Grandes lignes du projet

Prérequis: - Avoir un compte AWS
           - petites notions sur linux, nous travaillerons sur Ubuntu 24.04
           - connaissance de base docker, docker-compose


lien, si vous n'avez pas de compte AWS et vous voulez en créé : https://smart.dhgate.com/step-by-step-guide-to-creating-and-securing-your-aws-account-for-optimal-use/


un pipeline Airflow orchestrant des extractions de données depuis un bucket S3 de Amazon (weather(météo), COVID) et ensuite des transformations utilisant duckdb via dbt.

- Installation de airflow via le script docker compose que vous pouvez trouver sur le site officiel de airflow: https://airflow.apache.org/docs/apache-airflow/stable/howto/docker-compose/index.html#fetching-docker-compose-yaml,



 - cliquez sur le lien en bleu dans l'image ou directement ici : docker-compose.yml
pour telecharger le script docker compose.

- Créé un dossier airflow et mettre le script docker-compose.yml téléchargé, dans ce dossier

- créé cette structure de base à l'interieur du dossier airflow avec cette commande : mkdir -p dags plugins scripts config data/{raw,processed} tests


ensuite :

- on se deplace dans le dossier où se trouve le
script et on exécute la commande : docker-compose up -d
il installera tout les conteneurs liés à docker et va les lancer.

```
d-a-s@d-a-s-ThinkPad-X380-Yoga:~$ cd Desktop/
d-a-s@d-a-s-ThinkPad-X380-Yoga:~/Desktop$ cd airflow/
d-a-s@d-a-s-ThinkPad-X380-Yoga:~/Desktop/airflow$ docker-compose up
```



```
                d-a-s@d-a-s-ThinkPad-X380-Yoga: ~          ×          d-a-s@d-a-s-ThinkPad-X380-Yoga: ~/Desktop/airflow          ×     ⌄
example_passing_params_via_test_command to 2025-09-24 21:45:00+00:00, run_after=2025-09-24 21:45:00+00:
00
airflow-dag-processor_1  | [2025-09-24T21:45:43.387+0000] {dag.py:1622} INFO - Sync 1 DAGs
airflow-dag-processor_1  | [2025-09-24T21:45:43.395+0000] {dag.py:2236} INFO - Setting next_dagrun for
example_simplest_dag to None, run_after=None
airflow-dag-processor_1  | [2025-09-24T21:45:43.461+0000] {dag.py:1622} INFO - Sync 3 DAGs
airflow-dag-processor_1  | [2025-09-24T21:45:43.474+0000] {dag.py:2236} INFO - Setting next_dagrun for
example_branch_datetime_operator to 2025-09-24 00:00:00+00:00, run_after=2025-09-24 00:00:00+00:00
airflow-dag-processor_1  | [2025-09-24T21:45:43.479+0000] {dag.py:2236} INFO - Setting next_dagrun for
example_branch_datetime_operator_2 to 2025-09-24 00:00:00+00:00, run_after=2025-09-24 00:00:00+00:00
airflow-dag-processor_1  | [2025-09-24T21:45:43.483+0000] {dag.py:2236} INFO - Setting next_dagrun for
example_branch_datetime_operator_3 to 2025-09-24 00:00:00+00:00, run_after=2025-09-24 00:00:00+00:00
airflow-dag-processor_1  | [2025-09-24T21:45:43.536+0000] {dag.py:1622} INFO - Sync 2 DAGs
