

AWS Pyspark Development

Below is a detailed explanation of the project structure, full code implementation, test cases, and instructions on how to run or test the application.

Project Structure

```
1  pyspark_aws_app/
2  |─ src/
3  |   |─ __init__.py
4  |   |─ config_reader.py    # Reads configuration from S3
5  |   |─ athena_loader.py    # Queries Athena tables
6  |   |─ data_transformer.py # PySpark transformations
7  |   |─ s3_writer.py        # Writes data to S3
8  |   └─ main.py             # Main application script
9  |─ tests/
10 |   |─ __init__.py
11 |   |─ test_config_reader.py
12 |   |─ test_athena_loader.py
13 |   |─ test_data_transformer.py
14 |   |─ test_s3_writer.py
15 |   └─ test_main.py
16 |─ requirements.txt        # Dependencies
17 |─ README.md               # Documentation
18
```

Code Implementation

src/config_reader.py

```
1  import json
2  import boto3
3
4  def read_config_from_s3(bucket_name, key):
5      s3 = boto3.client("s3")
6      obj = s3.get_object(Bucket=bucket_name, Key=key)
7      return json.loads(obj["Body"].read())
8
```

src/athena_loader.py

```
1  from pyathena import connect
2
3  def query_athena(sql_query, database, work_group):
4      conn = connect(database=database, work_group=work_group)
5      cursor = conn.cursor()
6      cursor.execute(sql_query)
7      return cursor.fetchall()
8
```

src/data_transformer.py

```

1 from pyspark.sql import DataFrame
2
3 def transform_data(df: DataFrame) -> DataFrame:
4     return df.filter(df["Age"] > 20).select("Name", "Age")
5

```

src/s3_writer.py

```

1 def write_to_s3(df, bucket_name, key):
2     df.write.csv(f"s3a://{bucket_name}/{key}", mode="overwrite", header=True)
3

```

src/main.py

```

1 from pyspark.sql import SparkSession
2 from src.config_reader import read_config_from_s3
3 from src.athena_loader import query_athena
4 from src.data_transformer import transform_data
5 from src.s3_writer import write_to_s3
6
7 def main():
8     spark = SparkSession.builder.appName("PySparkAWSApp").getOrCreate()
9
10    # Read configuration
11    config = read_config_from_s3("test-bucket", "config.json")
12
13    # Query Athena
14    sql_query = config["athena_query"]
15    data = query_athena(sql_query, config["athena_database"], config["athena_work_group"])
16
17    # Load data into PySpark
18    df = spark.createDataFrame(data, schema=["Name", "Age"])
19
20    # Transform data
21    transformed_df = transform_data(df)
22
23    # Write data back to S3
24    write_to_s3(transformed_df, config["output_bucket"], "output/data.csv")
25
26 if __name__ == "__main__":
27     main()
28

```

Test Cases

tests/test_config_reader.py

```

1 import boto3
2 from moto import mock_s3
3 from src.config_reader import read_config_from_s3
4
5 @mock_s3
6 def test_read_config_from_s3():
7     s3 = boto3.client("s3", region_name="us-east-1")
8     s3.create_bucket(Bucket="test-bucket")

```

```

9     s3.put_object(Bucket="test-bucket", Key="config.json", Body='{ "athena_query": "SELECT * FROM test",
"output_bucket": "output-bucket"}')
10
11     config = read_config_from_s3("test-bucket", "config.json")
12     assert config["athena_query"] == "SELECT * FROM test"
13

```

tests/test_athena_loader.py

```

1 from src.athena_loader import query_athena
2
3 def test_query_athena(mocked):
4     mocked.patch("src.athena_loader.connect")
5     query_athena("SELECT * FROM test", "test_db", "primary")
6

```

tests/test_data_transformer.py

```

1 from pyspark.sql import SparkSession
2 from src.data_transformer import transform_data
3
4 def test_transform_data():
5     spark = SparkSession.builder.appName("Test").getOrCreate()
6     data = [("John", 30), ("Doe", 25), ("Jane", 15)]
7     df = spark.createDataFrame(data, ["Name", "Age"])
8
9     transformed_df = transform_data(df)
10    assert transformed_df.count() == 2
11

