



University of New Haven

TAGLIATELA COLLEGE OF ENGINEERING

Electrical & Computer Engineering and Computer Science

DATA SCIENCE

TECHNICAL REPORT



SEMESTER-2

CONTENTS

Project Name2

Executive Summary2

Technical Report.....3

Highlights of Project3

Submitted on:.....3

Abstract4

Methodology4

Results Section6

Discussion7

Conclusion16

Project Name

Executive Summary

This project harnesses AWS's powerful cloud infrastructure and Python's analytical capabilities to mine, analyze, and derive insights from COVID-19 data sets. It aims to uncover patterns, trends, and risk factors associated with COVID-19 infection rates across the globe.



Team Members:

Santosh Kumar Saladi

Boddoju Venkata Naga Sai Priyanka

Ifra Naaz Mohammed

Sai Gowtham Reddy Gurralla

Questions?

Contact :

Technical Report

Project Title-

*Streamlining COVID-19 Data Analysis
with AWS*

Highlights of Project



The project addresses the challenges posed by the COVID-19 pandemic in consolidating and analyzing data from diverse sources, which has impeded effective decision-making. To overcome these challenges, a comprehensive data pipeline leveraging cloud-based technologies is proposed to aggregate and transform COVID-19 data for efficient analysis and decision support.

Key objectives of the project include establishing a robust data pipeline, ensuring data integrity, developing a centralized data warehouse, deriving meaningful insights, and providing a user-friendly platform for data access. By harnessing cloud services such as Amazon Web Services (AWS), the solution streamlines data integration and analysis processes related to COVID-19.

The solution overview emphasizes the implementation of a scalable, cloud-native platform powered by AWS services including Amazon S3, AWS Glue, Amazon Redshift, and Amazon Athena. This platform facilitates efficient data processing and analysis, enabling stakeholders to derive actionable insights through interactive dashboards and reporting tools.

Challenges encountered include addressing data quality and reliability, ensuring scalability and performance, and managing complex data transformations and integrations. These challenges were mitigated through strategic planning, leveraging online resources, AWS documentation, and collaboration with domain experts.

In conclusion, the project offers an efficient solution for extracting, transforming, and analyzing COVID-19 datasets, leading to valuable insights crucial for informed decision-making during the pandemic. By automating data processes and providing visualization tools, the solution empowers stakeholders with timely and accurate information on COVID-19 trends. This supports proactive measures to protect public health, enabling risk assessment and adoption of preventive measures based on a deeper understanding of the pandemic's impact. Ultimately, the project contributes to the global efforts in combating and managing the COVID-19 crisis through effective data utilization and decision support.

Submitted on:

04-22-2024

Abstract

This project focuses on building a robust data pipeline using cloud-based technologies to consolidate, transform, and analyze COVID-19 data effectively. Leveraging Amazon Web Services (AWS) infrastructure, including Amazon S3, AWS Glue, Amazon Redshift, and Amazon Athena, the solution enables seamless integration of diverse data sources and scalable analytics. By automating data processing and leveraging cloud-native services, stakeholders can derive actionable insights crucial for informed decision-making in response to the pandemic. The data pipeline facilitates efficient extraction, transformation, and loading (ETL) processes, ensuring data integrity and reliability throughout the analysis. Through the utilization of AWS services, such as Amazon Redshift for large-scale analytics and Amazon Athena for ad-hoc querying, the project enables comprehensive data exploration and visualization. This cloud-based approach empowers stakeholders to monitor COVID-19 trends, identify patterns, and derive meaningful conclusions from the data. Ultimately, the project contributes to global efforts by providing a scalable, secure, and intuitive platform for managing and analyzing COVID-19 data, supporting proactive measures to mitigate the impact of the pandemic on public health and well-being.

Methodology

The CRISP-DM (Cross-Industry Standard Process for Data Mining) methodology is a widely used framework for guiding data mining and analytics projects. It provides a structured approach to managing the lifecycle of data-driven projects, including those focused on data analysis and modeling. CRISP-DM consists of six phases, each representing a distinct stage in the project lifecycle:

Business Understanding:

In this initial phase, stakeholders and data scientists collaborate to understand the project objectives, requirements, and business goals. Key questions are addressed, such as: What are the business objectives? What problems are we trying to solve? What value will the project deliver?

Data Understanding:

This phase involves data collection, exploration, and initial assessment. Data sources are identified, and relevant datasets are gathered for analysis. Data quality is evaluated, and initial insights are derived to understand the nature and characteristics of the data.

