

AWS Glue DataBrew

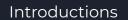














Overview on DataBrew



Demo



Q&A

Speaker



Sanchit Jain Lead Architect - AWS at Quantiphi **AWS APN Ambassador**

FOLLOW ME







Why do we need DataBrew?

- We all are fascinated with various analytical stuff like fancy visualization, BI report, Machine learning output
- But do we know, 50% of the time is just spend in cleansing, understanding & exploring your data
- We all like the elegant outcome but really no one want really to invest the time required to clean the data due to combustion process, multiple rounds of back and forth, and time consuming, etc.
- With this growing ask from the industry, AWS launched Glue DataBrew service in Nov 2020 with an idea to purely focus on business use-case rather than preparation work required for the same.



What is Glue DataBrew?

AWS Glue DataBrew is a visual data preparation tool that makes it easy for data analysts and data scientists to prepare data with an interactive, point-and-click visual interface without writing code.





NEED FOR DATABREW

"Upto 50% of data analysis time is spent on preparing data"

Time Consuming

- Multi-step process to extract, clean, normalize & load data at scale
- The right tools for the right persona must be integrated

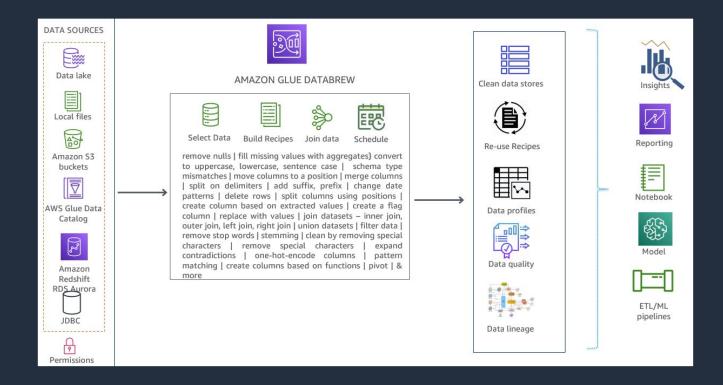
Expensive

- Costly user licenses & siloed tools that cause rework
- Often requires moving large amount of data into silos

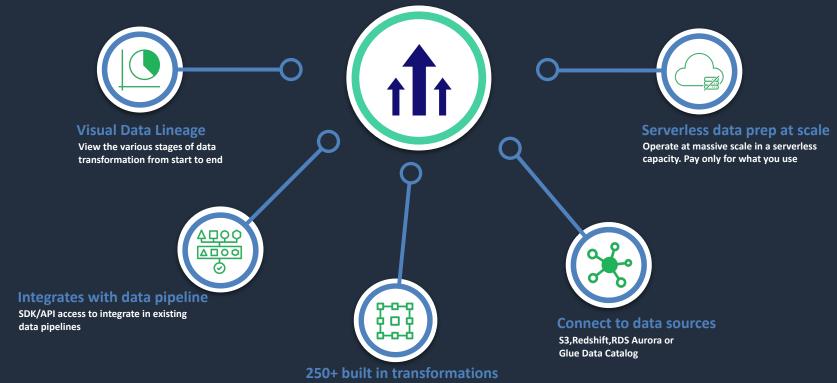
Manual

- Needs a lot of code-based heavy lifting to work at scale
- Hard to operationalize & build repeatable workflows

Overview of DataBrew



Advantages of DataBrew



Join,tokenize,split,merge,extract, remove, Group,pivot, normalize, label encode or more

Core Concepts

Dataset - a set of data—rows or records that are divided into columns or fields



Project - The interactive data preparation workspace in DataBrew



Recipe - is a set of instructions or steps for data that you want DataBrew to act on



Job - The process of running these instructions when you make the recipe is called a job.



Data Lineage - DataBrew tracks your data in a visual interface to determine its origin, called a data lineage.



