

Optimizing Retail Operations through Data Warehousing

Anshu Reddy Dhamana - 018172141
Brendan Chao - 014142726
Sri Gopi Sarath Gode - 018191537
Srithareddy Devireddy - 018175456
Vinuthna Papana - 018176106



Introduction to Data Warehousing

A centralized repository for storing structured data from various sources. AWS Redshift serves as the data warehouse for this project. It consolidates sales, customer, and product data into a single source of truth.

Objective is to create a robust and scalable data warehousing solution to store and analyze retail data, enabling actionable insights and informed decision-making.

S3 (Data Source) -> AWS Glue (ETL) -> Redshift (Data Warehouse) -> Visualizations.

Motivation

In today's competitive retail environment, businesses generate vast amounts of data from transactions, customer interactions, and product sales. However, traditional systems struggle to process and analyze this data, leading to missed opportunities and suboptimal decisions.

This project leverages a scalable data warehouse on Amazon Redshift to transform raw retail sales data into actionable insights. Key benefits include:

- Rapid handling of massive datasets.
- Identifying trends, customer behavior patterns, and product performance.
- Enhancing inventory control, marketing strategies, and sales forecasting.



Overview of ETL Process



Extract: Data Collection

Source: Online Retail II dataset (CSV format).

Tools: Python for data extraction and Amazon S3 for storage.

Transform: Data Cleaning and Structuring

Handle missing values and remove duplicates.

Standardize data formats (e.g., date-time fields).

Create calculated fields like totalprice (quantity × unit price).

Generate time-based dimensions (e.g., year, month, day).

Load: Data Storage

Use Amazon Redshift as the target data warehouse.

Load data into fact and dimension tables using Redshift's COPY commands.

Key Features:

Automated, scalable pipeline for processing large datasets.

Ensures data consistency and readiness for analysis.

Tools Used:

Python (Pandas, Boto3), AWS S3, AWS Glue, Amazon Redshift, SQL.

Role of Amazon Redshift as a Data Warehouse

Centralized Data Repository - Stores transformed and structured data in a star schema. Facilitates integration of sales, customer, product, and time-related data for analysis.

Performance and Scalability - Optimized for running complex analytical queries on large datasets. Supports scalability to accommodate growing data volumes.

Efficient Query Execution - Leverages columnar storage for faster query performance. Uses distribution keys and sort keys to optimize data retrieval.

Data Accessibility - Accessible via BI tools (Tableau, Power BI) for creating visualizations. Enables real-time insights through SQL queries.

Integration with AWS Ecosystem - Seamless integration with AWS Glue for ETL processes. Direct connection to Amazon S3 for fast data loading using COPY commands.

Key Contributions to the Project - Acts as the backbone for querying and analyzing retail sales data. Enables calculation of metrics like total sales, customer purchasing patterns, and best-selling products.

load customer data into redshift cluster

Data from the "customer dimension" table was extracted, transformed, and loaded into the Redshift cluster for centralized storage.

This process ensures integration with the star schema, enabling efficient querying and analysis of customer-related metrics. Redshift's SQL-based querying capabilities allow for detailed customer insights and streamlined data exploration.

The screenshot shows the AWS Redshift Query Editor interface. On the left, the sidebar displays the database schema with tables like customerdimension, productsdimension, salesfacts, and timedimension. The main area contains a query editor window with the following SQL code:

```
1 SELECT * FROM "dev"."public"."customerdimension";
2 WHERE customerid = 13685;
3
4 SELECT * FROM "dev"."public"."timedimension";
5
6 SELECT * FROM "dev"."public"."salesfacts";
7
8 SELECT a.customerid, c.customerid
9 FROM salesfacts s
10 INNER JOIN customerdimension c
11 ON s.customerid = c.customerid
12
13 LIMIT 100;
```

Below the code, five result sets are displayed, each showing 100 rows of data. The first result set, titled 'customer', shows columns 'customerid' and 'country'. The data includes rows for United Kingdom, France, and USA. The last row shown is customerid 14911 from the United Kingdom.

