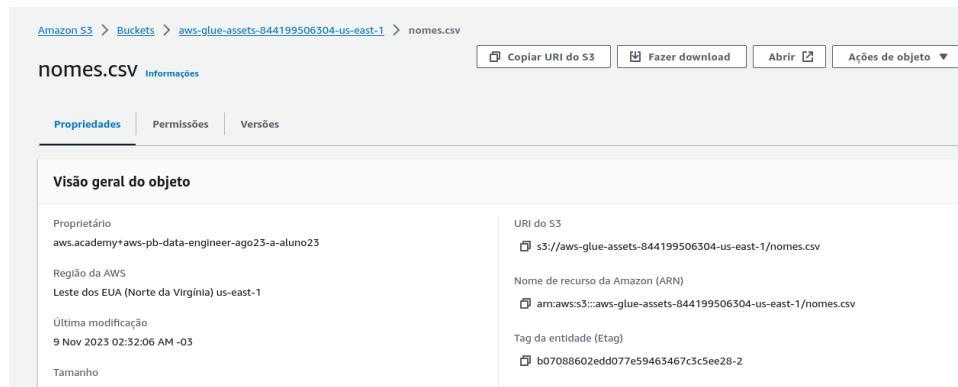


## LabAwsGlue

Objetivo: construir um job Glue.

Atenção: Esse material segue as orientações da documentação da documentação “LabAWSGlueSteps” presente nesse diretório.

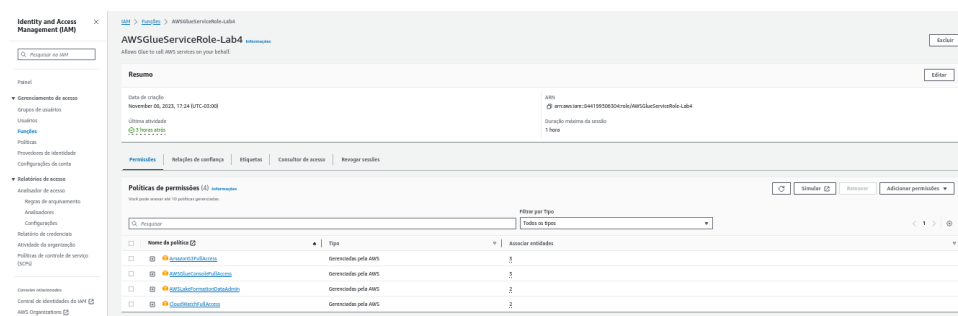
- Os dados de origem foram exportados para dentro do S3:



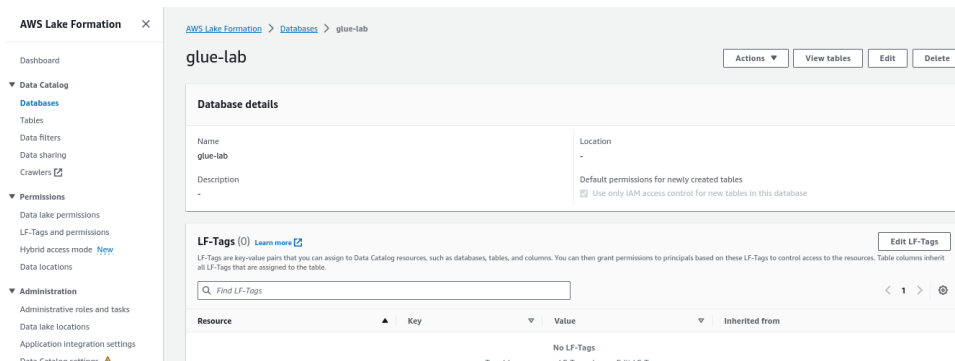
- Para a configuração da conta para utilizar o AWS glue foi necessário criar um usuário IAM novo:



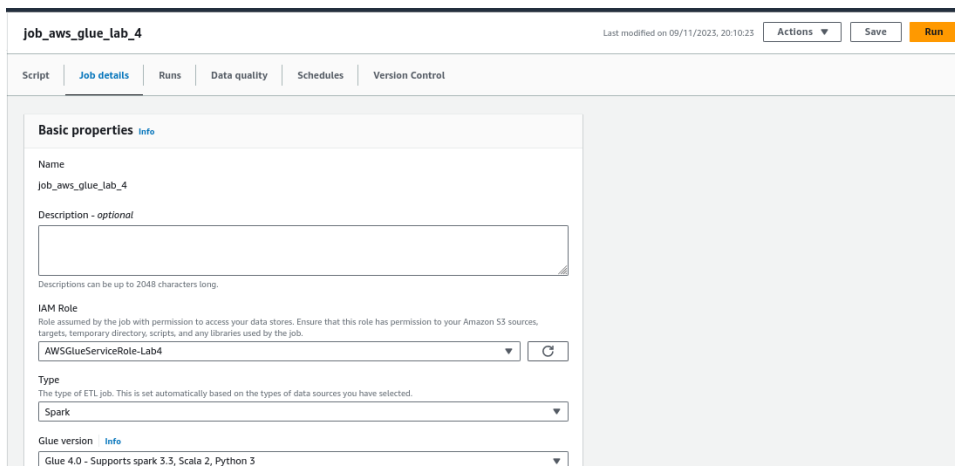
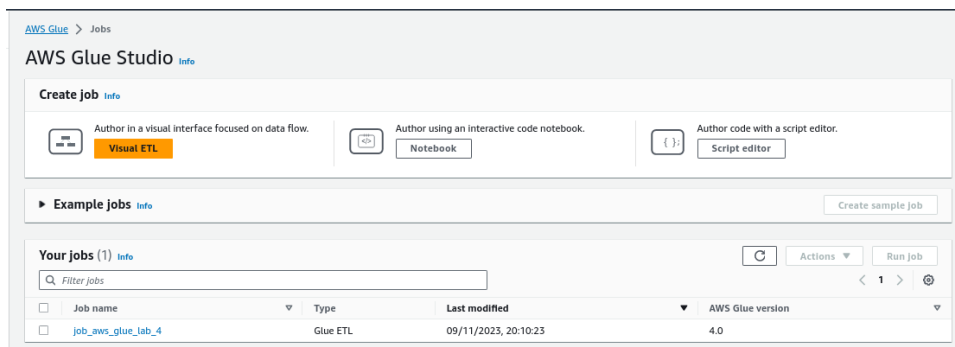
- A IAM Role para os jobs do aws glue foi criada:



- A configuração das permissões no AWS Lake Formation foi feita:

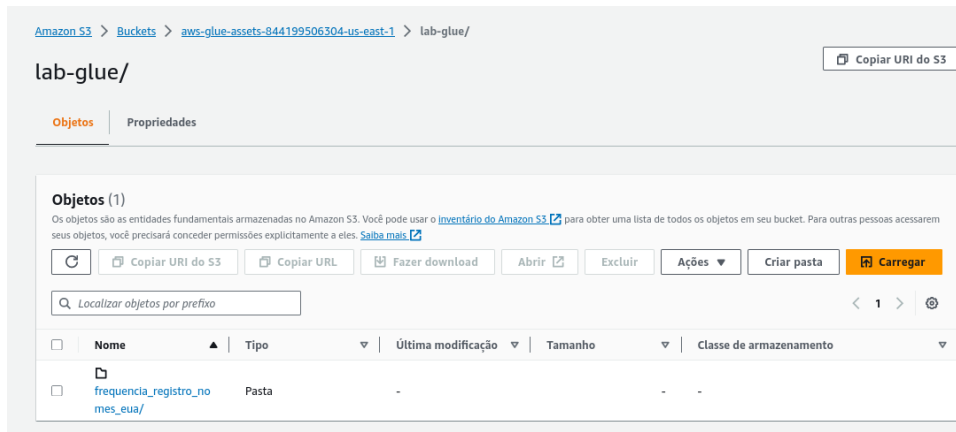


- O novo job no AWS Glue foi criado:



**Atenção:** os caminhos de entrada e saída já foram configurados como pode ser visto na imagem acima.

- O caminho de destino foi configurado:

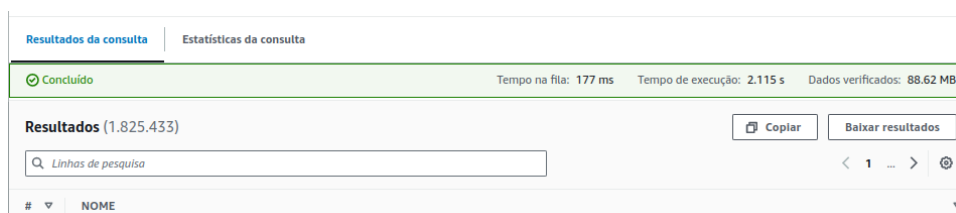


Agora vamos construir um job Glue nos moldes dos exemplos anteriores. Seguem os passos para desenvolver:

- Ler o arquivo nomes.csv no S3 (lembre-se de realizar upload do arquivo antes).
- Imprima o schema do dataframe gerado no passo anterior.



- Escrever o código necessário para alterar a caixa dos valores da coluna nome para
- MAIÚSCULO.



- Imprimir a contagem de linhas presentes no dataframe.



- Imprimir a contagem de nomes, agrupando os dados do dataframe pelas colunas ano e sexo.
- Ordene os dados de modo que o ano mais recente apareça como primeiro registro do dataframe.

