

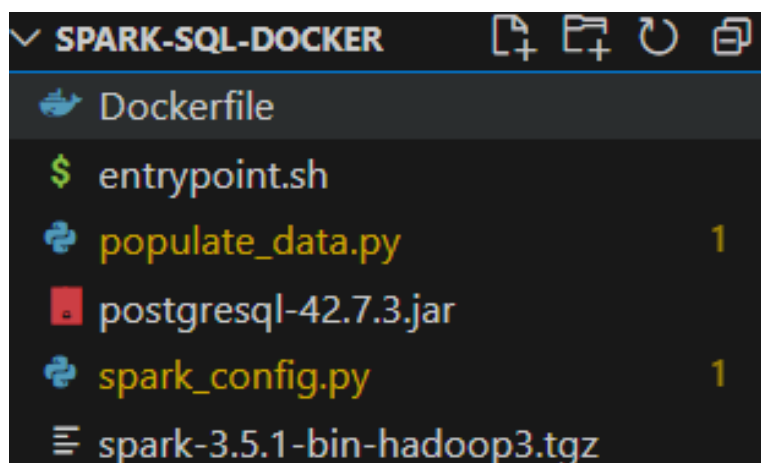
Hướng dẫn Thiết lập và Sử dụng Docker Container cho SparkSQL và PostgreSQL

Mục tiêu

Mục tiêu của dự án này là thiết lập và chạy SparkSQL trong một Docker container. Docker container sẽ chứa tất cả các phụ thuộc cần thiết để chạy Spark và một cơ sở dữ liệu PostgreSQL để lưu trữ dữ liệu. Người dùng sẽ có thể thực hiện các truy vấn SQL trên cơ sở dữ liệu này thông qua SparkSQL.

Yêu cầu

1. Docker: Đảm bảo Docker đã được cài đặt trên hệ thống. Có thể tải xuống và cài đặt Docker từ trang web chính thức của Docker.
2. Apache Spark: Truy cập trang web chính thức của Apache Spark và tải xuống phiên bản `spark-3.5.1-bin-hadoop3.tgz`. Đảm bảo rằng tệp này được lưu trong thư mục dự án
3. JDBC PostgreSQL: Có thể tải PostgreSQL từ trang web chính thức. Đảm bảo bạn đã tải xuống file JAR PostgreSQL JDBC driver và đặt nó trong cùng thư mục với Dockerfile
4. Python và PySpark: Đã cài đặt
5. Tệp và thư mục: Tạo một thư mục mới cho dự án và di chuyển tệp tin Spark đã tải xuống vào trong thư mục đó. Cấu trúc thư mục nên trông giống như sau



Thực hành

Bước 1: Tạo Dockerfile và Script Entry Point

1. Tạo Dockerfile

Tạo một tệp có tên là `Dockerfile` trong thư mục dự án.

Sử dụng `Dockerfile` đã cung cấp để xây dựng Docker image.

```
1  # Sử dụng image OpenJDK chính thức làm cơ sở
2  FROM openjdk:11-jre-slim
3
4  # Đặt biến môi trường cho Spark
5  ENV SPARK_HOME=/opt/spark
6  ENV PATH=$PATH:$SPARK_HOME/bin:$SPARK_HOME/sbin
7
8  # Cài đặt các gói phụ thuộc
9  RUN apt-get update && \
10     apt-get install -y curl python3 python3-pip postgresql postgresql-cont
11     apt-get clean
12
13  # Sao chép và cài đặt Spark
14  COPY spark-3.5.1-bin-hadoop3.tgz /opt/
15  RUN tar -xzf /opt/spark-3.5.1-bin-hadoop3.tgz -C /opt/ && \
16     rm /opt/spark-3.5.1-bin-hadoop3.tgz && \
17     mv /opt/spark-3.5.1-bin-hadoop3 /opt/spark
18
19  # Sao chép script entrypoint
20  COPY entrypoint.sh /usr/local/bin/
21
22  # Đặt quyền thực thi cho script entrypoint
23  RUN chmod +x /usr/local/bin/entrypoint.sh
24
25  # Đặt entrypoint
26  ENTRYPOINT ["/usr/local/bin/entrypoint.sh"]
27
28  # Mở các cổng cần thiết
29  EXPOSE 4040 8080 7077 5432
```

2. Tạo Script Entry Point

Tạo một tệp có tên là `entrypoint.sh` trong cùng thư mục với `Dockerfile`.

```
1  #!/bin/bash
2
3  # Khởi động dịch vụ PostgreSQL
4  service postgresql start
5
6  # Tạo người dùng và cơ sở dữ liệu PostgreSQL
7  su - postgres -c "psql -c \"CREATE USER sparkuser WITH PASSWORD 'sparkpassword';\""
8  su - postgres -c "psql -c \"CREATE DATABASE sparkdb;\""
9  su - postgres -c "psql -c \"GRANT ALL PRIVILEGES ON DATABASE sparkdb TO sparkuser;\""
10
11 # Khởi động Spark
12 $SPARK_HOME/sbin/start-master.sh
13 $SPARK_HOME/sbin/start-worker.sh spark://$(hostname):7077
14
15 # Giữ container chạy
16 tail -f /dev/null
```

Bước 2: Xây dựng và Chạy Docker Container

1. Xây dựng Docker Image

Mở terminal và điều hướng đến thư mục chứa `Dockerfile` và các tệp khác. Sau đó, chạy lệnh sau để xây dựng Docker image:

```
docker build -t spark-sql-docker .
```

Lệnh này sẽ tạo một Docker image có tên `spark-sql-docker` sử dụng các chỉ thị trong `Dockerfile`.

