



Data Warehouse

Analyse de Marché d'Hôtellerie à Marrakech

Encadré par Pr Imade Benelallam

AARABA ANASS

aaaraba@insea.ac.ma

DSE

OMAR FAROUQ LAFDIL

o.lafdil@insea.ac.ma

DSE

Année académique 2021/2022

Table des matières :

Introduction

- 1- Conception du datamart
- 2- Apify TripAdvisor Scraper
- 3- Préparation de l'environnement de Pycharm et Airflow sous WSL (Ubunto)
- 4- Création et exécution du DAG TripAdvisor Hotels dag sous Airflow
- 5- Dashbording sous Power BI

Conclusion

Introduction :

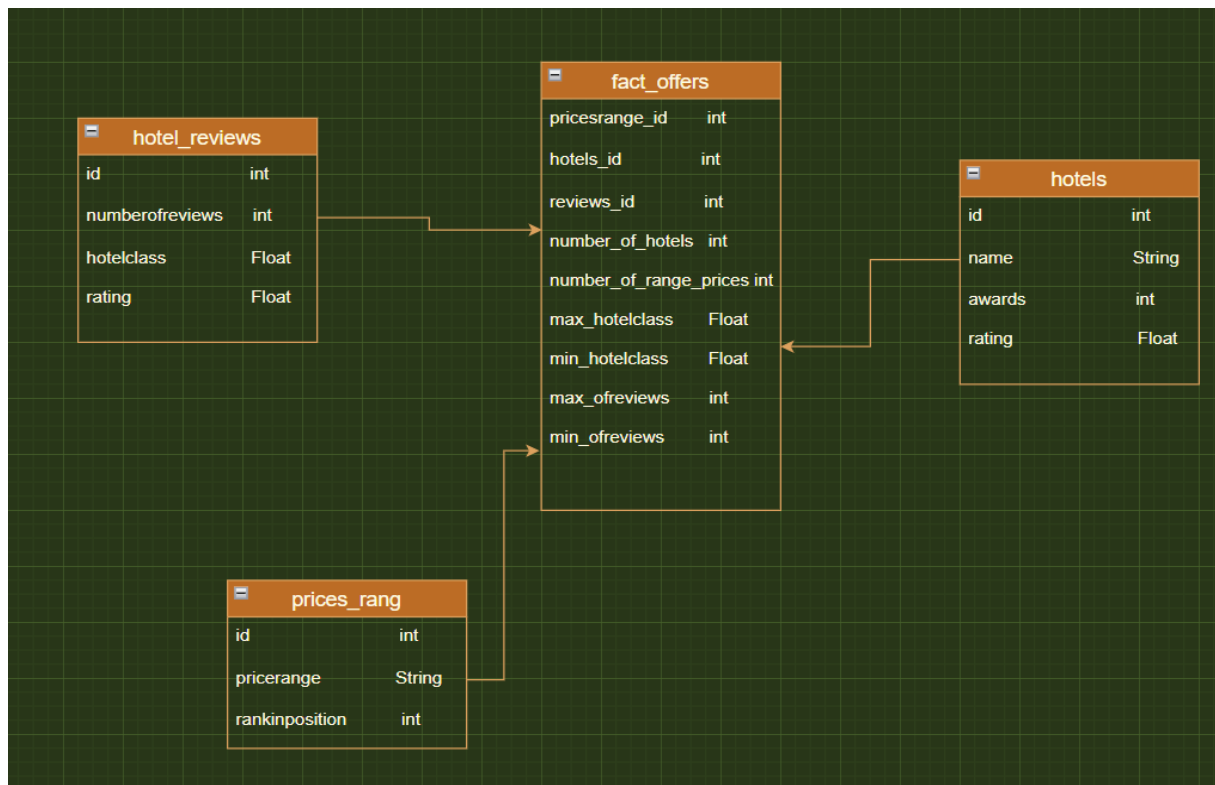
La ville de Marrakech est la capitale touristique du Maroc. L'hôtellerie dans Marrakech est le facteur d'étude de tourisme dans la ville. Après la réouverture des frontières connue pendant la pandémie de COVID-19, les responsables du secteur touristique, notamment le ministère du tourisme, pourrait demander l'acquisition un système décisionnel permettant de gérer les offres des hôtels à Marrakech.

