp\_cw2\_bucket

Location

eu (multiple regions in European Union)

Storage class

Standard

Public access

Not public

Protection

None

OBJECTS

CONFIGURATION

PERMISSION

PROTECTION

LIFECYCLE

OBSERVABILITY

Buckets > bdp\_cw2\_bucket

UPLOAD FILES

UPLOAD FOLDER

CREATE FOLDER

TRANSFER DATA

MANAGE HOLDS

DOWNLOAD

DELETE

Filter by name prefix only

Filter

Filter objects and folders

Show deleted data

<input type="checkbox"/>	Name	Size	Type	Created	Storage class	
<input type="checkbox"/>	bdp_cw2_partb.py	3.4 KB	text/x-python	18 Aug 2023, 01:15:24	Standard	<div><div>↓</div><div>⋮</div></div>
<input type="checkbox"/>	bdp_visual.py	4.8 KB	text/x-python	18 Aug 2023, 06:51:25	Standard	<div><div>↓</div><div>⋮</div></div>
<input type="checkbox"/>	data/	—	Folder	—	—	<div><div>⋮</div></div>
<input type="checkbox"/>	test.csv	62.4 KB	text/csv	17 Aug 2023, 19:47:03	Standard	<div><div>↓</div><div>⋮</div></div>
<input type="checkbox"/>	train.csv	119.5 KB	text/csv	17 Aug 2023, 19:46:45	Standard	<div><div>↓</div><div>⋮</div></div>

Get started with Cloud Storage

Getting bucket information

Help document

Get information on the size and metadata of your Cloud Storage buckets.

Uploading objects

Help document

Upload the objects containing your data to your Cloud Storage buckets.

Downloading objects

Help document

Download the objects from your Cloud Storage buckets.

Use cases for Cloud Storage

Help document

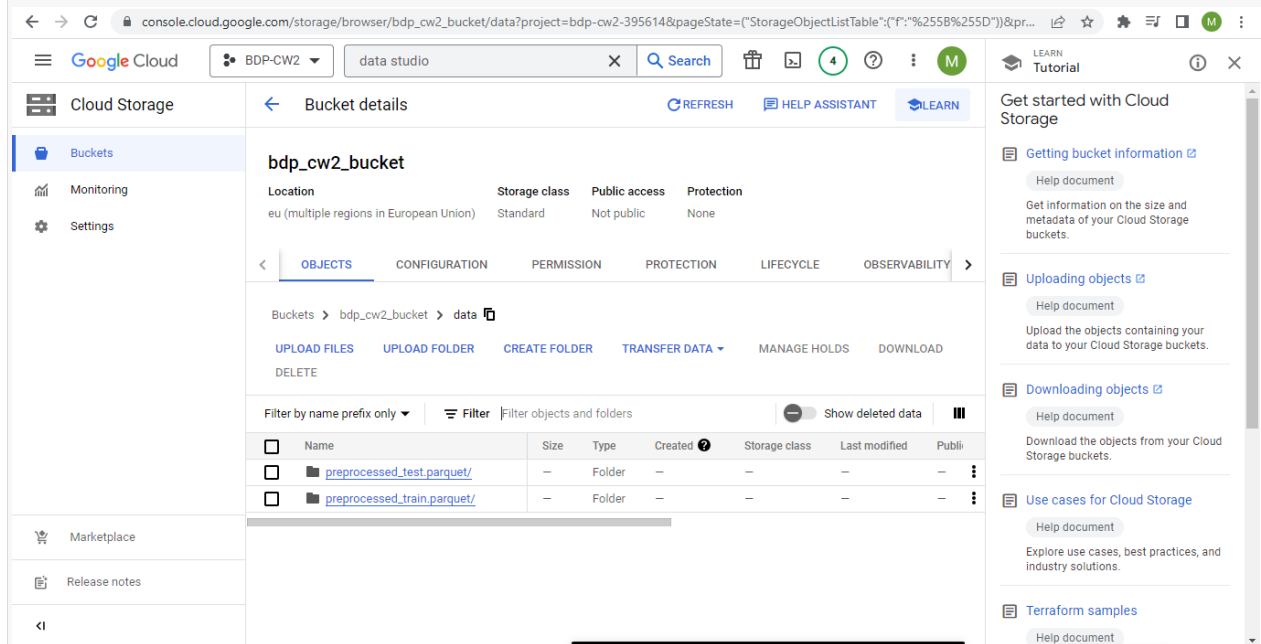
Explore use cases, best practices, and industry solutions.