At the bottom right of the editor, it says 'Query ID 1412987 Elapsed time: 16 ms Total rows: 100'.

load product data into redshift cluster

Successfully loaded product data into Amazon Redshift for efficient querying and analysis. Utilized SQL workbench for data integration, ensuring accurate representation of product

The screenshot shows the Redshift query editor v2 interface. The left sidebar includes sections for Editor, Queries, Notebooks, Divers, History, and Scheduled queries. The main area has tabs for Untitled 1 and Untitled 2. Untitled 1 is active, showing a query window with the following SQL code:

```
1 SELECT * FROM "dev"."public"."customerdimension";
2 SELECT * FROM "dev"."public"."productdimension";
3 SELECT * FROM "dev"."public"."timedimension";
4 SELECT * FROM "dev"."public"."salesfacts";
5
6 SELECT s.customerid, c.customerid
7 FROM salesfacts s
8 INNER JOIN customerdimension c
9 ON s.customerid = c.customerid
10
11 LIMIT 100;
```

The results tab shows 100 rows from Result 2 (100), displaying columns productid, description, and price, all showing the value 'HEARTS WRAPPING TA...' and a price of 0.05.

Untitled 2 is also visible, showing a table definition for 'timedimension' with columns dated, year, month, day, weekday, day_string, month_string, year_string, and dated_string, each with specific data types and constraints.

load salesfacts into redshift cluster

The screenshot shows the AWS Redshift Query Editor interface. On the left, the sidebar displays the database schema with tables like customerdimension, productdimension, salesfacts, and timedimension. The main area shows a SQL query being run against the 'dev' database. The query joins customerdimension and salesfacts tables based on customerid. The results are displayed in a table with columns: invoiceid, productid, customerid, dated, quantity, totalprice, and invoice time. The table contains 100 rows of transaction data. At the bottom, the status bar indicates the query ID, elapsed time (66 ms), and total rows (100).

```
1 SELECT * FROM "dev"."public"."customerdimension";
2 SELECT * FROM "dev"."public"."productdimension";
3 SELECT * FROM "dev"."public"."timedimension";
4 SELECT * FROM "dev"."public"."salesfacts";
5
6 SELECT s.customerid, c.customerid
7 FROM salesfacts s
8 INNER JOIN customerdimension c
9 ON s.customerid = c.customerid
10
11
```

invoiceid	productid	customerid	dated	quantity	totalprice	invoice time
489437	10002	15362	2009-12-01	12	10.2	09:08:00
490136	10002	17432	2009-12-03	1	0.85	19:13:00
490140	10002	12872	2009-12-03	4	3.4	20:03:00
490144	10002	13154	2009-12-04	12	10.2	08:46:00
490229	10002	18223	2009-12-04	12	10.2	12:20:00
490458	10002	12583	2009-12-06	48	40.8	11:57:00
491506	10002	16470	2009-12-11	9	7.649999999999999	12:21:00
491752	10002	14032	2009-12-14	12	10.2	12:02:00
491826	10002	16619	2009-12-14	24	20.4	14:12:00
492945	10002	14032	2009-12-21	12	10.2	13:29:00
493375	10002	15514	2009-12-23	1	0.85	12:07:00
493446	10002	17807	2010-01-04	3	2.55	14:00:00
493678	10002	14769	2010-01-08	12	10.2	09:08:00
494390	10002	18841	2010-01-14	110	93.5	10:54:00
495183	10002	15296	2010-01-21	12	10.2	13:48:00
495534	10002	14032	2010-01-25	12	10.2	15:36:00
496033	10002	12872	2010-01-28	5	4.25	14:19:00
496216	10002	16892	2010-01-29	1	0.85	13:44:00
496237	10002	17677	2010-01-29	24	20.4	14:37:00
496400	10002	17365	2010-02-01	3	2.55	11:56:00
496616	10002	12682	2010-02-03	12	10.2	09:03:00
498249	10002	14911	2010-02-17	12	10.2	14:39:00
498361	10002	13394	2010-02-18	12	10.2	13:34:00
498850	10002	16416	2010-02-23	36	38.6	12:00:00

Successfully ingested sales transaction data into the Redshift cluster for advanced querying and analysis.

Ensured accurate integration of key metrics like quantity, totalprice, and invoice time to enable precise sales trend analysis.

Leveraged SQL workbench to validate the integrity and completeness of the loaded data.

load timedimen into redshift cluster

Successfully integrated time-based dimension data into the Redshift cluster for temporal analysis.

