

Abstract geometric lines in the top-left corner of the slide, consisting of several overlapping, irregular polygons and lines in a dark gray color.

DESAFIO DE ANÁLISE DE DADOS: FILMES E SÉRIES

Jade Marinho Torres

ÍNDICE

Introdução

Principais objetivos

Etapas do Desafio

Dashboard

Áreas de Conhecimento

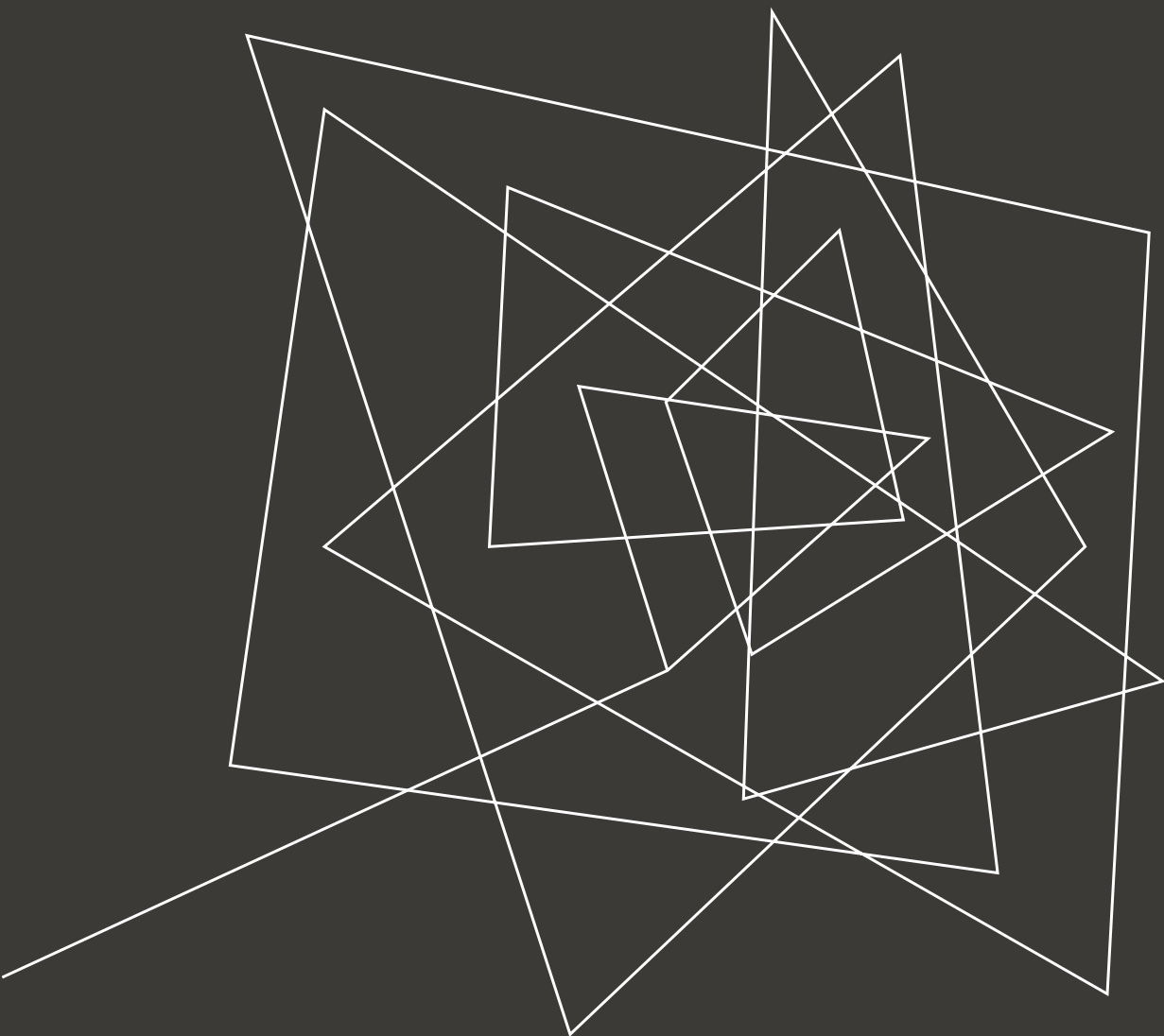
Geração de Valor

Resumo

INTRODUÇÃO

Olá a todos! Meu nome é Jade e estarei com vocês para apresentar minha análise.

