# IBM Data Engineering Capstone Project

## Project Overview

As part of the Capstone project, you will assume the role of an Associate Data Engineer who has recently joined an e-commerce organization. You will be presented with a business challenge that requires building a data platform for retailer data analytics.

In this Capstone project, you will:
Design a data platform that uses MySQL as an OLTP database and MongoDB as a NoSQL database.
Design and implement a data warehouse and generate reports from the data.
Design a reporting dashboard that reflects the key metrics of the business.
Extract data from OLTP and NoSQL databases, transform it and load it into the data warehouse, and then create an ETL pipeline.
And finally, create a Spark connection to the data warehouse, and then deploy a machine learning model.

In Module 1, you will design the OLTP database for an E-Commerce website, populate the OLTP Database with the data provided and automate the export of the daily incremental data into the data warehouse.

In Module 2, you will set up a NoSQL database to store the catalogue data for an E-Commerce website, load the E-Commerce catalogue data into the NoSQL database, and query the E-Commerce catalogue data in the NoSQL database.

In Module 3, you will design the schema for a data warehouse based on the schema of the OLTP and NoSQL databases. You'll then create the schema and load the data into fact and dimension tables, automate the daily incremental data insertion into the data warehouse, and create Cubes and Rollups to make the reporting easier.
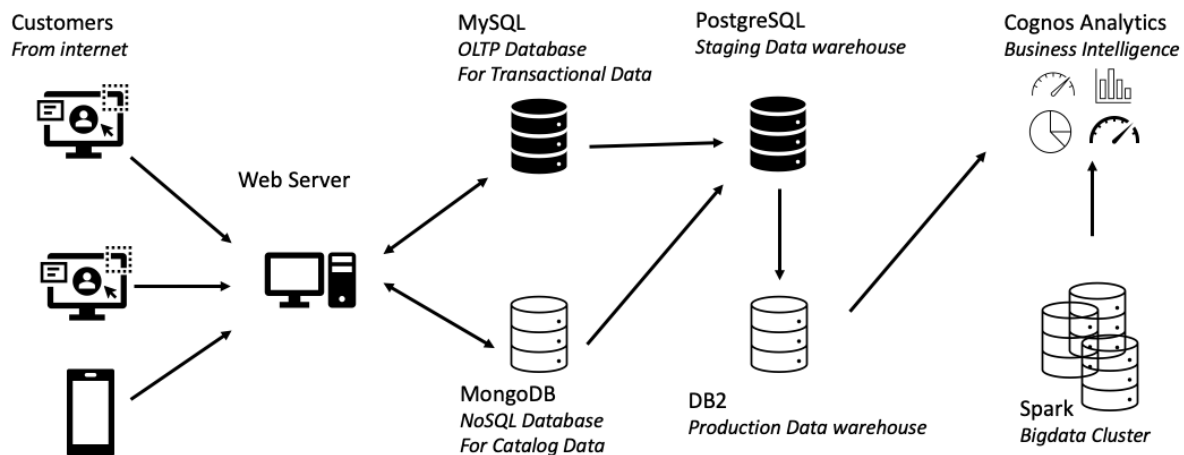
In Module 4, you will create a Cognos data source that points to a data warehouse table, create a bar chart of Quarterly sales of cell phones, create a pie chart of sales of electronic goods by category, and create a line chart of total sales per month for the year 2020.

In Module 5, you will extract data from OLTP, NoSQL, and MongoDB databases into CSV format. You will then transform the OLTP data to suit the data warehouse schema and then load the transformed data into the data warehouse. Finally, you will verify that the data is loaded properly.

In the sixth and final module, you will use your skills in Big Data Analytics to create a Spark connection to the data warehouse, and then deploy a machine learning model on SparkML for making sales projections.

# Data Platform Architecture

## Data Platform Architecture



This is the data platform architecture of an ecommerce company named SoftCart. SoftCart uses a hybrid architecture, with some of its databases on premises and some on cloud.

## Tools and Technologies:

- OLTP database - MySQL
- NoSQL database - MongoDB
- Production Data warehouse – DB2 on Cloud
- Staging - Data warehouse – PostgreSQL
- Big data platform - Hadoop
- Big data analytics platform – Spark
- Business Intelligence Dashboard - IBM Cognos Analytics
- Data Pipelines - Apache Airflow

## Process:

SoftCart's online presence is primarily through its website, which customers access using a variety of devices like laptops, mobiles and tablets. All the catalogue data of the products is stored in the MongoDB NoSQL server. All the transactional data like inventory and sales are stored in the MySQL database server. SoftCart's webserver is driven entirely by these two databases. Data is periodically extracted from these two databases and put into the staging data warehouse running on PostgreSQL. Production data warehouse is on the cloud instance of IBM DB2 server. BI teams connect to the IBM DB2 for operational dashboard creation. IBM Cognos Analytics is used to create dashboards. SoftCart uses Hadoop cluster as its big data platform where all the data collected for analytics purposes. Spark is used to analyse the data on the Hadoop cluster. To move data between OLTP, NoSQL and the data warehouse, ETL pipelines are used and these run on Apache Airflow.

# Module 1: with link to OLTP database exercise

## OLTP database requirements and design

### OLTP database:
OLTP database is generally used to handle everyday business transactions of an organization like a bank or a super market chain. OLTP databases can be write heavy or may have a balanced read/write load.

### OLTP database requirements:
An OLTP database is expected to handle a huge number of transactions per second. Each transaction usually involves accessing (read/write) a small portion of the database, in other words the payload per transaction is small.
The time taken to execute a transaction usually called latency needs to be very less.

### OLTP database design:
The schema of an OLTP database is highly normalized so as to achieve a very low latency. To further improve the latency an OLTP database stores only the recent data like the last few week's data. They are usually run on storage that is very fast like SSD.

## Scenario
You are a data engineer at an e-commerce company. Your company needs you to design a data platform that uses MySQL as an OLTP database. You will be using MySQL to store the OLTP data.

## Design the OLTP Database
Creating sales database and sales_data table in sales database.

```
Oracle is a registered trademark of Oracle Corporation and/or its
affiliates. Other names may be trademarks of their respective
owners.

Type 'help;' or '\h' for help. Type '\c' to clear the current input statement.

mysql> CREATE DATABASE sales;
Query OK, 1 row affected (0.01 sec)

mysql> USE sales;
Database changed
mysql> CREATE TABLE sales_data (
    -> product_id INT,
    -> customer_id INT,
    -> price INT,
    -> quantity INT,
    -> timestamp DATETIME);
Query OK, 0 rows affected (0.04 sec)

mysql>
```

# Load the given data(oltpdata.csv) into sales_data table

Importing data from oltpdata.csv file into sales_data table using phpMyAdmin GUI.



List the tables in the sales database and query the count of records in the sales_data table to make sure the data are imported correctly and successfully.