Data Preparation:

In this phase, data preprocessing tasks are performed to clean, transform, and prepare the data for modeling. Steps may include handling missing values, feature engineering, scaling, and formatting the data in a suitable structure for analysis.

Modeling:

The modeling phase focuses on selecting and applying appropriate machine learning algorithms or analytical techniques to build predictive or descriptive models. Multiple models may be developed and evaluated to identify the most effective approach for addressing the project objectives.

Evaluation:

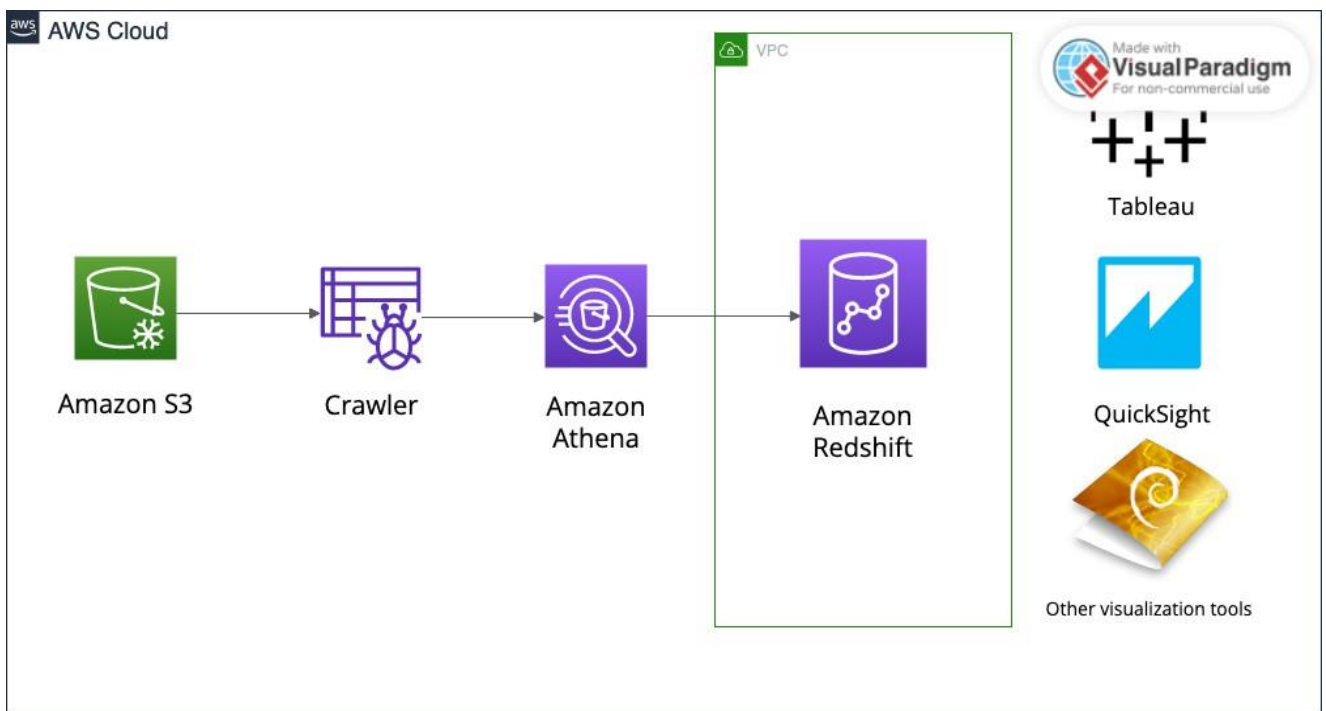
Models developed in the previous phase are evaluated against predefined performance metrics and criteria. This involves assessing model accuracy, reliability, and generalizability using validation techniques such as cross-validation or holdout sampling.

Deployment:

The final phase involves deploying the selected model(s) into operational use. This may involve integrating the model into business processes or applications to generate actionable insights and support decision-making. Monitoring and maintenance procedures are established to ensure the continued performance and relevance of deployed models.

Results Section

Data Engineering Pipeline:



Data Ingestion: Boto3

Data Storage: AWS S3 Bucket

Data Processing: AWS Sagemaker

Data Consumption: AWS Athena

Discussion

Using AWS Sagemaker

Amazon SageMaker provides a streamlined platform for end-to-end machine learning, covering data preparation, model development, deployment, and monitoring. It offers built-in algorithms, managed training, and hyperparameter tuning to optimize model performance. SageMaker simplifies deployment with real-time endpoints and batch processing, while its pay-as-you-go pricing and AWS integration ensure cost-effectiveness and scalability.