Data Profile - A report which summarizes your existing shape of your data



Demo



Prerequisites



Create Project & Dataset



Exploring & Preparing Dataset



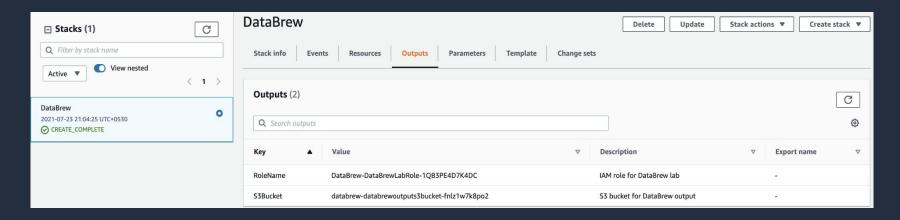
Creating a DataBrew job



Viewing data lineage

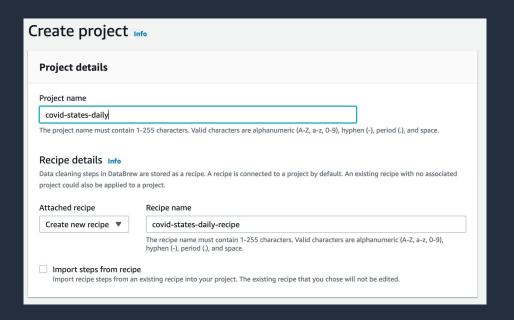
Prerequisites

- Download the <u>cloudformation template</u> from this link and Deploy it
- Download the <u>dataset</u> from this link and upload it to the S3 bucket
- Once the Cloudformation stack is deployed successfully please capture the values for RoleName and S3Bucket details



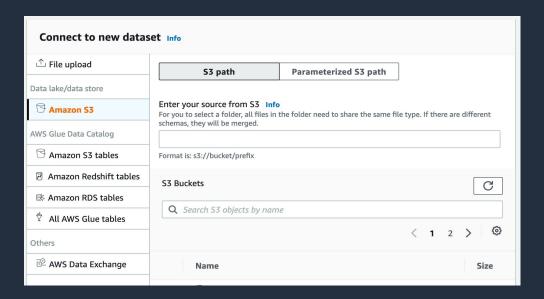
Creating a project

- Navigate to the AWS Glue DataBrew service
- On the DataBrew console, select Projects
- Click Create project
- In the Project details section, enter covid-states-daily as the project name



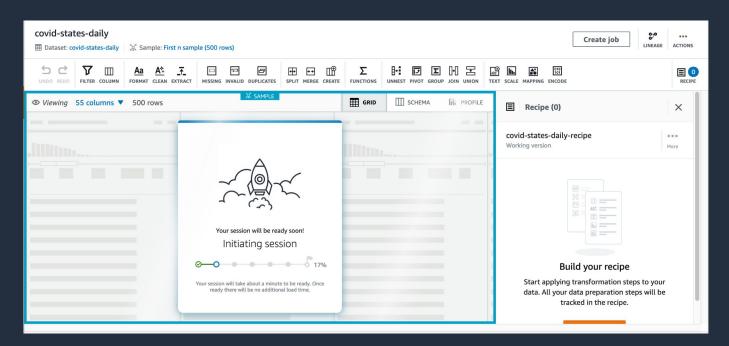
Creating a dataset

- In the Select a dataset section, select New dataset and enter covid-states-daily-stats
- In the Connect to a new dataset section, select Amazon S3 under "Data lake/data store" and Enter the S3 path. Leave the default configuration values
- In the Permissions section, select the role DataBrew-DataBrewLabRole--xxxxx from the drop-down list
- Click Create project

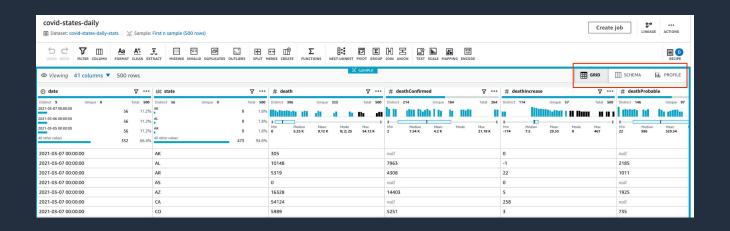


Creating a dataset

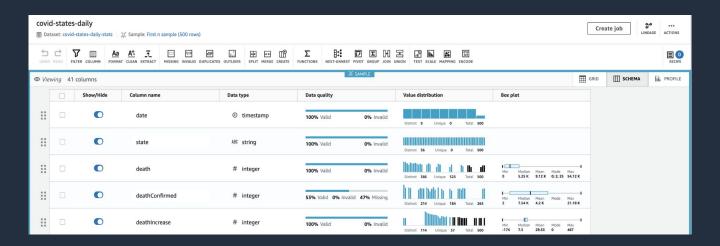
• Glue DataBrew will create the project, this may take a few minutes.