# Set up Admin tasks

Creating index on the timestamp field to speed up queries.

```
mysql> CREATE INDEX ts ON sales_data(timestamp);
Query OK, 0 rows affected (0.09 sec)
Records: 0  Duplicates: 0  Warnings: 0

mysql> SHOW INDEX FROM sales_data;
+------------+------------+----------+--------------+-------------+-----------+-------------+----------+
--------+------+------------+---------+---------------+---------+------------+
| Table      | Non_unique | Key_name | Seq_in_index | Column_name | Collation | Cardinality | Sub_part |
 Packed | Null | Index_type | Comment | Index_comment | Visible | Expression |
+------------+------------+----------+--------------+-------------+-----------+-------------+----------+
--------+------+------------+---------+---------------+---------+------------+
| sales_data |          1 | ts       |            1 | timestamp   | A         |        2605 |     NULL |
   NULL | YES  | BTREE      |         |               | YES     | NULL       |
+------------+------------+----------+--------------+-------------+-----------+-------------+----------+
--------+------+------------+---------+---------------+---------+------------+
1 row in set (0.04 sec)

mysql>
```

Writing a bash script(datadump.sh) that exports all the rows in the sales_data table to a file named sales_data.sql to automate the export of the daily incremental data into the data warehouse.

```
datadump.sh
 1   #!/bin/sh
 2   # The above line tells the interpreter this code needs to be run as a shell script.
 3
 4   # This will be printed on to the screen. In the case of cron job, it will be printed to the logs.
 5   echo "Pulling Database: This may take a few minutes"
 6   # Create a backup
 7   if mysqldump --user=root --password='NzQzNS12YXN0ZXXJz' sales sales_data > sales_data.sql ; then
 8     echo 'sales_data.sql created'
 9   else
10     echo 'Error sales_data.sql was not created!'
11     exit
12   fi
```

```
theia@theiadocker-vastersett:/home/project$ ls -l datadump.sh
-rw-r--r-- 1 theia users 450 Jul  8 05:29 datadump.sh
theia@theiadocker-vastersett:/home/project$ chmod +x datadump.sh
theia@theiadocker-vastersett:/home/project$ ls -l datadump.sh
-rwxr-xr-x 1 theia users 450 Jul  8 05:29 datadump.sh
theia@theiadocker-vastersett:/home/project$ ./datadump.sh
Pulling Database: This may take a few minutes
mysqldump: [Warning] Using a password on the command line interface can be insecure.
sales_data.sql created
theia@theiadocker-vastersett:/home/project$ ls -l
total 220
-rwxr-xr-x 1 theia users    450 Jul  8 05:29 datadump.sh
-rw-r--r-- 1 theia users 218964 Jul  8 05:34 sales_data.sql
theia@theiadocker-vastersett:/home/project$
```

## Module 2:

## Scenario

You are a data engineer at an e-commerce company. Your company needs you to design a data platform that uses MongoDB as a NoSQL database. You will be using MongoDB to store the e-commerce catalogue data.

## Importing data into MongoDB database

Installing mongoimport and mongoexport. Importing given data(catalog.json) into a database named 'catalog' and a collection named 'electronics' on mongodb server.

List out all databases and collections on mongodb server to check if catalog database and electronics collection are created successfully or not.

```
theia@theiadocker-vastersett:/home/project$ mongo -u root -p MTMwMzUtdmFzdGVy --authenticationDatab
ase admin local
MongoDB shell version v3.6.3
connecting to: mongodb://127.0.0.1:27017/local
MongoDB server version: 3.6.3
Server has startup warnings:
2022-07-11T12:01:11.984+0000 I STORAGE  [initandlisten]
2022-07-11T12:01:11.984+0000 I STORAGE  [initandlisten] ** WARNING: Using the XFS filesystem is str
ongly recommended with the WiredTiger storage engine
2022-07-11T12:01:11.984+0000 I STORAGE  [initandlisten] **          See http://dochub.mongodb.org/c
ore/prodnotes-filesystem
> show dbs
admin    0.000GB
catalog  0.000GB
config   0.000GB
local    0.000GB
>
```

```
> use catalog
switched to db catalog
> show collections
electronics
>
```

Creating index on the field 'type' for faster queries.

```
> db.electronics.createIndex({"type":1})
{
        "createdCollectionAutomatically" : false,
        "numIndexesBefore" : 1,
        "numIndexesAfter" : 2,
        "ok" : 1
}
```

# Trying out a few queries to check the imported data:

Finding the count of laptops.

```
> db.electronics.count({"type":"laptop"})
389
```

Finding the number of smart phones with screen size of 6 inches.

```
> db.electronics.count({"type":"smart phone","screen size":6})
8
>
```

Finding the average screen size of smart phones.

```
> db.electronics.aggregate([{"$group":{"_id":"$type","averge":{"$avg":"$screen size"}}}])
{ "_id" : "television", "averge" : 39.8 }
{ "_id" : "laptop", "averge" : 14.568123393316196 }
{ "_id" : "smart phone", "averge" : 6 }
> db.electronics.aggregate([{"$match":{"type":"smart phone"}},{"$group":{"_id":"$type","averge":{"$
avg":"$screen size"}}}])
{ "_id" : "smart phone", "averge" : 6 }
>
```

# Export data from MongoDB database as csv file

Exporting the fields _id, "type", "model", from the 'electronics' collection into a file named electronics.csv

```
> exit
bye
theia@theiadocker-vastersett:/home/project$ mongoexport -u root -p MTMwMzUtdmFzdGVy --authenticatio
nDatabase admin --db catalog --collection electronics --out electronics.csv --type=csv --fields _id
,type,model
2022-07-11T12:49:29.772+0000    connected to: mongodb://localhost/
2022-07-11T12:49:29.786+0000    exported 438 records
theia@theiadocker-vastersett:/home/project$ 
```

# Module 3.1:

## Scenario

You are a data engineer hired by an ecommerce company named SoftCart.com. The company retails download only items like E-Books, Movies, Songs etc. The company has international presence and customers from all over the world. The company would like to create a data warehouse so that it can create reports like:

- total sales per year per country
- total sales per month per category
- total sales per quarter per country
- total sales per category per country

You will use your data warehousing skills to design and implement a data warehouse for the company.

## Designing a Data Warehouse

The ecommerce company has provided you the sample data.

| OrderID | Item | Category | Price | Country | Date |
|---------|------|----------|-------|---------|------|
| 2123 | The Matrix | Movie | 9.99 | USA | 20-Feb-21 |
| 3254 | The Alchemist | Ebook | 5.99 | Canada | 20-Feb-21 |
| 4901 | Baby Shark | Song | 2.49 | Japan | 20-Feb-21 |
| 5679 | The Lord of the Rings | Ebook | 6.99 | Cyprus | 20-Feb-21 |

You will start your project by designing a Star Schema for the warehouse by identifying the columns for the various dimension and fact tables in the schema. Name your database as softcart.