Ce projet s'inscrit dans le cadre de l'élément de module Data Warehouse & Business Analytics. Il est réalisé en utilisant Airflow, MySQL Workbench, Mysql server, Apify tripadvisor scraper et Power BI.

I- Conception du datamart :

On veut produire un entrepôt de données pour étudier les offres des hôtels dans la ville de Marrakech à partir de trois tables de dimensions :
hotel_reviews, hotels et prices_rang.

Le schéma en étoile pour l'étude des offres d'hôtels est :



Sous Mysql-python, son implémentation est donnée par :

```

2
3 import mysql.connector as mysql
4
5 conn = mysql.connect(host='127.0.0.1', user='root', password='Anasaar123@', database="tp2_mysql")
6 cursor = conn.cursor()
7
8 #tables des dimensions
9 cursor.execute("CREATE TABLE IF NOT EXISTS prices_rang(id INT PRIMARY KEY NOT NULL AUTO_INCREMENT, pricerange VARCHAR "
10                "(30), "
11                "rankingposition INT)")
12 cursor.execute("CREATE TABLE IF NOT EXISTS hotels(id INT PRIMARY KEY NOT NULL AUTO_INCREMENT, name VARCHAR (100), "
13                "awards INT, "
14                "rating DOUBLE)")
15 cursor.execute("CREATE TABLE IF NOT EXISTS hotel_reviews(id INT PRIMARY KEY NOT NULL AUTO_INCREMENT, numberofreviews "
16                  "INT, "
17                  "hotelclass DOUBLE, "
18                  "rating DOUBLE)")
19 #Table de fait
20 cursor.execute("CREATE TABLE IF NOT EXISTS fact_offers(pricesrange_id INT, hotels_id INT, reviews_id INT, "
21                "number_of_hotels INT, number_ofrange_prices INT, max_hotelclass INT, min_hotelclass INT, "
22                "max_ofreviews INT, min_ofreviews INT, KEY pricesrange_id (pricesrange_id), KEY hotels_id (hotels_id), "
23                "KEY reviews_id (reviews_id), CONSTRAINT fact_offers_ibfk_1 FOREIGN KEY (pricesrange_id) REFERENCES "
24                "prices_rang (id), CONSTRAINT fact_offers_ibfk_2 FOREIGN KEY (hotels_id) REFERENCES hotels (id), "
25                "CONSTRAINT fact_offers_ibfk_3 FOREIGN KEY (reviews_id) REFERENCES hotel_reviews (id))")
26
27

```

Résultats dans MYSQL SERVER :

```

mysql> show tables;
+-----+
| Tables_in_tp2_mysql |
+-----+
| fact_offers          |
| hotel_reviews        |
| hotels               |
| prices_rang          |
+-----+
4 rows in set (0.00 sec)

```

II- Apify TripAdvisor Scraper :

C'est un robot web hébergé qui permet à toute personne possédant des compétences élémentaires en programmation d'extraire des données structurées de n'importe quel site web. Contrairement au scraping web par pointer-cliquer, Apify fonctionne sur des sites web modernes de plus en plus complexes et dynamiques. Il peut être utilisé par une large gamme d'utilisateurs, quiconque a besoin de données provenant du web, qu'il s'agisse d'un étudiant, un journaliste, une start-up ou une grande entreprise.

Il s'agit d'explorer des sites Web, d'en extraire des données structurées et de les exporter dans des formats tels qu'Excel, CSV ou JSON.

Pour notre cas on utilise les hôtels dans la ville de Marrakech et on l'exporte sous JSON :

The screenshot displays the Apify web interface. On the left is a sidebar with the Apify logo, a user profile for 'Anas', and navigation links: Home, Actors (selected), Tasks, Runs, Schedules, Storage, Proxy, Custom solutions, and Settings. The main area shows the 'Actors / maxcopell/tripadvisor' page. It features the 'Tripadvisor scraper' actor icon and description: 'Scrape Tripadvisor restaurants and hotels. Get reviews, pricing, contact details, amenities, awards. Download your data as HTML table, JSON, CSV, Excel, XML, and RSS feed.' Below this, it states 'Actor is rented and the trial is expiring at May 11, 2022 11:39 AM. Afterwards, \$40.00 will be charged from your credits monthly.' A tabbed interface at the top includes 'Input' (selected), Info, Runs, Builds, Tasks, Issues, and Integrations. A link 'Switch to JSON editor' is visible. The 'Input' tab contains a text area with the URL 'https://www.tripadvisor.com/Hotels-g293734-Marrakech_Marrakech_Safi-Hotels.html', a 'Start URLs (optional)' section with an '+ Add' button and a 'Text file' dropdown, and a 'Max search results (optional)' section with a value of '100' and '+'/'-' buttons.