airflow-dag-processor_1  | [2025-09-24T21:45:43.542+0000] {dag.py:2236} INFO - Setting next_dagrun for
example_dynamic_task_mapping to None, run_after=None
airflow-dag-processor_1  | [2025-09-24T21:45:43.544+0000] {dag.py:2236} INFO - Setting next_dagrun for
example_task_mapping_second_order to None, run_after=None
airflow-dag-processor_1  | [2025-09-24T21:45:43.595+0000] {dag.py:1622} INFO - Sync 1 DAGs
airflow-dag-processor_1  | [2025-09-24T21:45:43.615+0000] {dag.py:2236} INFO - Setting next_dagrun for
example_branch_dop_operator_v3 to 2025-09-24 21:45:00+00:00, run_after=2025-09-24 21:45:00+00:00
airflow-dag-processor_1  | [2025-09-24T21:45:43.672+0000] {dag.py:1622} INFO - Sync 1 DAGs
airflow-dag-processor_1  | [2025-09-24T21:45:43.682+0000] {dag.py:2236} INFO - Setting next_dagrun for
tutorial to 2025-09-24 21:45:43.681904+00:00, run_after=2025-09-24 21:45:43.681904+00:00
```



```
        d-a-s@d-a-s-ThinkPad-X380-Yoga: ~     ×     d-a-s@d-a-s-ThinkPad-X380-Yoga: ~/Des... ×     d-a-s@d-a-s-ThinkPad-X380-Yoga: ~/Des... ×     ⌄
d-a-s@d-a-s-ThinkPad-X380-Yoga:~/Desktop/airflow$ docker ps
CONTAINER ID   IMAGE                         COMMAND                  CREATED         STATUS
            PORTS                                          NAMES
47e35a0c0722   airflow_airflow-worker        "/usr/bin/dumb-init …"   24 hours ago    Up About a minut
e (healthy)    8080/tcp                                       airflow_airflow-worker_1
6c04087dd2e0   airflow_airflow-triggerer     "/usr/bin/dumb-init …"   24 hours ago    Up About a minut
e (healthy)    8080/tcp                                       airflow_airflow-triggerer_1
d79c5bdb9c5d   airflow_airflow-dag-processor "/usr/bin/dumb-init …"   24 hours ago    Up About a minut
e (healthy)    8080/tcp                                       airflow_airflow-dag-processor_1
409bfccb5407   airflow_airflow-apiserver     "/usr/bin/dumb-init …"   24 hours ago    Up About a minut
e (healthy)    0.0.0.0:8080->8080/tcp, [::]:8080->8080/tcp    airflow_airflow-apiserver_1
88bc4d21f31e   airflow_airflow-scheduler     "/usr/bin/dumb-init …"   24 hours ago    Up About a minut
e (healthy)    8080/tcp                                       airflow_airflow-scheduler_1
cde6b682c55d   postgres:13                   "docker-entrypoint.s…"   24 hours ago    Up About a minut
e (healthy)    5432/tcp                                       airflow_postgres_1
e3754f863416   redis:7.2-bookworm            "docker-entrypoint.s…"   24 hours ago    Up About a minut
e (healthy)    6379/tcp                                       airflow_redis_1
d-a-s@d-a-s-ThinkPad-X380-Yoga:~/Desktop/airflow$
```
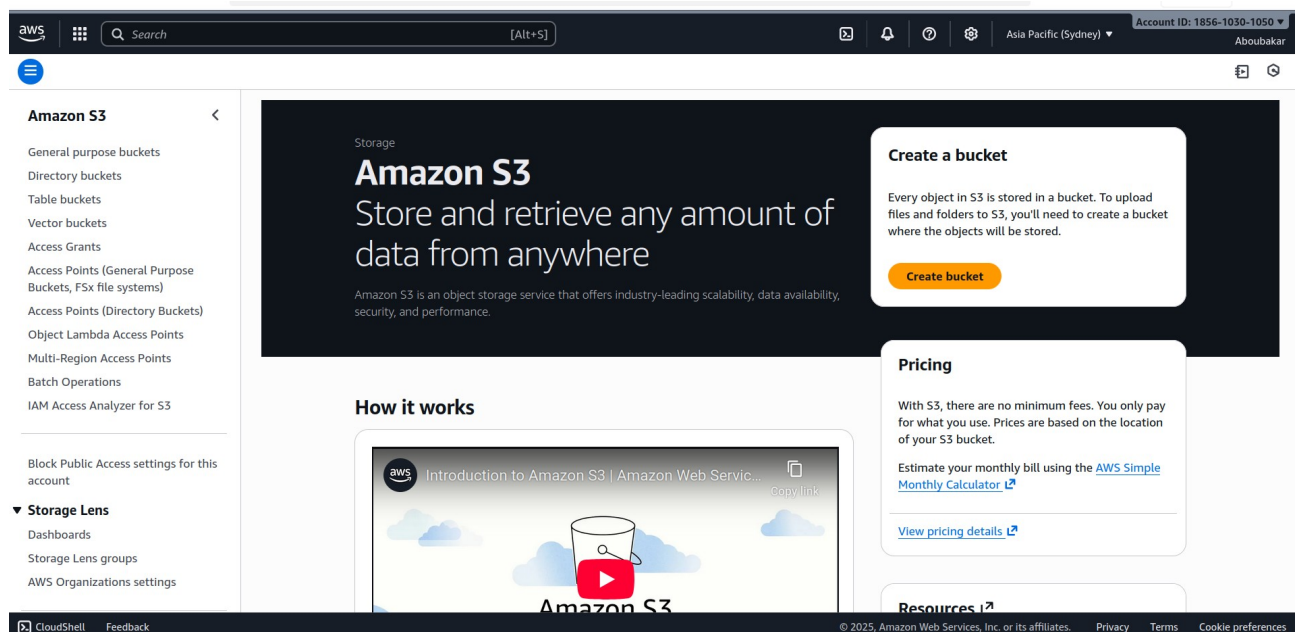
- Installer dbt-core + l'adaptateur DuckDB car on ne va pas utiliser Postgres mais DuckDB vu que