Creating tables and relationships amongst created tables, using the ERD design Tool of pgAdmin.

# Create the schema

After creating tables and relationships using ERD design tool, generate the schema sql using the tool. Then use the created schema sql to create the schema in a database named staging.



# Module 3.2:

## Scenario

You are a data engineer hired by an ecommerce company named SoftCart.com. The company retails download only items like E-Books, Movies, Songs etc. The company has international presence and customers from all over the world. You have designed the schema for the data warehouse in the previous assignment. Data engineering is a team game. Your senior data engineer reviewed your design. Your schema design was improvised to suit the production needs of the company. In this assignment you will generate reports out of the data in the data warehouse.

# Load data into the Data Warehouse (IBM DB2)

Loading the data provided by the company in csv format into the tables (DimDate, DimCategory, DimCountry, FactSales).

Load Data    Load History    Tables    Views    Indexes    Aliases    MQTs    Sequences    Application objects

HHM70408.DIMCOUNTRY

Back

Export to CSV

| | COUNTRYID SMALLINT | COUNTRY VARCHAR(20) |
|---|---|---|
| 1 | 1 | Argentina |
| 2 | 2 | Australia |
| 3 | 3 | Austria |
| 4 | 4 | Azerbaijan |
| 5 | 5 | Belgium |
| 6 | 6 | Brazil |
| 7 | 7 | Bulgaria |
| 8 | 8 | Canada |
| 9 | 9 | Cyprus |
| 10 | 10 | Czech Republic |
| 11 | 11 | Denmark |
| 12 | 12 | Egypt |
| 13 | 13 | Estonia |
| 14 | 14 | Finland |
| 15 | 15 | France |
| 16 | 16 | Germany |

Load Data    Load History    Tables    Views    Indexes    Aliases    MQTs    Sequences    Application objects

HHM70408.FACTSALES

Back

Export to CSV

| | ORDERID INTEGER | DATEID SMALLINT | COUNTRYID SMALLINT | CATEGORYID SMALLINT | AMOUNT SMALLINT |
|---|---|---|---|---|---|
| 1 | 1 | 1 | 6 | 4 | 5190 |
| 2 | 2 | 1 | 25 | 2 | 1205 |
| 3 | 3 | 1 | 22 | 5 | 3155 |
| 4 | 4 | 1 | 40 | 2 | 268 |
| 5 | 5 | 1 | 28 | 3 | 3199 |
| 6 | 6 | 1 | 30 | 4 | 6643 |
| 7 | 7 | 1 | 11 | 2 | 3995 |
| 8 | 8 | 1 | 14 | 4 | 2318 |
| 9 | 9 | 1 | 18 | 2 | 3956 |
| 10 | 10 | 1 | 37 | 1 | 3681 |
| 11 | 11 | 1 | 10 | 3 | 1309 |
| 12 | 12 | 1 | 1 | 1 | 5937 |
| 13 | 13 | 1 | 9 | 2 | 6487 |
| 14 | 14 | 1 | 42 | 1 | 7745 |
| 15 | 15 | 1 | 6 | 1 | 3597 |
| 16 | 16 | 1 | 11 | 5 | 399 |

# Queries for data analytics

Creating and running queries to check if the created data warehouse can create reports.

Grouping sets query using the columns country, category, totalsales. (Total sales per country and Total sales per category)



Rollup query using the columns year, country, and totalsales. (Total sales per year per country and Total sales per year)

Cube query using the columns year, country, and average sales. (Average sales per year per country, Average sales per year, Average sales per country)



Creating a materialized query table (MQT) named total_sales_per_country that has the columns country and total_sales to improve the performance of complex queries that operate on very large amounts of data.

Db2 uses a materialized query table to precompute the results of data that is derived from one or more tables. When you submit a query, Db2 can use the results that are stored in a materialized query table rather than compute the results from the underlying source tables on which the materialized query table is defined.
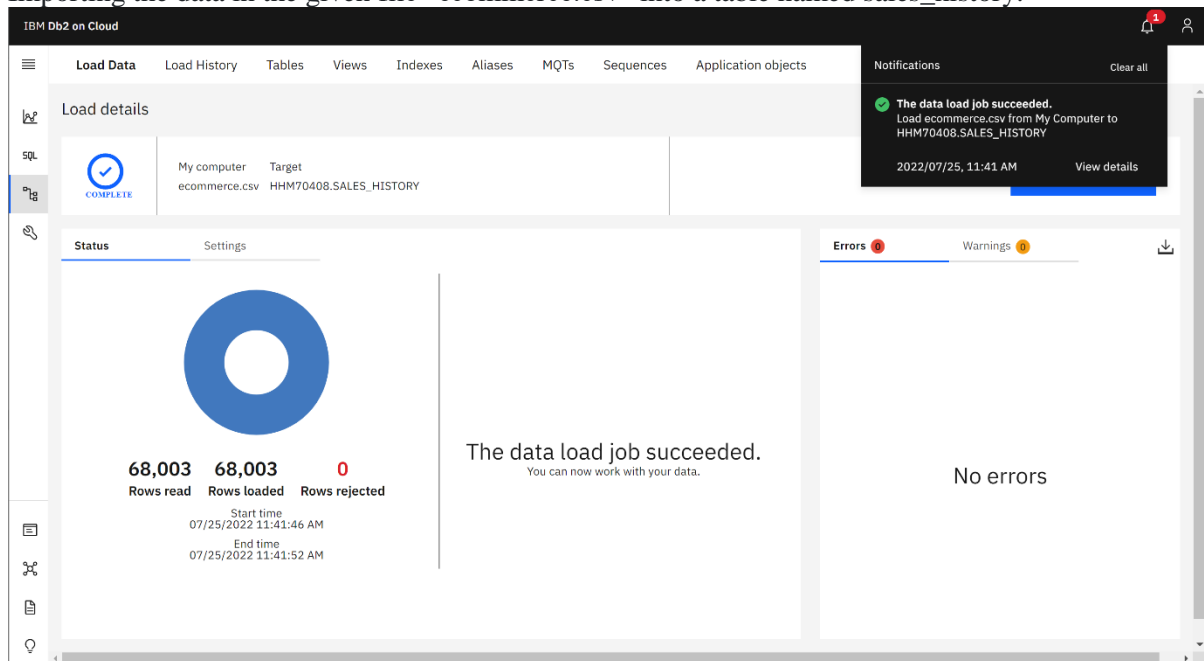
# Module 4:

## Scenario

You are a data engineer at an e-commerce company. Your company has finished setting up a data warehouse. Now you are assigned the responsibility to design a reporting dashboard that reflects the key metrics of the business.