Les résultats sont :

SUCCEEDED

100 results

3 of 3 requests handled

0 of 0 webhooks finished

Log Info Input Key-value store Dataset Request queue Live view Webhooks

Showing only the latest lines of the log.

```

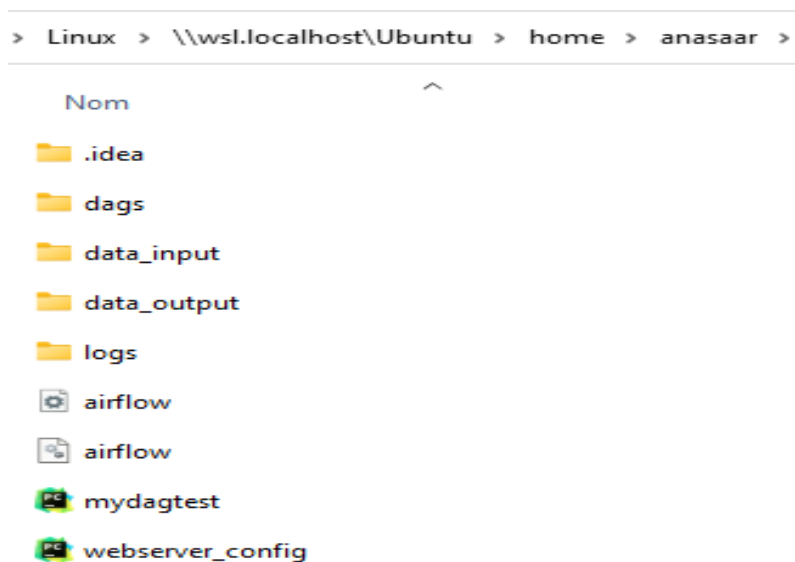
2022-05-07T16:14:46.348Z ACTOR: Pulling Docker image from repository.
2022-05-07T16:14:46.988Z ACTOR: Creating Docker container.
2022-05-07T16:14:47.176Z ACTOR: Starting Docker container.
2022-05-07T16:14:47.761Z Starting X virtual framebuffer using: Xvfb :99 -ac -screen 0 1280x720x16 -nolisten tcp
2022-05-07T16:14:47.761Z Executing main command
2022-05-07T16:14:49.056Z INFO System info {"apifyVersion":"2.3.2","apifyClientVersion":"2.3.1","osType":"Linux","nodeVersion":"v16.15.0"}
2022-05-07T16:14:49.349Z ACTOR: Actor run will metamorph
2022-05-07T16:14:49.352Z ACTOR: Sending Docker container SIGTERM signal.
2022-05-07T16:14:49.377Z ACTOR: Pulling Docker image from repository.
2022-05-07T16:14:50.268Z ACTOR: Creating Docker container.
2022-05-07T16:14:50.306Z ACTOR: Starting Docker container.
2022-05-07T16:14:53.556Z INFO System info {"apifyVersion":"2.3.2","apifyClientVersion":"2.3.1","osType":"Linux","nodeVersion":"v16.15.0"}
2022-05-07T16:14:53.675Z INFO Input validation OK
2022-05-07T16:14:58.331Z INFO Fetched locationId: 293734 for location: https://www.tripadvisor.com/Hotels-g293734-Marrakech_Marrakech_Safi-Hotels.html
2022-05-07T16:14:58.756Z INFO Starting the crawler...
2022-05-07T16:14:58.847Z INFO BasicCrawler:AutoscaledPool: state {"currentConcurrency":0,"desiredConcurrency":2,"systemStatus":{"isSystemIdle":true,"memInfo":{"isOverloaded":false,"limitRatio":0.2,"actualRatio":0.1,"freeMemory":10737408,"totalMemory":10737408}}
2022-05-07T16:15:00.048Z INFO Processing postprocessors with last data offset: 100

