# Database Project Report

**Student 1 : Fahadh Mohamed – 62720**

**Student 2 : Kuiyu Ding – 62716**
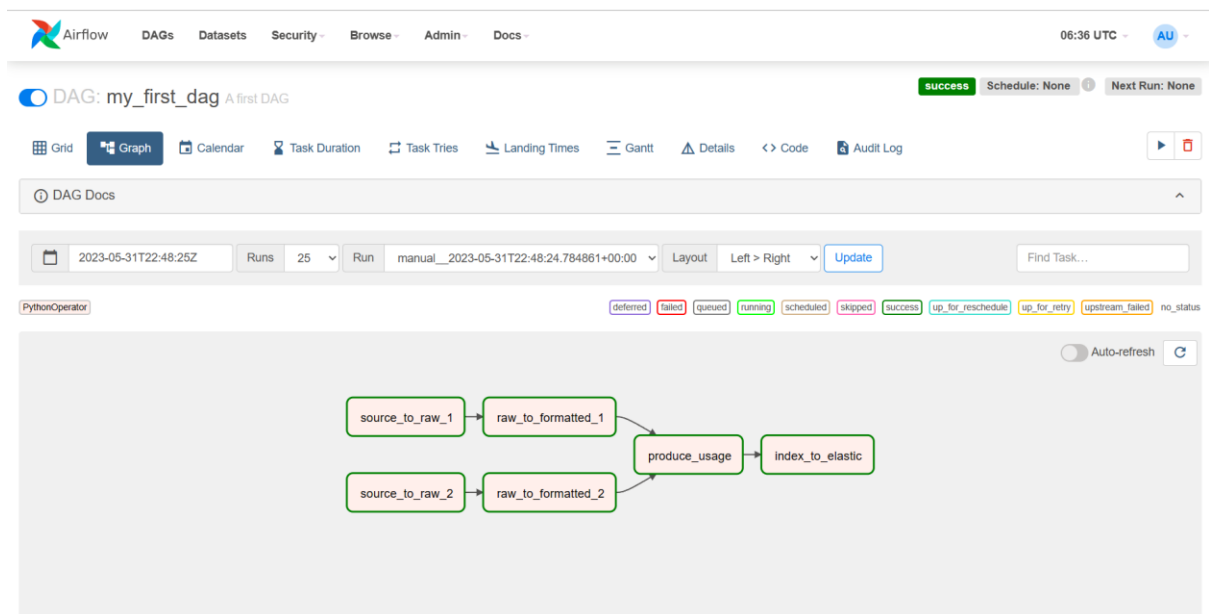
## Project Overview :

In this project, we analyse the data of movies , we have showed two usages, First usage is to analyse the average rating for each year and the second usage is to recommend the movies based on highest numVotes for the movies came after 2000(year) and has rating above 8.

## Airflow :

Airflow is an open-source platform designed to programmatically author, schedule, and monitor workflows. It allows users to define complex data pipelines as a collection of tasks, each with its dependencies and scheduling properties.

## Dag:

We created a dag to perform set of tasks which converts from source to formatted and combine the results, finally it will show in the website using elasticsearch and kibana. Each task will run automatically because of the airflow.

## Source_to_raw_1 :

For the source 1, we have to write a code to get the data from api, the api is taken from rapidapi website which consists of top 100 imdb movies.

```
url = "https://imdb-top-100-movies1.p.rapidapi.com/"

headers = {
    "X-RapidAPI-Key": "b4c2b91fcemshec2addb6175f22ep122986jsn03535c95625d",
    "X-RapidAPI-Host": "imdb-top-100-movies1.p.rapidapi.com"
}
```

