**Guided Lab: Introduction to AWS Glue Using Python Shell**

**Description**

In this guided lab, we'll work with a CSV file containing transaction records and use Python Shell to redact sensitive information.

**AWS Glue** is a fully managed extract, transform, and load (ETL) service that makes it easy for customers to prepare and load their data for analytics. It provides a Python shell script that can be used to perform data transformations, including redacting sensitive data.

**CSV** stands for **Comma Separated Values**. It is a file format used to store tabular data, such as in spreadsheets or databases. In a CSV file, each line represents a record, and the fields in each record are separated by commas. The first line of the file typically contains the names of the fields. CSV files are commonly used for data exchange between different applications and are supported by many software programs.

**Prerequisites**

This lab assumes you have basic knowledge of Python programming and experience creating an Amazon S3 bucket and are familiar with its basic components.

If you find any gaps in your knowledge, consider taking the following labs:

- Creating an Amazon S3 bucket.

**Objectives**

In this lab, you will:

- Learn how to use AWS Glue Python Shell to process and transform data stored in an S3 bucket using Python scripts.

[Subscribe to access AWS PlayCloud Labs](#)

**Lab Steps**

**Prepare Your Environment**

1. **Log into the AWS Management Console**: Navigate to the console and sign in.

2. **Create an S3 Bucket**:

    o Go to the S3 service in the AWS Console.

    o Click **"Create bucket"**.

    o Give your bucket a unique name.

    o Leave the rest of the settings as default and click **Create bucket**

    o

3. **Create a two folders:**

- 
  - Navigate to your newly created bucket.
  - Create two folders
    - **input**
    - **output**



4. Download the CSV Data:

https://media.tutorialsdojo.com/public/transactions.csv

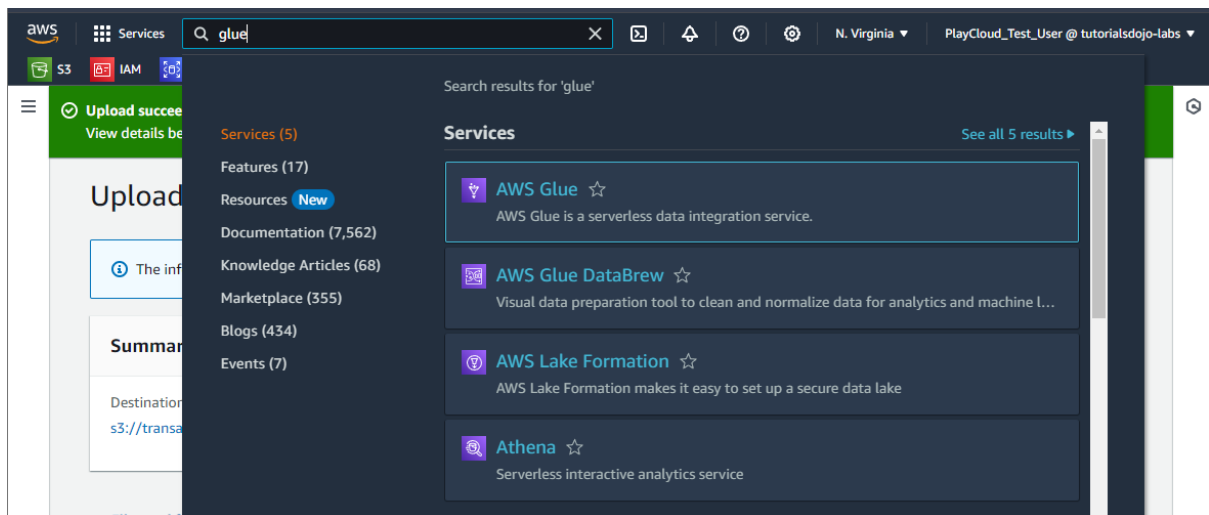*Take a look at the provided CSV file by opening the CSV file using Excel or any other spreadsheet application to visualize the data. When you are done you can continue and upload the file.*