## Load data into the data warehouse (IBM DB2)

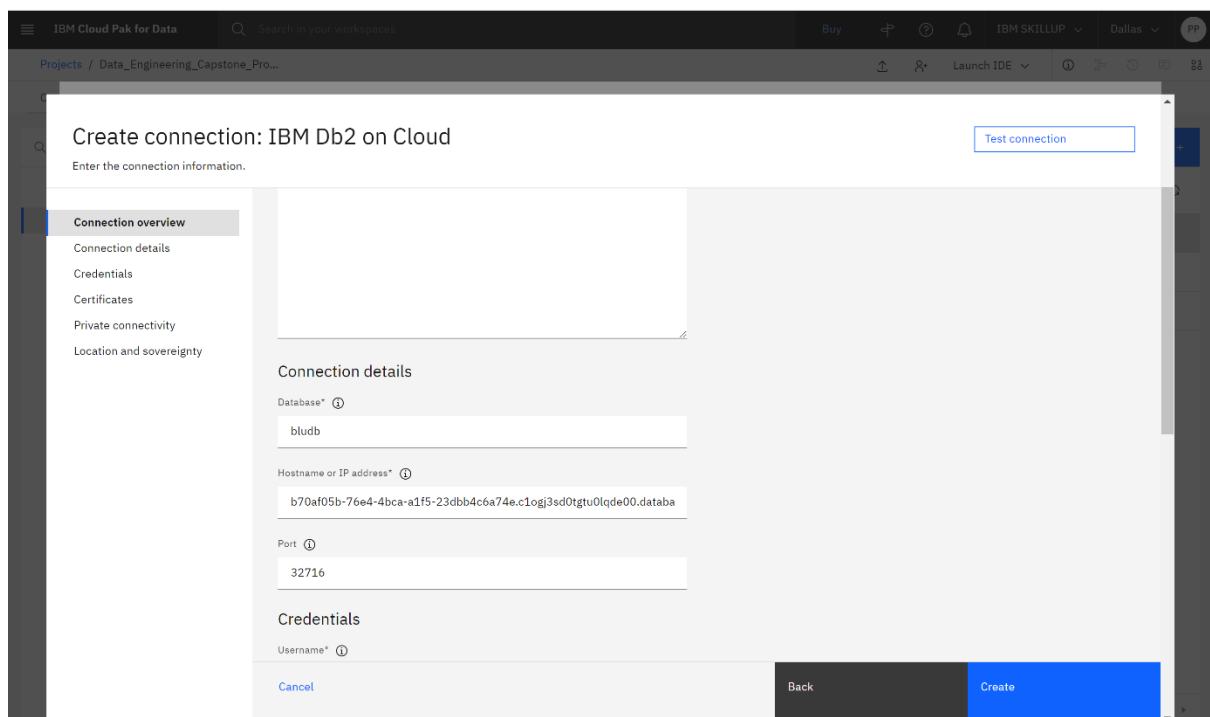Importing the data in the given file "ecommerce.csv" into a table named sales_history.
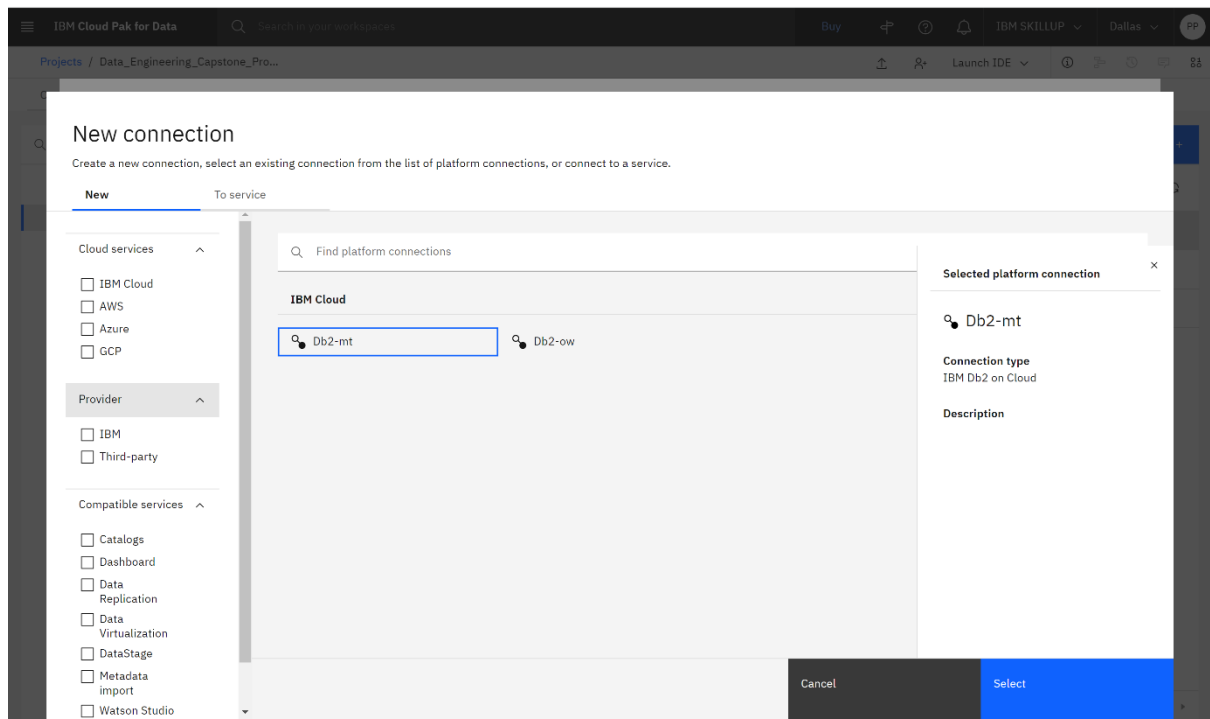