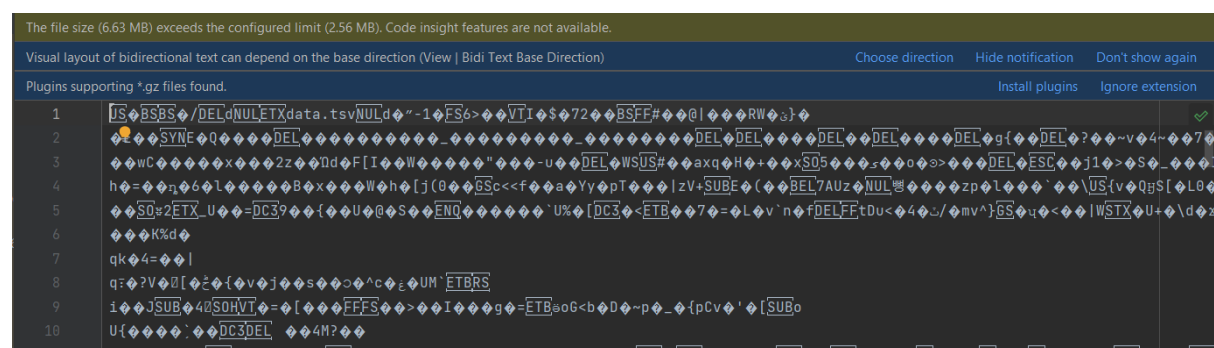I get the data as Json format and I saved it in the raw folder. You can see the data below,

**Source_to_raw_2 :**

Similarly for the source 2, we have to write a code to get the data from api, the api is taken from imdb website, which shows ratings and numVotes of the movies.

```python
url = 'https://datasets.imdbws.com/title.ratings.tsv.gz'



path =  DATALAKE_ROOT_FOLDER + "raw/imdb/" + current_day + "/"
os.makedirs(path, exist_ok=True)
r = requests.get(url, allow_redirects=True)
open(path + 'title.ratings.tsv.gz', 'wb').write(r.content)
```

I get the data as tsv.gz file and I saved it in raw folder. You can see the data below.

The file size (6.63 MB) exceeds the configured limit (2.56 MB). Code insight features are not available.

Visual layout of bidirectional text can depend on the base direction (View | Bidi Text Base Direction)    Choose direction   Hide notification   Don't show again

Plugins supporting *.gz files found.    Install plugins   Ignore extension

1    US�BSBS�/DELdNULETXdata.tsvNULd�"-1�FS6>��VTI�$�72��BSFF#��@|���RW�₃}�
2    �²�₂��SYNE�Q���⚬DEL����������⚬-����������-��������DEL⚬DEL����DEL⚬⚬DEL����DEL⚬g{⚬⚬DEL⚬?⚬⚬~v⚬4~⚬⚬7⚬
3    ��⚬WC�����⚬x⚬��2z⚬⚬Dd⚬F[I⚬⚬W�����⚬"��⚬-u⚬⚬DEL⚬WSUS#⚬⚬axq⚬H⚬+⚬⚬xSO5⚬⚬⚬s⚬⚬o⚬⊃>⚬⚬⚬DEL⚬ESC⚬⚬j1⚬>⚬S⚬_⚬⚬⚬I
4    h⚬=⚬⚬η⚬6⚬l⚬⚬⚬⚬⚬B⚬x⚬⚬⚬W⚬h⚬[j(0⚬⚬GSc<<f⚬⚬a⚬Yy⚬pT⚬⚬⚬|zV+SUBE⚬(⚬⚬BEL7AUz⚬NUL嶝⚬⚬⚬⚬zp⚬l⚬⚬⚬`⚬⚬\US{v⚬Q♁S[⚬L0⚬
5    ⚬⚬SO≠2ETX_U⚬⚬=DC39⚬⚬{⚬⚬U⚬@⚬S⚬⚬ENQ⚬⚬⚬⚬⚬⚬`U%⚬[DC3⚬<ETB⚬⚬7⚬=⚬L⚬v`n⚬fDELFFtDu<⚬4⚬⌁/⚬mv^}GS⚬q⚬<⚬⚬|WSTX⚬U+⚬\d⚬x
6    ⚬⚬⚬K%d⚬
7    qk⚬4=⚬⚬|
8    q₮⚬?V⚬Ø[⚬ċ⚬{⚬v⚬j⚬⚬s⚬⚬⊃⚬^c⚬¿⚬UM`ETBRS
9    i⚬⚬JSUB⚬4ØSOHVT⚬=⚬[⚬⚬⚬FFFS⚬⚬>⚬⚬I⚬⚬⚬g⚬=ETBăoG<b⚬D⚬~p⚬_⚬{pCv⚬'⚬[SUBo
10   U{⚬⚬⚬⚬`⚬⚬DC3DEL ⚬⚬4M?⚬⚬