É um prazer estar aqui hoje para compartilhar os resultados do desafio de análise de dados sobre filmes. Durante esse processo, explorei diversos aspectos dos gêneros Ação e Aventura e suas características, a fim de entregar um produto refinado e relevante para o Programa de Bolsas da Compasso.



PRINCIPAIS OBJETIVOS

Explorar a correlação entre a duração dos filmes e suas notas médias pelas décadas dos filmes mais populares.



GÊNERO DE FILMES E REFINAMENTO

Em primeiro lugar, nossa Squad 5 – Five It – foi atribuído nossa análise nos gêneros de Ação e Aventura. Esses gêneros são amplamente apreciados pelo público.

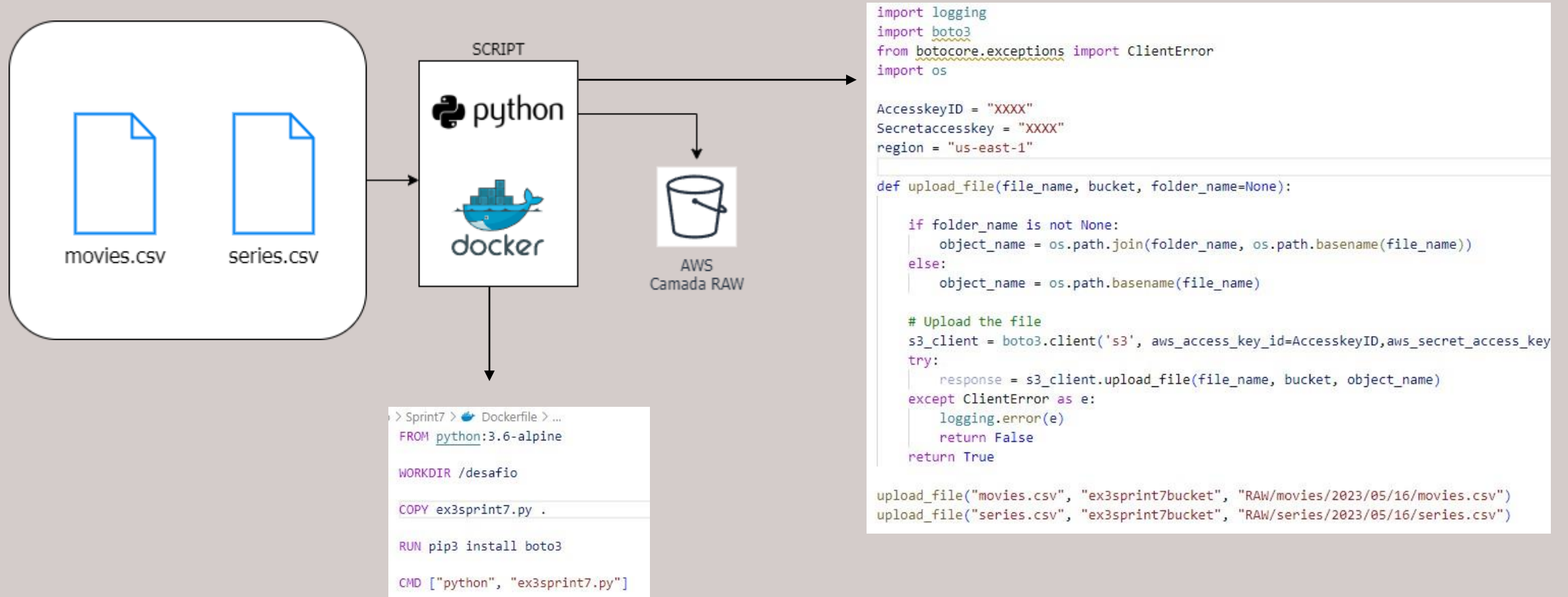
Além disso, decidi por um refinamento no meu estudo, considerando apenas os filmes de Ação e Aventura que possuíam mais de 15 mil votos. Esse critério me permitiu trabalhar com filmes verdadeiramente populares.



PRIMEIRA ETAPA DO DESAFIO

Para alcançar nossos resultados, passamos por diferentes etapas. Primeiro, coletamos todos os dados sobre filmes e series do arquivo CSV. Em seguida, os aloquei dentro do Bucket AWS via Docker com Boto3.