2. Chạy Docker Container

Sau khi xây dựng xong, chạy lệnh sau để khởi động Docker container:

```
docker run -it -p 4040:4040 -p 8080:8080 -p 7077:7077 -p 5432:5432 spark-sql-docker
```

- `-it`: Chạy container trong chế độ tương tác với terminal .
- `-p 4040:4040 -p 8080:8080 -p 7077:7077 -p 5432:5432`: Ánh xạ các cổng của container với các cổng trên máy chủ để có thể truy cập các dịch vụ bên trong container từ bên ngoài.

Sau khi chạy 2 lệnh trên, ta sẽ được đưa vào bên trong docker có dạng như sau

```
C:\Users\lamtu>docker ps
CONTAINER ID   IMAGE          NAMES      COMMAND                  CREATED        STATUS        PORTS
96bc1c7a19f9   spark-sql-docker   "/usr/local/bin/entr..."  2 hours ago   Up 2 hours   0.0.0.0:4040->4040/tcp, 0.0.0.0:5432->5432/tcp, 0.0.0.0:7077->7077/tcp, 0.0.0.0:8080->8080/tcp
```

Bước 3: Thêm Dữ liệu và Thực hiện Truy vấn

1. Thêm Dữ liệu vào Cơ sở Dữ liệu

Tạo một tệp có tên là `populate_data.py` để nhập dữ liệu mẫu vào cơ sở dữ liệu PostgreSQL:

```
1 import psycopg2
2 import random
3
4 conn = psycopg2.connect(
5     dbname="sparkdb",
6     user="sparkuser",
7     password="sparkpassword",
8     host="localhost"
9 )
10 cur = conn.cursor()
11
12 cur.execute("""
13     CREATE TABLE sample_data (
14         id SERIAL PRIMARY KEY,
15         value INTEGER
16     );
17 """)
18
19 for _ in range(1000000):
20     cur.execute("INSERT INTO sample_data (value) VALUES (%s)", (random.randint(1, 100),))
21
22 conn.commit()
23 cur.close()
24 conn.close()
```

Chạy các lệnh sau trong một terminal mới để sao chép script Python vào container và thực thi nó để nhập dữ liệu mẫu vào cơ sở dữ liệu:

```
docker cp populate_data.py <container_id>:/populate_data.py
```

```
docker exec -it <container_id> python3 /populate_data.py
```

2. Thực hiện Truy vấn SparkSQL

Tạo một tệp có tên là `spark_config.py` để cấu hình Spark và thực hiện truy vấn SparkSQL:

```
1  from pyspark.sql import SparkSession
2
3  # Cấu hình Spark session
4  spark = SparkSession.builder \
5      .appName("Spark PostgreSQL") \
6      .config("spark.jars", "C:\Users\lamtu\spark-sql-docker\postgresql-42.7.3.jar") \
7      .getOrCreate()
8
9  # Đọc dữ liệu từ PostgreSQL
10 df = spark.read \
11     .format("jdbc") \
12     .option("url", "jdbc:postgresql://localhost:5432/sparkdb") \
13     .option("dbtable", "sample_table") \
14     .option("user", "sparkuser") \
15     .option("password", "sparkpassword") \
16     .option("driver", "org.postgresql.Driver") \
17     .load()
18
19 df.show()
```

Chạy các lệnh sau trong một terminal mới để sao chép file `spark_config.py` vào container và cấu hình Spark để kết nối với cơ sở dữ liệu PostgreSQL và thực hiện truy vấn:

```
docker cp spark_config.py 96bc1c7a19f9:/spark_config.py
```

```
docker cp C:\Users\lamtu\spark-sql-docker\postgresql-42.7.3.jar
96bc1c7a19f9: C:\Users\lamtu\spark-sql-docker\postgresql-42.7.3.jar
```

```
docker exec -it 96bc1c7a19f9 python3 /spark_config.py
```

Khi thấy thông báo này, database đã được hoàn thành, ta có thể kiểm tra bằng cách đếm số lượng dòng trong database:

```
root@b32d1cef2829:/# sqlite3 /opt/spark/examples/src/main/resources/people.db
SQLite version 3.34.1 2021-01-20 14:10:07
Enter ".help" for usage hints.
sqlite> SELECT COUNT (*) FROM people;
1000000
```

Chúng ta có thể thấy số lượng dòng trong database.

Sau đây là kết quả khi chạy các file ở trên:

```
INFO:__main__:Database connection successful.  
INFO:__main__:Temporary view created.
```

Thông báo về việc đã kết nối database và tạo temporary view.

```
INFO:__main__:Performance test with large dataset...  
INFO:__main__:Count all records: 1000000  
INFO:__main__:Time to count all records: 0.042037010192871094 seconds  
INFO:__main__:Count records with age > 50: 500130  
INFO:__main__:Time to count records with age > 50: 0.05854320526123047 seconds  
INFO:py4j.clientserver:Closing down clientserver connection
```

Kiểm tra thời gian count toàn bộ record và count có điều kiện WHERE và đóng kết nối server với database.