Key fields like dateid, year, month, weekday, and formatted strings were loaded to support seasonal and trend analysis.

The screenshot shows the AWS Redshift Query Editor interface. On the left, the sidebar displays the schema structure under the 'dev' database, specifically the 'public' schema, where the 'timedimension' table is selected. The main area contains a query editor window with the following SQL code:

```
1 SELECT * FROM "dev"."public"."customerdimension";
2 SELECT * FROM "dev"."public"."productdimension";
3 SELECT * FROM "dev"."public"."timedimension";
4 SELECT * FROM "dev"."public"."salesfacts";
5
6 SELECT s.customerid, c.customerid
7 FROM salesfacts s
8 INNER JOIN customerdimension c
9 ON s.customerid = c.customerid
10 LIMIT 100;
```

The results of the query are displayed in a table titled 'Result 1 (100)'. The table has columns: dated, year, month, day, and weekday. The data shows the first 100 rows of the joined tables. At the bottom of the results pane, it says 'Query ID 1412993 Elapsed time: 72 ms Total rows: 100'.

dated	year	month	day	weekday
2009-01-01	2009	1	1	Thursday
2009-01-02	2009	1	2	Friday
2009-01-03	2009	1	3	Saturday
2009-01-04	2009	1	4	Sunday
2009-01-05	2009	1	5	Monday
2009-01-06	2009	1	6	Tuesday
2009-01-07	2009	1	7	Wednesday
2009-01-08	2009	1	8	Thursday
2009-01-09	2009	1	9	Friday
2009-01-10	2009	1	10	Saturday
2009-01-11	2009	1	11	Sunday
2009-01-12	2009	1	12	Monday
2009-01-13	2009	1	13	Tuesday
2009-01-14	2009	1	14	Wednesday
2009-01-15	2009	1	15	Thursday
2009-01-16	2009	1	16	Friday
2009-01-17	2009	1	17	Saturday
2009-01-18	2009	1	18	Sunday
2009-01-19	2009	1	19	Monday
2009-01-20	2009	1	20	Tuesday
2009-01-21	2009	1	21	Wednesday
2009-01-22	2009	1	22	Thursday
2009-01-23	2009	1	23	Friday
2009-01-24	2009	1	24	Saturday

IAM Roles and Security

In Amazon Web Services (AWS), Identity and Access Management (IAM) roles serve as fundamental tools for securely delegating permissions within AWS resources. These roles define a set of permissions that dictate actions authorized for entities, be it an AWS service or a user, on AWS resources. They eliminate the necessity for long-term credentials like usernames and passwords when granting access to resources.

Introduction to IAM Roles in AWS



AmazonS3FullAccess: Grants complete access to Amazon S3 (Simple Storage Service), enabling any action on S3 buckets within the AWS account.

AmazonRedshiftFullAccess: Provides full access to Amazon Redshift, empowering the IAM role to manage Redshift clusters comprehensively. This includes tasks such as creating, modifying, deleting clusters, and handling administrative functions.

AWSGlueServiceRole: Grants the essential permissions for executing actions within AWS Glue, a fully managed ETL service facilitating data preparation and loading for analytics.



Benefits of IAM Roles for Security

Simplified Access: Enables secure and temporary access to AWS services like S3, Redshift, and Glue without hardcoding credentials.

Enhanced Security: Enforces least privilege, limiting access to only necessary resources.

Streamlined Integration: Facilitates seamless, secure interactions between AWS services in the ETL pipeline.

Audit and Compliance: Tracks and logs resource access with CloudTrail for monitoring and compliance.

Project Impact: Ensures secure data handling during ETL and protects sensitive information.