DIAGRAMA E CÓDIGOS DA PRIMEIRA ETAPA

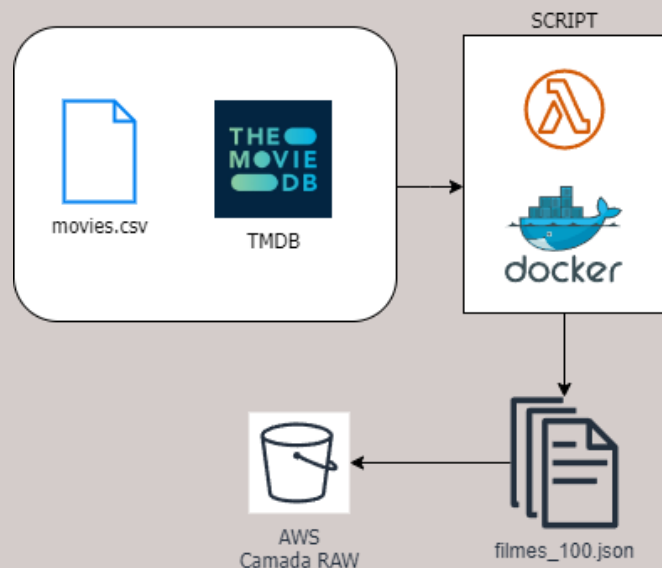




SEGUNDA ETAPA DO DESAFIO

Após obtermos nosso CSV na camada RAW, avançamos para a importação dos dados JSON do TMDb. Focando especificamente em filmes de Ação e Aventura, procedi ao processamento desses dados a partir do CSV, aplicando um filtro para selecionar somente os filmes populares, com mais de 15 mil votos e que não tenha valores nulos no tempo em minutos. Em seguida, colocamos os dados JSON no Bucket da AWS.

DIAGRAMA E CÓDIGOS DA SEGUNDA ETAPA



```
# Carregar o arquivo CSV 'movies.csv' em um DataFrame
df_movies = pd.read_csv('movies.csv', delimiter='|')

# Filtrar apenas os filmes de Ação/Aventura
df_filtered = df_movies[(df_movies['tempoMinutos'] != 'N') &
                        (df_movies['numeroVotos'] >= 15000) &
                        (df_movies['genero'].str.contains('action', case=False) &
                         df_movies['genero'].str.contains('adventure', case=False))]

df_filtered = df_filtered[['id', 'tituloPrincipal', 'tempoMinutos', 'anoLancamento', 'genero', 'notaMedia', 'numeroVotos']]
df_filtered = df_filtered.groupby('tituloPrincipal').first()
df_filtered = df_filtered.sort_values(by=['notaMedia', 'tempoMinutos'], ascending=[False, True])
```