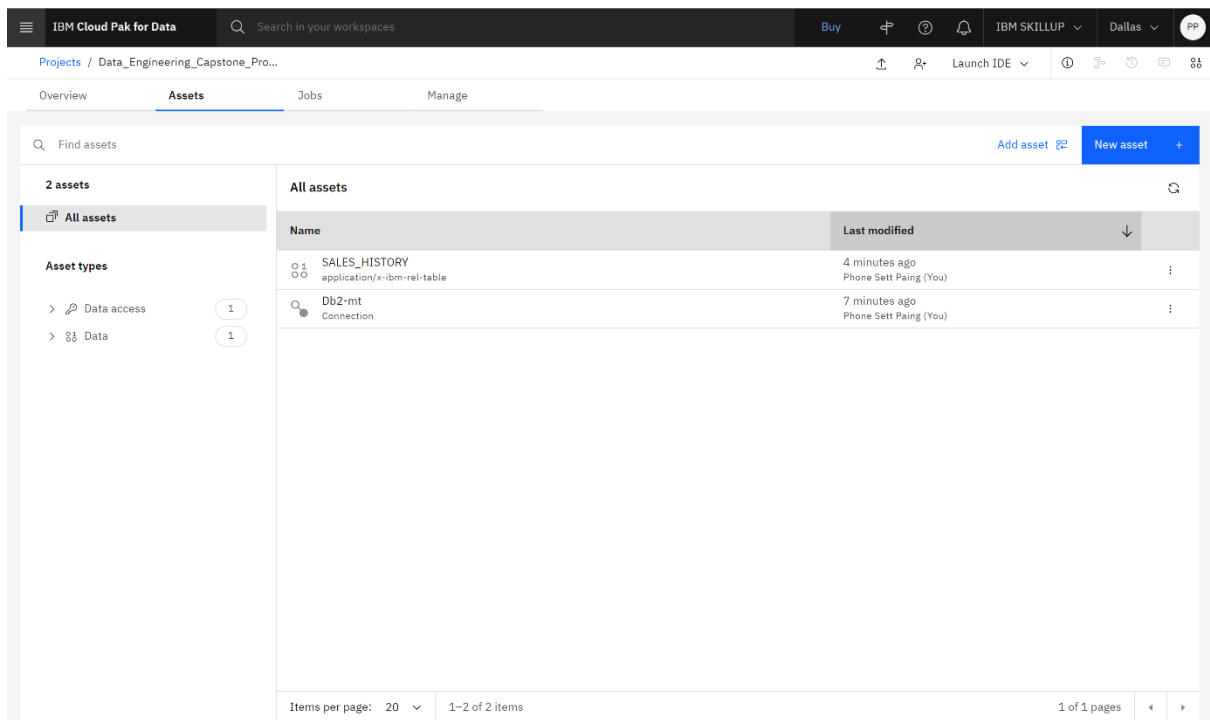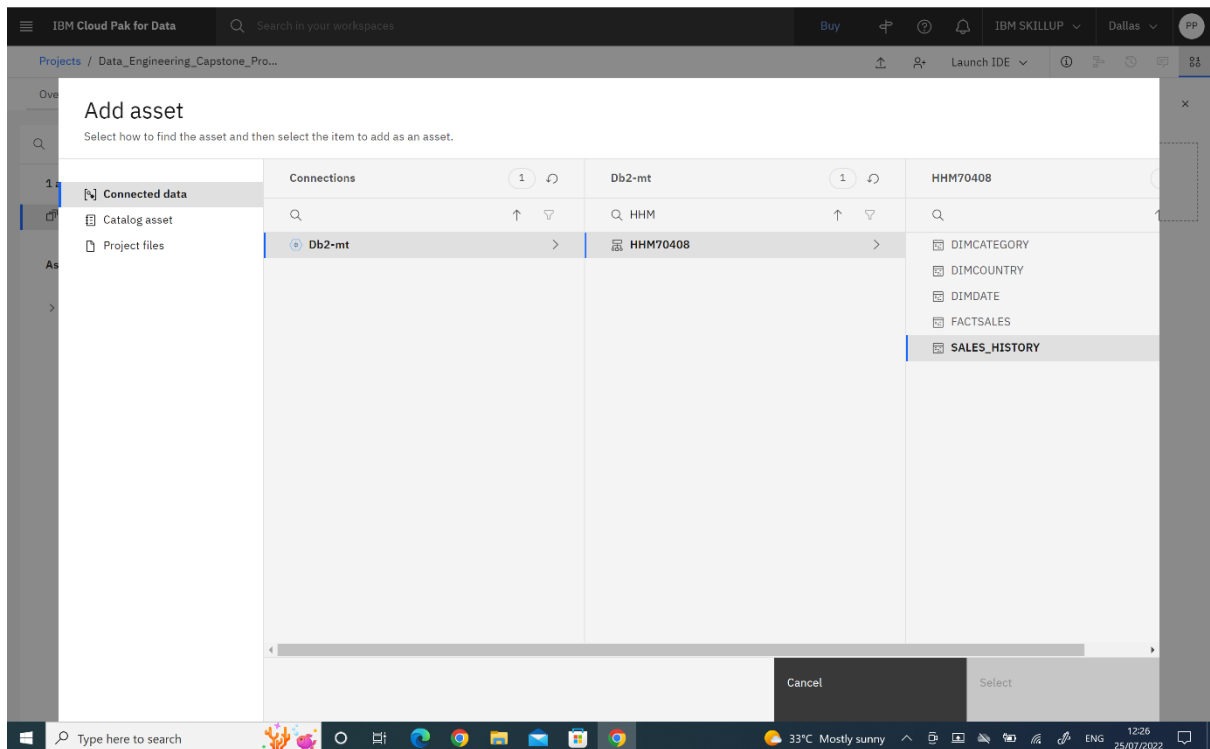Listing the first 10 rows in the sales_history table to check the loaded data.
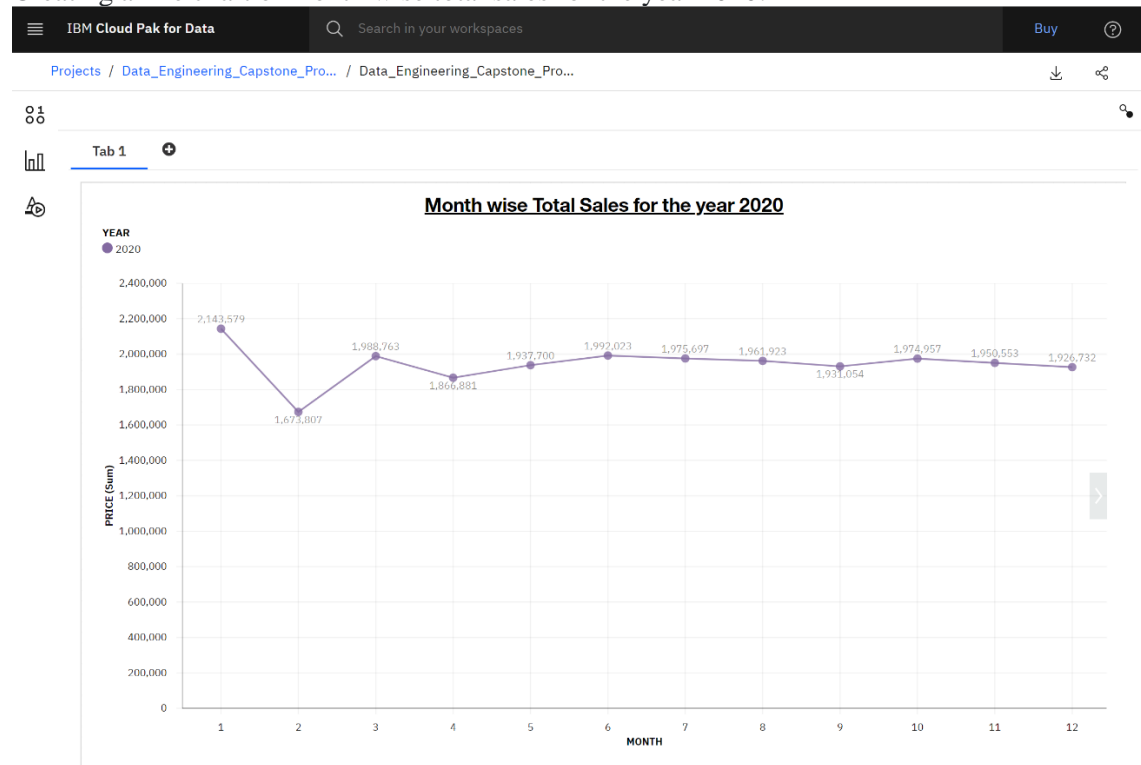
# Create data source in IBM Cognos