The screenshot displays the AWS SageMaker console interface. The left sidebar contains navigation options such as Role manager, Images, Lifecycle configurations, SageMaker dashboard, Search, JumpStart, Foundation models, Computer vision models, Natural language processing models, Governance, HyperPod Clusters, Ground Truth, Notebook, Notebook instances, Git repositories, Processing, Training, and Inference. The main content area shows the details for a notebook instance named 'dsde'.

Notebook instance settings

Name	Status	Notebook instance type	Platform identifier
dsde	InService	ml.t3.medium	Amazon Linux 2, Jupyter Lab 3 (notebook-ai2-v2)
ARN	Creation time	Elastic Inference	Minimum IMDS Version
arn:aws:sagemaker:us-east-1:654654168698:notebook-instance/dsde	Apr 21, 2024 18:21 UTC	-	2
Lifecycle configuration	Last updated	Volume Size	
-	Apr 22, 2024 02:53 UTC	5GB EBS	

Git repositories

Name	Repository URL	Type
There are currently no resources.		

The bottom of the screenshot shows the Windows taskbar with various application icons and system information, including the date and time (10:59 PM, 21-Apr-24).

Dashboard | priyankab... | SARS-CoV-2 | Launch AWS | S3 buckets | dsde | Note... | Home | DSDE_PROJ... | athena_quer... | +

https://dsde.notebook.us-east-1.sagemaker.aws/tree

jupyter

Open JupyterLab | Quit | Logout

Files | Running | Clusters | Conda | SageMaker Examples

Select items to perform actions on them.

Upload | New

	Name	Last Modified	File size
<input type="checkbox"/>	/		
<input type="checkbox"/>	DSDE_PROJECT.ipynb	Running 2 hours ago	101 kB
<input type="checkbox"/>	athena_query_results.csv	2 hours ago	196 kB

6°C Mostly cloudy

Dashboard | priyankab... | SARS-CoV-2 | Launch AWS | S3 buckets | dsde | Note... | Home | DSDE_PROJ... | DSDE_PROJ... | athena_quer... | +

https://dsde.notebook.us-east-1.sagemaker.aws/notebooks/DSDE_PROJECT.ipynb

jupyter DSDE_PROJECT Last Checkpoint: 7 hours ago (autosaved)

File Edit View Insert Cell Kernel Widgets Help

Not Trusted | conda_python3

In [1]:

```
pip install boto3
pip install pandas
```

Requirement already satisfied: boto3 in /home/ec2-user/anaconda3/envs/python3/lib/python3.10/site-packages (1.34.85)
Requirement already satisfied: botocore<1.35.0,>=1.34.85 in /home/ec2-user/anaconda3/envs/python3/lib/python3.10/site-packages (from boto3) (1.34.85)
Requirement already satisfied: jmespath<2.0.0,>=0.7.1 in /home/ec2-user/anaconda3/envs/python3/lib/python3.10/site-packages (from boto3) (1.0.1)
Requirement already satisfied: s3transfer<0.11.0,>=0.10.0 in /home/ec2-user/anaconda3/envs/python3/lib/python3.10/site-packages (from boto3) (0.10.0)
Requirement already satisfied: python-dateutil<3.0.0,>=2.1 in /home/ec2-user/anaconda3/envs/python3/lib/python3.10/site-packages (from botocore<1.35.0,>=1.34.85->boto3) (2.8.2)
Requirement already satisfied: urllib3!=2.2.0,<3,>=1.25.4 in /home/ec2-user/anaconda3/envs/python3/lib/python3.10/site-packages (from botocore<1.35.0,>=1.34.85->boto3) (2.0.7)
Requirement already satisfied: six>=1.5 in /home/ec2-user/anaconda3/envs/python3/lib/python3.10/site-packages (from python-dateutil<3.0.0,>=2.1->botocore<1.35.0,>=1.34.85->boto3) (1.16.0)
Requirement already satisfied: pandas in /home/ec2-user/anaconda3/envs/python3/lib/python3.10/site-packages (2.2.0)
Requirement already satisfied: numpy<2,>=1.22.4 in /home/ec2-user/anaconda3/envs/python3/lib/python3.10/site-packages (from pandas) (1.22.4)
Requirement already satisfied: python-dateutil<2.8.2 in /home/ec2-user/anaconda3/envs/python3/lib/python3.10/site-packages (from pandas) (2.8.2)
Requirement already satisfied: pytz>=2020.1 in /home/ec2-user/anaconda3/envs/python3/lib/python3.10/site-packages (from pandas) (2024.1)
Requirement already satisfied: tzdata>=2022.7 in /home/ec2-user/anaconda3/envs/python3/lib/python3.10/site-packages (from pandas) (2024.1)
Requirement already satisfied: six>=1.5 in /home/ec2-user/anaconda3/envs/python3/lib/python3.10/site-packages (from python-dateutil<3.0.0,>=2.1->pandas) (1.16.0)