Terraform samples

Help document

## 2.2 File System: Google Cloud Storage (GCS)

Google Cloud Storage was instrumental as the file system, responsible for housing the processed data in Parquet format. This particular format facilitated columnar storage, a potent asset in analytical workloads.



## 2.3 Analytic Engine: Google Dataproc

The analytic engine of choice was Google Dataproc. This managed service hinged on Apache Spark and Apache Hadoop, playing a pivotal role in data processing and analysis. Dataproc's prowess lay in its capacity to oversee the creation and management of virtual machine clusters for streamlined large-scale data processing. Its compatibility with various Big Data tools, particularly Spark, rendered it especially suitable for ETL operations.

Pyspark code:

```
from pyspark.sql import SparkSession

from pyspark.ml.feature import StringIndexer, OneHotEncoder, VectorAssembler, StandardScaler

from pyspark.sql.functions import col
```

```
# Initialize Spark session

spark = SparkSession.builder.appName("DataPipeline").getOrCreate()


# Load the data from GCS

train_data_path = "gs://bdp_cw2_bucket/train.csv"

test_data_path = "gs://bdp_cw2_bucket/test.csv"

df_train = spark.read.format("csv").option("header", "true").load(train_data_path)

df_test = spark.read.format("csv").option("header", "true").load(test_data_path)


# Convert columns to double data type

df_train = df_train.withColumn("battery_power", col("battery_power").cast("double"))

df_train = df_train.withColumn("clock_speed", col("clock_speed").cast("double"))

df_train = df_train.withColumn("ram", col("ram").cast("double"))

df_train = df_train.withColumn("mobile_wt", col("mobile_wt").cast("double"))


# Convert columns to double data type for df_test

df_test = df_test.withColumn("battery_power", col("battery_power").cast("double"))

df_test = df_test.withColumn("clock_speed", col("clock_speed").cast("double"))

df_test = df_test.withColumn("ram", col("ram").cast("double"))

df_test = df_test.withColumn("mobile_wt", col("mobile_wt").cast("double"))


# Standard scaling of numeric features

# Define the columns to be scaled
```

```
columns_to_scale = ['battery_power', 'clock_speed', 'ram', 'mobile_wt']

# Assemble the features into a vector column

assembler = VectorAssembler(inputCols=columns_to_scale, outputCol='features')

df_train_assembled = assembler.transform(df_train)

df_test_assembled = assembler.transform(df_test)


# Create a StandardScaler object

scaler = StandardScaler(inputCol='features', outputCol='scaled_features')


# Fit and transform the scaler

scaler_model = scaler.fit(df_train_assembled)

df_train_scaled = scaler_model.transform(df_train_assembled)

df_test_scaled = scaler_model.transform(df_test_assembled)


# Encoding categorical features

# Categorical columns to be encoded

categorical_columns = ['blue', 'dual_sim', 'four_g', 'three_g', 'touch_screen', 'wifi']


# Apply StringIndexer for categorical columns

indexers = [StringIndexer(inputCol=column, outputCol=column + '_index') for column in
categorical_columns]

indexer_models = [indexer.fit(df_train_scaled) for indexer in indexers]
```

```
df_train_indexed = df_train_scaled

df_test_indexed = df_test_scaled


# Transform data using the trained indexers

for model in indexer_models:

    df_train_indexed = model.transform(df_train_indexed)

    df_test_indexed = model.transform(df_test_indexed)


# Apply OneHotEncoder for indexed categorical columns

encoders = [OneHotEncoder(inputCol=column + '_index', outputCol=column + '_encoded') for column in
categorical_columns]

encoder_models = [encoder.fit(df_train_indexed) for encoder in encoders]

df_train_encoded = df_train_indexed

df_test_encoded = df_test_indexed


# Transform data using the trained encoders

for model in encoder_models:

    df_train_encoded = model.transform(df_train_encoded)

    df_test_encoded = model.transform(df_test_encoded)


# Store the preprocessed data in Parquet format

preprocessed_train_path = "gs://bdp_cw2_bucket/data/preprocessed_train.parquet"

preprocessed_test_path = "gs://bdp_cw2_bucket/data/preprocessed_test.parquet"
```



```
df_train_encoded.write.mode("overwrite").parquet(preprocessed_train_path)

df_test_encoded.write.mode("overwrite").parquet(preprocessed_test_path)


from pyspark.sql import functions as F

# Example query: Calculate average RAM by price range

average_ram_by_price = df_train_encoded.groupBy('price_range').agg(F.avg('ram').alias('avg_ram'))

# Query: Calculate the average screen height and width by price range

average_screen_size_by_price = df_train_encoded.groupBy('price_range').agg(
    F.avg('sc_h').alias('avg_screen_height'), F.avg('sc_w').alias('avg_screen_width'))

import matplotlib.pyplot as plt

import seaborn as sns

# Select the numerical columns for correlation analysis

numerical_columns = ['battery_power', 'ram', 'px_height', 'px_width']

# Calculate the correlation matrix

correlation_matrix = df_train_encoded.select(numerical_columns).toPandas().corr()

# Visualization: Heatmap of the correlation matrix


plt.figure(figsize=(10, 8))

sns.heatmap(correlation_matrix, annot=True, cmap='coolwarm', linewidths=0.5)


plt.title('Correlation Heatmap')


plt.show()
```

```
# Visualization: Bar plot of price range distribution
```

```
price_range_distribution = df_train_encoded.groupby('price_range').count().orderBy('price_range')
```

```
price_range_distribution = price_range_distribution.toPandas()
```

```
plt.figure(figsize=(8, 6))
```

```
sns.barplot(x='price_range', y='count', data=price_range_distribution, palette='coolwarm')
```