S3 est un outil de stockage objet qui permet de faire des laques de données qui n' ont pas une structure relationnelle pour pouvoir effectuer des requetes sur S3, on ne peut pas utiliser Postgres car gérant les db relationnelles, on va utiliser DuckDB qui peut requeté sur des fichiers Json, CSV et à des connecteurs pour S3 qu'on va utiliser.
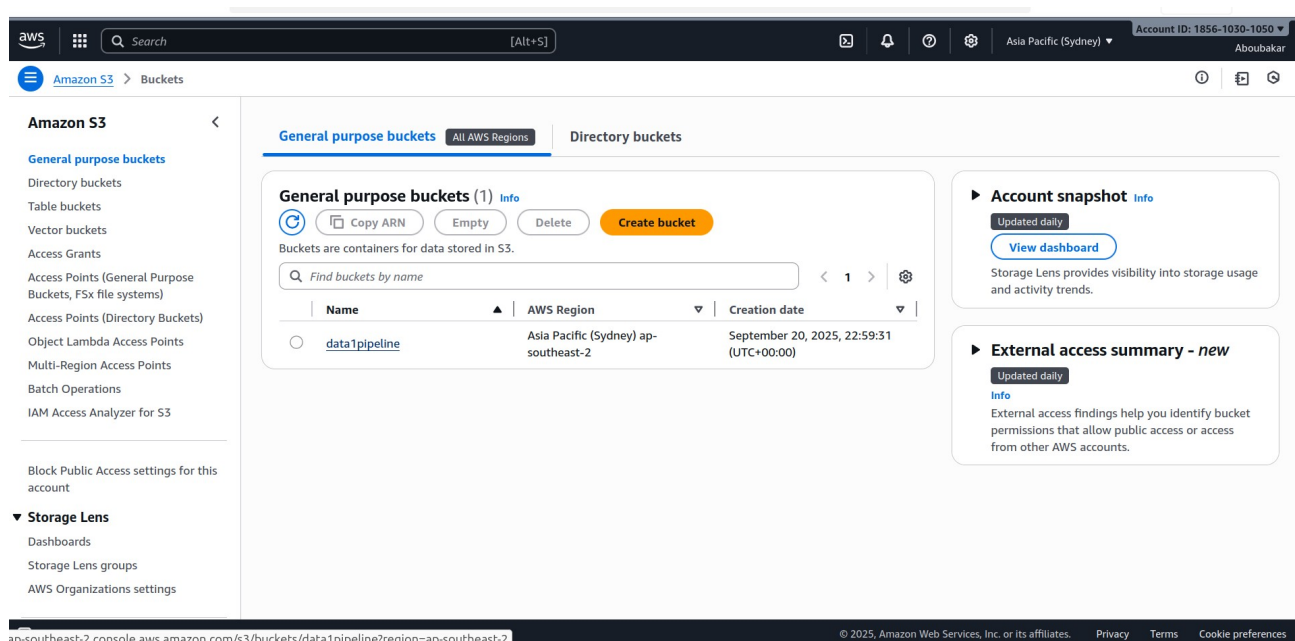On le fera plus tard vers la fin dans un conteneur airflow pour qu'il soit visible à l'interieur d' airflow.

