Data Analytics with Hadoop Project.

- 1. Tools
- a. Hadoop Server
- b. Visual Studio Code
- c. Pyhton Version 3.10 or Upper
- d. Install the Package (psycopg2-binary, SQLAlchemy, sqlparse, pandas, hdfs, pywebhdfs, mrjob)
- e. DBeaver

2. Steps by Steps

a. config.json

```
"marketplace prod": {
    "host": "103.150.197.96",
    "db": "project 4",
    "user": "postgres",
    "password": "dsbatch13!",
    "port": 5431
},
"dwh": {
    "host": "103.150.197.96",
    "db": "dwh project 4",
    "user": "postgres",
    "password": "dsbatch13!",
    "port": 5431
},
"hadoop": {
    "client": "http://hadoop-server:9870"
}
```

1. config.json

Script ini berfungsi untuk mengonfigurasi 3 koneksi serta pengaturan akses antara dua database (data source marketplace_prod dan destination dwh) dengan host, nama database, pengguna, kata sandi, dan port yang telah dibuat yang kemudian data tersebut tersimpan didalam Hadoop (http://hadoop-server:9870).

b. Design DWH

```
DROP TABLE IF EXISTS dim_orders_alif;

CREATE TABLE dim_orders_alif (
    order_id INT NOT NULL,
    order_date DATE NOT NULL,
    user_id INT NOT NULL,
    payment_name VARCHAR(255),
    shipper_name VARCHAR(255),
    order_price INT,
    order_discount INT,
    voucher_name VARCHAR(255),
    voucher_price INT,
    order_total INT,
    rating_status VARCHAR(255)
);
```

2. Design DWH

Query diatas berfungsi untuk membuat Table dim_orders_alif yang akan digunakan sebagai tempat penyimpanan data dari data source.

c. Koneksi ke Hadoop

```
import os
import json
import psycopg2
import hdfs
from sqlalchemy import create engine
def config(connection db):
   path = os.getcwd()
   with open(path+'/'+'config.json') as file:
       conf =json.load(file)[connection db]
   return conf
def psql conn(conf, name conn):
   try:
        conn = psycopg2.connect(
           host=conf['host'],
            database=conf['db'],
            user=conf['user'],
            password=conf['password'],
            port=conf['port']
        )
        print(f'[INFO] Success connect PostgreSQL {name conn}')
        engine =
create engine(f"postgresql+psycopg2://{conf['user']}:{conf['password']}@{conf[
'host']}:{conf['port']}/{conf['db']}")
        return conn, engine
    except Exception as e:
       print(f"[INFO] Can't connect PostgreSQL {name conn}")
        print(str(e))
def hadoop conn(conf):
   client = conf['client']
   try:
        conn = hdfs.InsecureClient(client)
        print("[INFO] Success to connect HADOOP ...")
        return conn
    except:
        print("[INFO] Can't connect HADOOP ...")
```

3. Koneksi ke Hadoop

Script berfungsi untuk mengatur koneksi dan akses ke database PostgreSQL dan Hadoop. Pertama fungsi "config" digunakan untuk membaca file konfigurasi JSON ("config.json") yang telah dibuat sebelumnya. Kemudian fungsi "psql_conn" digunakan untuk membuat koneksi ke database PostgreSQL dengan menggunakan package psycopg2. Terakhir fungsi "hadoop_conn" digunakan untuk membuat koneksi ke Hadoop dengan menggunakan package hdfs.