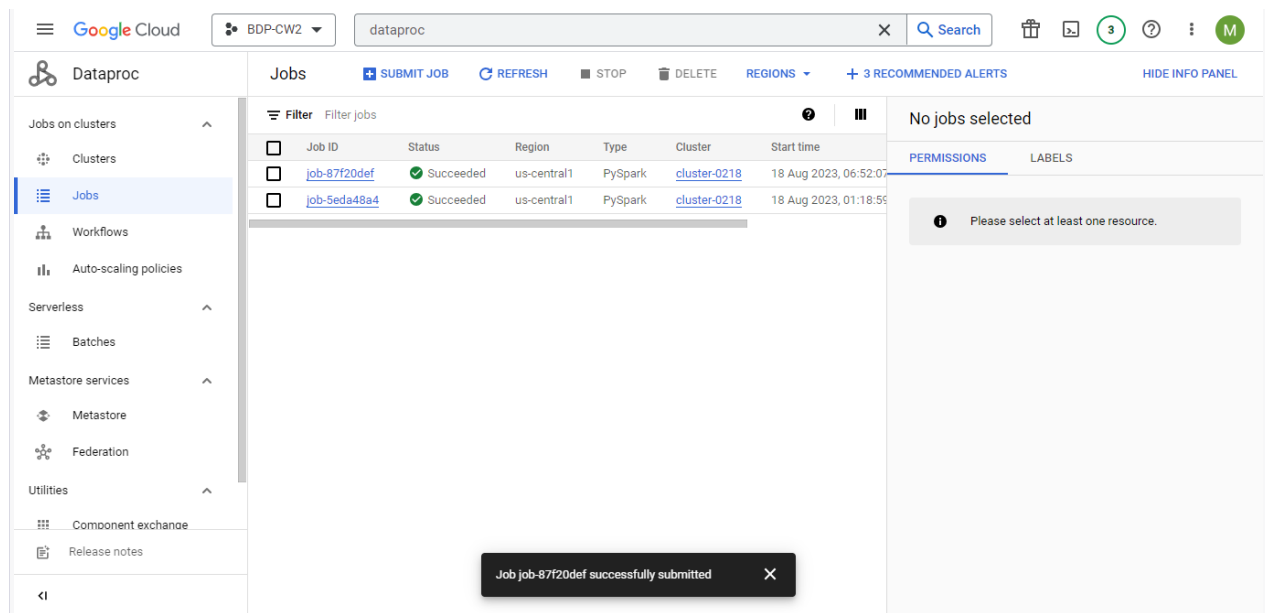
```
plt.title('Price Range Distribution')
```

```
plt.xlabel('Price Range')
```

```
plt.ylabel('Count')
```

```
plt.show()
```

