

Project Architecture

- **Data Ingestion:** *Data is downloaded from the internet and stored in an Amazon S3 bucket.*

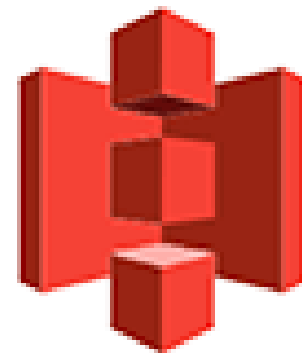
- **Data Cataloging:** *AWS Glue Crawler is used to infer the schema and create a data catalog from the raw data stored in S3.*

- **Data Querying:** *Amazon Athena is employed to run SQL queries on the data stored.*

- **Data Loading:** *AWS Glue ETL jobs load the cataloged data into an Amazon Redshift data warehouse.*

- **Data Visualization:** *Power BI is connected to the Redshift cluster for visualization and reporting.*

- **Automation:** *Apache Airflow is used to automate the entire ETL process.*



amazon
S3

☰

Amazon S3 > Buckets > kaggle---dataset

kaggle---dataset

Info

Objects

Properties

Permissions

Metrics

Management

Access Points

Objects (1)

Info

🔄

📄 Copy S3 URI

📄 Copy URL

📄 Download

🔗 Open

Delete

Actions ▼

Create folder

📄 Upload

Objects are the fundamental entities stored in Amazon S3. You can use [Amazon S3 inventory](#) to get a list of all objects in your bucket. For others to access your objects, you'll need to explicitly grant them permissions. [Learn more](#)

🔍 Find objects by prefix

< 1 > ⚙️

<input type="checkbox"/>	Name	Type	Last modified	Size	Storage class
<input type="checkbox"/>	📄 Telco_customer_churn.csv	csv	October 15, 2024, 01:40:42 (UTC+03:00)	1.7 MB	Standard



AWS Glue



AWS Glue

✕

Getting started

ETL jobs

Visual ETL

Notebooks

Job run monitoring

Data Catalog tables

Data connections

Workflows (orchestration)

▼ Data Catalog

Databases

Tables

Stream schema registries

Schemas

Connections

Crawlers

Classifiers

Catalog settings

► Data Integration and ETL

► Legacy pages

What's New

Documentation

Crawlers

A crawler connects to a data store, progresses through a prioritized list of classifiers to determine the schema for your data, and then creates metadata tables in your data catalog.

Crawlers (2)

Info

Last updated (UTC)

October 15, 2024 at 00:37:51

🔄

Action

Run

Create crawler

View and manage all available crawlers.

🔍 Filter crawlers

< 1 > ⚙️

<input type="checkbox"/>	Name	State	Schedule	Last run	Last run timest...	Log	Table changes fro...
<input type="checkbox"/>	redshift-crawler	🟢 Ready		🟢 Succeeded	October 14, 2024 ...	View log	1 created
<input type="checkbox"/>	s3-glue-crawler	🟢 Ready		🟢 Succeeded	October 14, 2024 ...	View log	1 updated

us-east-1.console.aws.amazon.com/gluestudio/home?region=us-east-1#/editor/job/s3_upload_to_redshift/runs

aws Services Search [Alt+S]

AWS Glue

Getting started

ETL jobs

Visual ETL

Notebooks

Job run monitoring

Data Catalog tables

Data connections

Workflows (orchestration)

Data Catalog

Databases

Tables

Stream schema registries

Schemas

Connections

Crawlers

Classifiers

Catalog settings

Data Integration and ETL

Legacy pages

What's New

Documentation

s3_upload_to_redshift

Last modified on 10/15/2024, 3:37:52 AM

Actions Save Run

Visual Script Job details Runs Data quality Schedules Version Control

Job runs (1/3) Info

Last updated (UTC) October 15, 2024 at 00:39:13

View details Stop job run Table View Card View

Filter job runs by property

	Run status	Retries	Start time (Local)	End time (Local)	Duration	Capacity (D...	Worker type	Glue versio
	Succeeded	0	10/15/2024 03:32:22	10/15/2024 03:34:00	1 m 23 s	2 DPUs	G.1X	4.0
	Succeeded	0	10/15/2024 03:26:37	10/15/2024 03:28:32	1 m 37 s	2 DPUs	G.1X	4.0
	Succeeded	0	10/15/2024 03:00:05	10/15/2024 03:03:24	2 m 59 s	2 DPUs	G.1X	4.0

Run details

Input arguments (10)

Continuous logs

Run insights

Metrics

Spark UI

Job name	Start time (Local)	Glue version	Last modified on (Local)
s3_upload_to_redshift	10/15/2024 03:32:22	4.0	10/15/2024 03:34:00
Id	End time (Local)	Worker type	Log group name
jr_ec61ed9e0e90cdf2fa63904a00ca2b76f19cf4cb8510f72cfde0b95bbc3cdf8	10/15/2024 03:34:00	G.1X	/aws-glue/jobs
Run status	Start-up time	Max capacity	Number of workers
Succeeded	15 seconds	2 DPUs	2
Retry attempt number	Execution time	Execution class	Timeout
Initial run	1 minute 23 seconds	Standard	2880 minutes
Trigger name	Security configuration	Cloudwatch logs	Usage profile

CloudShell Feedback

© 2024, Amazon Web Services, Inc. or its affiliates. Privacy Terms Cookie preferences