```

III- Préparation de l'environnement de Pycharm et Airflow sous WSL (Ubuntu) :

1- Préparation de Airflow sous WSL (Ubuntu) :

Organisation du dossier Airflow sous WSL :

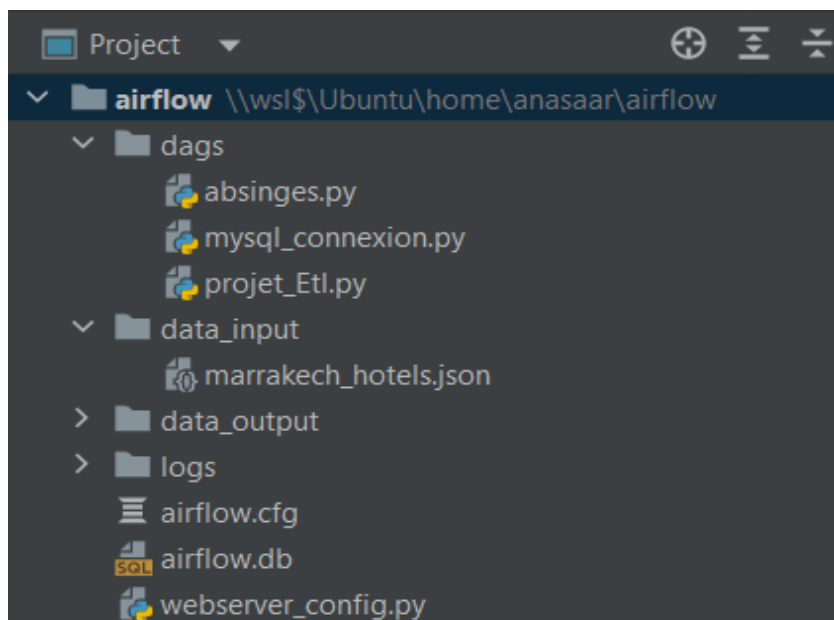


Connexion d'airflow avec Mysql server :

Edit Connection	
Connection Id *	tp2_mysql
Connection Type *	MySQL <small>Connection Type missing? Make sure you've installed the corresponding Airflow Provider Package.</small>
Description	
Host	mysql
Schema	
Login	root
Password

2- Préparation de l'environnement du Pycharm :

L'environnement Pycharm jouera un rôle fondamental dans ce projet. Il contiendra dans un dossier un code python qui traitera les différentes bases de données et le Dag utilisée pour nos objectifs. Sa structure est la suivante :



Les dossiers et fichiers peuvent se décrire comme suit :

- **Projet_Etl.py** : Noyau du projet contenant le dag utilisé.

- `Mysql_connexion.py` : Celui responsable de créer les tables de dimension ainsi que la table de fait.
- `Marrakech_hotels.json` : Le fichier contenant les données extrait de TripAdvisor.
- `Data_output` : Le dossier qui va contenir les fichiers CSV tirés du fichier JSON.

IV- Création et exécution du DAG TripAdvisor Hotels dag sous Airflow :

Airflow est une plateforme qui permet de créer, de planifier et de surveiller des workflows (flux de travail) par le biais de la programmation informatique. Un workflow peut se définir comme une suite de tâches mis en place selon un calendrier ou déclenché par un événement. La mise en œuvre d'un workflow passe par la création d'un DAG.

1- Création du DAG TripAdvisor Hotels dag :

- `Load_file` :

Permet d'importer les données depuis le dossier input :

```
def load_file():
    with open(input, 'r') as f:
        data = json.load(f)
    return data
```