- Grid view When the project has been created, you will be presented with the Grid view. This is the default view, where a sample of the data is shown in tabular format. The Grid view shows
 - Columns in the dataset
 - Data type of each column
 - Summary of the range of values that have been found
 - Statistical distribution for numerical columns

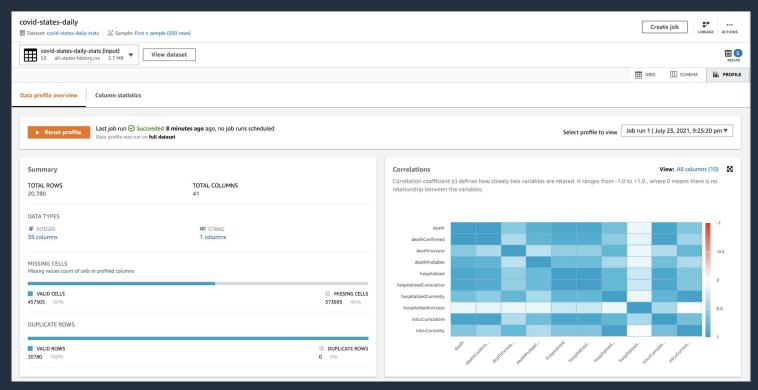


- Schema view The Schema view shows the schema that has been inferred from the dataset. In schema view, you can see statistics about the data values in each column. In the Schema view, you can
 - Select the checkbox next to a column to view the summary of statistics for the column values
 - Show/Hide columns
 - Rename columns
 - Change the data type of columns
 - Rearrange the column order by dragging and dropping the columns

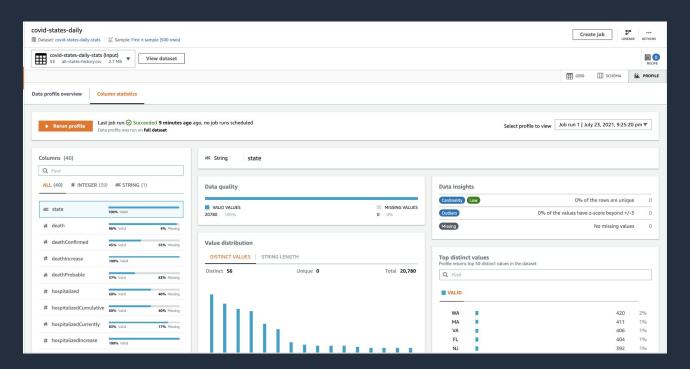


- Profile view In the Profile view, you can run a data profile job to examine and collect statistical summaries about the data. A data profile is an assessment in terms of structure, content, relationships, and derivation.
 - Job Name
 - Click on Run data profile
 - In the job details and job run sample panels, leave the default values
 - Job Output Setting
 - In the Job output settings section, select the S3 bucket with the name DataBrew-DataBrewLabRole--xxxxx and a folder name (eg. data-profile)
 - In the Permissions section, select the IAM role with the name databrew-lab-DataBrewLabRole-xxxxx
 - Leave all other settings as the default values
 - Click Create and run job
 - When the profile job has successfully completed, click on View data profile under Jobs from the menu on the left hand side of the DataBrew console

Profiling Output

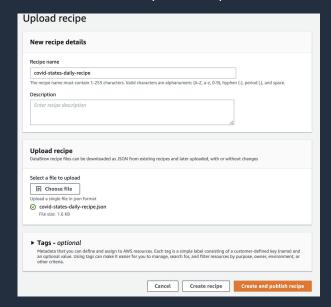


• Click on the Column statistics tab to view a column-by-column breakdown of the data values.



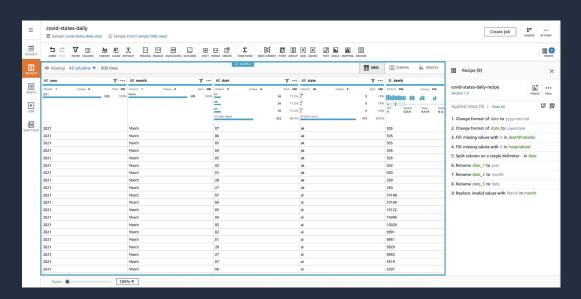
Preparing the dataset

- In this section, we will apply the different transformations to the dataset.
 - Rename columns
 - Change the data type of columns
 - Filled with the most frequent value
- Download the recipe from this <u>link</u>
- Select on Recipe from the menu on the left-hand side of the DataBrew console and click Upload Recipe
- Provide the below details
 - o Recipe Name
 - Upload Recipe json script downloaded under Step 1
 - Select Create and publish recipe



Preparing the dataset

- Now let's apply this recipe, click the project we configured now and the right side, click Import recipe
- Under Import recipe, select the recipe we configured and click Next
- Select the Append option from the right side and click Next
- Now let's validate the recipe and wait for all validation to be successful
- Once the validation is successful, click Import
- Now we will be back to the project screen and with the recipe implied on the dataset

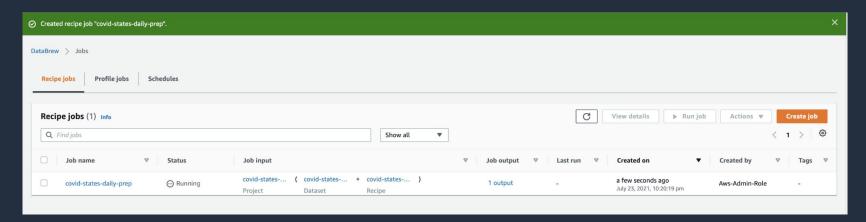


Creating a DataBrew job

- Click on Jobs from the menu on the left hand side of the DataBrew console
 - On the Recipe jobs tab, click on Create job
 - Enter covid-states-daily-prep for the job name
 - Select Create a recipe job
 - Select the covid-states-daily dataset
 - Select the 'covid-states-daily-recipe'
 - o In the Job output settings section, enter the S3 location s3://bucket-name/job-outputs/.
 - In the Permissions section, select the role DataBrew-DataBrewLabRole--xxxxx
 - Click Create and run job

Creating a DataBrew job

• DataBrew job is created and the job status is Running

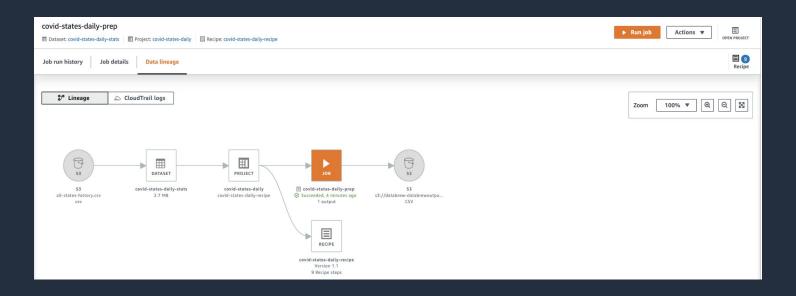


• Click on the link to the job output, and verify that the output files in the S3 bucket

Name	•	Type ▽	Last modified	▽	Size ▽	Storage class	▽
covid-states-daily-prep_23Jul2021_1627059065181_part00000.csv		csv	July 23, 2021, 22:21:19 (UTC+05:30)		2.6 MB	Standard	

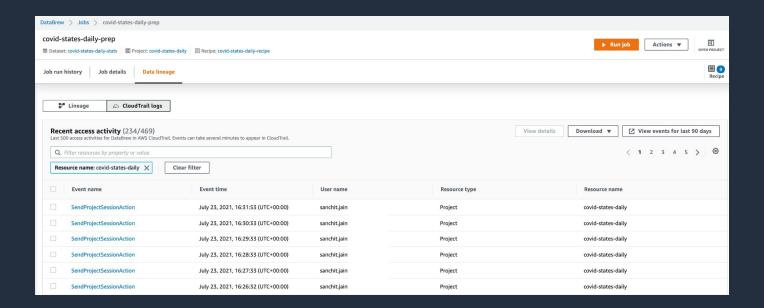
Viewing data lineage

- In DataBrew, navigate back to the covid-states-daily project
- Click on Lineage at the top right



Viewing data lineage

Select CloudTrail logs to view all the action on this dataset.



THANK YOU