d. App.py

```
# connection data source
    conf = connection.config('marketplace prod')
    conn, engine = connection.psql conn(conf, 'DataSource')
    cursor = conn.cursor()
    # connection dwh
    conf dwh = connection.config('dwh')
    conn dwh, engine_dwh = connection.psql_conn(conf_dwh, 'DataWarehouse')
    cursor dwh = conn dwh.cursor()
    # connection dwh to hadoop
    conf dwh hadoop = connection.config('hadoop')
    client hadoop = connection.hadoop conn(conf dwh hadoop)
    # get query string
   path query = os.getcwd()+'/query/'
    query = sqlparse.format(
        open(path query+'query.sql', 'r').read(), strip comments=True
    ).strip()
    # get schema dwh design
   path dwh design = os.getcwd()+'/query/'
    dwh design = sqlparse.format(
        open(path_dwh_design+'dwh_design.sql', 'r').read(),
strip comments=True
   ).strip()
    try:
        # get data
        print('[INFO] Service ETL is Running ...')
        df = pd.read sql(query, engine)
        # upload file to hadoop as DWH
        filetime = datetime.now().strftime('%Y%m%d')
        my file = f'dim orders {filetime} alif.csv'
       my file with path = f'/digitalskola/project4/{my file}'
       with client hadoop.write(my file with path, encoding='utf-8') as
writer:
            df.to csv(writer, index=False)
        print(f"[INFO] Upload Data in HADOOP Success ...")
        # get data from hadoop for create data mart
        print(f"[INFO] Get Data in HADOOP ....")
        hdfs=PyWebHdfsClient(host='hadoop-server',port='9870',
user name='hduser')
        filetime = datetime.now().strftime('%Y%m%d')
        data = hdfs.read file(str(my file with path))
        data = data.decode().split('\n')
        data list = []
        for item in data:
            item = item.replace('\r', '')
```

4b. App.py

Pada script diatas diatas menjelasakan Connection data source: Membuat koneksi ke data source "marketplace_prod" menggunakan fungsi "config" dan "psql_conn" dari modul "connection". Connection dwh: Membuat koneksi ke data warehouse "dwh" menggunakan fungsi "config" dan "psql_conn". Connection dwh to hadoop: Membuat koneksi ke Hadoop menggunakan fungsi "config" dan "hadoop_conn". Get query string: Membaca dan memformat string kueri dari file "query.sql" yang terletak di direktori "query/" menggunakan fungsi "sqlparse.format". Get schema dwh design: Membaca dan memformat schema desain data warehouse dari file "dwh_design.sql" yang terletak di direktori "query/". Get data: Mengeksekusi query SQL pada data source menggunakan Pandas dan menyimpan hasilnya dalam DataFrame "df". Upload file to Hadoop: Mengunggah DataFrame "df" ke Hadoop sebagai data warehouse. Nama file yang diunggah disesuaikan dengan timestamp. Get data from Hadoop: Membaca file yang diunggah dari Hadoop, melakukan manipulasi data, dan menyimpannya dalam file "dim_orders_{timestamp}_alif.csv". Perform MapReduce: Menjalankan skrip "mapReduce.py" untuk melakukan operasi MapReduce pada file hasil sebelumnya dan menyimpan hasilnya dalam file "output/Wordercount_output_hadoop_map.txt".

e. MapReduce.py

```
cols =
'order id, order date, user id, payment name, shipper name, order price, orde
r discount, voucher name, voucher price, order total, rating status'.split(
',')
def csv readline(line):
    """Given a sting CSV line, return a list of strings."""
    for row in csv.reader([line]):
        return row
class OrderDateCount(MRJob):
    def steps(self):
        return [
            MRStep (mapper=self.mapper, reducer=self.reducer),
            MRStep(reducer=self.sort)
        ]
    def mapper(self, , line):
        # Convert each line into a dictionary
        row = dict(zip(cols, csv readline(line)))
        #skip first row as header
        if row['order id'] != 'order id':
            # Yield the order date
```

5a. MapReduce.py

```
yield row['order_date'][0:7], 1

def reducer(self, key, values):
    #for 'order_date' compute
    yield None, (key,sum(values))

def sort(self, key, values):
    data = []
    for order_date, order_count in values:
        data.append((order_date, order_count))
        data.sort()

for order_date, order_count in data:
        yield order_date, order_count
```