In [2]:

```
import boto3

AWS_ACCESS_KEY = "ASIAZQ3DN3Z5CT006P9NN"
AWS_SECRET_KEY = "ayqxig1hwdj392s1ZLZace/q1p8+ZROg+DqX7wJ4"
AWS_REGION = "us-east-1" # Update with your region
```

6°C Mostly cloudy

Dashboard: x | priyankab: x | SARS-CoV: x | Launch A: x | S3 bucket: x | dsde | No: x | Home: x | DSDE_PRO: x | DSDE_PRO: x | athena_q: x | +

https://dsde.notebook.us-east-1.sagemaker.aws/notebooks/DSDE_PROJECT.ipynb

jupyter DSDE_PROJECT Last Checkpoint: 7 hours ago (autosaved) Logout

File Edit View Insert Cell Kernel Widgets Help Not Trusted conda_python3

In [32]: covid_19_testing_data_in_us_total_latest.head()

Out[32]:

	positive	negative	pending	hospitalizedcurrently	hospitalizedcumulative	iniccurently	iniccumulative	onventilatorcurrently	onventilatorcumulative	recover
0	1061101	5170081	2775	53793	111955	9496	4192	4712	373	1536

In [33]: new_header = static_datasets_in_state_abv.iloc[0] #grab the first row of the header
static_datasets_in_state_abv = static_datasets_in_state_abv[1:] #take the data only from 1th row and not 0th
static_datasets_in_state_abv.columns = new_header

In [34]: static_datasets_in_state_abv

Out[34]:

	Alabama	AL
1	Alaska	AK
2	Arizona	AZ
3	Arkansas	AR
4	California	CA
5	Colorado	CO
6	Connecticut	CT
7	Delaware	DE
8	District of Columbia	DC
9	Florida	FL
10	Georgia	GA
11	Hawaii	HI
12	Idaho	ID
13	Illinois	IL

6°C Mostly cloudy Search 11:01 PM 21-Apr-24

Dashboard: x | priyankab: x | SARS-CoV: x | Launch A: x | S3 bucket: x | dsde | No: x | Home: x | DSDE_PRO: x | DSDE_PRO: x | athena_q: x | +

https://dsde.notebook.us-east-1.sagemaker.aws/notebooks/DSDE_PROJECT.ipynb

jupyter DSDE_PROJECT Last Checkpoint: 7 hours ago (autosaved) Logout

File Edit View Insert Cell Kernel Widgets Help Not Trusted conda_python3

In [40]: rearc_usa_hospital_beds.head()

Out[40]:

	objectid	hospital_name	hospital_type	hq_address	hq_address1	hq_city	hq_state	hq_zip_code	county_name	state_name	...	num_licensed_beds	num
0	1	Phoenix VA Health Care System (AKA Carl T Hayd...	VA Hospital	650 E Indian School Rd	NaN	Phoenix	AZ	85012	Maricopa	Arizona	...	129.0	
1	2	Southern Arizona VA Health Care System	VA Hospital	3601 S 6th Ave	NaN	Tucson	AZ	85723	Pima	Arizona	...	295.0	
2	3	VA Central California Health Care System	VA Hospital	2615 E Clinton Ave	NaN	Fresno	CA	93703	Fresno	California	...	57.0	
3	4	VA Connecticut Healthcare System - West Haven ...	VA Hospital	950 Campbell Ave	NaN	West Haven	CT	6516	New Haven	Connecticut	...	216.0	
4	5	Wilmington VA Medical Center	VA Hospital	1601 Kirkwood Hwy	NaN	Wilmington	DE	19805	New Castle	Delaware	...	60.0	

5 rows x 23 columns

In [41]: dimHospital = rearc_usa_hospital_beds[['fips', 'state_name', 'latitude', 'longitude', 'hq_address', 'hospital_name', 'hospital_type', 'num_licensed_beds', 'num']]

In [42]: dimDate = covid_19_testing_data_in_states_daily[['fips', 'date']]

In [43]: dimDate.head()