5. Upload the file to the **input/** folder

**Create Your First Python Shell Job in AWS Glue**

1. Navigate to the **AWS Glue Console**.



2. You can either click on **ETL Jobs** on the left bar or click on the **Author and edit ETL jobs**

3. Click **Script Editor**



4. Under Engine select **Python shell** and select **Start fresh** as the Options

Click **Create script**

5. Provide a name for the job, e.g., "DataProcessingJob".



6. In the script section, paste the provided Python script.

import sys

import pandas as pd

from awsglue.utils import getResolvedOptions


# Fetch job parameters

args = getResolvedOptions(sys.argv, [])

input_path = "s3://{bucket_name}/input/transactions.csv"

output_path = "s3://{bucket_name}/output/"


# Load data from S3

transactions = pd.read_csv(input_path)

```python
# Function to redact card numbers

def redact_card_number(card_number):

    # Convert the card number to a string, just in case it's in numeric format

    card_number = str(card_number)

    # Replace all but the last 4 digits with asterisks

    return '*' * (len(card_number) - 4) + card_number[-4:]


# Apply the redaction function to the CardNumber column

transactions['CardNumber'] = transactions['CardNumber'].apply(redact_card_number)


# Save processed data

transactions.to_csv(output_path + "transactions_redacted.csv", index=False)


print("Data processing complete. Card numbers redacted.")
```

*DO NOT FORGET TO CHANGE THE PLACEHOLDER, {bucket_name}*

- CODE EXPLANATION:
  **Importing Libraries**:

  - import sys: This module provides access to some variables used or maintained by the Python interpreter and to functions that interact strongly with the interpreter.

  - import pandas as pd: This imports the pandas library, which is a powerful data manipulation tool built on top of the Python programming language.

  - from awsglue.utils import getResolvedOptions: This imports a function getResolvedOptions from the awsglue.utils module. This function is used to fetch job parameters when running an AWS Glue job.

**Fetching Job Parameters**:

  - args = getResolvedOptions(sys.argv, []): This line fetches the job parameters using the getResolvedOptions function. It retrieves the parameters passed to the Glue job at runtime and stores them in the args dictionary.

**Defining Input and Output Paths**:

  - input_path = "s3://{bucket_name}/input/transactions.csv": This specifies the S3 path where the input CSV file (transactions.csv) is located.

  - output_path = "s3://{bucket_name}/output/": This specifies the S3 path where the processed data will be saved.