- Transformation du data de fichier JSON au fichiers CSV :

```

def transform_apifydata(path: str):
    apifydata = load_file()
    transformed_data = []
    try:
        for DataRow in apifydata:
            transformed_data.append({
                'ID': DataRow['id'],
                'Name': DataRow['name'],
                'Type': DataRow['type'],
                'rating': DataRow['rating'],
                'Awards': len(DataRow['awards']),
                'RankingPosition': DataRow['rankingPosition'],
                'HotelClass': DataRow['hotelClass'],
                'NumberOfReviews': DataRow['numberOfReviews'],
                'priceRange': DataRow['priceRange']
            })
    except:
        pass
    marakech_df1 = pd.DataFrame(transformed_data)
    marakech_df1.to_csv(path, index=False)

def prices_rang_data(path: str):
    apifydata = load_file()
    transformed_data = []
    try:
        for DataRow in apifydata:
            transformed_data.append({
                'priceRange': DataRow['priceRange'],

```

```

        for DataRow in apifydata:
            transformed_data.append({
                'priceRange': DataRow['priceRange'],
                'RankingPosition': DataRow['rankingPosition']
            })
        except:
            pass
    prices_rang_df1 = pd.DataFrame(transformed_data)
    prices_rang_df1.to_csv(path, index=False)

def hotels_data(path: str):
    apifydata = load_file()
    transformed_data = []
    try:
        for DataRow in apifydata:
            transformed_data.append({
                'Name': DataRow['name'],
                'Awards': len(DataRow['awards']),
                'rating': DataRow['rating']
            })
        except:
            pass
    hotels_data_df1 = pd.DataFrame(transformed_data)
    hotels_data_df1.to_csv(path, index=False)

def hotel_reviews_data(path: str):
    apifydata = load_file()
    transformed_data = []

```

```

def hotel_reviews_data(path: str):
    apifydata = load_file()
    transformed_data = []
    try:
        for DataRow in apifydata:
            transformed_data.append({
                'NumberOfReviews': DataRow['numberOfReviews'],
                'HotelClass': DataRow['hotelClass'],
                'rating': DataRow['rating']
            })
        except:
            pass
    hotel_reviews_df1 = pd.DataFrame(transformed_data)
    hotel_reviews_df1.to_csv(path, index=False)

```

- Insertion des données dans la base donnée crée :

```

def data_dimensions_insertion():
    csvdata = pd.read_csv('/home/anasaar/airflow/data_output/marrakech1.csv', index_col=False, delimiter=',')
    csvdata.head()

    for i, row in csvdata.iterrows():
        sql = "INSERT INTO prices_rang(pricerange, rankingposition) VALUES (%s,%s)"
        ms.cursor.execute(sql, tuple(row))
        ms.conn.commit()

    csvdata1 = pd.read_csv('/home/anasaar/airflow/data_output/marrakech2.csv', index_col=False, delimiter=',')
    csvdata1.head()

    for i, row in csvdata1.iterrows():
        sql = "INSERT INTO hotels( name, awards, rating) VALUES (%s,%s,%s)"
        ms.cursor.execute(sql, tuple(row))
        ms.conn.commit()

    csvdata2 = pd.read_csv('/home/anasaar/airflow/data_output/marrakech3.csv', index_col=False, delimiter=',')
    csvdata2.head()

    for i, row in csvdata2.iterrows():
        sql = "INSERT INTO hotel_reviews(numberofreviews, hotelclass, rating) VALUES (%s,%s,%s)"
        ms.cursor.execute(sql, tuple(row))
        ms.conn.commit()

def data_fact_insertion():
    sql="insert into fact_offers(number_of_hotels, number_ofrange_prices, max_hotelclass, min_hotelclass, " \
        "max_ofreviews, min_ofreviews) select count(name), count(pricerange), max(hotelclass), min(hotelclass), " \
        "max(numberofreviews), min(numberofreviews) from hotels, prices_rang, hotel_reviews;"
    ms.cursor.execute(sql)
    ms.conn.commit()

```

- Définition du DAG :

```

with DAG(
    dag_id='TripAdvisor_Hotels_dag',
    start_date=pendulum.datetime(2022, 1, 1, tz="UTC"),
    schedule_interval='@Daily',
    catchup=False
) as dag:

```

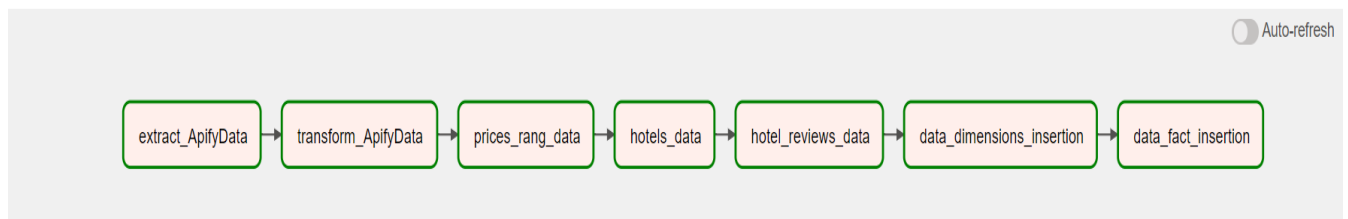
- Définition des taches utilisées :

```
task_extract_ApifyData = PythonOperator(
    task_id='extract_ApifyData',
    python_callable=load_file,
    dag=dag
)
task_transform_ApifyData = PythonOperator(
    task_id='transform_ApifyData',
    python_callable=transform_apifydata,
    op_kwargs={'path': '/home/anasaar/airflow/data_output/marrakech.csv'},
    dag=dag
)
task_prices_rang_data = PythonOperator(
    task_id='prices_rang_data',
    python_callable=prices_rang_data,
    op_kwargs={'path': '/home/anasaar/airflow/data_output/marrakech1.csv'},
    dag=dag
)
task_hotels_data = PythonOperator(
    task_id='hotels_data',
    python_callable=hotels_data,
    op_kwargs={'path': '/home/anasaar/airflow/data_output/marrakech2.csv'},
    dag=dag
)
task_hotel_reviews_data = PythonOperator(
    task_id='hotel_reviews_data',
    python_callable=hotel_reviews_data,
    op_kwargs={'path': '/home/anasaar/airflow/data_output/marrakech3.csv'},
    dag=dag
)
```

```
task_insertion_dimension_tables = PythonOperator(
    task_id='data_dimensions_insertion',
    python_callable=data_dimensions_insertion,
    dag=dag
)
task_insertion_fact_table = PythonOperator(
    task_id='data_fact_insertion',
    python_callable=data_fact_insertion,
    dag=dag
)
```

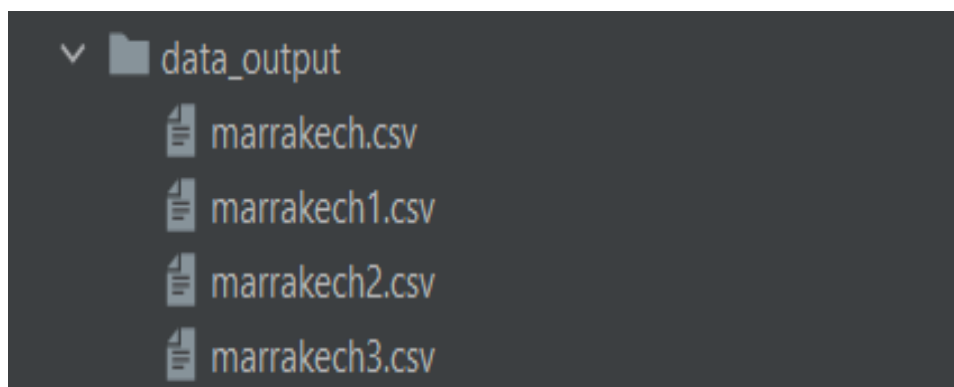
```
task_extract_ApifyData >> task_transform_ApifyData >> task_prices_rang_data >> task_hotels_data >>\
task_hotel_reviews_data >> task_insertion_dimension_tables >> task_insertion_fact_table
```

2- Exécution du DAG TripAdvisor Hotels dag :



- Visualisation des résultats du DAG :

Pour les fichiers CSV :



Pour la base données Mysql (Dans le Workbench) :

Table prices_rang :

	id	pricerange	rankingposition
▶	11	\$82 - \$155	52
	12	\$126 - \$201	51
	13	\$99 - \$124	56
	14	\$354 - \$817	59
	15	\$126 - \$327	57
✱	NULL	NULL	NULL

Table hotel_reviews :

	id	numberofreviews	hotelclass	rating
▶	101	1241	4	5
	102	582	4	5
	103	272	3,5	5
	104	864	5	5
	105	531	3	5
	106	1133	5	4,5
	107	1023	3,5	5
	108	955	5	5
	109	310	5	5
	110	331	4	5
	111	914	4	5