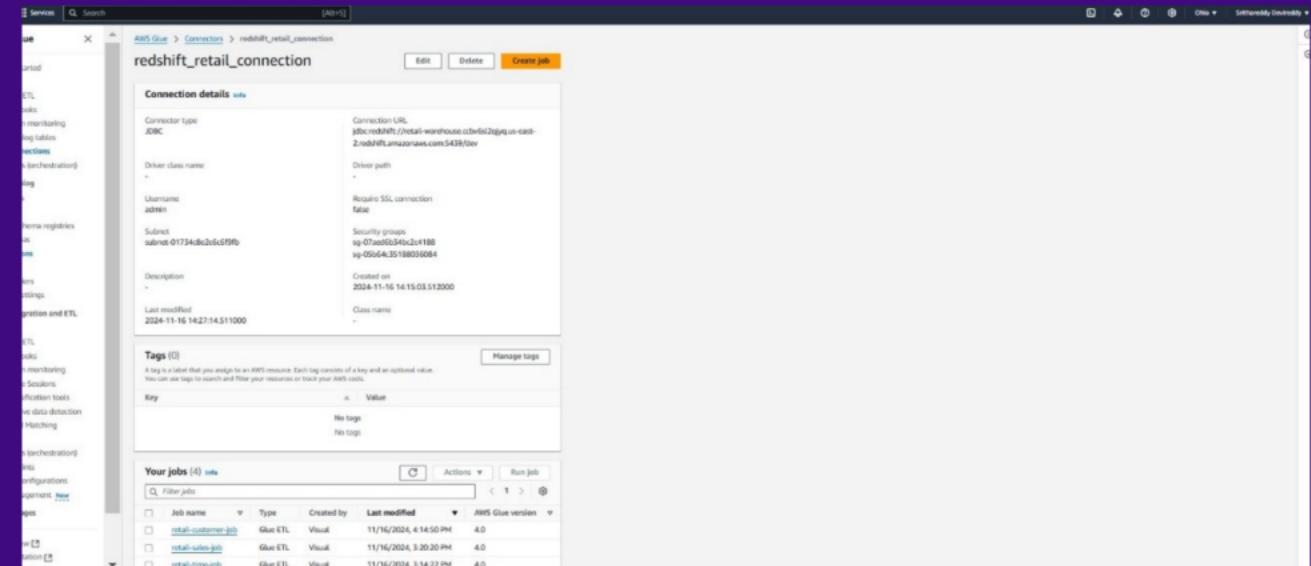
Centralized Permission Management

By consolidating these permissions within a single IAM role, entities assuming this role gain the combined access and capabilities across Amazon S3, Amazon Redshift, and AWS Glue. This practice aligns with security best practices, adhering to the principle of least privilege by granting only the necessary permissions for specific tasks



JDBC Connection to Temporary Database:

In AWS Glue, a JDBC connection establishes access to a temporary database, essential for the ETL process. JDBC, as a Java Database Connectivity tool, enables data interaction and transfer between the Glue environment and external databases. This connectivity supports AWS Glue in extracting, transforming, and loading data across diverse databases, ensuring smooth data processing throughout the ETL operations.



Contents of Temporary Database:

Contents of Temporary Database:

The screenshot shows the AWS Glue Data Catalog interface. On the left, a sidebar menu includes: Getting started, ETL jobs, Visual ETL, Notebooks, Job run monitoring, Data Catalog tables, Data connections, Workflows (orchestration), Data Catalog (expanded), Databases (selected), Tables, Stream schema registries, Schemas, Connections, Crawlers, Classifiers, Catalog settings, and Data Integration and ETL (expanded), ETL jobs, Visual ETL, and Monitoring.

The main content area displays the **retail_metadata_db** database properties. It shows the database name is **retail_metadata_db**, there is no description, the location is unspecified, and it was created on November 16, 2024 at 10:32:37 UTC. There are buttons for Edit and Delete.

Below the properties, the **Tables (6)** section is shown. It lists six tables: cleaned_customer_dat, customer_data_csv, modified_sales_data_c, product_data_csv, and sales_data_csv. All tables are located in the **retail_metadata_db** database and have an S3 location of **s3://retail-store-data-l**. They are all classified as CSV files and are not deprecated. Each table has a "View data" and "View data quality" link.

Name	Database	Location	Classification	Deprecated	View data	Data quality
cleaned_customer_dat	retail_metadata_db	s3://retail-store-data-l	CSV	-	Table data	View data quality
customer_data_csv	retail_metadata_db	s3://retail-store-data-l	CSV	-	Table data	View data quality
modified_sales_data_c	retail_metadata_db	s3://retail-store-data-l	CSV	-	Table data	View data quality
product_data_csv	retail_metadata_db	s3://retail-store-data-l	CSV	-	Table data	View data quality
sales_data_csv	retail_metadata_db	s3://retail-store-data-l	CSV	-	Table data	View data quality

Overview of the ETL Process

Extract:

Data extracted from S3 bucket containing raw retail sales data.