2023-11-09T20:12:18.671-03:00

ano	sexo	count
2014	M	13977
2014	F	19067
2013	F	19191
2013	M	14012
2012	F	19468
2012	M	14216
2011	F	19540
2011	M	14329
2010	M	14241
2010	F	19800
2009	F	20165
2009	M	14519
2008	F	20439
2008	M	14606
2007	F	20918
2007	M	14383
2006	F	20043
2006	M	14026
2005	M	13358
2005	F	19175

only showing top 20 rows

Copiar

- Apresentar qual foi o nome feminino com mais registros e em que ano ocorreu.

2023-11-09T20:12:22.774-03:00

NOME	sexo	total	ano
Linda	F	99680	1947

Copiar

- Apresentar qual foi o nome masculino com mais registros e em que ano ocorreu.

2023-11-09T20:12:26.446-03:00

NOME	sexo	total	ano
James	M	94755	1947

Copiar

- Apresentar o total de registros (masculinos e femininos) para cada ano presente no dataframe.
- Considere apenas as primeiras 10 linhas, ordenadas pelo ano, de forma crescente.

2023-11-09T20:12:30.921-03:00

ano	registros
1880	201484
1881	192699
1882	221538
1883	216950
1884	243467
1885	240855
1886	255319
1887	247396
1888	299480
1889	288950

only showing top 10 rows

Copiar

- Escrever o conteúdo do dataframe com os valores de nome em maiúsculo no S3.



```
dynamicFrame = glueContext.create_dynamic_frame.from_options(  
    connection_type="s3",  
    connection_options={"paths": [source_file]},  
    format="csv",  
    format_options={  
        "withHeader": True,  
        # "optimizePerformance": True,  
    },  
)
```

```
spark_df = dynamicFrame.toDF()
```

```
spark_df = spark_df.withColumnRenamed("nome", "NOME")
```

```
spark_df = spark_df.withColumn("ano", spark_df["ano"].cast(IntegerType()))
```

```
spark_df = spark_df.withColumn("total", spark_df["total"].cast(IntegerType()))
```

```
dynamicFrame.printSchema()
```

```
row_count = spark_df.count()
```

```
print("Numero de linhas no DataFrame: ", row_count)
```

```
result=spark_df.groupBy("ano",  
    "sexo").agg(count("NOME").alias("count")).orderBy(col("ano").desc())  
result.show()
```

```
spark_df = spark_df.orderBy(col("ano").desc())
```

```
f_data = spark_df.filter(col("sexo") == "F").orderBy(col("total").desc()).limit(1)  
f_data.show()
```

```
m_data = spark_df.filter(col("sexo") == "M").orderBy(col("total").desc()).limit(1)
m_data.show()
```

```
agg_result = spark_df.groupBy("ano").agg(sum("total").alias("registros")).orderBy("ano")
agg_result.show(10)
```

```
spark_df = spark_df.withColumn("NOME", upper(spark_df["NOME"]))
```

```
dynamicFrame = DynamicFrame.fromDF(spark_df.repartition(1), glueContext, "dynamic_frame")
```

```
glueContext.write_dynamic_frame.from_options(
    frame=dynamicFrame,
    connection_type="s3",
    connection_options={"path": target_path},
    format="json",
)
```

```
job.commit()
```

- Criando crawler:

[AWS Glue](#) > Crawlers

### Crawlers

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.

**Crawlers (1)** [Info](#)

Last updated (UTC)  
November 10, 2023 at 03:35:55

Action ▾

Run

Create crawler

View and manage all available crawlers.

Q

Filter crawlers

< 1 > ⚙

<input type="checkbox"/>	Name ▾	State ▾	Schedule	Last run ▾	Last run timestamp ▾	Log	Table changes from L...
<input type="checkbox"/>	<a href="#">FrequenciaRegistroNo...</a>	Ready		Succeeded	November 9, 2023 at ...	<a href="#">View log</a>	1 updated

[AWS Glue](#) > [Crawlers](#) > FrequenciaRegistroNomesCrawle

**FrequenciaRegistroNomesCrawle**

Last updated (UTC)  
November 10, 2023 at 03:36:07

Run crawler

Edit

Delete

**Crawler properties**

Name  
FrequenciaRegistroNomesCrawle

Description  
-

Maximum table threshold  
-

IAM role  
[AWSGlueServiceRole-Lab4](#)

Security configuration  
-

Database  
glue-lab

Lake Formation configuration  
-

State  
READY

Table prefix  
-

► Advanced settings

[Crawler runs](#) | [Schedule](#) | [Data sources](#) | [Classifiers](#) | [Tags](#)

**Crawler runs (24)**

Stop run

[View CloudWatch logs](#)

[View run details](#)

The list of crawler runs for this crawler.

Q

Filter data

📅

Filter by a date and time range

< 1 2 > ⚙

	Start time (UTC) ▲	End time (UTC) ▾	Current/last duration ▾	Status ▾	DPU hours ▾	Table changes
<input type="radio"/>	November 9, 2023 at 23:13:12	November 9, 2023 at 23:14:06	54 s	Completed	0.079	1 table change, 0 partition changes
-	-	-	-	-	-	-