**Raw_to_formatted_1:**

Now the data is in json format, it is not suitable for analysis. So, we have to convert into parquet format. We used spark to convert the data into parquet format.

```python
df = spark.read.json(RAW_FOLDER)
parquet_file_name = file_name.replace(".json", ".snappy.parquet")
df.write.format("parquet").option("compression", "snappy").save(FORMATTED_FOLDER + parquet_file_name)
```

After that I saved the data in the formatted folder. You can see the data below, which shows id, title, description of the movies, imdbid, rating, release year, director name, genre etc...

## Raw_to_formatted_2:

Similarly we have to do this for raw_to_formatted_2, the data is in tsv.gz format, it is not suitable for analysis. So, we have to convert into parquet format. We used spark to convert the data into parquet format.

```python
# Read TSV data using Spark
df = spark.read.option("header", True).option("delimiter", "\t").csv(RAW_FOLDER)

# Write data as Parquet with Snappy compression
parquet_file_name = file_name.replace(".tsv.gz", ".snappy.parquet")
df.write.format("parquet").option("compression", "snappy").save(FORMATTED_FOLDER + parquet_file_name)
```

And we saved the parquet data in formatted folder. You can see the data below,

## Produce_usage :

After converting this data to parquet, we have to combine this data. We used spark to combine this data, we have to write a SQL query to join the data by using (JOIN), here we joined the data using imdbid from both of the data.

```
averageRating_year = sqlContext.sql(" SELECT api2.year, "
" AVG(api1.averageRating) AS average_rating "
" FROM data AS api1 "
" JOIN ratings AS api2 ON api1.tconst = api2.imdbid "
" GROUP BY api2.year "
" ORDER BY api2.year; " )


numVotes_movies = sqlContext.sql(" SELECT api1.numVotes, "
" api1.averageRating AS rating, "
" api2.title, "
" api2.year "
" FROM data AS api1 "
" JOIN ratings AS api2 ON api1.tconst = api2.imdbid "
" WHERE api2.year > 2000 AND api1.averageRating > 8.0 "
)
```

**1ˢᵗ usage :** we joined a sql query using imdbid from both of the formatted data and calculated average rating by each year.

**2ⁿᵈ usage :** Similarly, we joined a sql query using imdbid from both of the formatted data and we displayed the data that has above average rating 8 and the movies above the year 2000. It shows numVotes, rating, title and year.

And then, we converted the data into parquet format and saved the data in usage folder.
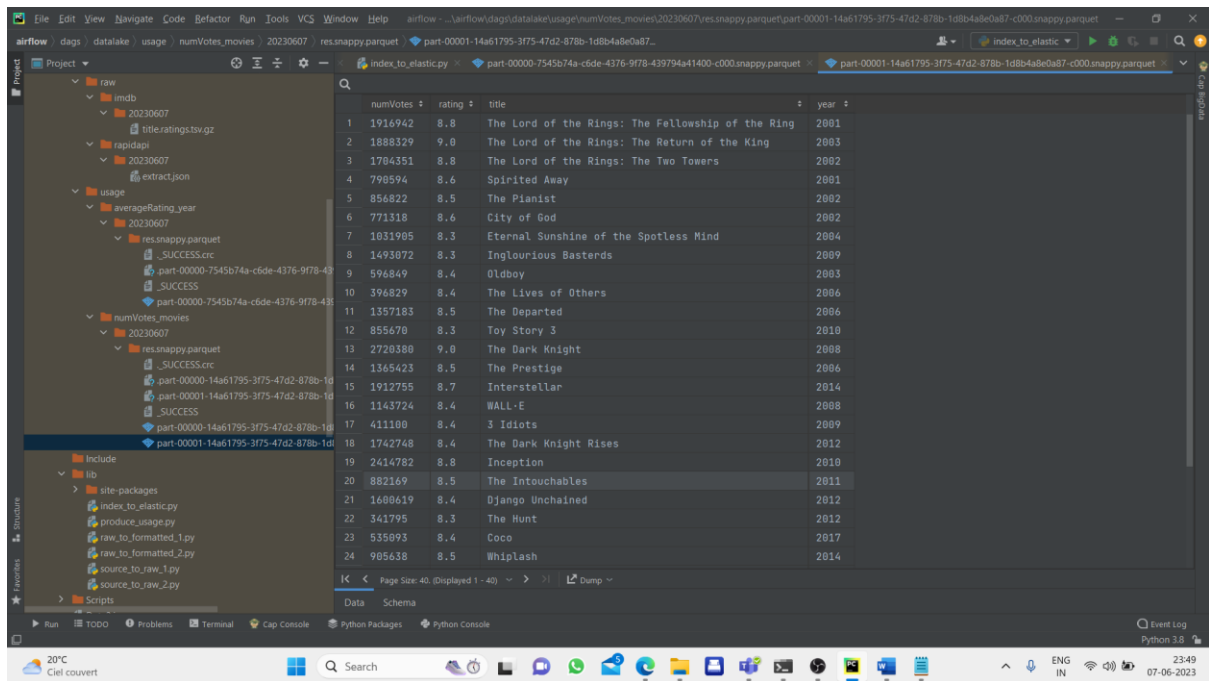
```python
# Check content of the DataFrame stats_df and save it:
print(averageRating_year.show())
averageRating_year.write.save(USAGE_FOLDER1 + "res.snappy.parquet", mode="overwrite")

# Check content of the DataFrame top10_df  and save it:
print(numVotes_movies.show())
numVotes_movies.write.save(USAGE_FOLDER2 + "res.snappy.parquet", mode="overwrite")
```

You can see the usage data in parqueted format.

## Index_to_elastic:

After getting the final data from combing the data, We used kibana, elasticsearch to display the data in website. Initially we have to create a index in elasticsearch, we created index for both of the usages.

```
C:\Users\fahad\Downloads\elasticsearch-7.8.1-windows-x86_64\elasticsearch-7.8.1\bin>curl -X GET http://localhost:9200/_cat/indices?v
health status index          uuid                    pri rep docs.count docs.deleted store.size pri.store.size
green  open   movies         0MlNvC2GQ52jegm1j77pBA   1   0      55           0         6.5kb       6.5kb
green  open   moviesvotes    U4tcMAN8R5OoMtNgK46trw   1   0      31           0         7.7kb       7.7kb
```

