# Using AWS Cloud Computing Services
## Sonam Dawani

This project is to present the use of Amazon Web Services. The following Amazon Web Services are used in this project:

- S3
- AWS Glue
  - Crawler
  - ETL Job
- AWS Lambda
- AWS CloudWatch
- AWS Athena
- AWS SageMaker

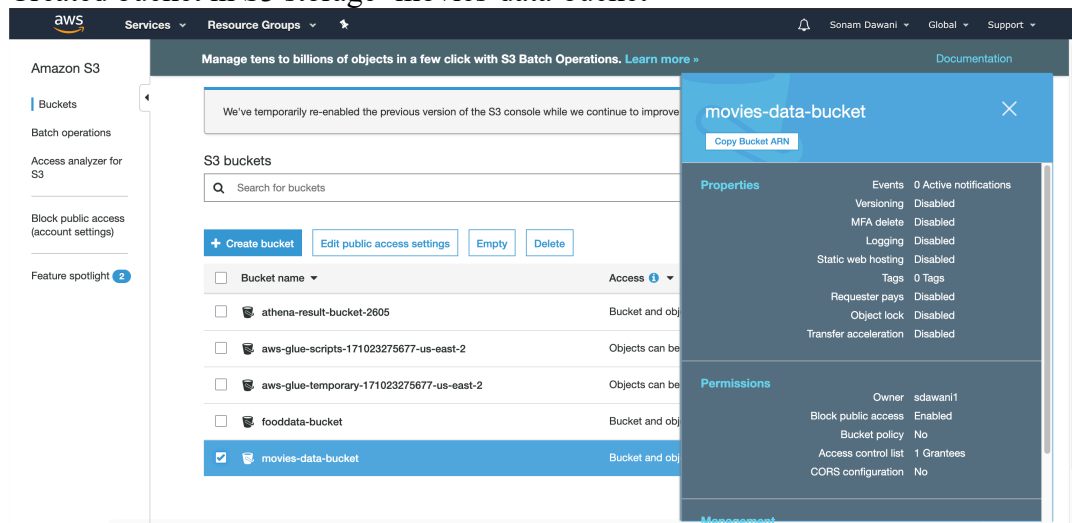The code and scripts implemented in this project are uploaded on Github.
https://github.com/SONAMDAWANI/UsingAWS_DataAnalysis

For this project simple data is used: movies.csv

movies

| | A | B | C | D | E | F | G | H |
|---|---|---|---|---|---|---|---|---|
| 1 | **Film** | **Genre** | **Lead Studio** | **Audience score %** | **Profitability** | **Rotten Tomatoes %** | **Worldwide Gross** | **Year** |
| 2 | Zack and Miri Make a Porno | Romance | The Weinstein Company | 70 | 1.747541667 | 64 | $41.94 | 2008 |
| 3 | Youth in Revolt | Comedy | The Weinstein Company | 52 | 1.09 | 68 | $19.62 | 2010 |
| 4 | You Will Meet a Tall Dark Stranger | Comedy | Independent | 35 | 1.211818182 | 43 | $26.66 | 2010 |
| 5 | When in Rome | Comedy | Disney | 44 | 0 | 15 | $43.04 | 2010 |
| 6 | What Happens in Vegas | Comedy | Fox | 72 | 6.267647029 | 28 | $219.37 | 2008 |
| 7 | Water For Elephants | Drama | 20th Century Fox | 72 | 3.081421053 | 60 | $117.09 | 2011 |
| 8 | WALL-E | Animation | Disney | 89 | 2.896019067 | 96 | $521.28 | 2008 |
| 9 | Waitress | Romance | Independent | 67 | 11.0897415 | 89 | $22.18 | 2007 |
| 10 | Waiting For Forever | Romance | Independent | 53 | 0.005 | 6 | $0.03 | 2011 |
| 11 | Valentine's Day | Comedy | Warner Bros. | 54 | 4.184038462 | 17 | $217.57 | 2010 |

## Using S3 Buckets
Created bucket in S3 storage 'movies-data-bucket'

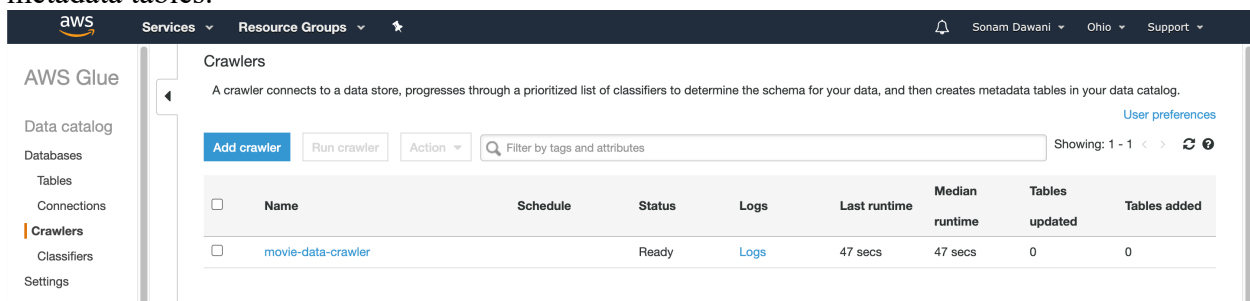Created read and write folders in the 'movies-data-bucket' bucket.