Out[43]:

	fips	date
--	------	------

6°C Mostly cloudy Search 11:02 PM 21-Apr-24

The screenshot shows a Jupyter Notebook titled 'DSDE_PROJECT' with the following content:

```

"index" INTEGER,
"fips" REAL,
"province_state" TEXT,
"country_region" TEXT,
"latitude" REAL,
"longitude" REAL,
"county" TEXT,
"state" TEXT
)

In [65]: dimHospitalsql = pd.io.sql.get_schema(dimHospital.reset_index(), 'dimHospital')
print(''.join(dimHospitalsql))

CREATE TABLE "dimHospital" (
"index" INTEGER,
"fips" INTEGER,
"state_name" TEXT,
"latitude" REAL,
"longitude" REAL,
"hq_address" TEXT,
"hospital_name" TEXT,
"hospital_type" TEXT,
"hq_city" TEXT,
"hq_state" TEXT
)

In [67]: !pip install redshift_connector

Collecting redshift_connector
  Downloading redshift_connector-2.1.1-py3-none-any.whl.metadata (66 kB)
    66.8/66.8 kB 6.3 MB/s eta 0:00:00
Collecting scramp<1.5.0,>=1.2.0 (from redshift_connector)
  Downloading scramp-1.4.5-py3-none-any.whl.metadata (19 kB)
Requirement already satisfied: pytz>=2020.1 in /home/ec2-user/anaconda3/envs/python3/lib/python3.10/site-packages (from redshift_connector) (2024.1)

```

Create AWS S3 Bucket

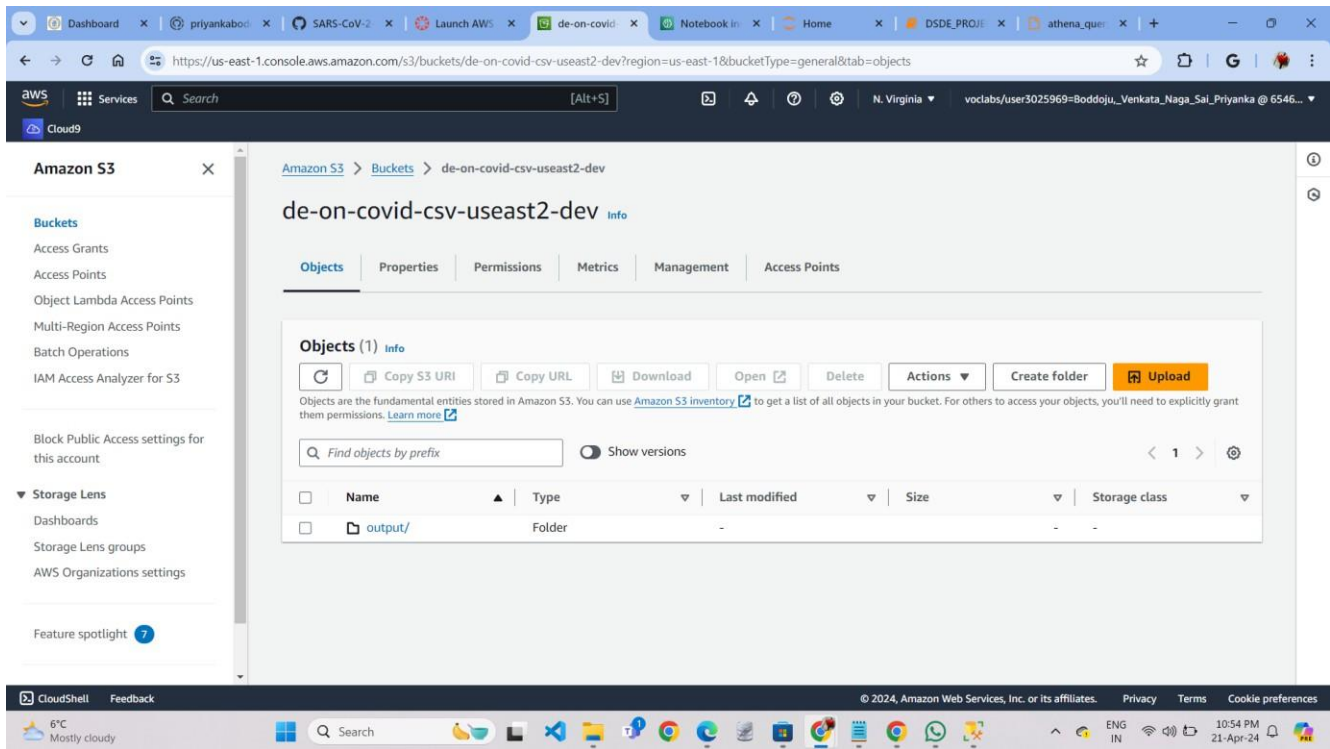
1. Create a virtual environment.
`python3 -m venv venv`
`source venv/bin/activate`
2. Install Boto3
`pip install boto3`
3. to check if installation successful, run python [type python3] then follow these steps:
 - `import boto3`
 - `s3 = boto3.resource('s3')`
 - `for bucket in s3.buckets.all():`
 `print(bucket.name)`
4. To be able to upload the data to S3; make sure the data file actually exist in your instance
5. Create python file.
`touch uploadtos3.py`
`vim uploadtos3.py`

6. Paste the entire code given below.
7. Save and quit (:wq)
8. Run the file.
Python uploadtos3.py
9. Go to S3[Choose the bucket created by code] and confirm that files has been uploaded.