Table hotels :

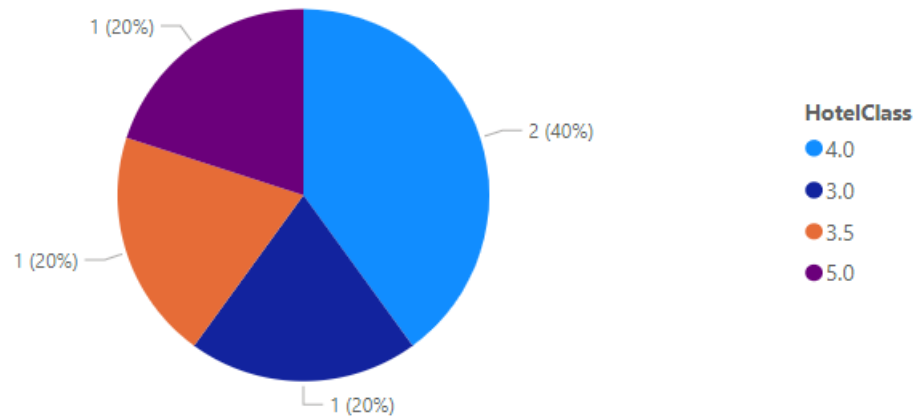
	id	name	awards	rating
▶	331	Riad Dar Najat	10	5
	332	Zamzam Riad	7	5
	333	Riad Tzarra	5	5
	334	Villa des Orangers	11	5
	335	Palais Khum	7	5
	336	Club Med Marrakech le Riad	10	4,5
	337	Riad Miski	10	5
	338	Villa Makassar	10	5
	339	Ksar Char-Bagh Small Luxury hotel	7	5
	340	Riad La Terrasse des Oliviers	8	5
	341	Riad 58 Blu	11	5

Table fact_offers :

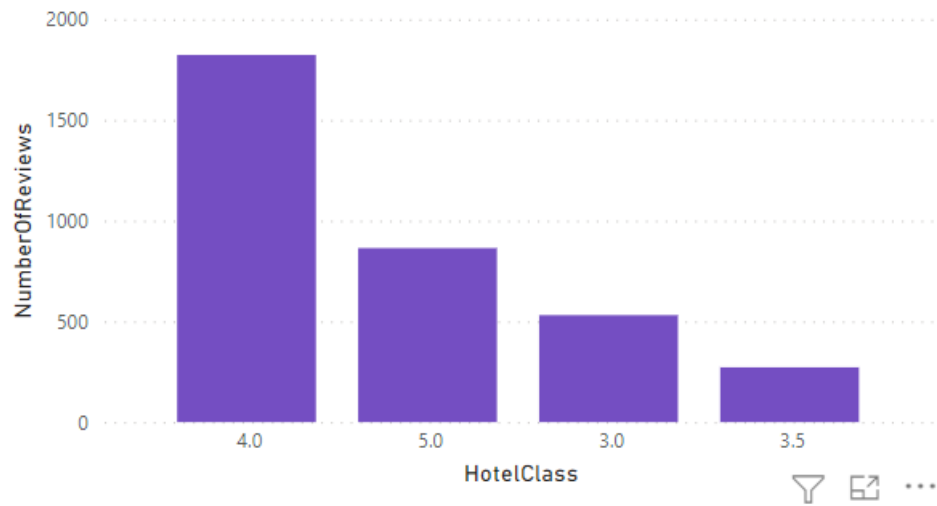
	views_id	number_of_hotels	number_ofrange_prices	max_hotelclass	min_hotelclass	max_ofreviews	min_ofreviews
▶		27500	27500	5	0	2117	141

V- Dashbording sous Power BI :