Uses Python scripts and AWS Glue Crawlers to identify data structure and schema.

Transform:

Cleaning data by addressing missing values and inconsistencies.

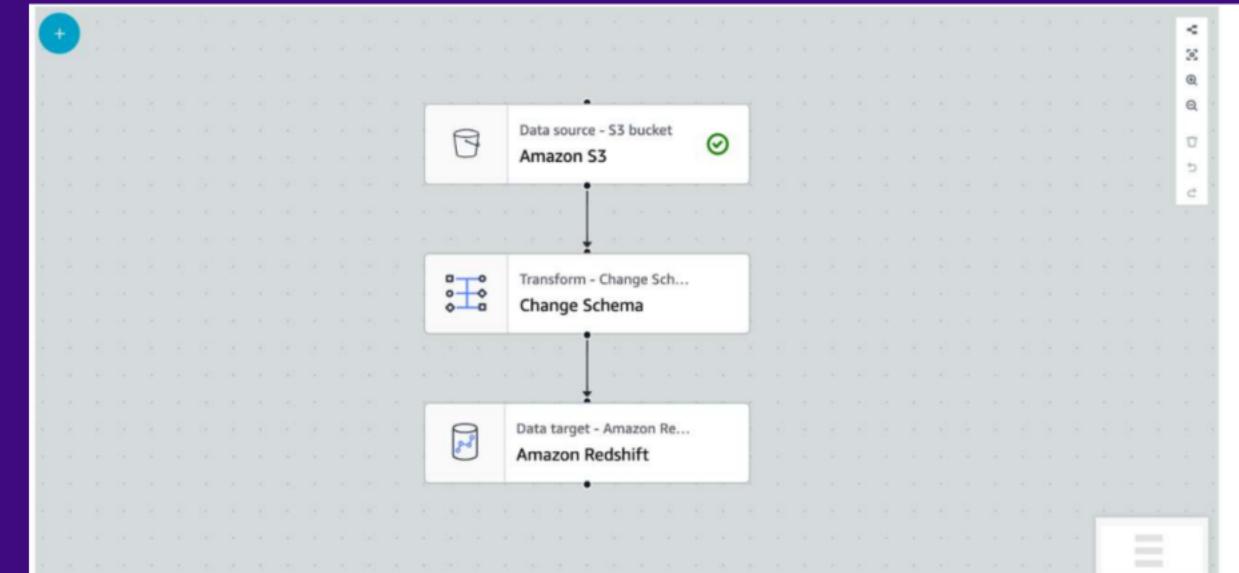
Applying transformations like creating time dimensions, calculating metrics (e.g., total sales).

Restructuring data into a star schema with fact and dimension tables.

Load:

Processed data loaded into Amazon Redshift for warehousing.

Utilizes SQL COPY commands for efficient bulk loading.



Execution of ETL job:

The screenshot shows the AWS Glue Job Editor interface for the job **retail-etl-fact-job**. The top navigation bar includes the VS icon, Services, Search, [Alt+S], and account information for Ohio and Srithareddy Devireddy.

The main view displays the **Runs** tab, which shows one successful run. The run details are as follows:

Run status	Retries	Start time (Local)	End time (Local)	Duration	Capacity...	Worker type	Glue version
Success Succeeded	0	11/17/2024 02:17:25	11/17/2024 02:20:04	2 m 25 s	10 DPUs	G.1X	4.0

The **Run details** section provides more specific information about the run:

Run details	Input arguments (10)	Continuous logs	Run insights	Metrics	Spark UI
Job name		Start time (Local)		Glue version	Last modified on (Local)
retail-etl-fact-job		11/17/2024 02:17:25		4.0	11/17/2024 02:20:04
Id		End time (Local)		Worker type	Log group name
jr_1fc86b1c07ca6f25d5f543ddb064d979cc84ff9	11/17/2024 02:20:04			G.1X	/aws-glue/jobs
7b487cb1d09dc9bebd2db736c	7b487cb1d09dc9bebd2db736c			Max capacity	Number of workers
Run status		Start-up time		10 DPUs	10
Success Succeeded		13 seconds		Execution class	Timeout
Retry attempt number		Execution time			

At the bottom, there are links for CloudWatch, Feedback, © 2024 Amazon Web Services, Inc. or its affiliates, Privacy, Terms, and Cookie preferences.