```

tests/test_s3_writer.py

```

1 from pyspark.sql import SparkSession
2 from moto import mock_s3
3 import boto3
4 from src.s3_writer import write_to_s3
5
6 @mock_s3
7 def test_write_to_s3():
8     spark = SparkSession.builder.appName("Test").getOrCreate()
9     s3 = boto3.client("s3", region_name="us-east-1")
10    s3.create_bucket(Bucket="output-bucket")
11
12    data = [("John", 30), ("Doe", 25)]
13    df = spark.createDataFrame(data, ["Name", "Age"])
14    write_to_s3(df, "output-bucket", "output/data.csv")
15

```

How to Run or Test the Application

1. Install Dependencies:

```

1 pip install -r requirements.txt
2

```

2. Run Tests:

```
1 pytest tests/
2
```

3. Run the Application:

```
1 python src/main.py
2
```

Key Notes

- Use the `moto` library to mock AWS services for local testing.
- Use `pytest-mock` for mocking Athena queries.
- Configure `spark-defaults.conf` for local S3 access via `s3a://`.
- For production, switch mock configurations to actual AWS settings.

Let me know if you need further clarifications or enhancements!

update the above code, to include reading, transforming and writing to redshift

To include functionality for reading, transforming, and writing data to Amazon Redshift in your PySpark application, we'll use the `JDBC` connector provided by Spark. Below is the updated project structure, full code, and how to test the application.

Updated Project Structure

```
1 pyspark_redshift_app/
2 |─ src/
3 |   |─ __init__.py
4 |   |─ config_reader.py      # Reads configuration from S3
5 |   |─ redshift_loader.py    # Reads data from Redshift
6 |   |─ data_transformer.py   # PySpark transformations
7 |   |─ redshift_writer.py    # Writes data to Redshift
8 |   |─ main.py               # Main application script
9 |─ tests/
10 |   |─ __init__.py
11 |   |─ test_config_reader.py
12 |   |─ test_redshift_loader.py
13 |   |─ test_data_transformer.py
14 |   |─ test_redshift_writer.py
15 |   |─ test_main.py
16 |─ requirements.txt         # Dependencies
17 |─ README.md                # Documentation
18
```

Code Implementation

`src/config_reader.py`

```
1 import json
2 import boto3
```

```

3
4 def read_config_from_s3(bucket_name, key):
5     s3 = boto3.client("s3")
6     obj = s3.get_object(Bucket=bucket_name, Key=key)
7     return json.loads(obj["Body"].read())
8

```

src/redshift_loader.py

```

1 from pyspark.sql import SparkSession
2
3 def read_from_redshift(spark: SparkSession, jdbc_url: str, table: str, user: str, password: str):
4     return spark.read.format("jdbc").option("url", jdbc_url) \
5         .option("dbtable", table) \
6         .option("user", user) \
7         .option("password", password) \
8         .load()
9

```

src/data_transformer.py

```

1 from pyspark.sql import DataFrame
2
3 def transform_data(df: DataFrame) -> DataFrame:
4     return df.filter(df["age"] > 20).select("name", "age")
5

```

src/redshift_writer.py

```

1 def write_to_redshift(df, jdbc_url, table, user, password):
2     df.write.format("jdbc").option("url", jdbc_url) \
3         .option("dbtable", table) \
4         .option("user", user) \
5         .option("password", password) \
6         .mode("overwrite") \
7         .save()
8

```

src/main.py

```

1 from pyspark.sql import SparkSession
2 from src.config_reader import read_config_from_s3
3 from src.redshift_loader import read_from_redshift
4 from src.data_transformer import transform_data
5 from src.redshift_writer import write_to_redshift
6
7 def main():
8     spark = SparkSession.builder.appName("PySparkRedshiftApp").getOrCreate()
9
10    # Read configuration
11    config = read_config_from_s3("test-bucket", "config.json")
12
13    # Load data from Redshift
14    df = read_from_redshift(
15        spark,
16        config["redshift_jdbc_url"],
17        config["input_table"],
18        config["redshift_user"],

```

```

19     config["redshift_password"]
20 )
21
22 # Transform data
23 transformed_df = transform_data(df)
24
25 # Write transformed data to Redshift
26 write_to_redshift(
27     transformed_df,
28     config["redshift_jdbc_url"],
29     config["output_table"],
30     config["redshift_user"],
31     config["redshift_password"]
32 )
33
34 if __name__ == "__main__":
35     main()
36