Creating connection between sales_history table in IBM DB2 and Watson Studio with added Cognos Dashboard Embedded (CDE) service to use it as a data source.
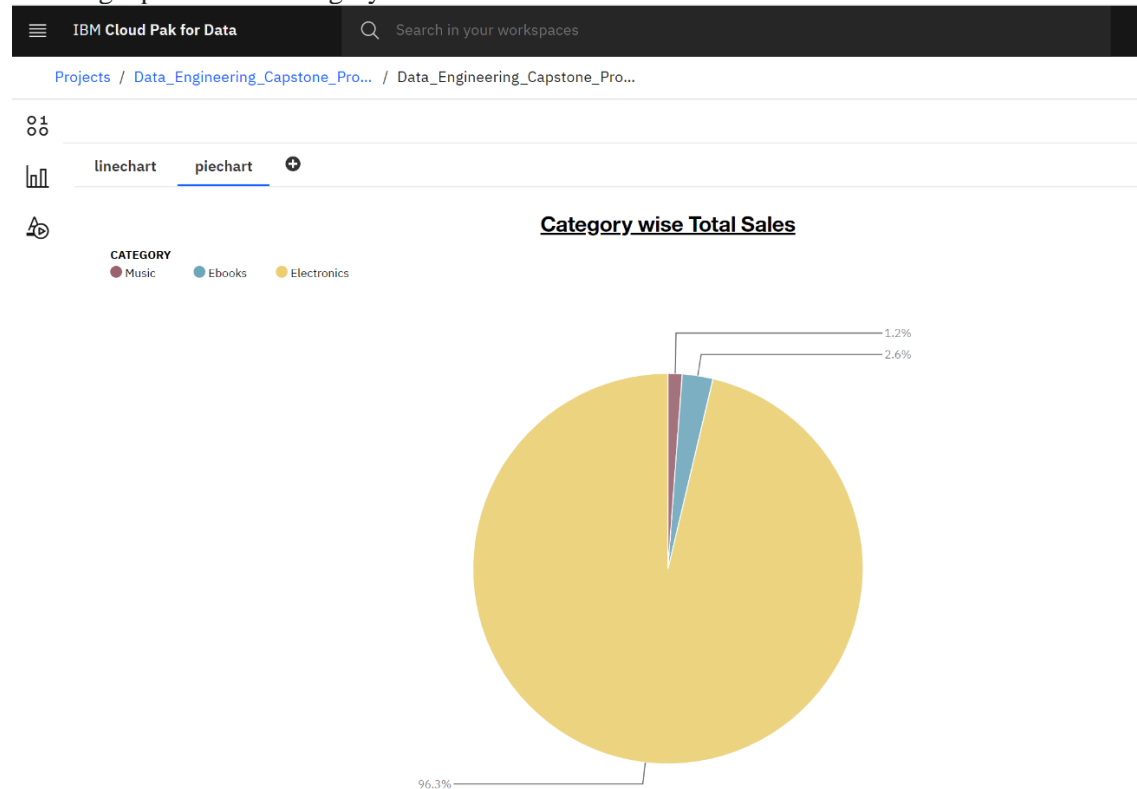
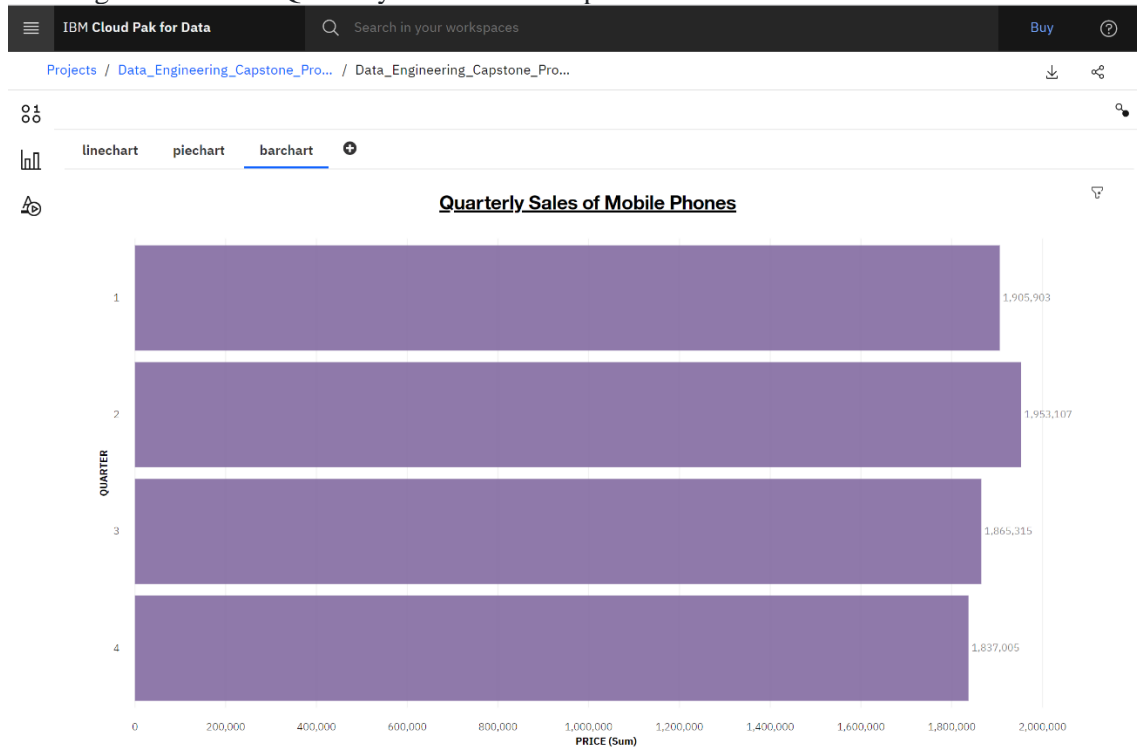# Create dashboard using IBM Cognos

Creating a line chart of month wise total sales for the year 2020.