**Loading Data from S3**:

- o transactions = pd.read_csv(input_path): This line uses the pandas read_csv function to load the data from the CSV file located at input_path into a pandas DataFrame named transactions.
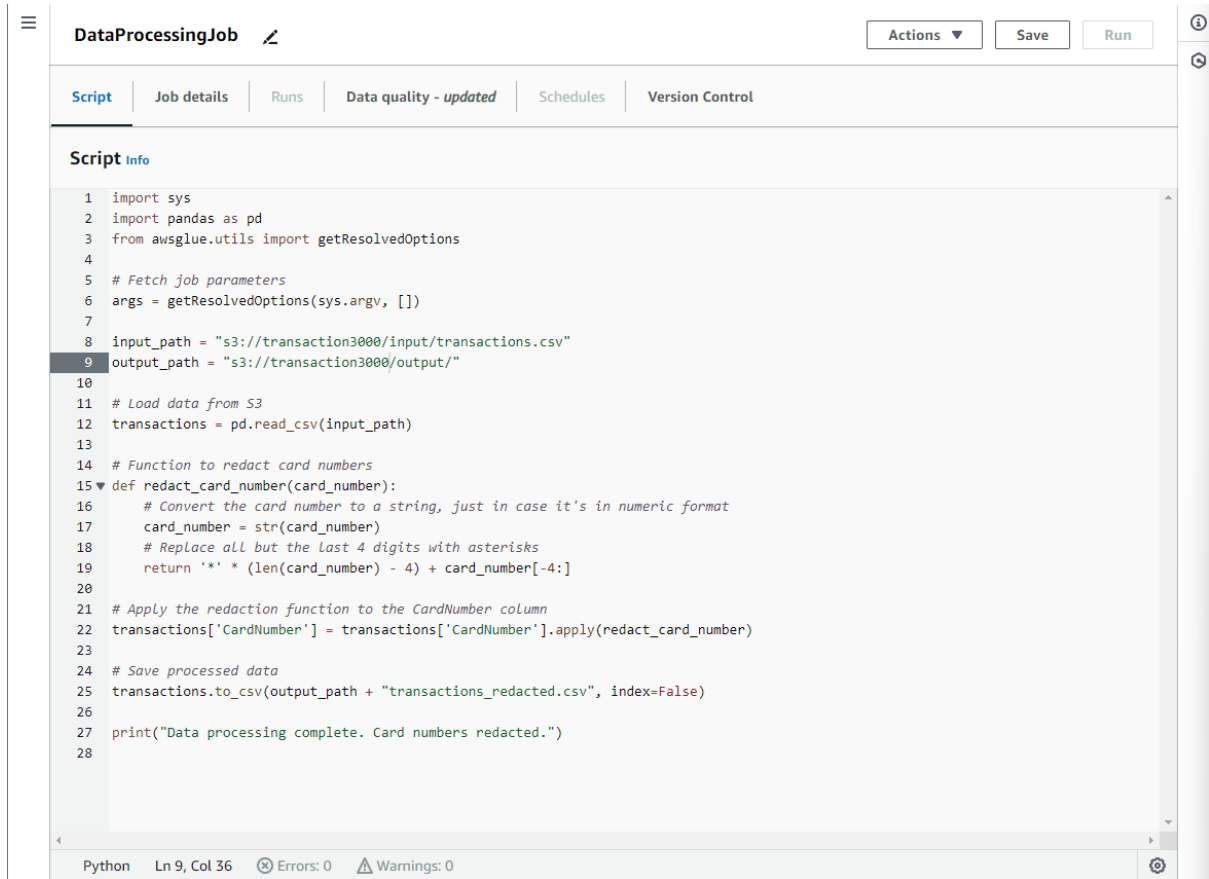
**Redacting Card Numbers**:

- o def redact_card_number(card_number): This defines a function named redact_card_number that takes a card_number as input.

- o card_number = str(card_number): This line converts the card_number to a string, ensuring it's in string format.

- o return '*' * (len(card_number) - 4) + card_number[-4:]: This line replaces all characters of the card_number except the last four digits with asterisks and returns the redacted card number.

- o transactions['CardNumber'] = transactions['CardNumber'].apply(redact_card_number): This applies the redact_card_number function to the 'CardNumber' column of the transactions DataFrame, replacing the original card numbers with the redacted versions.

**Saving Processed Data**:

- o transactions.to_csv(output_path + "transactions_redacted.csv", index=False): This line saves the processed data (with redacted card numbers) to a new CSV file named transactions_redacted.csv (You can changed the name as desired) at the specified output_path.

**Printing Message**:

- o print("Data processing complete. Card numbers redacted."): This line prints a message indicating that the data processing is complete and the card numbers have been redacted.
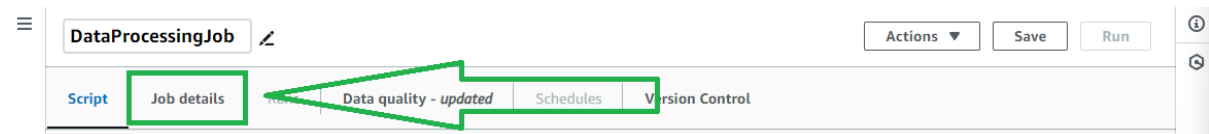
```
DataProcessingJob

Script    Job details    Runs    Data quality - updated    Schedules    Version Control

Script Info

1   import sys
2   import pandas as pd
3   from awsglue.utils import getResolvedOptions
4
5   # Fetch job parameters
6   args = getResolvedOptions(sys.argv, [])
7
8   input_path = "s3://transaction3000/input/transactions.csv"
9   output_path = "s3://transaction3000/output/"
10
11  # Load data from S3
12  transactions = pd.read_csv(input_path)
13
14  # Function to redact card numbers
15 ▼ def redact_card_number(card_number):
16      # Convert the card number to a string, just in case it's in numeric format
17      card_number = str(card_number)
18      # Replace all but the last 4 digits with asterisks
19      return '*' * (len(card_number) - 4) + card_number[-4:]
20
21  # Apply the redaction function to the CardNumber column
22  transactions['CardNumber'] = transactions['CardNumber'].apply(redact_card_number)
23
24  # Save processed data
25  transactions.to_csv(output_path + "transactions_redacted.csv", index=False)
26
27  print("Data processing complete. Card numbers redacted.")
28

Python    Ln 9, Col 36    ⊗ Errors: 0    ⚠ Warnings: 0
```