- Créé un bucket S3, on l'a nommé dans notre cas data1pipeline vous pouvez le changé

- Allez sur votre console AWS, connectez vous, tapez dans la barre de recherche S3, vous allez être ici



-Cliquez sur general purpose  bucket

- Cliquez sur create bucket et entrez le nom laissez les parametre ar defaut et cliquez tout en bas sur create bucket

- Créé S3 Access Key et Secret Key
1- Ces identifiants permettent à dbt (ou tout autre outil) d'accéder à votre bucket S3.
2- Ils sont générés depuis la console AWS IAM (Identity and Access Management) :
3- Créer un utilisateur IAM avec accès programmatique (Programmatic access).
4- Attacher une politique (Policy) permettant l'accès à S3 (ex: AmazonS3FullAccess pour test ou une policy restreinte pour la production).
6- Cliquez sur votre utilisateur et cherchez : create access key et generé les
7- La console génère alors une Access Key ID et une Secret Access Key.
Ces clés peuvent être utilisées dans dbt, Airflow, ou tout script pour se connecter à S3.

- Placer les fichiers data_extraction_dag.py et dbt_pipeline_dag.py dans le dossier airflow créé par docker compose

- Créé dans le dossier airflow, le dossier dbt et à l interieur de ce dossier, le dossier "votre projet"

# Aller dans ton dossier airflow (si tu es sur ton Desktop par exemple)
cd ~/Desktop/airflow

# Créer le dossier dbt à l'intérieur du dossier airflow
mkdir -p dbt

# Aller dans ce dossier dbt
cd dbt

# Créer le dossier pour ton projet dbt (exemple : my_dbt_project)
mkdir -p my_dbt_project

# Initialiser un projet dbt à l'intérieur
cd my_dbt_project
dbt init my_dbt_project

Volumes montés dans le conteneur Airflow via Docker Compose :
- /home/d-a-s/Desktop/airflow/dbt:/opt/airflow/dbt: on le met dans la partie volume de notre fichier docker compose

→ Le dossier local ~/airflow/dbt est monté dans le conteneur à /opt/airflow/dbt.

→ Sert à stocker les projets dbt (models, macros, analyses, seeds, etc.) de manière persistante.

→ Toute modification dans le conteneur sera visible sur la machine hôte et vice-versa.

- /home/d-a-s/.dbt:/opt/airflow/.dbt : on le met dans la partie volume de notre fichier docker compose

→ Le dossier local ~/.dbt, il contient le fichier profile, vous pouvez le mettre ailleurs mais il faut specifier le dossier et monté dans le conteneur à /opt/airflow/.dbt.

→ Contient les fichiers de configuration dbt, notamment profiles.yml qui définit les connexions aux bases de données.

→ Permet au conteneur d'utiliser les mêmes configurations que celles définies sur la machine hôte.

Configuration dbt pour le projet my_dbt_project (profiles.yml) :

- target: dev

→ Définit l'environnement de travail par défaut pour dbt.

- outputs:

dev:

type: duckdb

path: /opt/airflow/dbt/my_dbt_project/my_duckdb.duckdb

→ Fichier DuckDB stocké dans le dossier du projet dbt monté dans le conteneur Airflow.