Amazon Athena

Data

Data source

AwsDataCatalog

Database

churn_data_db

Tables and views

Create

Filter tables and views

Tables (1)

kaggle__dataset

customerid string

count bigint

country string

state string

city string

zip code bigint

lat long string

latitude double

longitude double

Query 4

+

⌵

1 SELECT *

2 FROM "churn_data_db"."kaggle__dataset";

SQL

Ln 1, Col 9

Run again

Explain

Cancel

Clear

Create

Reuse query results

up to 60 minutes ago

Query results

Query stats

Completed

Time in queue: 105 ms

Run time: 792 ms

Data scanned: 1.66 MB

Results (7,043)

Copy

Download results

Search rows

#	customerid	count	country	state	city	zip code	lat long	latitude	longitude
1	3668-QPYBK	1	United States	California	Los Angeles	90003	"33.964131	33.964131	
2	9237-HQITU	1	United States	California	Los Angeles	90005	"34.059281	34.059281	
3	9305-CDSKC	1	United States	California	Los Angeles	90006	"34.048013	34.048013	

EC2 | us-east-1

EC2 Instance

kaggle---dataset

Crawlers

CloudWatch

Query editor

Policies | IAM

DAGs - Airflow

us-east-1.console.aws.amazon.com/athena/home?region=us-east-1#/query-editor/history/1e1fec74-3f51-41e1-b896-e337b4a891c9

aws Services Search [Alt+S]

N. Virginia Ziad_Ahmed

Amazon Athena > Query editor

Editor Recent queries Saved queries Settings

Workgroup primary

Data

Data source
AwsDataCatalog

Database
churn_data_db

Tables and views
Create

Filter tables and views

Tables (1)
kaggle__dataset
customerid string
count bigint
country string
state string
city string

Query 4

1 SELECT count(*)
2 FROM "churn_data_db"."kaggle__dataset";

SQL Ln 1, Col 1

Run Explain Cancel Clear Create

Reuse query results
up to 60 minutes ago

Query results Query stats

Completed Time in queue: 101 ms Run time: 489 ms Data scanned: 1.66 MB

Results (1)
Search rows

_col0
1 7043

CloudShell Feedback

© 2024, Amazon Web Services, Inc. or its affiliates. Privacy Terms Cookie preferences



Amazon Redshift

▼

Conn

×

EC2 Ir

×

EC2 Ir

×

Creat

×

Consc

×

S3 bu

×

Tables

×

Redsh

×

Editor

×

custo

×

+

—

📄

×

←→↻🔍us-east-1.console.aws.amazon.com/sqlworkbench/home?region=us-east-1#/client🔍☆📄Z⋮

awsServices🔍Search[Alt+S]

☰

Editor

📁

Queries

📖

Notebooks

📈

Charts

🕒

History

📅

Scheduled queries

🌙

⚙️

Redshift query editor v2

⊕ Create

⬅️ Load data

⏪

🔍 Filter resources

↻

▼🔗 customer-churn-redshift...

> 📁 awsdatalog

> 📁 dev

▼ 📁 public

▼ 📄 Tables1

custo...

> 👁 Views0

📄 customer_churn

↻

×

	Field	Type
A	customerid	character varying(255)
A	city	character varying(255)
#	zip_code	integer
A	gender	character varying(255)
A	senior_citizen	character

+

☰

Untitled 1

×

▶ Run

🔴

🔘 Limit 100

🔘 Explain

🔘 Isolated session

📄

customer-chu...

▼

dev

▼

📅 Schedule

💾

🔄

⋮

1 SELECT COUNT(*)

2 FROM customer_churn;

Result 1 (1)

Export

🔘 Chart

🔄

⌵

☐	count	
☐	7043	

Row 2, Col 21, Chr 37

Query ID 2375

Elapsed time: 148 ms

Total rows: 1

CloudShellFeedback

© 2024, Amazon Web Services, Inc. or its affiliates. PrivacyTermsCookie preferences

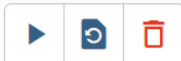


Apache
Airflow

DAG: customer_churn_dag

Schedule: @weekly ⓘ

Next Run ID: 2024-10-20, 00:00:00 UTC



2024-10-15

12:26:56 AM ⌚

All Run Types ▾

All Run States ▾

Clear Filters

Auto-refresh ☐

25 ▾

Press **shift** + **/** for Shortcuts

deferred failed queued removed restarting running scheduled shutdown skipped success up_for_reschedule up_for_retry upstream_failed no_status



DAG

Run

Task

customer_churn_dag / ▶ 2024-10-06, 00:00:00 UTC / tsk_is_glue_job_finish_running

Clear task

Mark state as... ▾

Filter DAG by task ▾

Details

Graph

Gantt

<> Code

Event Log

Logs

XCom

Task Duration

Layout:

Left -> Right ▾

tsk_glue_job_trigger

tsk_grab_glue_job_run_id

tsk_is_glue_job_finish_running

tsk_glue_job_trigger

■ success

PythonOperator

tsk_grab_glue_job_run_id

■ success

PythonOperator

tsk_is_glue_job_finish_running

■ success

GlueJobSensor



Power BI

FileHomeInsertModelingViewOptimizeHelp

Paste

Cut

Copy

Format painter

Clipboard

Get data

Excel workbook

OneLake data hub

SQL Server

Enter data

Data

Dataverse

Recent sources

Transform data

Refresh data

Queries

New visual

Text box

More visuals

Insert

New visual calculation

New measure

Quick measure

Calculations

Sensitivity

Sensitivity

Publish

Share

Copilot

Copilot