5a. MapReduce.py

Script ini merupakan implementasi dari MapReduce menggunakan framework MRJob. Tujuannya adalah untuk menghitung jumlah pesanan berdasarkan tanggal pesanan (order_date). Variabel "cols": Menyimpan daftar nama kolom yang digunakan dalam file CSV yang akan diproses. Kemudian fungsi "csv readline": Membaca baris CSV sebagai string dan mengembalikan daftar string yang berisi nilai-nilai kolom. Lalu kelas "OrderDateCount": Merupakan kelas utama yang mengimplementasikan algoritma MapReduce. Method "steps" menentukan langkah-langkah MapReduce yang akan dilakukan, yaitu "mapper", "reducer", dan "sort". Method "mapper": Menerima baris input dan mengonversinya menjadi sebuah kamus (dictionary) menggunakan nama kolom dari "cols". Kemudian, untuk setiap baris yang bukan merupakan header, menghasilkan pasangan kunci-nilai dengan kunci berupa bulan dan tahun dari tanggal pesanan (order date) dan nilai 1. Method "reducer": Menerima pasangan kunci-nilai dari mapper, yang memiliki kunci yang sama, yaitu bulan dan tahun pesanan. Menghasilkan pasangan kunci-nilai di mana kunci adalah None (tidak digunakan) dan nilai berupa pasangan kunci-nilai dari bulan dan jumlah pesanan pada bulan tersebut. Method "sort": Menerima pasangan kunci-nilai dari reducer dan mengumpulkannya dalam sebuah list. Kemudian, list tersebut diurutkan berdasarkan tanggal pesanan (order_date). Akhirnya, menghasilkan pasangan kunci-nilai dengan kunci berupa tanggal pesanan dan nilai berupa jumlah pesanan pada tanggal tersebut.

3. Hasil

a. App.py

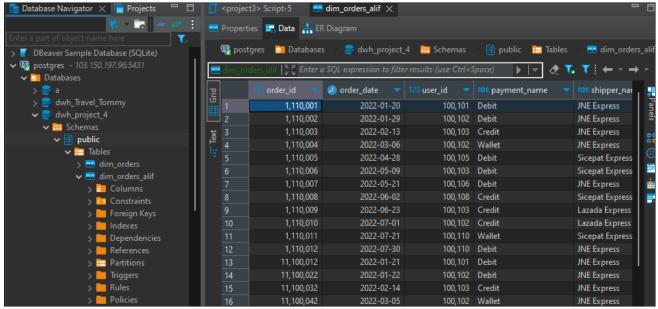
Setelah semua script tidak ada peringatan error pada Visual Studio Code jalankan command "python app.py". Kemudian akan muncul hasil seperti dibawah.

```
@LAPTOP-HMKJLOOC MINGW64 /d/BOOTCAMP/bahan/PROJECT4/digitaskola_de_10_batch_processing-main/digitaskola_de_10_batch_processing-main
$ python app.py
[INFO] Service ETL is Starting ...
[INFO] Success connect PostgreSQL DataSource
[INFO] Success connect PostgreSQL DataWarehouse
[INFO] Success to connect HADOOP ...
[INFO] Service ETL is Running ...
[INFO] Upload Data in HADOOP Success ...
[INFO] Get Data in HADOOP ...
No configs found; falling back on auto-configuration
No configs specified for inline runner
Creating temp directory C:\Users\Cerv\AppData\Local\Temp\mapReduce.Cerv.20230630.070003.839655
Running step 1 of 2...
Running step 2 of 2...
iob output is in C:\Users\Cerv\AppData\Local\Temp\mapReduce.Cerv.20230630.070003.839655\output
Streaming final output from C:\Users\Cerv\AppData\Local\Temp\mapReduce.Cerv.20230630.070003.839655\output...
 emoving temp directory C:\Users\Cerv\AppData\Local\Temp\mapReduce.Cerv.20230630.070003.839655...
 INFO] Download Data in HADOOP Success and Created file mart....
[INFO] Service ETL is Success ...
```

6. Python app.py

Pada gambar diatas menampilkan bahwa Service ETL sudah mulai hingga Service ETL berhasil.

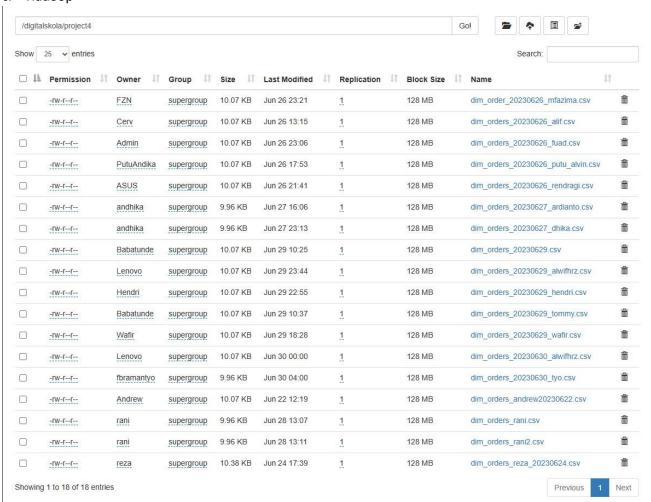
b. Data di Dbeaver



7. Data Dbeaver

Gambar diatas menampilan bahwa berhasil mengtransform data ke DBeaver.