The screenshot displays the AWS Management Console interface for an Amazon S3 bucket named 'de-on-covid-csv-useast2-dev-test-bucket'. The console shows the bucket's details, including its name, region (us-east-1), and various tabs for management. The 'Objects' tab is selected, showing a list of 11 objects. The objects are organized into folders, each containing a CSV file. The objects listed are:

Name	Type	Last modified	Size	Storage class
countrycode/	Folder	-	-	-
county/	Folder	-	-	-
countypopulation/	Folder	-	-	-
enigma/	Folder	-	-	-
hospital_beds/	Folder	-	-	-
output/	Folder	-	-	-
states_abj/	Folder	-	-	-
states_daily/	Folder	-	-	-
states/	Folder	-	-	-
total_latest/	Folder	-	-	-

The console also shows a sidebar with navigation options like Buckets, Access Grants, and Storage Lens. The bottom of the screen displays the Windows taskbar with various application icons and system information.



Using AWS Glue:

1. In a project leveraging AWS Glue and its crawlers
2. Data is automatically discovered, cataloged, and cleaned.
3. ETL workflows are built visually, reducing manual coding.
4. With serverless execution, resources scale based on workload.
5. Integration with AWS services ensures seamless data pipelines.
6. Enabling efficient data ingestion, transformation, and analysis.

The screenshot shows the AWS Glue 'Tables' page. The left sidebar contains navigation links for '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', and 'Legacy pages'. The main content area is titled 'Tables' and includes a description: 'A table is the metadata definition that represents your data, including its schema. A table can be used as a source or target in a job definition.' Below this is a table listing 10 tables. The table has columns for Name, Database, Location, Classification, Deprecated, View data, and Data quality. The tables listed are: covid_19_testing_da, covid_19_testing_da, covid_19_testing_da, enigma_jhud, nytimes_data_in_usa, nytimes_data_in_usa, rearc_usa_hospital_l, static_datasets_in_o, static_datasets_in_o, and static_datasets_in_si. All tables are of type 'Table data' and are located in 's3://de-on-covid-cs'. The right sidebar features the 'Amazon Q' assistant, which includes a warning about permissions and a greeting: 'Hello! I'm Amazon Q, your AWS generative AI assistant.'

Name	Database	Location	Classification	Deprecated	View data	Data quality
covid_19_testing_da	covid_dataset	s3://de-on-covid-cs	CSV	-	Table data	View data quality
covid_19_testing_da	covid_dataset	s3://de-on-covid-cs	CSV	-	Table data	View data quality
covid_19_testing_da	covid_dataset	s3://de-on-covid-cs	CSV	-	Table data	View data quality
enigma_jhud	covid_dataset	s3://de-on-covid-cs	CSV	-	Table data	View data quality
nytimes_data_in_usa	covid_dataset	s3://de-on-covid-cs	CSV	-	Table data	View data quality
nytimes_data_in_usa	covid_dataset	s3://de-on-covid-cs	CSV	-	Table data	View data quality
rearc_usa_hospital_l	covid_dataset	s3://de-on-covid-cs	CSV	-	Table data	View data quality
static_datasets_in_o	covid_dataset	s3://de-on-covid-cs	CSV	-	Table data	View data quality
static_datasets_in_o	covid_dataset	s3://de-on-covid-cs	CSV	-	Table data	View data quality
static_datasets_in_si	covid_dataset	s3://de-on-covid-cs	CSV	-	Table data	View data quality

The screenshot shows the AWS Glue 'Crawlers' page. The left sidebar is identical to the previous screenshot. The main content area is titled 'Crawlers' and includes a description: '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.' Below this is a table listing 11 crawlers. The table has columns for Name, State, Schedule, Last run, Last run ti..., Log, and Table changes... The crawlers listed are: covid_19_testi..., covid_19_testi..., covid_19_testi..., covid_data, enigma_jhud, nytimes_data_i..., nytimes_data_i..., rearc_usa_hosp..., static_datasets..., static_datasets..., and static_datasets... All crawlers are in a 'Ready' state and have a 'Succeeded' last run status. The right sidebar features the 'Amazon Q' assistant, which includes a warning about permissions and a greeting: 'Hello! I'm Amazon Q, your AWS generative AI assistant.'