The read folder contains the original csv file, whereas write data is initially empty.
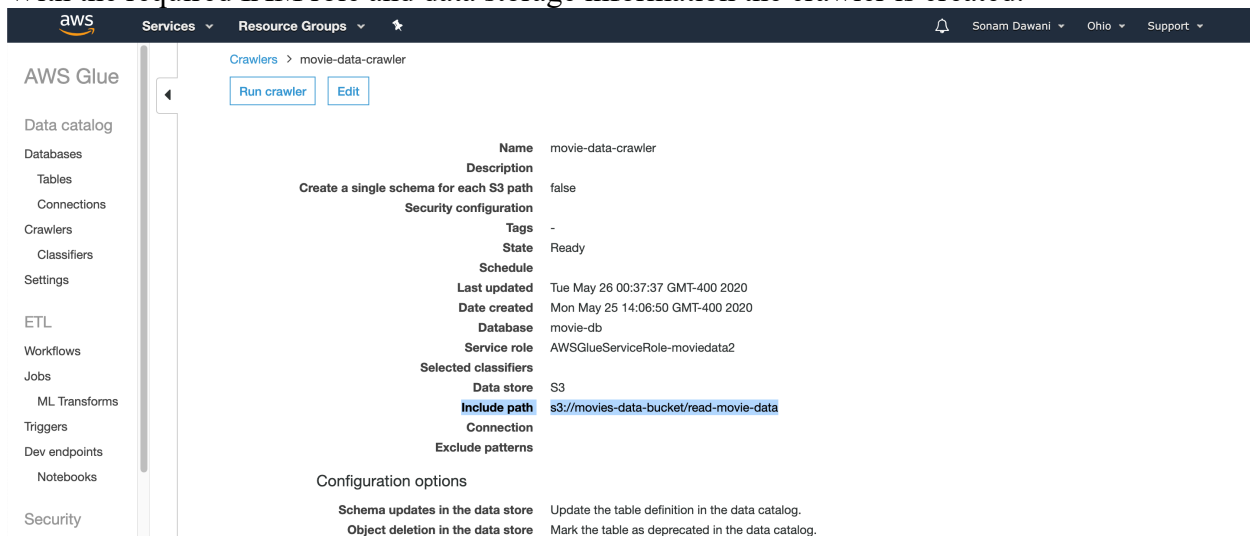The result after ETL transformation will be stored in the write folder.

## Using AWS Glue

- ### Crawler

First created a crawler which connects to the data storage (here the S3 bucket) and then creates metadata tables.



With the required IAM role and data storage information the crawler is created.

The execution of crawler above created the specified database 'movies-db' with 'read_movie_data' table in it.



- **ETL Job**

Created run On-Demand ETL job.



The script written for this job extract the 'read_movie_data' in a dataframe, performs transformation steps and load the resultant dataframe as CSV in the 'write-movie-data' folder of the S3 bucket.

In the transformation steps the movies are group by decades and mean rating of each decade is calculated.



The transformed data:

run-1590446627509-part-r-00000

| decade | movie_count | rating mean |
|--------|-------------|-------------|
| 2010 | 34 | 42.73529411764706 |
| 2000 | 43 | 49.83720930232558 |