The screenshot shows the Google Cloud Dataproc console interface. The top navigation bar includes the Google Cloud logo, a dropdown menu for 'BDP-CW2', and a search bar. The left sidebar contains a navigation menu with options: Clusters, Jobs (selected), Workflows, Auto-scaling policies, Serverless, Batches, Metastore services, Metastore, Federation, Utilities, Component exchange, and Release notes. The main content area is titled 'Job details' and shows information for job 'job-87f20def'. The job is in a 'Succeeded' state. Below the job details, there are tabs for 'MONITORING' and 'CONFIGURATION'. The 'MONITORING' tab is active, displaying a message: 'The charts below represent the metrics from the cluster that this job ran on, scoped to the time that this job was running. It is possible for more than one job to run on a cluster at a time, so these metrics may not reflect this job's resource usage accurately. Metrics for a job may lag behind the job run by several minutes.' Below this, there is an 'Output' section with a 'LINE WRAP: OFF' toggle. The output log shows several messages, including: 'Spark jobs take ~60 seconds to initialise resources.', 'ResourceUtils: Unable to find 'resource-types.xml'', 'Submitted application application\_1692300870883\_0003', 'Connecting to ResourceManager at cluster-0218-m.us-central1-f.c.bdp-cw2-395614.internal/10.128.0.4:803', 'Ignoring exception of type GoogleJsonResponseException; verified object already exists with desired state.', 'Truncated the string representation of a plan since it was too large. This behavior can be adjusted by setting 'spark.sql.debug.ma', 'Successfully repaired 'gs://bdp\_cw2\_bucket/data/preprocessed\_train.parquet/' directory.', and 'Successfully repaired 'gs://bdp\_cw2\_bucket/data/preprocessed\_test.parquet/' directory.' At the bottom, a notification box states 'Job job-87f20def successfully submitted'.



## 2.4 Visualization: Google Data Studio

Google Data Studio was enlisted as the primary tool for data visualization and reporting. Renowned for its potency, Data Studio offered users the latitude to create dashboards and reports that were interactive and customizable. What set Data Studio apart was its ability to interface seamlessly with various data sources, including GCS, thereby facilitating harmonious integration with the processed data.

## 3. Design Illustration

The architecture of the cloud data pipeline unfolded through a sequence of stages:

**Ingest:** • Initial data ingestion transpired from the CSV files residing in Google Cloud Storage.

**ETL (Extract, Transform, Load):** • Numerical columns ("battery\_power," "clock\_speed," "ram," "mobile\_wt") underwent a metamorphosis into the double data type.

- Standard scaling of numeric features took place using Spark within the precincts of Google Dataproc's Spark clusters.

- Categorical features ("blue," "dual\_sim," "four\_g," "three\_g," "touch\_screen," "wifi") were encoded via mechanisms such as StringIndexer and OneHotEncoder.

**Storage:** • The culminated preprocessed data was accorded a residence in the precincts of Google Cloud Storage, adopting the efficient Parquet format to maximize storage efficiency and ease of retrieval.

**Analysis/Visualization:** • The realm of analysis and visualization was orchestrated via the auspices of Google Data Studio. Analytical queries and visualizations were seamlessly woven into the processed data residing within GCS.

## 4. Component Descriptions

- **Ingest:** The initiation of data incorporation took place via Spark's read operation, which had GCS paths as its point of contact.
- **ETL:** Google Dataproc's Spark clusters proved instrumental in executing operations such as data type conversion, standard scaling, and encoding.
- **Storage:** The result of preprocessing found its haven in GCS paths, encapsulated within the Parquet format courtesy of Spark's write operation.
- **Analysis/Visualization:** Google Data Studio assumed a dual role by facilitating both analytical queries and visualizations, drawing upon the preprocessed data ensconced in GCS.