Redshift cluster

Details of Redshift cluster:

Amazon Redshift Serverless > Namespace configuration > default-namespace

default-namespace [Info](#)

General information

Namespace default-namespace	Status Available	Admin user name admin
Namespace ID dfdc940c-e435-4b2b-a6fa-7f0696aa1b3b	Date created November 16, 2024, 16:16 (UTC-08:00)	Database name dev
Namespace ARN arn:aws:redshift-serverless:us-east-2:783764579792:namespace/dfdc940c-e435-4b2b-a6fa-7f0696aa1b3b	Storage used 861 MB	Total table count 2

Workgroup [Actions](#) [Query data](#)

Workgroup name

Workgroup default-workgroup	Status Available
--------------------------------	--

[Actions](#)

Connection details of Redshift Cluster:

Amazon Redshift Serverless > Workgroup configuration > default-workgroup

default-workgroup [Info](#)

General information

Workgroup default-workgroup	Date created November 16, 2024, 16:16 (UTC-08:00)
Namespace default-namespace	Status Available
Namespace ARN arn:aws:redshift-serverless:us-east-2:783764579792:workgroup/b6510ad-0d41-4fd1-90f1-46532d7740f7	Configuration Production
Workgroup version 1.0.76178	Base capacity 128 RPU:s
	Custom domain name -
	Patch version Patch 166

[Actions](#) [Query data](#)

Data access [Lists](#) [Performance](#) [Tags](#)

Network and security [Info](#)

VPC security group sg-05644c35188030084	Enhanced VPC routing On
VPC endpoint ID spca-04a2a6ab6205a8ff	Publicly accessible On
	IP address type IPv4



Amazon S3 Bucket and Data Storage

Amazon S3 serves as a highly scalable storage solution for various types of data, including CSV files that are pivotal in ETL processes. S3 buckets act as data repositories, providing a reliable means to store raw data before it undergoes transformation and loading.



Data Transformation using AWS Glue

AWS Glue facilitates the transformation of raw data through automated ETL processes, allowing users to clean, reorganize, and prepare data for analysis. It efficiently handles complex transformation logic and integrates seamlessly with other AWS services such as Amazon Redshift.



Data Crawler Functionality

Crawlers are vital components in the ETL process, automating metadata discovery and organization, which significantly enhances data processing efficiency.

Importance of Crawlers in ETL

In AWS Glue, crawlers serve as automated processes used for discovering and cataloging metadata from diverse data sources. These processes analyze data structure and schema, facilitating streamlined data processing within the ETL (Extract, Transform, Load) workflow.

The S3 crawler within AWS Glue inspects data stored in Amazon S3 buckets. It identifies file formats and schema, organizing the data for subsequent use in the ETL process. By scanning S3 data, this crawler comprehends its structure, aiding in the preparation of data for transformation.

Similarly, the Redshift crawler in AWS Glue targets Redshift clusters. It assesses the data within Redshift tables, capturing their structure and schema. This acquired information becomes crucial for mapping and efficiently transforming data during ETL operations. It ensures consistency and accuracy in data processing within the Redshift environment.



Redshift Crawler:

The screenshot shows the AWS Glue console with the path: AWS Glue > Crawlers > retail-redshift-crawler. The crawler properties section includes:

- Name: retail-redshift-crawler
- IAM role: AWSGlueServiceRole-AmazonRedshiftFullAccess
- Description: None
- Database: retail_metadata_db
- Security configuration: None
- Status: READY
- Table prefix: None

The Crawler runs tab shows one run listed:

Start time (UTC)	End time (UTC)	Current/last duration	Status	DPU hours	Table changes
November 16, 2024 at 22:45:48	November 16, 2024 at 22:45:48	55 s	Completed	0.071	-



Retail S3 Crawler:

The screenshot shows the AWS Glue console with the path: AWS Glue > Crawlers > retail-s3-crawler. The crawler properties section includes:

- Name: retail-s3-crawler
- IAM role: AWSGlueServiceRole-AmazonS3FullAccess
- Description: None
- Database: retail_metadata_db
- Security configuration: None
- Status: READY
- Table prefix: None