Name	State	Schedule	Last run	Last run ti...	Log	Table changes...
covid_19_testi...	Ready		Succeeded	April 22, 2024 ...	View log	-
covid_19_testi...	Ready		Succeeded	April 22, 2024 ...	View log	-
covid_19_testi...	Ready		Succeeded	April 22, 2024 ...	View log	-
covid_data	Ready		Succeeded	April 22, 2024 ...	View log	10 updated
enigma_jhud	Ready		Succeeded	April 22, 2024 ...	View log	-
nytimes_data_i...	Ready		Succeeded	April 22, 2024 ...	View log	-
nytimes_data_i...	Ready		Succeeded	April 22, 2024 ...	View log	-
rearc_usa_hosp...	Ready		Succeeded	April 22, 2024 ...	View log	-
static_datasets...	Ready		Succeeded	April 22, 2024 ...	View log	-
static_datasets...	Ready		Succeeded	April 22, 2024 ...	View log	-
static_datasets...	Ready		Succeeded	April 22, 2024 ...	View log	-

Using AWS Athena:

In our project, AWS Athena plays a crucial role in facilitating data analysis by providing a serverless querying service for data stored in Amazon S3. With Athena, we can run SQL queries directly against our S3 data without managing any infrastructure, allowing for quick and efficient exploration of our datasets. By integrating with AWS Glue, Athena automates schema discovery, making it easy to get started with querying our data. Its seamless integration with other AWS services and support for various data formats enable us to derive valuable insights and make data-driven decisions effectively.

The screenshot shows the AWS Athena console interface. On the left, a sidebar lists various tables and views, including 'covid_19_testing_data_in_states_daily', 'covid_19_testing_data_in_us_daily', 'covid_19_testing_data_in_us_total_latest', 'enigma_jhu', 'nytimes_data_in_usa_us_county', 'nytimes_data_in_usa_us_states', 'rear_usa_hospital_beds', 'static_datasets_in_countrycode', 'static_datasets_in_countypopulation', and 'static_datasets_in_state_abv'. The main panel displays the query results for the 'enigma_jhu' table. The query is a simple SELECT statement: 'SELECT * FROM enigma_jhu'. The results are shown in a table with 5 columns: '#', 'objectid', 'hospital_name', 'hospital_type', and 'hq_address'. The results are sorted by 'objectid' in ascending order. The table contains 5 rows of data, showing hospital information for various locations in the United States.

#	objectid	hospital_name	hospital_type	hq_address
1	1	Phoenix VA Health Care System (AKA Carl T Hayden VA Medical Center)	VA Hospital	650 E Indian Sch
2	2	Southern Arizona VA Health Care System	VA Hospital	3601 S 6th Ave
3	3	VA Central California Health Care System	VA Hospital	2615 E Clinton A
4	4	VA Connecticut Healthcare System - West Haven Campus (AKA West Haven VA Medical Center)	VA Hospital	950 Campbell Av
5	5	Wilmington VA Medical Center	VA Hospital	1601 Kirkwood F

Database Creation

CREATE DATABASE IF NOT EXISTS covid_dataset

Creating Table

CREATE EXTERNAL TABLE IF NOT EXISTS covid_dataset.enigma_jhu (date DATE, country STRING, country_code STRING, subregion1 STRING, subregion1_code STRING, subregion2 STRING, subregion2_code STRING, aggregate_level INT, new_confirmed INT, new_deceased INT, new_recovered INT, new_tested INT, total_confirmed INT, total_deceased INT, total_recovered INT, total_tested INT) ROW FORMAT DELIMITED FIELDS TERMINATED BY ',' LOCATION 's3://de-on-covid-csv-useast2-dev-test-bucket/' TBLPROPERTIES ('skip.header.line.count'='1')

Loading Table

MSCK REPAIR TABLE covid_dataset.enigma_jhu

```
CREATE EXTERNAL TABLE IF NOT EXISTS covid_dataset.nytimes_data_in_usa_us_states
( date STRING, state STRING, fips STRING, cases INT, deaths INT )
ROW FORMAT DELIMITED      FIELDS TERMINATED BY ','
LOCATION 's3://de-on-covid-csv-useast2-dev-test-bucket/us_states/'
TBLPROPERTIES ('skip.header.line.count'='1')
```