→ Persistant : toutes les modifications restent sur la machine hôte via le volume Docker.

extensions:

- httpfs

→ Permet à DuckDB d'accéder à des fichiers stockés sur S3 directement depuis dbt.

threads: 1

→ Nombre de threads utilisés pour l'exécution des modèles dbt.

schema: main

→ Schéma par défaut utilisé dans DuckDB.

s3_region: ap-southeast-2

s3_access_key_id: AKIASWNZ4NZ5JMMRMLOV

s3_secret_access_key: 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

→ Informations de connexion pour accéder au bucket S3 depuis dbt/DuckDB.

Fichier sources.yml (dans models/staging)

Rôle :

• Définit les sources de données externes que dbt va utiliser dans les modèles.

• Permet à dbt de référencer les fichiers S3 (ou d'autres sources) de manière centralisée et documentée.

• Facilite la traçabilité et la gestion des données avant toute transformation (staging).

Contenu et explication :

version: 2

sources:
- name: raw

```
# Nom logique de la source de données
tables:
- name: weather_data # Table ou fichier météo
description: "Fichiers météo depuis S3"
external:
# Données externes (non transformées)
location: "s3://data1pipeline/raw/weather/" # Chemin S3
format: json
# Format des fichiers
pattern: "*.json" # Pattern pour filtrer les fichiers
- name: covid_data # Table ou fichier COVID
description: "Fichiers COVID depuis S3"
external:
location: "s3://data1pipeline/raw/covid/"
format: csv
pattern: "*.csv"
```

Résumé :

• Ce fichier n'est pas strictement obligatoire pour faire fonctionner dbt, mais il est fortement recommandé.

• Il sert à documenter les sources et permettre aux modèles dbt de les référencer facilement via source('raw', 'weather_data') ou source('raw', 'covid_data').

• Cela centralise la configuration S3, donc si le chemin change, il suffit de modifier ici.

- Étapes pour ajouter tes modèles

1. Aller dans ton projet dbt :
cd ~/Desktop/airflow/dbt/my_dbt_project

2. Créer les dossiers nécessaires :

mkdir -p models/staging

mkdir -p models/marts

mettez y les fichiers
-- models/staging/stg_weather.sql
-- models/staging/stg_covid.sql
-- models/mart/daily_weather_metrics.sqlCréation et récupération de la clé API OpenWeatherMap

1. Créer un compte :

Aller sur https://home.openweathermap.org/users/sign_up

2. Connexion :

Se connecter à https://home.openweathermap.org

3. Générer une clé API :

• Aller dans le menu API keys.

• Cliquer sur "Generate" ou utiliser la clé par défaut fournie.

• Donner un nom à la clé (exemple : airflow_project_key).
Vous utiliserez cette clé dans le data_extraction_dag.py pour recupérer les données de meteo(weather)

4.
1️⃣ Accéder au conteneur Airflow
Identifie d'abord le nom du conteneur Airflow en cours d'exécution :
avec docker ps
Ensuite, connecte-toi au conteneur (par ex. le scheduler ou le webserver)
avec la commande : docker exec -it airflow-scheduler-1 bash

5.
2️⃣ Installer dbt et dbt avec l adaptateur duckdb dans le
conteneur avec : pip install --no-cache-dir dbt-core dbt-
duckdb

6. Vérifier l'installation avec
dbt –version
vous dois voir quelque chose comme :
installed version: 1.8.x
plugins:
- duckdb: 1.8.x

7. ouvrir la page web de airflow http://localhost:8080/
ET executé les dag : data_extraction_dag et dbt_pipeline_dag


**Resultats**

**On verifie qu'il y a des données dans le bucket**



**Le pipeline d'extraction des données des API et le stockage vers S3**

**On lance le pipeline d'extraction de données des API de Covid et Météo et on envoie dans le bucket S3 Selon le script data_extraction_dag.py_py qui dans dags du dossier airflow**