7. Review the script to ensure it performs data redaction as required.

8. Set the **IAM role** you created earlier.

- Click on Job Details



- Under IAM role, Select **PlayCLoud-Sandbox.**

**Basic properties** Info

**Name**

DataProcessingJob

**Description - *optional***

Descriptions can be up to 2048 characters long.

**IAM Role**
Role assumed by the job with permission to access your data stores. Ensure that this role has permission to your Amazon S3 sources, targets, temporary directory, scripts, and any libraries used by the job.

PlayCLoud-Sandbox ▼

**Type**
The type of ETL job. This is set automatically based on the types of data sources you have selected.

Python Shell

**Python version**

Python 3.9 ▼

☑ Load common analytics libraries (recommended)
Include common Python libraries from pypi.org ⧉ such as pandas, numpy and s3fs. Uncheck if you are loading your own libraries, and wish to avoid version conflicts.

**Automatically scale the number of workers**
☐ AWS Glue will optimize costs and resource usage by dynamically scaling the number of workers up and down throughout the job run. Requires Glue 3.0 or later.

- Click **Save**



| DataProcessingJob ✎ | | | | Actions ▼ | Save | Run |

| Script | Job details | Runs | Data quality - *updated* | Schedules | Version Control |

9. You should be able to see the **Runs Tab** now and the **Run** button in the upper right corner.



DataProcessingJob          Last modified on 4/29/2024, 6:02:27 PM    Actions ▼    Save    **Run**

| Script | Job details | **Runs** | Data quality - *updated* | Schedules | Version Control |

**Basic properties** Info

**Name**
DataProcessingJob

**Description - *optional***

Descriptions can be up to 2048 characters long.

10. Go to the **Runs Tab,** Click either **Run job** or **Run** in the upper right corner.

11. You will be prompted with a **Successfully Started Job.**



12. Wait for it to finished and you should see similar output like the image below.



13. Now, navigate back to your S3 bucket's **output/ folder**

- You should see a csv file named **transactions_redacted.csv**
  *Remember, in our code we intentionally named this **transactions_redacted.csv.** You can changed this as desired in the code editor.*

14. Download the object and open the CSV file using Excel or any other spreadsheet application to visualize the data. When you are done you can continue and upload the file.

| | A | B | C | D |
|---|---|---|---|---|
| 1 | Date | Amount | CardNumber | Merchant |
| 2 | 2024-01-12 | $150.00 | ************3456 | Coffee Shop |
| 3 | 2024-01-13 | $45.00 | ************4321 | Bookstore |
| 4 | 2024-01-14 | $78.25 | ************4677 | Electronics Store |
| 5 | 2024-01-15 | $22.50 | ************4321 | Grocery Store |
| 6 | 2024-01-16 | $130.00 | ************4321 | Online Marketplace |

*Take your time to look at the data and observe what are the difference from the original CSV file data.*

Congratulations! You just learn how to use AWS Glue Python Shell to process and transform data stored in an S3 bucket using Python scripts. This is only an introductory of AWS Glue can do. Remember redacting is just one of the data transformation you can do with AWS Glue. You can also drop columns, join and relationalize data, etc. From this point you can go and try and experiment different python shell script as you see fit.

One last thing! It is a good practice to clean up the resources created during this lab. Not only will it make you a better professional, but you will also become a more organized person. Happy learning!