ID par HotelClass

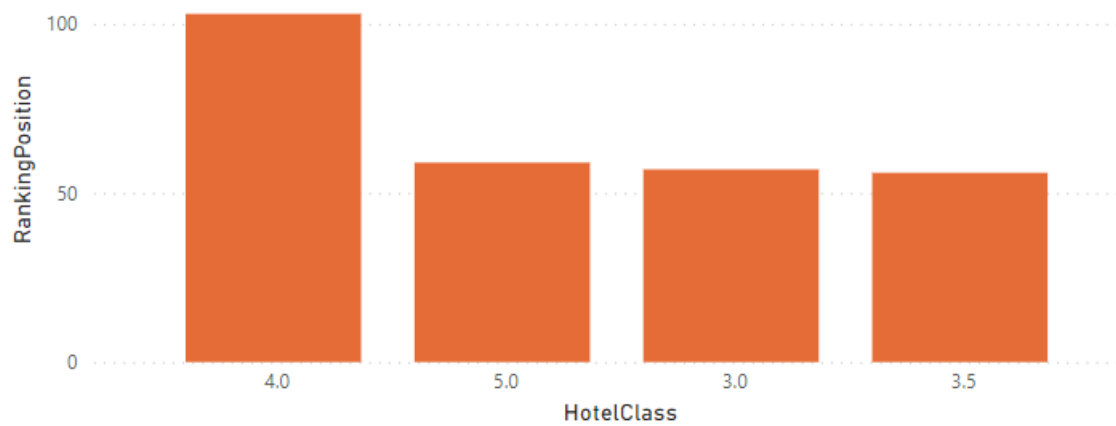


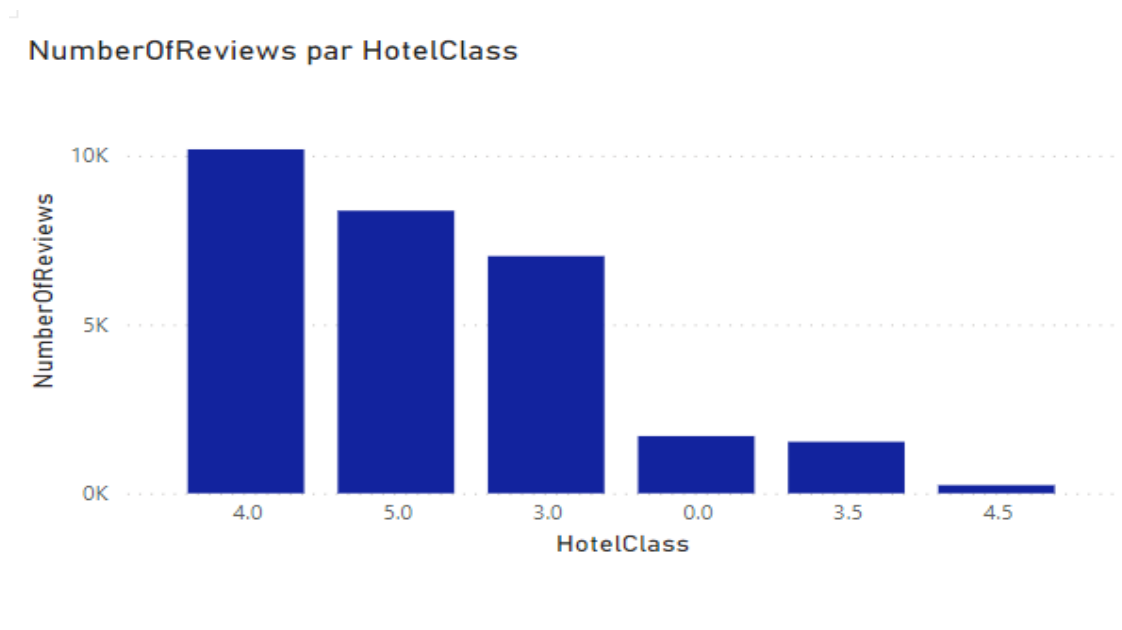
NumberOfReviews par HotelClass



NumberOfReviews par HotelClass

RankingPosition par HotelClass





Conclusion :

Ce travail nous a permis de bien saisir des notions pratiques relatives aux technologies tel que Airflow, Apify, MySQL, MySQL server et Power BI. Aussi ce travail est conçu pour permettre aux gens dans les postes de décision dans les hôtels à Marrakech et dans le secteur du tourisme d'essayer d'améliorer sa qualité et arriver à mieux gérer les hôtels et l'industrie d'hôtellerie au Maroc.