**Le pipeline éxécuté avec succès**

**Sous bucket créés**



**données dans le sous bucket weather**

**données dans le sous bucket covid**



**Execution du pipeline de transformation dbt des données stocker dans le sous bucket covid et weather sur S3**

**Tout les modèles sont éxécutés**



**On se connecte au conteneur airflow scheduler où on a installé dbt et on éxécute la commande pour voir que la connexion de DuckDB à S3 est fait**

Sep 24 22:04
d-a-s@d-a-s-ThinkPad-X380-Yoga: ~/Desktop/airflow
d-a-s@d-a-s-ThinkPad-X380-Yoga: ~/Desktop/airflow          d-a-s@d-a-s-ThinkPad-X380-Yoga: ~/Desktop/airflow

```
22:03:27  Running with dbt=1.10.11
22:03:28  dbt version: 1.10.11
22:03:28  python version: 3.12.11
22:03:28  python path: /home/airflow/.local/bin/python
22:03:28  os info: Linux-6.14.0-24-generic-x86_64-with-glibc2.36
22:03:28  Using profiles dir at /opt/airflow/.dbt
22:03:28  Using profiles.yml file at /opt/airflow/.dbt/profiles.yml
22:03:28  Using dbt_project.yml file at /opt/airflow/dbt/my_dbt_project/dbt_project.yml
22:03:28  adapter type: duckdb
22:03:28  adapter version: 1.9.6
22:03:28  Configuration:
22:03:28    profiles.yml file [OK found and valid]
22:03:28    dbt_project.yml file [OK found and valid]
22:03:28  Required dependencies:
22:03:28   - git [OK found]

22:03:28  Connection:
22:03:28    database: my_duckdb
22:03:28    schema: main
22:03:28    path: /opt/airflow/dbt/my_dbt_project/my_duckdb.duckdb
22:03:28    config_options: None
22:03:28    extensions: ['httpfs']
22:03:28    settings: {}
22:03:28    external_root: .
22:03:28    use_credential_provider: None
22:03:28    attach: None
22:03:28    filesystems: None
22:03:28    remote: None
22:03:28    plugins: None
22:03:28    disable_transactions: False
22:03:28  Registered adapter: duckdb=1.9.6
22:03:54    Connection test: [OK connection ok]

22:03:54  All checks passed!
root@88bc4d21f31e:/opt/airflow/dbt/my_dbt_project#
```

**Connexion réussi**

```
root@88bc4d21f31e:/opt/airflow/dbt/my_dbt_project# curl -L https://github.com/duckdb/duckdb/releases/download/v1.1.3/duckdb_cli-linux-amd64.zip -o duckdb.zip
  % Total    % Received % Xferd  Average Speed   Time    Time     Time  Current
                                 Dload  Upload   Total   Spent    Left  Speed
  0     0    0     0    0     0      0      0 --:--:-- --:--:-- --:--:--     0
100 15.9M  100 15.9M    0     0  1466k      0  0:00:11  0:00:11 --:--:-- 2206k
root@88bc4d21f31e:/opt/airflow/dbt/my_dbt_project# apt-get install -y unzip
Reading package lists... Done
Building dependency tree... Done
Reading state information... Done
unzip is already the newest version (6.0-28).
0 upgraded, 0 newly installed, 0 to remove and 66 not upgraded.
root@88bc4d21f31e:/opt/airflow/dbt/my_dbt_project# unzip duckdb.zip -d /usr/local/bin
chmod +x /usr/local/bin/duckdb
Archive:  duckdb.zip
  inflating: /usr/local/bin/duckdb
root@88bc4d21f31e:/opt/airflow/dbt/my_dbt_project# duckdb --version
v1.1.3 19864453f7
root@88bc4d21f31e:/opt/airflow/dbt/my_dbt_project# duckdb my_duckdb.duckdb
Error: unable to open database "my_duckdb.duckdb": IO Error: Could not set lock on file "my_duckdb.duckdb": Conflicting lock is held in /usr/local/bin/python3
12 (PID 963). See also https://duckdb.org/docs/connect/concurrency
root@88bc4d21f31e:/opt/airflow/dbt/my_dbt_project# ps -p 963 -f
```