```

Test Cases

tests/test_redshift_loader.py

```

1 from pyspark.sql import SparkSession
2 from src.redshift_loader import read_from_redshift
3
4 def test_read_from_redshift(mock):
5     spark = SparkSession.builder.appName("Test").getOrCreate()
6     mock.patch("pyspark.sql.DataFrameReader.load", return_value=spark.createDataFrame([("John", 30)],
7 ["name", "age"]))
8     df = read_from_redshift(spark, "jdbc:redshift://test", "test_table", "user", "password")
9     assert df.count() == 1
10    assert df.columns == ["name", "age"]
11

```

tests/test_redshift_writer.py

```

1 from pyspark.sql import SparkSession
2 from src.redshift_writer import write_to_redshift
3
4 def test_write_to_redshift(mock):
5     spark = SparkSession.builder.appName("Test").getOrCreate()
6     data = [("John", 30), ("Doe", 25)]
7     df = spark.createDataFrame(data, ["name", "age"])
8
9     mock.patch("pyspark.sql.DataFrameWriter.save")
10    write_to_redshift(df, "jdbc:redshift://test", "test_table", "user", "password")
11

```

tests/test_data_transformer.py

```

1 from pyspark.sql import SparkSession
2 from src.data_transformer import transform_data
3
4 def test_transform_data():
5     spark = SparkSession.builder.appName("Test").getOrCreate()
6     data = [("John", 30), ("Doe", 25), ("Jane", 15)]

```

```
7 df = spark.createDataFrame(data, ["name", "age"])
8
9 transformed_df = transform_data(df)
10 assert transformed_df.count() == 2
11
```

How to Run the Application

1. **Install Dependencies:** Add the following to `requirements.txt`:

```
1 boto3
2 pyspark
3 pytest
4 pytest-mock
5
```

Install using:

```
1 pip install -r requirements.txt
2
```

2. **Set Up Redshift JDBC Driver:**

- Download the Redshift JDBC driver from AWS.
- Add the driver to your Spark classpath:

```
1 export SPARK_CLASSPATH=/path/to/redshift-jdbc-driver.jar
2
```

3. **Run Tests:**

```
1 pytest tests/
2
```

4. **Run the Application:**

```
1 python src/main.py
2
```

Key Notes

1. Replace `jdbc:redshift://...` with your actual Redshift JDBC URL in the configuration file.
2. Ensure the IAM role or credentials used have appropriate access to the Redshift cluster.
3. Use `pytest-mock` to mock external dependencies for testing Redshift interactions.

Let me know if you need further assistance!

how to package the above code and submit it to the emr cluster

Packaging your PySpark application and submitting it to an EMR cluster involves the following steps:

Steps to Package and Submit the Application

1. Organize the Project Structure

Ensure your project is structured properly:

```
1 pyspark_redshift_app/
2 |─ src/
3 |   |─ __init__.py
4 |   |─ config_reader.py
5 |   |─ redshift_loader.py
6 |   |─ data_transformer.py
7 |   |─ redshift_writer.py
8 |   └─ main.py
9 |─ tests/                # Optional, not required for deployment
10 |─ requirements.txt      # Python dependencies
11 |─ README.md             # Documentation
12 |─ setup.py              # For packaging the project
13 |─ config.json           # Configuration file (optional)
14
```

2. Create a `setup.py` File for Packaging

The `setup.py` file allows you to package the code as a Python module.

```
1 from setuptools import setup, find_packages
2
3 setup(
4     name="pyspark_redshift_app",
5     version="1.0.0",
6     packages=find_packages(),
7     install_requires=[
8         "boto3",
9         "pyspark",
10        "pyathena"
11    ],
12    entry_points={
13        "console_scripts": [
14            "run_app=src.main:main"
15        ]
16    },
17 )
18
```

3. Package the Application

1. Create a `.tar.gz` or `.whl` package:

```
1 python setup.py sdist
2
```

This will create a `dist/` directory containing the packaged application, e.g., `pyspark_redshift_app-1.0.0.tar.gz`.