Creating a pie chart of category wise total sales.

Creating a bar chart of Quarterly sales of mobile phones.



# Module 5.1:

## Scenario

You are a data engineer at an e-commerce company. You need to keep data synchronized between different databases/data warehouses as a part of your daily routine. One task that is routinely performed is the sync up of staging data warehouse and production data warehouse. Automating this sync up will save you a lot of time and standardize your process. You will be given a set of python scripts to start with. You will use/modify them to perform the incremental data load from MySQL server which acts as a staging warehouse to the IBM DB2 which is a production data warehouse. This script will be scheduled by the data engineers to sync up the data between the staging and production data warehouse.

# Automate loading of incremental data into the data warehouse

Creating the function get_last_rowid() that connect to the DB2 data warehouse and return the last rowid of the table sales_data.

```python
automation.py > ...
36
37    # Find out the last rowid from DB2 data warehouse
38    # The function get_last_rowid must return the last rowid of the table sales_data or
39
40  ∨ def get_last_rowid():
41        SQL = "SELECT MAX(rowid) FROM SALES_DATA"
42        stmt = ibm_db.exec_immediate(conn, SQL)
43        tuple = ibm_db.fetch_tuple(stmt)
44  ∨     while tuple != False:
45            last_row = tuple[0]
46            return last_row
47
48    last_row_id = get_last_rowid()
49    print("\nLast row id on production datawarehouse = ", last_row_id)
50
```

Creating the function get_latest_records() that connect to MySQL database and return all records later than the given last_rowid from DB2 data warehouse.
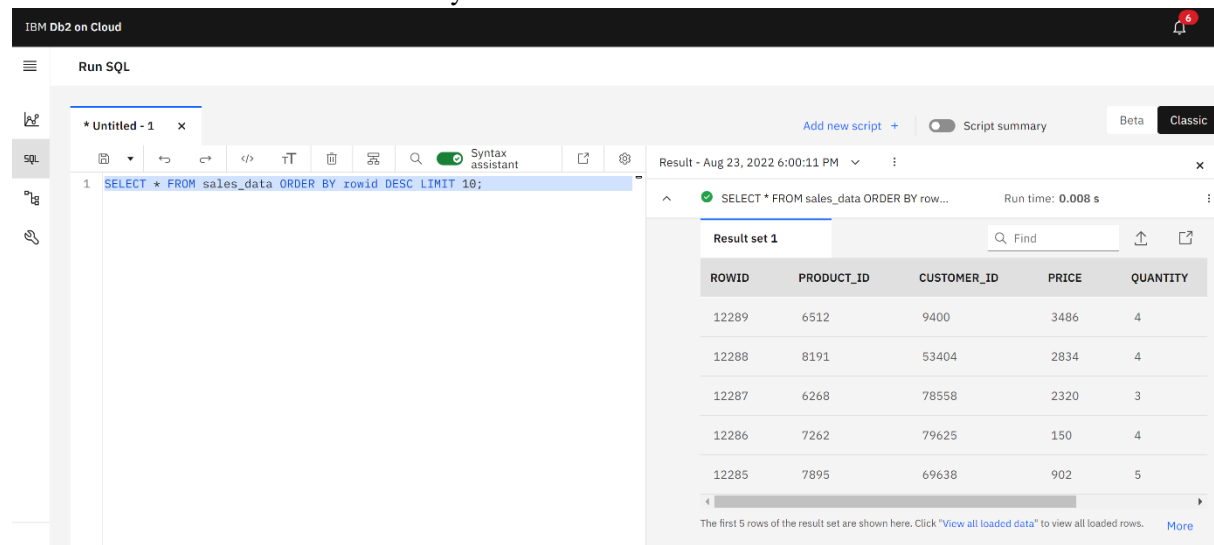
```python
51    # List out all records in MySQL database with rowid greater than the one on the Data wa
52    # The function get_latest_records must return a list of all records that have a rowid g
53
54  ∨ def get_latest_records(rowid):
55        latest_row = []
56        SQL = f"""SELECT * FROM sales_data WHERE rowid > {rowid}"""
57        cursor.execute(SQL)
58  ∨     for row in cursor.fetchall():
59            latest_row.append(row)
60        return latest_row
61
62    new_records = get_latest_records(last_row_id)
63
64    print("\nNew rows on staging datawarehouse = ", len(new_records))
65
```

Creating the function insert_records() that connect to the DB2 data warehouse and insert all the given records (latest records from MySQL database).

```python
65
66    # Insert the additional records from MySQL into DB2 data warehouse.
67    # The function insert_records must insert all the records passed to it into the sales_data table in IBM
68
69  ∨ def insert_records(records):
70  ∨     for record in records:
71            SQL = f"INSERT INTO SALES_DATA(rowid,product_id,customer_id,quantity) VALUES {record}"
72            ibm_db.exec_immediate(conn,SQL)
73
74    insert_records(new_records)
75    print("\nNew rows inserted into production datawarehouse = ", len(new_records))
76
77    # disconnect from mysql warehouse
78    connection.close()
79    # disconnect from DB2 data warehouse
80    ibm_db.close(conn)
81    # End of program
82
```

# Testing the data synchronization

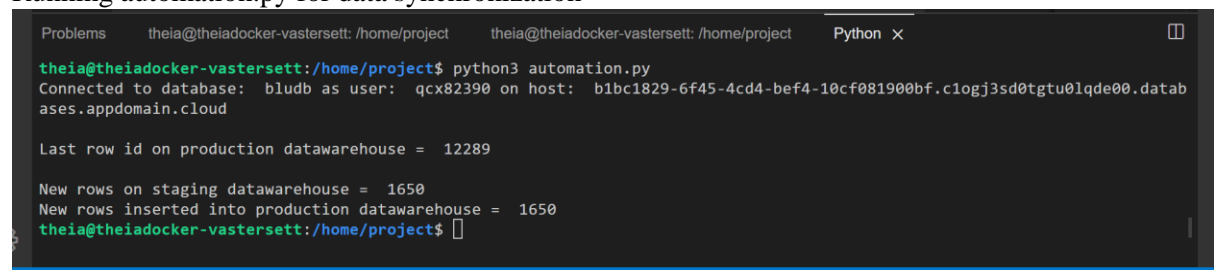Data in DB2 data warehouse before synchronization



Running automation.py for data synchronization



Data in DB2 data warehouse after synchronization



Data have been successfully synchronized.