c. Hadoop



8. Hadoop

Gambar diatas menunjukkan berhasil menyimpan data di Hadoop

d. Output CSV

```
order\_id, order\_date, user\_id, payment\_name, shipper\_name, order\_price, order\_discount, voucher\_name, voucher\_price, order\_total, rating\_status
1110001,2022-01-20,100101,Debit,JNE Express,250000,15000,New User,5000.0,230000,Medium Impact 1110002,2022-01-29,100102,Debit,JNE Express,620000,40000,New User,5000.0,575000,Medium Impact
1110003,2022-02-13,100103,Credit,JNE Express,6000000,1000000,New User,5000.0,4995000,Very Low Impact 1110004,2022-03-06,100102,Wallet,JNE Express,3150000,45000,,3105000,High Impact
1110005,2022-04-28,100105,Debit,Sicepat Express,4000000,1000000,New User,5000.0,2995000,Very Low Impact
1110006,2022-05-09,100103,Debit,Sicepat Express,4500000,1030000,,3470000,Medium High Impact 1110007,2022-05-21,100106,Debit,JNE Express,870000,25000,,845000,High Impact
1110008,2022-06-02,100108,Credit,Sicepat Express,2000000,0,New User,5000.0,1995000,Medium High Impact
1110009,2022-06-23,100103,Credit,Lazada Express,2000000,0,,,2000000,High Impact
1110010,2022-07-01,100102,Credit,Lazada Express,1050000,45000,,,1005000,Low Impact
1110011,2022-07-21,100110,Wallet,Sicepat Express,550000,15000,,,535000,High Impact
1110012,2022-07-30,100110,Debit,JNE Express,490000,35000,Body Soap Promo,10000.0,445000,Medium High Impact
11100012,2022-01-21,100101,Debit,JNE Express,250000,15000,New User,5000.0,230000,Medium Impact
11100022,2022-01-22,100102,Debit,JNE Express,620000,40000,New User,5000.0,575000,Medium Impact
11100032,2022-02-14,100103,Credit,JNE Express,6000000,1000000,New User,5000.0,4995000,Very Low Impact
11100042,2022-03-05,100102,Wallet,JNE Express,3150000,45000,,,3105000,High Impact
11100052,2022-04-28,100105,Debit,Sicepat Express,4000000,1000000,New User,5000.0,2995000,Very Low Impact
11100062,2022-05-07,100103,Debit,Sicepat Express,4500000,1030000,,,3470000,Medium High Impact
11100072,2022-05-22,100106,Debit,JNE Express,870000,25000,,,845000,High Impact
11100082,2022-06-03,100108,Credit,Sicepat Express,2000000,0,New User,5000.0,1995000,Medium High Impact
11100092,2022-06-24,100103,Credit,Lazada Express,2000000,0,,,20000000,High Impact
11100102,2022-07-04,100102,Credit,Lazada Express,1050000,45000,,,1005000,Low Impact
11100112,2022-07-23,100110,Wallet,Sicepat Express,550000,15000,,,535000,High Impact
11100122,2022-07-23,100110,Debit,JNE Express,490000,35000,Body Soap Promo,10000.0,445000,Medium High Impact 11100013,2022-01-23,100101,Debit,JNE Express,250000,15000,New User,5000.0,230000,Medium Impact
11100023,2022-01-29,100102,Debit,JNE Express,620000,40000,New User,5000.0,575000,Medium Impact
11100033.2022-02-14.100103.Credit.JNF Fxnress.6000000.1000000.New User.5000.0.4995000.Verv Low Tmnact
```

9. Output CSV

Output pada data yang tersimpan.

e. Output MapReduce

```
output > \( \bigsize \) Wordercount_output_hadoop_map.txt
        "2022-01"
  1
                       18
        "2022-02"
                       9
        "2022-03"
                       9
        "2022-04"
                       9
        "2022-05"
                       18
        "2022-06"
                       18
        "2022-07"
                       27
```

10. Output MapReduce

Output hasil dari MapReduce.