

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

job_aws_glue_lab_4

Last modified on 21/05/2023, 00:21:46

Try new UI

Actions

Save

Run

Successfully updated job

Successfully updated job job_aws_glue_lab_4. To run the job choose the Run Job button.

Script

Job details

Runs

Data quality New

Schedules

Version Control

Script Info

```
1 import sys
2 from awsglue.transforms import *
3 from awsglue.utils import getResolvedOptions
4 from pyspark.context import SparkContext
5 from awsglue.context import GlueContext
6 from awsglue.job import Job
7 from pyspark.sql import functions as F
8
9 ## @params: [JOB_NAME]
10 args = getResolvedOptions(sys.argv, ['job_aws_glue_lab_4','S3_INPUT_PATH','S3_TARGET_PATH'])
11 sc = SparkContext()
12 glueContext = GlueContext(sc)
13 spark = glueContext.spark_session
14 job = Job(glueContext)
15 job.init(args['job_aws_glue_lab_4'], args)
16
17 # Ler o arquivo CSV
18 df = glueContext.create_dynamic_frame.from_options(
19     "s3",
20     {
21         "paths": [args['S3_INPUT_PATH']],
22     },
23     "csv",
24     {
25         "withHeader": True,
26         "separator": ",",
27     },
28 )
29
30 # Imprimir o schema do dataframe
31 df.printSchema()
32
33 # Alterar a caixa dos valores da coluna nome para MAIÚSCULO
34 df = df.toDF()
35 df = df.withColumn('nome', F.upper(F.col('nome')))
36 df = DynamicFrame.fromDF(df, glueContext, "df")
37
38 # Imprimir a contagem de Linhas presentes no dataframe
39 print("Número de linhas: ", df.count())
40
41 # Imprimir a contagem de nomes, agrupando os dados do dataframe pelas colunas ano e sexo.
42 df.groupBy("ano", "sexo").agg(F.sum("total").alias("contagem")).show()
43
44 # Ordene os dados de modo que o ano mais recente apareça como primeiro registro do dataframe.
45 df = df.orderBy(df.ano.desc())
46
47 # Apresentar qual foi o nome feminino com mais registros e em que ano ocorreu.
48 nome_feminino_mais_registrado = df.filter(df.sexo == 'F').groupBy('nome', 'ano').sum('total').orderBy('sum(total)',
49     ascending=False).first()
50 print(f"Nome feminino mais registrado: {nome_feminino_mais_registrado['nome']} no ano
51     {nome_feminino_mais_registrado['ano']}")
52
53 # Apresentar qual foi o nome masculino com mais registros e em que ano ocorreu.
54 nome_masculino_mais_registrado = df.filter(df.sexo == 'M').groupBy('nome', 'ano').sum('total').orderBy('sum(total)',
55     ascending=False).first()
56 print(f"Nome masculino mais registrado: {nome_masculino_mais_registrado['nome']} no ano
57     {nome_masculino_mais_registrado['ano']}")
58
59 # Apresentar o total de registros (masculinos e femininos) para cada ano presente no dataframe.
60 df.groupBy('ano').sum('total').show()
61
62 # Considere apenas as primeiras 10 linhas, ordenadas pelo ano, de forma crescente.
63 df = df.orderBy(df.ano.asc())
64 df = DynamicFrame.fromDF(df.limit(10), glueContext, "df")
65
66 # Escrever o conteúdo do dataframe com os valores de nome em maiúsculo no S3.
67 glueContext.write_dynamic_frame.from_options(
68     frame = df,
69     connection_type = "s3",
70     connection_options = {
71         "path": args['S3_TARGET_PATH'] + "/frequencia_registro_nomes_eua",
72         "partitionKeys": ["sexo", "ano"]
73     },
74     format = "json"
```

Connections
Crawlers
Classifiers

```
71 )  
72 job.commit()  
73
```

PythonLn 67, Col 22

Errors: 0Warnings: 0