As you can see, we created two indexes, movies and moviesdata. And we updated the data in the elastic search.

```
settings = {
    "settings" : {
        "number_of_shards" : 1,
        "number_of_replicas" : 0
    },
    "mappings" : {
        "properties": {
            "Number" : {"type":"integer"},
            "Year" : {"type":"integer"},
            "Average_rating" : {"type" : "float"},
        }
    }
}
client.indices.create(index="movies",body=settings)
```

```
settings1 = {
    "settings" : {
        "number_of_shards" : 1,
        "number_of_replicas" : 0
    },
    "mappings" : {
        "properties": {
            "Number" : {"type":"integer"},
            "Year" : {"type":"integer"},
            "Average_rating" : {"type" : "float"},
            "Votes" : {"type" : "integer"},
            "Title" : {"type" : "text"},
        }
    }
}
client.indices.create(index="moviesvotes",body=settings1)
```
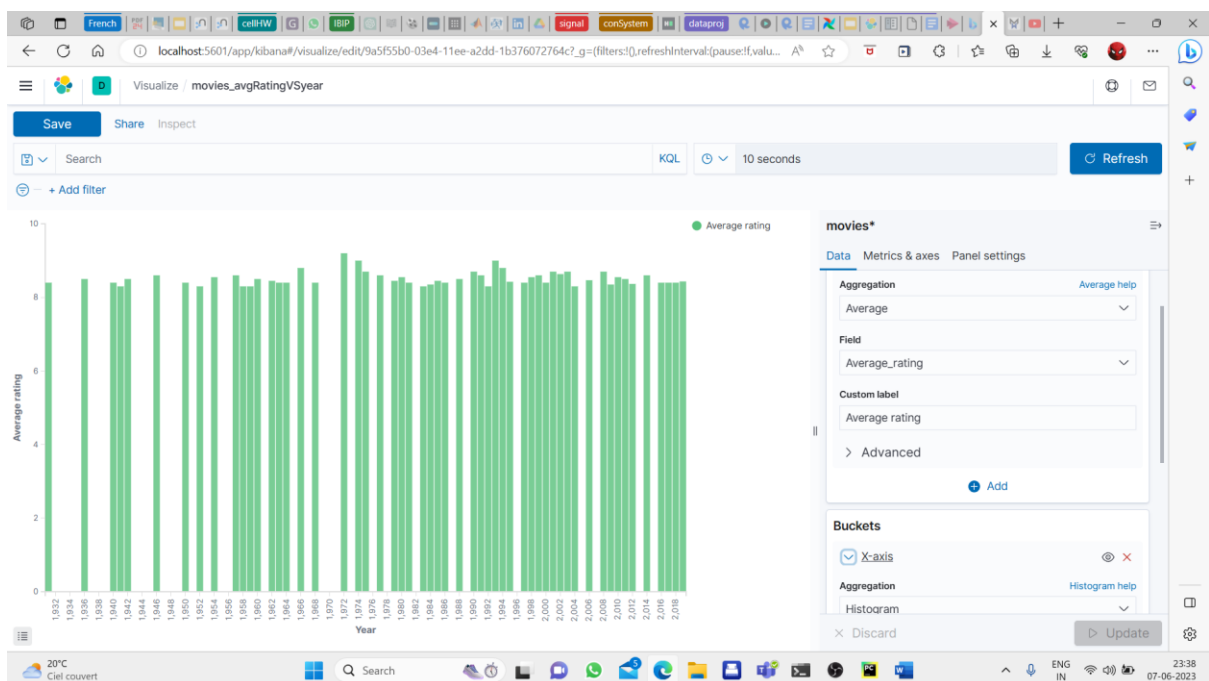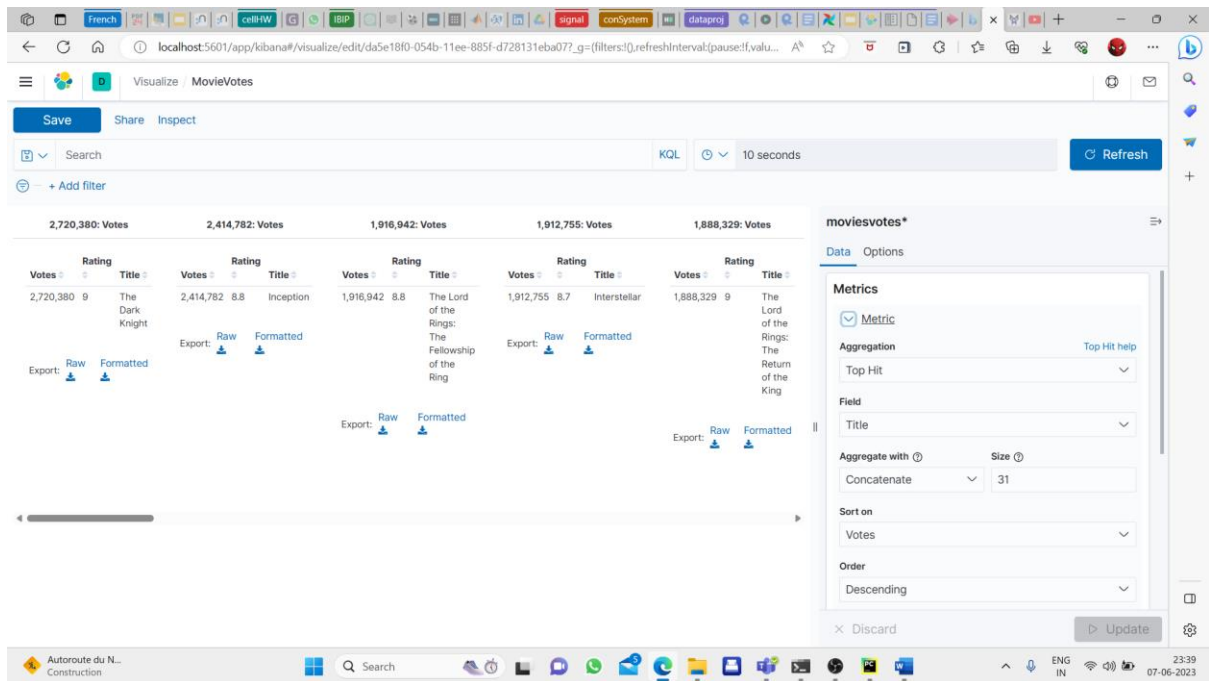
**Elasticsearch and Kibana :**