**Installation de DuckDB dans le conteneur pour verifier qu'il y a effectiovement les tables liés à l'éxécution de dbt run**

```
bash: ps: command not found
root@88bc4d21f31e:/opt/airflow/dbt/my_dbt_project# kill -9 963
root@88bc4d21f31e:/opt/airflow/dbt/my_dbt_project# duckdb my_duckdb.duckdb -readonly
v1.1.3 19864453f7
Enter ".help" for usage hints.
D .show tables
Usage: .show
D SHOW TABLES;
┌─────────────────┐
│      name       │
│     varchar     │
├─────────────────┤
│ my_first_dbt_model │
│ my_second_dbt_model │
│ stg_covid       │
└─────────────────┘

D -- Voir tous les schémas
D .show schemas;
Usage: .show
D
D -- Voir toutes les tables et vues dans tous les schémas
D SELECT table_schema, table_name, table_type
  FROM information_schema.tables
  WHERE table_schema IN ('main', 'main_staging', 'main_mart');
┌──────────────┬──────────────────────┬────────────┐
│ table_schema │      table_name      │ table_type │
│   varchar    │       varchar        │  varchar   │
├──────────────┼──────────────────────┼────────────┤
│ main         │ my_first_dbt_model   │ BASE TABLE │
│ main         │ stg_covid            │ BASE TABLE │
│ main_mart    │ daily_weather_metrics│ BASE TABLE │
│ main         │ my_second_dbt_model  │ VIEW       │
│ main_staging │ stg_weather          │ VIEW       │
└──────────────┴──────────────────────┴────────────┘
```

**Verification qu'on a toutes les tables**

```
D SELECT * FROM main.stg_covid LIMIT 10;
```

| date<br>date | country<br>varchar | province<br>varchar | confirmed<br>int64 | deaths<br>int64 | recovered<br>double |
|---|---|---|---|---|---|
| 2020-01-22 | Afghanistan | | 0 | 0 | 0.0 |
| 2020-01-23 | Afghanistan | | 0 | 0 | 0.0 |
| 2020-01-24 | Afghanistan | | 0 | 0 | 0.0 |
| 2020-01-25 | Afghanistan | | 0 | 0 | 0.0 |
| 2020-01-26 | Afghanistan | | 0 | 0 | 0.0 |
| 2020-01-27 | Afghanistan | | 0 | 0 | 0.0 |
| 2020-01-28 | Afghanistan | | 0 | 0 | 0.0 |
| 2020-01-29 | Afghanistan | | 0 | 0 | 0.0 |
| 2020-01-30 | Afghanistan | | 0 | 0 | 0.0 |
| 2020-01-31 | Afghanistan | | 0 | 0 | 0.0 |

| 10 rows | | | | | 6 columns |

```
D SELECT * FROM main_mart.daily_weather_metrics LIMIT 10;
```

| city<br>json | date<br>date | avg_temperature<br>double | max_temperature<br>double | min_temperature<br>double | avg_humidity<br>double | observations_count<br>int64 |
|---|---|---|---|---|---|---|
| "Paris" | 2025-09-24 | 10.55 | 10.55 | 10.55 | 93.0 | 1 |

**données dans stagging covid et metrics de météo**

```
D SELECT * FROM main_staging.stg_weather LIMIT 2;
100%
```

| city<br>json | observation_time<br>timestamp with time zone | temperature<br>double | humidity<br>double | weather_description<br>json | wind_speed<br>double | loaded_at<br>timestamp with time zone |
|---|---|---|---|---|---|---|
| "Paris" | 2025-09-24 21:50:41+00 | 10.55 | 93.0 | "overcast clouds" | 4.63 | 2025-09-24 22:56:08.078+00 |
| "Paris" | 2025-09-24 22:36:56+00 | 10.55 | 93.0 | "overcast clouds" | 3.6 | 2025-09-24 22:56:08.078+00 |

**données dans stg_weather**