## Using AWS Lambda

AWS Lambda is used there to automate the process of running ETL job after the crawler execution is completed.
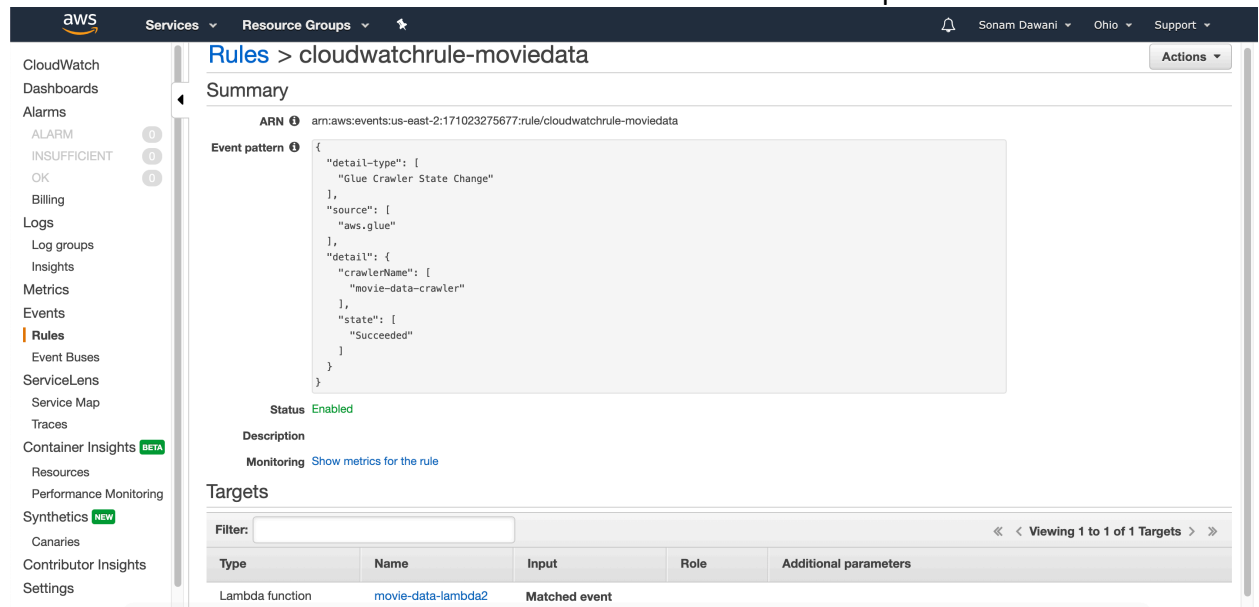
This is done in two steps:
- Creating the lambda function to start the ETL job
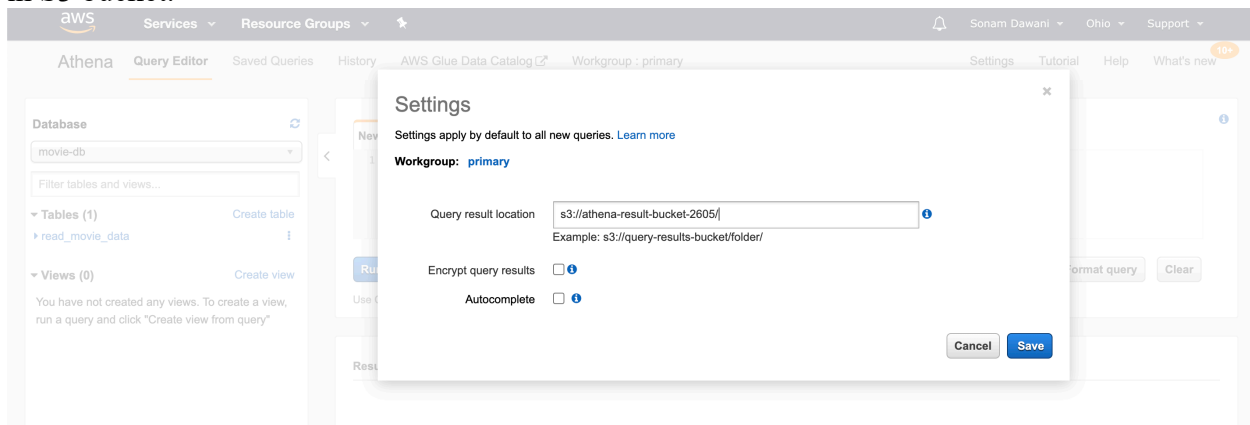


- Adding rule in CloudWatch

Created a rule in CloudWatch for the event of crawler. This rule is set to execute the lambda function once the execution of 'movie-data-crawler' is completed.



The steps above required to have **IAM roles** with some specific policies. Some of the IAM roles required to add Inline policies.

## Using AWS Athena

For analysis purpose the database table was queries using AWS Athena. The results were stored in S3 bucket.



Sample of querying data:

## Using AWS SageMaker

For SageMaker activity I used different dataset.
The dataset used here is food data containing recipe name, time to cook, nutritional values, ingredients etc. This data is collected from Food.com



The notebook contains:
- **Data reading from S3 bucket**
- Preprocessing
- Data Exploration
- Adding cuisine data using transfer learning
- Clustering of recipes using PCA and k-means
- Regression to predict nutritional value using Gradient Boosting
- Apriori / Market Basket Analysis

The github repo contains the PDF version of the notebook.
https://github.com/SONAMDAWANI/UsingAWS_DataAnalysis