

We completed various essential tasks related to the management and exploration of data on AWS. Initially, we uploaded the default credit card dataset from the UCI archives "<http://archive.ics.uci.edu/dataset/350/default+of+credit+card+clients>" to an AWS S3 bucket, utilizing Amazon S3 as a highly scalable and dependable storage solution for our datasets. With the data securely stored in S3, we utilized AWS Athena for data exploration. Athena, a serverless query service, allowed us to directly analyze data in S3 using SQL queries, offering a user-friendly approach for ad-hoc analysis without the necessity of overseeing a traditional database.

Furthermore, we experimented with AWS Glue, a fully managed ETL service. Specifically, we employed the Glue crawler to automatically discover and catalog metadata from our raw dataset, aiding in the organization and structuring of the data. Following this, we established a database and tables in the AWS Glue Data Catalog, enabling us to efficiently query and analyze the data. These tables seamlessly integrate into the broader AWS ecosystem, facilitating advanced analytics through smooth interaction with other services.

Through querying the tables within the database, we extracted valuable insights from the dataset. This comprehensive workflow, incorporating S3 for storage, Athena for query-driven exploration, and Glue for ETL and cataloging, constitutes a robust data pipeline on AWS. The integration of these services streamlines scalable and efficient data processing, empowering us to uncover meaningful information from the dataset and make well-informed decisions based on the insights derived.

Below are screenshots of the queries and outputs:

The screenshot displays the AWS Athena console interface. On the left, the 'Data' sidebar shows the 'Data source' as 'AwsDataCatalog', the 'Database' as 'projdb', and a list of 'Tables (1)' including 'projtable'. The main area shows a SQL query editor with the following query:

```
1 -- Get a sample of data
2 SELECT * FROM projtable LIMIT 10;
3
4
```

Below the query editor, the 'Query results' section shows the query status as 'Completed'. The 'Results (10)' section displays a table with 14 columns: #, Id, limit_bal, sex, education, marriage, age, pay_0, pay_2, pay_3, pay_4, pay_5, pay_6, and bill_amt. The first two rows of data are visible:

#	Id	limit_bal	sex	education	marriage	age	pay_0	pay_2	pay_3	pay_4	pay_5	pay_6	bill_amt
1	1	20000	2	2	1	24	2	2	-1	-1	-2	-2	3913
2	2	120000	2	2	2	26	-1	2	0	0	0	2	2682

The footer of the console shows the AWS logo, 'CloudShell', 'Feedback', and copyright information: '© 2023, Amazon Web Services, Inc. or its affiliates. Privacy Terms Cookie preferences'.

AWS Services Search [Option+S]

Query 2 Query 1 Query 3 Query 4 Query 5

Data source: AwsDataCatalog Database: projdb

Tables and views: Create Filter tables and views

Tables (1): projtable Views (0)

SQL Ln 1, Col 56

```
1 --Count of Default and Non-Default Cases in our dataset
2 SELECT "default", COUNT(*) AS count FROM projtable GROUP BY "default";
3
```

Run Explain Cancel Clear Create

Query results Query stats

Completed Time in queue: 64 ms Run time: 425 ms Data scanned: 2.76 MB

Results (2) Copy Download results

Search rows

#	default	count
1	1	6636
2	0	23364

© 2023, Amazon Web Services, Inc. or its affiliates. Privacy Terms Cookie preferences

AWS Services Search [Option+S]

Query 2 Query 1 Query 3 Query 4 Query 5

Data source: AwsDataCatalog Database: projdb

Tables and views: Create Filter tables and views

Tables (1): projtable Views (0)

SQL Ln 1, Col 67

```
1 --Average Amount of Given Credit for Default and Non-Default Cases
2 SELECT "default", AVG(limit_bal) AS avg_credit_amount FROM projtable GROUP BY "default";
3
```

Run Explain Cancel Clear Create

Query results Query stats

Completed Time in queue: 135 ms Run time: 583 ms Data scanned: 2.76 MB

Results (2) Copy Download results

Search rows

#	default	avg_credit_amount
1	0	178099.72607430234
2	1	130109.65641952984

© 2023, Amazon Web Services, Inc. or its affiliates. Privacy Terms Cookie preferences

we extended our efforts by working on creating an AWS pipeline:

