Prerequisites : Data catalog should be created in Glue. We have already done that in GlueDemo based on sample movies data (movies_sample.json).

1. Go to Athena. Click on 3 vertical dots besides table name and click on preview table. It will generate a query and will execute and show the result.



2. You can also download the result by clicking on highlighted yellow button in below snapshot.



3. Run below query to see how many movies were released year wise.

*SELECT year, COUNT(title) from "aug12-movies-database"."aug12moviestabledata"*
*GROUP BY year*
*ORDER BY COUNT(title);*

4. Modify above query to see percentage of movies released year wise.

<span style="color:red">SELECT year, COUNT(title), (COUNT(title) * 100.0/(SELECT COUNT(*) from "aug12-movies-database"."aug12moviestabledata")) from "aug12-movies-database"."aug12moviestabledata"
GROUP BY year
ORDER BY COUNT(title)</span>

*ASSIGNMENT:*

5. Let's now look at NestedJSON.json file.Create a new bucket with the NestedJSON.json file uploaded in it.

6. Create a crawler and name it as *nested-crawler.* While creating the crawler, give the database name as *nesteddb.*
7. Once the crawler is created, Run the crawler.
8. Crawler will create database *nesteddb* and within *nesteddb*, there will be a table created.
9. After clicking on table name, you will get below schema:

Schema

Showing: 1 - 1 of 1  ‹  ›

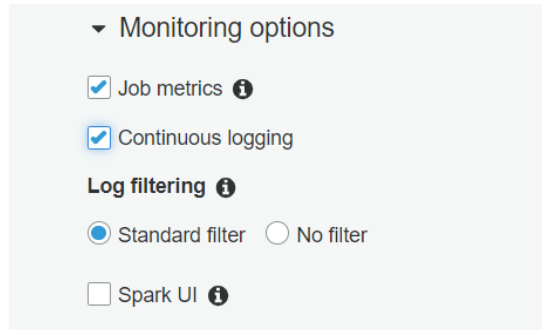| | Column name | Data type | Partition key | Comment |
|---|---|---|---|---|
| 1 | player | struct | | |

This table cannot be queried by Athena as when we click on struct in above snapshot, we get below details:

player schema details

```
        username:string
    ▼   characteristics:struct
            race:string
            class:string
            subclass:string
            power:int
            playercountry:string
    ▼   arsenal:struct
      ▼   kinetic:struct
              name:string
              type:string
              power:int
              element:string
      ▼   energy:struct
              name:string
              type:string
              power:int
```

10. We will now have to create Glue Job, which will convert NestedJSON file into a csv file.
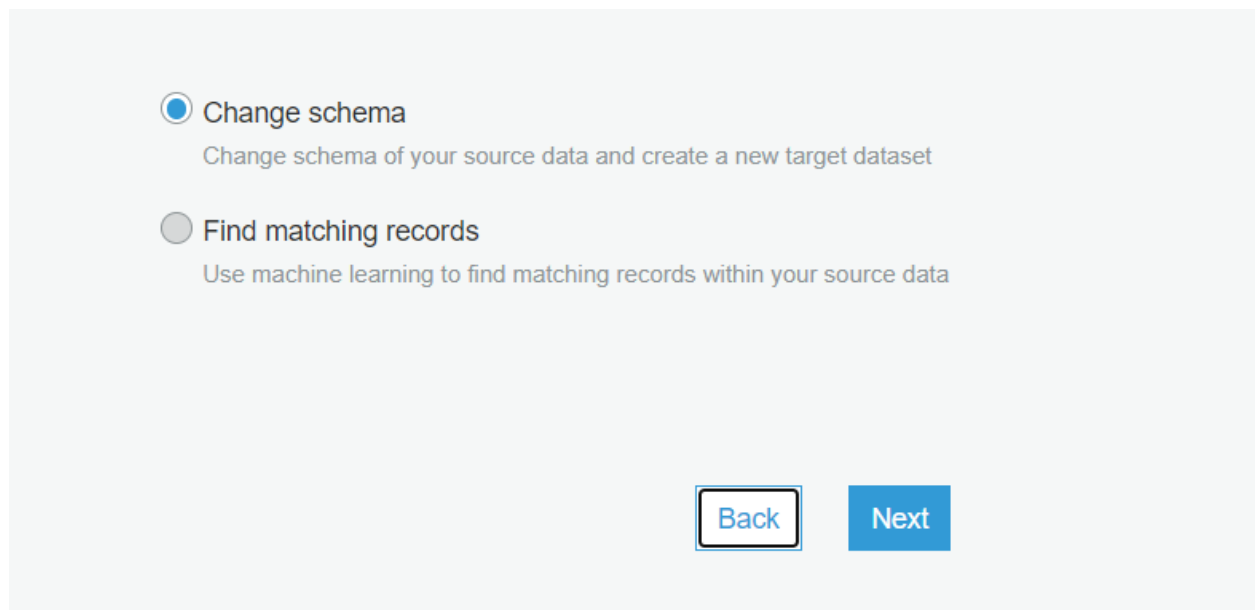11. Go to Glue jobs and click on Add Job. Name the job as flatten-data-job.

12. Expand Monitoring Options and select as shown below:



13. Select Change schema as shown below:Click Next.



14. Choose Create tables in your data target and select options as shown in below snapshot. Ensure that *transformeddata* folder is created in S3 bucket. Click Next

15. Output Schema definition will be created as shown below:



16. Click on Save Job and Edit Script.
17. Save Job and Run Job.
18. Go to transformeddata folder in S3 Bucket and download the zip file. Extract the zip and view the file data.
19. In order to query this data using Athena, create a crawler to construct Data Catalog from the data in compressed form available in *transformeddata* folder.
20. Click on Add crawler and give the name as csv-crawler. Select transformeddata folder on S3 as include path.
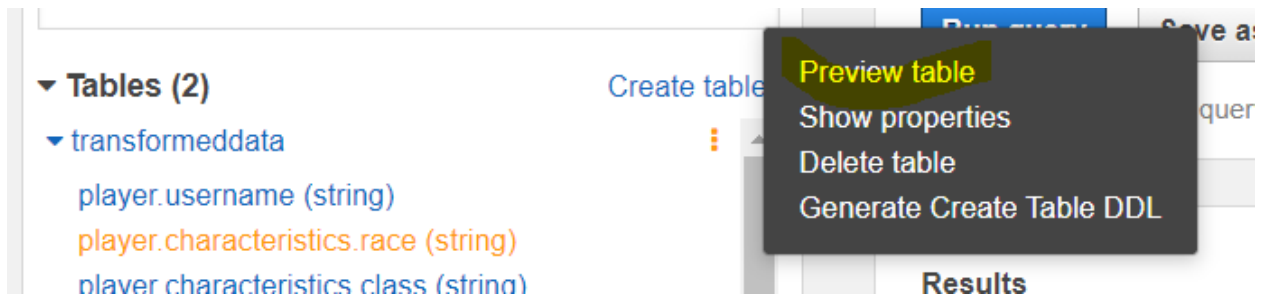
21. Choose the same database as created in earlier step i.e. nesteddb
22. Run the crawler.
23. Go to nesteddb and view tables. Click on transformeddata table and view schema.
24. Go to Athena.
25. Expand transformeddata table.

26. Click on Preview Table as shown below.

Tables (2)    Create table

▼ transformeddata

  player.username (string)
  player.characteristics.race (string)
  player.characteristics.class (string)

Preview table
Show properties
Delete table
Generate Create Table DDL

Results

27. You can now see the result as shown below :

```
1 SELECT * FROM "nesteddb"."transformeddata" limit 10;
```

Run query    Save as    Create ∨    (Run time: 0.41 seconds, Data scanned: 0.36 KB)        Format query    Clear

Use Ctrl + Enter to run query, Ctrl + Space to autocomplete        Athena engine version 1    Release versions

Results

| player.username | player.characteristics.race | player.characteristics.class | player.characteristics.subclass | player.characteristics.power | player.character |
|---|---|---|---|---|---|
| 1  user1 | Human | Warlock | Dawnblade | 300 | USA |

28. Navigate to S3 and all query results are written to S3. Observe the folder structure. Data will be available in csv format. Download the csv file and have a look at query result.