**Usage 1 :** To display the data in elasticsearch, we have to create a index and create a visualization, For the usage 1, we used vertical bar to display the data which shows average rating by year wise.
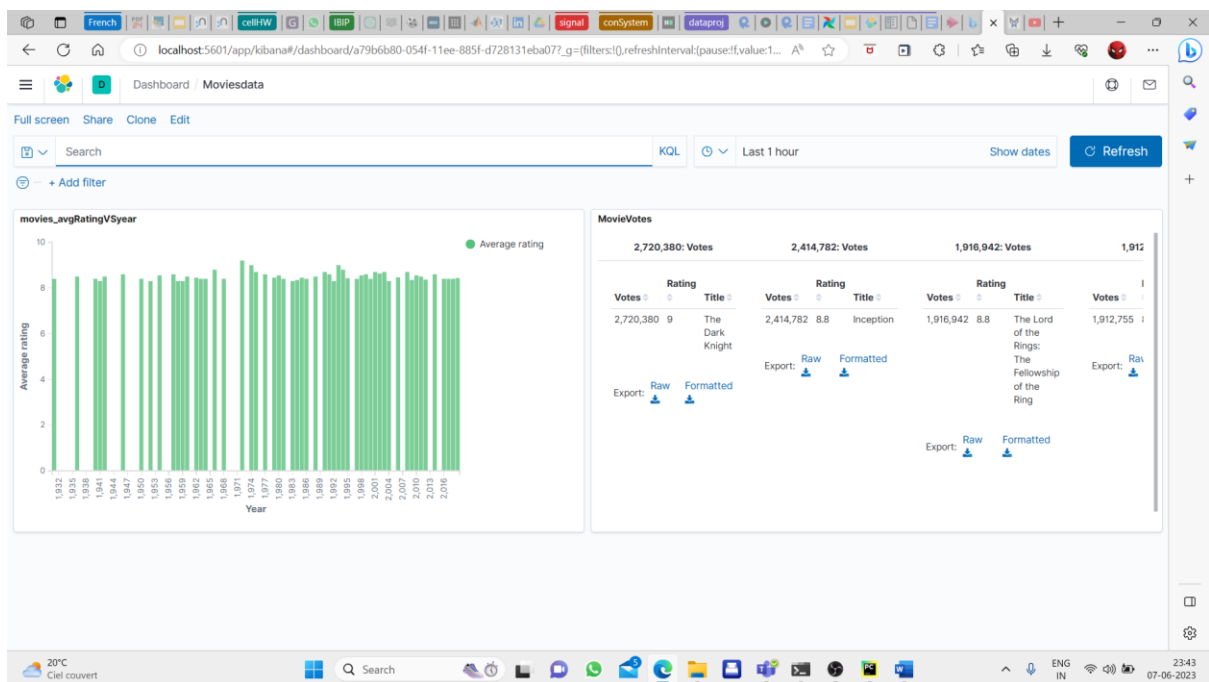


**Usage 2 :** For the usage 2, we used data model, which displays the numvotes of the movies by descending order, we also displayed the rating with this result.

## Elasticsearch Dashboard :

After getting this two data usages, we showed this data in the dashboard.

We have created specific folder to save the data.