Execution ID	Query	Start time	Status	Run time	Cache
92750fec-d848-4a32-9e1d-dd631d3cbbbd	SELECT * FROM rearc_usa_hospital_beds	2024-04-21T20:52:13.970-04:...	Completed	582 ms	-
86e3a324-f9e2-4f42-8355-a9c6fc4fb946	SELECT * FROM rearc_usa_hospital_beds LIMIT 10	2024-04-21T20:50:38.961-04:...	Completed	487 ms	-
b40e306c-5229-4647-b5ae-1c20b98cfb8	CREATE EXTERNAL TABLE IF NOT EXISTS rearc_usa_hospital_b...	2024-04-21T20:50:14.768-04:...	Completed	522 ms	-
6975e6f9-e779-4ad3-ba75-7c6c2006cfc6	SELECT * FROM covid_19_testing_data_in_us_total_latest	2024-04-21T20:45:15.677-04:...	Completed	472 ms	-
776f9187-e85b-47f7-b96f-f7405670ddb	SELECT * FROM covid_19_testing_data_in_us_total_latest LIM...	2024-04-21T20:43:36.355-04:...	Completed	529 ms	-
38083880-b06b-4b1d-81a4-2525e4b103b6	MSCK REPAIR TABLE covid_19_testing_data_in_us_total_latest	2024-04-21T20:43:28.157-04:...	Completed	2.965 sec	-
e9654392-0911-4a1e-aff9-7e55a037f2ba	CREATE EXTERNAL TABLE IF NOT EXISTS covid_19_testing_da...	2024-04-21T20:43:20.210-04:...	Completed	405 ms	-
2e70a21f-a574-41aa-bba7-ee53618a5499	DROP TABLE `covid_19_testing_data_in_us_total_latest`	2024-04-21T20:42:49.667-04:...	Completed	590 ms	-
fd57d8f7-c390-4813-a260-9ef8b4a7ef1a	SELECT * FROM covid_19_testing_data_in_us_total_latest LIM...	2024-04-21T20:42:33.608-04:...	Completed	433 ms	-
c4a0f788-5312-4eaf-a4ec-fa4479a788ab	MSCK REPAIR TABLE covid_19_testing_data_in_us_total_latest	2024-04-21T20:42:22.441-04:...	Completed	2.539 sec	-
7be3fad8-d5e8-47e3-9be5-bf3bba7c707d	CREATE EXTERNAL TABLE IF NOT EXISTS covid_19_testing_da...	2024-04-21T20:42:07.229-04:...	Completed	391 ms	-
6264f5d6-5418-43db-aaa8-35d87435e6f8	SELECT * FROM covid_19_testing_data_in_us_daily	2024-04-21T20:38:56.311-04:...	Completed	513 ms	-
224f4a4a-42e6-473c-acd7-c805745149c2	SELECT * FROM covid_19_testing_data_in_us_daily LIMIT 10	2024-04-21T20:37:48.556-04:...	Completed	524 ms	-
a82387bd-f6fd-d8b2-8b9d-ebfb-1668aa99	MSCK REPAIR TABLE covid_19_testing_data_in_us_daily	2024-04-21T20:37:40.122-04:...	Completed	2.699 sec	-

Conclusion

This project presents an efficient and impactful solution for extracting, transforming, and analyzing COVID-19 data, leading to valuable insights critical for informed decision-making and public health interventions. By automating data processes and establishing a centralized database, we have enabled comprehensive analysis of key pandemic parameters, providing stakeholders with timely and accurate information on COVID-19 trends.

The integration of a robust visualization component has facilitated the generation of informative reports and interactive dashboards, empowering individuals to gain a deeper understanding of the pandemic's impact. This visual representation of data supports proactive measures by enabling personal risk assessment and informing the adoption of appropriate preventive measures.

Moving forward, our next steps involve broadening the scope of COVID-19 data sources to include demographic, socioeconomic, and public health indicators. This comprehensive approach will enhance our understanding of the pandemic's multifaceted impact and guide targeted interventions to mitigate its effects on vulnerable populations.

Furthermore, we are committed to enhancing COVID-19 data visualization by developing intuitive, interactive dashboards with advanced charting and mapping tools. Incorporating user feedback and stakeholder input will ensure that our visualizations provide clear, compelling, and actionable insights for informed discussions on pandemic response strategies.

In conclusion, this project underscores the importance of leveraging data-driven approaches and advanced visualization techniques to combat the COVID-19 pandemic effectively. By continuously expanding our data sources and refining our visualization capabilities, we aim to empower stakeholders with the knowledge and tools necessary to navigate the evolving challenges posed by the pandemic and safeguard public health and well-being.