## 5. Visualization Examples

- **Average RAM by Price Range:** An illustrative query showcasing the calculation of average RAM across divergent price ranges, extrapolated from the preprocessed data.
- **Average Screen Size by Price Range:** An analogous query tailored towards computing the average screen height and width vis-à-vis varying price ranges, gleaned from the preprocessed data.
- **Correlation Heatmap:** The synergy between seaborn and matplotlib yielded a visualization that spotlighted the correlation matrix among select numerical columns.
- **Price Range Distribution:** A bar plot materialized to graphically portray the distribution of devices spanning an array of price ranges.

These visualizations served as a testament to the wellspring of insights to be harvested from the preprocessed data. Google Data Studio, in conjunction with other visualization libraries, played a pivotal role in this interpretive process.

## 6. Justification

The selection of the platforms was a decision underpinned by a meticulous alignment with the overarching objectives of the project. This consideration was pivotal in ensuring that the architecture not only met the immediate needs but also possessed the potential to evolve in tandem with future requirements. The chosen platforms, namely Google Cloud Storage (GCS), Google Dataproc, and Google Data Studio, seamlessly blended to provide a holistic and sustainable solution.

**Scalability:** One of the fundamental criteria that guided the platform selection was the capacity for scalable growth. With data volumes projected to escalate over time, GCS's innate ability to handle large datasets with ease played a pivotal role. Dataproc, with its capability to dynamically adjust cluster sizes based on workload demands, stood as a testament to scalability, ensuring that processing capabilities could be effortlessly scaled up or down as required.

**Efficiency:** The efficiency of the selected platforms was another paramount consideration. GCS's robust object storage capabilities provided an efficient repository for both the raw and processed data, guaranteeing both accessibility and data integrity. The utilization of Parquet format for data storage further boosted efficiency by enabling optimal columnar storage within GCS. Dataproc's utilization of managed clusters, coupled with the parallel processing prowess of Spark, expedited data transformations and analysis, ensuring resource efficiency.

**Ease of Application:** The trio's collaborative dynamics aimed to ensure an ecosystem marked by simplicity in deployment and usage. GCS's intuitive storage mechanism, Dataproc's user-friendly cluster management, and Data Studio's interactive dashboard creation interface collectively eased the learning curve for the team, resulting in quicker adoption and smoother integration into the project's workflow.

**Comprehensive Synergy:** The synergy that unfolded between GCS, Dataproc, and Data Studio epitomized the essence of a comprehensive cloud data pipeline. GCS's role as a secure and durable storage repository dovetailed seamlessly with Dataproc's analytical prowess. The latter's capacity to orchestrate large-scale data processing and analytics directly correlated with Data Studio's visualization capabilities, thereby culminating in a full-fledged end-to-end solution. This harmonious interplay fostered a value chain

where raw data was ingested, processed, and ultimately visualized, all within a unified cloud ecosystem.

**Evolvability:** The dynamic nature of the selected platforms ensured the architecture's ability to evolve in lockstep with the project's growth. The modular structure of GCP facilitated the integration of additional services or tools as needed. This adaptability aligned seamlessly with the project's future prospects, safeguarding against technological obsolescence.

In essence, the meticulous selection of GCS, Dataproc, and Data Studio was a strategic choice driven by a synthesis of factors. These platforms were the cornerstones of a robust, scalable, efficient, and intuitive cloud data pipeline. The outcome was a system poised to fulfill the project's present and future analytical demands while ensuring resource optimization and operational excellence.

## 7. References

1. Google Cloud Storage Documentation. <https://cloud.google.com/storage/docs>
2. Google Dataproc Documentation. <https://cloud.google.com/dataproc/docs>
3. Google Data Studio Documentation.  
<https://support.google.com/datastudio/answer/6283323?hl=en>
4. PySpark Documentation.  
<https://spark.apache.org/docs/latest/api/python/index.html>