2. If you prefer a wheel:

```
1 python setup.py bdist_wheel
2
```

4. Prepare Dependencies

1. List Python dependencies in `requirements.txt` :

```
1 boto3
2 pyspark
3 pyathena
4
```

2. Ensure the Redshift JDBC driver is available on the cluster. Either:

- Place the `.jar` file in an S3 bucket.
- Or directly reference it in your Spark configuration.

5. Upload Files to S3

Upload the following files to an S3 bucket accessible by your EMR cluster:

- Application package (`.tar.gz` or `.whl` file from `dist/`)
- `requirements.txt`
- `config.json` (if applicable)
- Redshift JDBC driver `.jar` file (if not already on the cluster).

Example command:

```
1 aws s3 cp dist/pyspark_redshift_app-1.0.0.tar.gz s3://your-bucket/path/
2 aws s3 cp requirements.txt s3://your-bucket/path/
3 aws s3 cp config.json s3://your-bucket/path/ # If used
4
```

6. Submit the Application to EMR

Option 1: Submit Using `spark-submit`

1. SSH into the EMR master node:

```
1 ssh -i your-key.pem hadoop@<emr-master-node-dns>
2
```

2. Run `spark-submit` :

```
1 spark-submit \
2   --deploy-mode cluster \
3   --master yarn \
4   --py-files s3://your-bucket/path/pyspark_redshift_app-1.0.0.tar.gz \
5   --jars s3://your-bucket/path/RedshiftJDBC42-no-awssdk-1.2.55.1083.jar \
6   s3://your-bucket/path/pyspark_redshift_app-1.0.0.tar.gz/src/main.py
7
```

Option 2: Submit Using EMR Steps

1. Add a Spark step to your EMR cluster:

```
1 aws emr add-steps --cluster-id <CLUSTER_ID> \
2   --steps Type=Spark,Name="PySparkApp",ActionOnFailure=CONTINUE,\
3
```

```
Args=[--deploy-mode,cluster,--master,yarn,--py-files,s3://your-bucket/path/pyspark_redshift_app-1.0.0.tar.gz,--jars,s3://your-bucket/path/RedshiftJDBC42-no-awssdk-1.2.55.1083.jar,s3://your-bucket/path/pyspark_redshift_app-1.0.0.tar.gz/src/main.py]
```

```
1
2 ---
3
4 ### **7. Debugging & Logs**
5 1. Monitor the EMR console for application status.
6 2. Review logs in the S3 bucket configured for EMR logs (found under EMR cluster settings).
7 3. SSH into the master node and check logs:
8 ```bash
9 yarn logs -applicationId <APPLICATION_ID>
10
```

Additional Notes

- Ensure the EMR cluster has appropriate IAM roles to access S3, Redshift, and other AWS services.
- Use the `config.json` file to dynamically configure settings like Redshift JDBC URL, user, password, input table, and output table.

Let me know if you need further clarifications!

Explain every file in the above code and what does each file do

Here's a detailed explanation of each file in the provided project structure, its role, and how it contributes to the functionality of the PySpark application.

Project Structure Overview

```
1 pyspark_redshift_app/
2 |─ src/
3 |   |─ __init__.py
4 |   |─ config_reader.py
5 |   |─ redshift_loader.py
6 |   |─ data_transformer.py
7 |   |─ redshift_writer.py
8 |   └─ main.py
9 |─ tests/
10 |─ requirements.txt
11 |─ setup.py
12 |─ config.json
13 └─ README.md
14
```

Detailed Explanation of Each File

1. `src/__init__.py`

- **Purpose:** Marks the `src/` directory as a Python package.
- **Contents:** Typically an empty file, but can optionally include initialization code for the package.

- **Usage:** Ensures other scripts can import modules from this directory, e.g., `from src.config_reader import read_config_from_s3`.
-

2. `src/config_reader.py`

- **Purpose:** Handles the reading of configurations stored in an Amazon S3 bucket.
 - **Key Functionality:**
 - Connects to S3 using the `boto3` library.
 - Retrieves the configuration file (JSON format) and parses it into a Python dictionary.
 - **Why It's Important:**
 - Centralized management of configuration values such as database connection details, table names, and other parameters.
 - Makes the application adaptable to different environments by simply changing the configuration file.
-