# Module 5.2:

## Scenario
Write a pipeline that analyses the web server log file, extracts the required fields, transforms and load.

## Creating Data Pipelines using Apache Airflow
Defining the DAG arguments.

```python
process_web_log.py ✕

airflow > dags > capstone > process_web_log.py > ...
1    # import the libraries
2    from datetime import timedelta
3    # The DAG object; we'll need this to instantiate a DAG
4    from airflow import DAG
5    # Operators; we need this to write tasks!
6    from airflow.operators.bash_operator import BashOperator
7    # This makes scheduling easy
8    from airflow.utils.dates import days_ago
9
10   #defining DAG arguments
11   # You can override them on a per-task basis during operator initialization
12   default_args = {
13    'owner': 'Phone Sett Paing',
14    'start_date': days_ago(0),
15    'email': ['vastersett@gmail.com'],
16    'email_on_failure': True,
17    'email_on_retry': True,
18    'retries': 1,
19    'retry_delay': timedelta(minutes=5),
20   }
21
22   # defining the DAG
```

Defining the DAG.

```python
process_web_log.py ✕

airflow > dags > capstone > process_web_log.py > ...
20   }
21
22   # defining the DAG
23   # define the DAG
24   dag = DAG(
25    dag_id='process_web_log',
26    default_args=default_args,
27    description='Analyzes the web server log file',
28    schedule_interval=timedelta(days=1),
29   )
30
31   # define the tasks
```

Creating a task to extract the ip address field from the web server log file into extracted_data.txt file.

```python
process_web_log.py ●
airflow > dags > capstone > process_web_log.py > ...
29    )
30
31    # define the tasks
32    # define the first task (extract)
33    extract_data = BashOperator(
34      task_id='extract_data',
35      bash_command='cut -d" " -f1 /home/project/airflow/dags/capstone/accesslog.txt > extracted_data.txt',
36      dag=dag,
37    )
38
39    # define the second task (transform)
```

Creating a task to filter out all the occurrences of ip address "198.46.149.143" from extracted_data.txt file into transformed_data.txt file.

```python
process_web_log.py ●
airflow > dags > capstone > process_web_log.py > ...
37    )
38
39    # define the second task (transform)
40    transform_data = BashOperator(
41      task_id='transform_data',
42      bash_command='grep -o "198.46.149.143" /home/project/airflow/dags/capstone/extracted_data.txt > transformed_data.txt',
43      dag=dag,
44    )
45
46    # define the third task (load)
47    load_data    BashOperator(
```

Creating a task to archive the file transformed_data.txt file into weblog.tar file.

```python
process_web_log.py ●
airflow > dags > capstone > process_web_log.py > ...
44    )
45
46    # define the third task (load)
47    load_data = BashOperator(
48      task_id='load_data',
49      bash_command='tar -cvf weblog.tar /home/project/airflow/dags/capstone/transformed_data.txt',
50      dag=dag,
51    )
52
53    # task pipeline
```

Defining the task pipeline.

```python
process_web_log.py ●
airflow > dags > capstone > process_web_log.py > ...
51    )
52
53    # task pipeline
54    extract_data >> transform_data >> load_data
55    |
```

# Getting the DAG operational

Submitting the DAG



Unpausing the DAG



Monitoring the DAG



# Module 6:

# Analyse search terms on the e-commerce web server

Download the search term data set for the e-commerce web server and analyse the data set, load the sales forecast model and predict the sales for the year 2023.

Printing the number of rows and columns in the data set.

```
[7]:  # Load the csv into a spark dataframe
```

```
[8]:  df = spark.read.csv('searchterms.csv',header=True).cache()
```

```
[9]:  # Print the number of rows and columns
      # Take a screenshot of the code and name it as shape.jpg)
```

```
[10]:  print((df.count(),len(df.columns)))
```

```
[Stage 1:>                                                        (0 + 1) / 1]
(10000, 4)
```

Printing the top 5 rows in the data set.

```
[11]:  # Print the top 5 rows
       # Take a screenshot of the code and name it as top5rows.jpg)
```

```
[12]:  df.show(5)
```

```
+---+-----+----+-------------+
|day|month|year|   searchterm|
+---+-----+----+-------------+
| 12|   11|2021| mobile 6 inch|
| 12|   11|2021| mobile latest|
| 12|   11|2021|   tablet wifi|
| 12|   11|2021|laptop 14 inch|
| 12|   11|2021|     mobile 5g|
+---+-----+----+-------------+
only showing top 5 rows
```

Finding the data type of the column "searchterm".

```
[13]:  # Find out the datatype of the column searchterm?
       # Take a screenshot of the code and name it as datatype.jpg)
```

```
[14]:  df.printSchema()
```

```
root
 |-- day: string (nullable = true)
 |-- month: string (nullable = true)
 |-- year: string (nullable = true)
 |-- searchterm: string (nullable = true)
```

Finding the number of times, the term "gaming laptop" was searched.

```
[15]: # How many times was the term `gaming laptop` searched?
      # Take a screenshot of the code and name it as gaminglaptop.jpg)
```

```
[16]: df.groupby("searchterm").count().filter("searchterm='gaming laptop'").show()
```

```
[Stage 13:======================================>          (58 + 8) / 75]
+-------------+-----+
|   searchterm|count|
+-------------+-----+
|gaming laptop|  499|
+-------------+-----+
```

Printing the top 5 frequently used search terms.

```
[17]: # Print the top 5 most frequently used search terms?
      # Take a screenshot of the code and name it as top5terms.jpg)
```

```
[18]: df.groupby("searchterm").count().sort('count',ascending=False).show(5)
```

```
[Stage 15:=========================================>        (176 + 11) / 200]
+-------------+-----+
|   searchterm|count|
+-------------+-----+
|mobile 6 inch| 2312|
|    mobile 5g| 2301|
|mobile latest| 1327|
|       laptop|  935|
|  tablet wifi|  896|
+-------------+-----+
only showing top 5 rows
```

Loading the sales forecast model.

```
sales_prediction.model/data/.part-00000-1db9fe2f-4d93-4b1f-966b-3b09e72d664e·
sales_prediction.model/data/._SUCCESS.crc
```

```
[22]: # Load the sales forecast model.
      # Take a screenshot of the code and name it as loadmodel.jpg)
      from pyspark.ml.regression import LinearRegressionModel

      model=LinearRegressionModel.load('sales_prediction.model')
```

```
[23]: # Using the sales forecast model, predict the sales for the year of 2023.
      # Take a screenshot of the code and name it as forecast.jpg
```

Predicting the sales for the year 2023, using the sales forecast model.

```
[23]: # Using the sales forecast model, predict the sales for the year of 2023.
      # Take a screenshot of the code and name it as forecast.jpg
```

```
[24]: from pyspark.ml.feature import VectorAssembler

      def predict(year):
          assembler = VectorAssembler(inputCols=["year"],outputCol="features")
          data=[[year,0]]
          columns=["year","sales"]
          _ = spark.createDataFrame(data,columns)
          __ = assembler.transform(_).select('features','sales')
          predictions = model.transform(__)
          predictions.select('prediction').show()

      predict(2023)
```

```
+------------------+
|        prediction|
+------------------+
|175.16564294006457|
+------------------+
```