The Crawler runs tab shows one run listed:

Start time (UTC)	End time (UTC)	Current/last duration	Status	DPU hours	Table changes
November 16, 2024 at 10:33:45	November 16, 2024 at 10:34:48	01 min 02 s	Completed	0.065	6 table changes, 0 partition changes

S3 Crawler: Functionality and Use Cases

Functionality:

AWS Glue S3 Crawler automates the discovery and cataloging of metadata from datasets stored in Amazon S3. Inspects data formats (CSV, JSON, Parquet, etc.) and schemas to generate a cataloged structure in the AWS Glue Data Catalog.

Identifies partition keys and organizes metadata for seamless querying and analysis.

Use Cases:

Data Preparation for ETL: Facilitates schema discovery and prepares datasets for ETL workflows in AWS Glue.

Data Lake Enablement: Organizes and catalogs S3 data for use in data lakes, enabling efficient querying and analytics.

Schema Evolution: Automatically updates metadata when the underlying S3 data structure changes.

Integration with Redshift: Provides metadata to map and load S3 datasets into Redshift for advanced analytics.

Seamless Data Exploration: Supports Athena and QuickSight by generating a queryable structure for large datasets.

Optimizing Retail Operations through Data Warehousing

Anshu Reddy Dhamana - 018172141
Brendan Chao - 014142726
Sri Gopi Sarath Gode - 018191537
Srithareddy Devireddy - 018175456
Vinuthna Papana - 018176106



Insights and Lessons Learned

Importance of a Robust ETL Pipeline:

Effective data cleaning, transformation, and loading are critical for ensuring accurate and efficient data warehouse operations.

Value of Amazon Redshift:

Leveraging Redshift's scalability and query optimization capabilities enabled fast and complex analytical processing of retail data.

Role of Data Modeling:

Implementing a star schema improved query performance and clarity, highlighting the significance of well-structured data models in decision support systems.

Significance of Proper Data Management

Enhanced Data Accuracy and Consistency:

Proper data management ensures data is clean, consistent, and reliable, reducing errors and improving the quality of insights.

Improved Decision-Making:

Organized and well-maintained data allows for faster, more informed business decisions, driving strategic advantages.

Scalability and Efficiency:

Effective data management enables seamless handling of large datasets, optimizing performance for analytics and future growth.

Key Lessons from the Project

Importance of ETL Processes:

Clean and structured data is crucial for reliable analytics; robust ETL pipelines ensure accuracy and efficiency in transforming raw data.

Significance of Data Modeling:

Applying a star schema streamlined querying and analysis, emphasizing the value of proper database design in supporting analytical goals.

Role of Automation and Scalability:

Leveraging AWS services like Glue and Redshift highlighted how cloud-based tools enhance scalability, reduce manual effort, and accelerate data processing.

Future Recommendations for Data Warehousing Best Practices

Enhance Scalability with Cloud Services:

Continuously leverage cloud-native tools like AWS Redshift and Glue to ensure scalability, flexibility, and high performance as data volume grows.

Implement Robust Data Governance:

Establish clear policies for data quality, access control, and security, ensuring compliance and maintaining data integrity.

Integrate Advanced Analytics:

Incorporate machine learning and predictive analytics capabilities into the data warehouse to derive deeper insights and actionable business intelligence.

Final Presentation Slides

The image displays three screenshots illustrating the data pipeline and analysis for a retail store.

Amazon S3 Bucket Overview: Shows the `retail-store-data-bucket` containing five CSV files: `cleaned_customer_data.csv`, `customer_data.csv`, `modified_sales_data.csv`, and `product_data.csv`.

AWS Glue Database Configuration: Displays the `retail_metadata_db` database properties and tables. The database contains six tables: `cleaned_customer_dat`, `customer_data_csv`, `modified_sales_data_c`, `product_data_cv`, and two unnamed tables.

Data Visualization Dashboard: A dashboard showing various metrics and charts. Key figures include 8.80M Sum of totalprice, 6M Sum of quantity, 4314 Count of customerid, and 4017 Count of productid. It also includes a Sales trend over time chart, a Top Selling Products bar chart, a customer Distribution by country map, and a line chart for Total Quantity by Year and Month.

Thank You