3. `src/redshift_loader.py`

- **Purpose:** Reads data from an Amazon Redshift table into a PySpark DataFrame.
 - **Key Functionality:**
 - Uses Spark's JDBC connector to fetch data from Redshift.
 - Takes parameters like the JDBC URL, table name, and credentials (user and password).
 - **Why It's Important:**
 - Facilitates reading large datasets efficiently from Redshift into Spark for processing.
-

4. `src/data_transformer.py`

- **Purpose:** Contains the data transformation logic for processing the DataFrame.
 - **Key Functionality:**
 - Defines functions to filter, clean, and reshape the data based on business requirements.
 - For example, in the provided code, it filters records where `age > 20` and selects only `name` and `age`.
 - **Why It's Important:**
 - Ensures the data meets the requirements of downstream systems or applications.
 - Encapsulates the transformation logic to keep the code modular and testable.
-

5. `src/redshift_writer.py`

- **Purpose:** Writes the transformed DataFrame back to an Amazon Redshift table.
 - **Key Functionality:**
 - Uses Spark's JDBC connector to write data to Redshift.
 - Supports different write modes (`overwrite`, `append`, etc.).
 - **Why It's Important:**
 - Handles writing processed data to the destination database for further use or analysis.
 - Ensures data integrity and proper handling of Redshift table schemas.
-

6. `src/main.py`

- **Purpose:** The entry point for the PySpark application.
- **Key Functionality:**
 - Initializes the Spark session.

- Orchestrates the workflow:
 - i. Reads the configuration using `config_reader`.
 - ii. Fetches data from Redshift using `redshift_loader`.
 - iii. Transforms the data using `data_transformer`.
 - iv. Writes the data back to Redshift using `redshift_writer`.

- **Why It's Important:**

- Provides a high-level overview of the application's workflow.
 - Acts as the glue that integrates all the other components.
-

7. `tests/`

- **Purpose:** Contains unit tests for the application.

- **Key Functionality:**

- Validates individual components (`config_reader`, `redshift_loader`, etc.) in isolation.
- Uses mock data to ensure functions work as expected without needing actual connections to S3 or Redshift.

- **Why It's Important:**

- Promotes test-driven development (TDD).
 - Ensures the application behaves correctly and prevents regressions.
-

8. `requirements.txt`

- **Purpose:** Lists all Python dependencies required for the project.

- **Contents:**

- `boto3`: AWS SDK for Python, used to interact with S3.
- `pyspark`: Required for Spark operations.
- `pyathena`: (Optional) If interacting with Athena.
- `pytest` and `pytest-mock`: For unit testing.

- **Why It's Important:**

- Makes it easy to set up the project environment by installing all dependencies with a single command (`pip install -r requirements.txt`).
-

9. `setup.py`

- **Purpose:** Configures the project for packaging and distribution.

- **Key Functionality:**

- Defines the package name, version, and dependencies.
- Allows the project to be packaged as a `.tar.gz` or `.whl` file.

- **Why It's Important:**

- Enables easy deployment of the project to clusters like EMR.
 - Standardizes the project structure for Python applications.
-

10. `config.json`

- **Purpose:** Stores configuration values such as database connection details, table names, and other parameters in JSON format.

- **Example Contents:**

```
1 {
```

```
2     "redshift_jdbc_url": "jdbc:redshift://your-redshift-cluster:5439/yourdb",
3     "input_table": "source_table",
4     "output_table": "target_table",
5     "redshift_user": "your_user",
6     "redshift_password": "your_password"
7 }
8
```

- **Why It's Important:**

- Provides a simple way to manage environment-specific configurations without hardcoding them in the code.
-

11. README.md

- **Purpose:** Provides documentation for the project.

- **Contents:**

- Overview of the application.
- How to set up the project.
- How to run the application.
- Examples of configuration and testing.

- **Why It's Important:**

- Helps new developers quickly understand the project.
 - Serves as a reference for maintaining and extending the application.
-

How These Files Work Together

1. `main.py` drives the workflow by calling functions from the other modules.
2. `config_reader.py` provides the necessary configuration for interacting with Redshift.
3. `redshift_loader.py` retrieves data from Redshift.
4. `data_transformer.py` applies transformations to the retrieved data.
5. `redshift_writer.py` writes the transformed data back to Redshift.
6. **Tests** validate the correctness of each component independently.

This modular design ensures the application is maintainable, testable, and reusable. Let me know if you need further details!