O código completo pode ser encontrado nesse link no meu github:
<https://github.com/jademarinho/Sprint1/blob/main/Desafio/Sprint8/nb%20com%20informacoes%20de%20api%20final.ipynb>

```
# Iterar sobre os IDs dos filmes filtrados
for movie_id in df_filtered['id']:
    # Construir a URL da API para obter as informações do país de produção do filme
    url = f'{base_url}{movie_id}?api_key={api_key}&language=pt-BR'

    # Fazer solicitação para a URL e obter os dados do filme
    response = requests.get(url)

    if response.status_code == 200:
        filme_data = response.json()

        # Extrair o nome do país de produção do filme
        countries = filme_data.get('production_countries', [])
        production_country = countries[0]['name'] if countries else 'N/A'

        # Adicionar a informação do país de produção ao dataframe
        df_filtered.loc[df_filtered['id'] == movie_id, 'production_country'] = production_country

        # Adicionar os dados do filme à lista de resultados
        filmes.append(filme_data)

# Atualizar a barra de progresso
pbar.update(1)

# Incrementar o contador de filmes
contador += 1

# Verificar se chegou a 100 filmes ou se atingiu o total de 754 filmes
if contador % 100 == 0 or contador == len(df_filtered):
    # Salvar o JSON em um arquivo
    json_data = json.dumps(filmes, indent=4)

    # Gerar o caminho do arquivo no S3
    folder = "Raw/TMDB/JSON/2023/06/02/"
    nome_arquivo = f"{folder}filmes_{contador}.json"

    # Enviar o JSON para o S3
    s3_client.put_object(Body=json_data, Bucket=bucket, Key=nome_arquivo)

    # Limpar a lista de filmes
    filmes = []

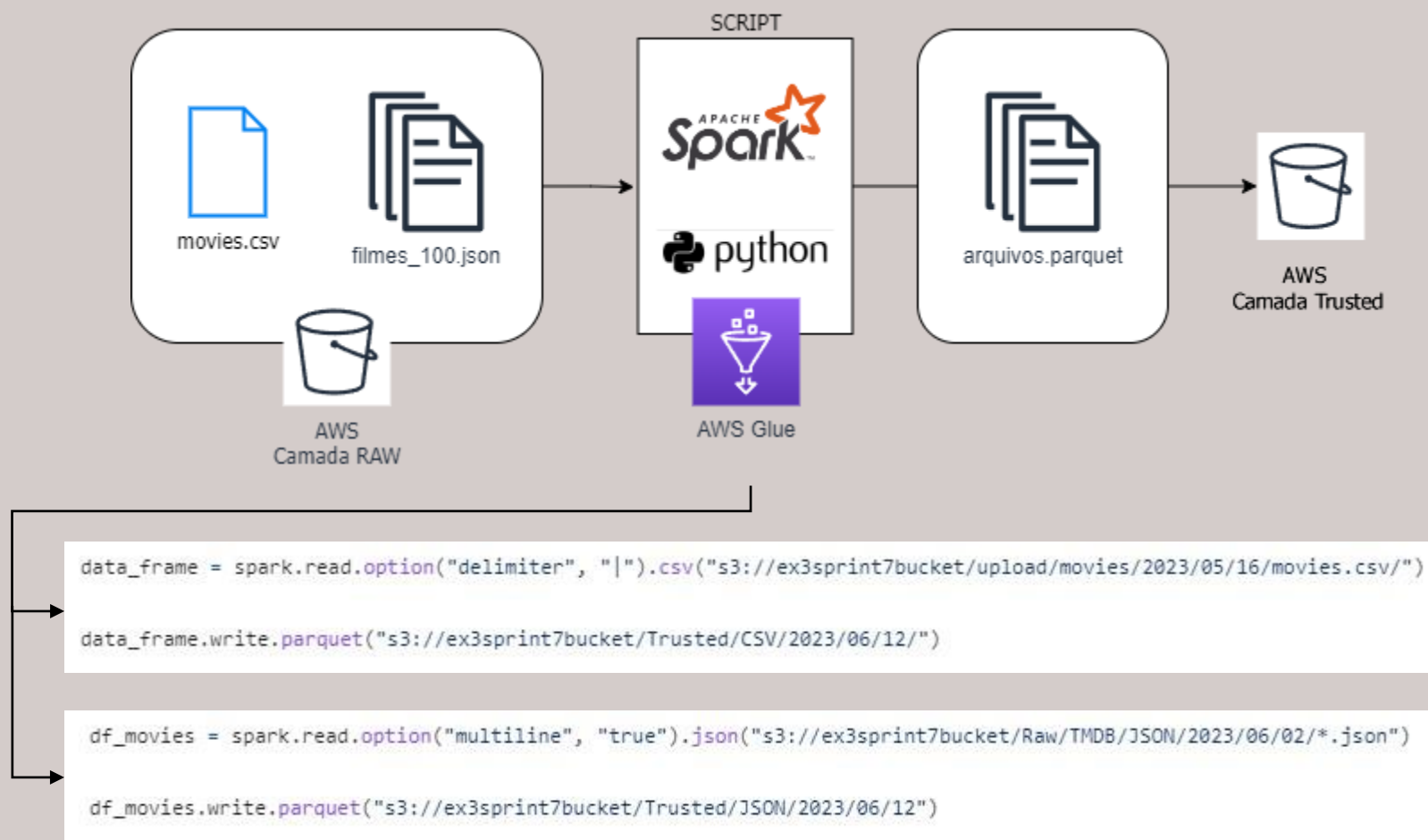
# Fechar a barra de progresso
pbar.close()
```



TERCEIRA ETAPA DO DESAFIO – PARTE 1

Persistir dados da RAW na camada Trusted do data lake no formato PARQUET através de 2 Jobs em PySpark no AWS Glue, um para o CSV e outro para o JSON.

DIAGRAMA E CÓDIGOS DA TERCEIRA ETAPA

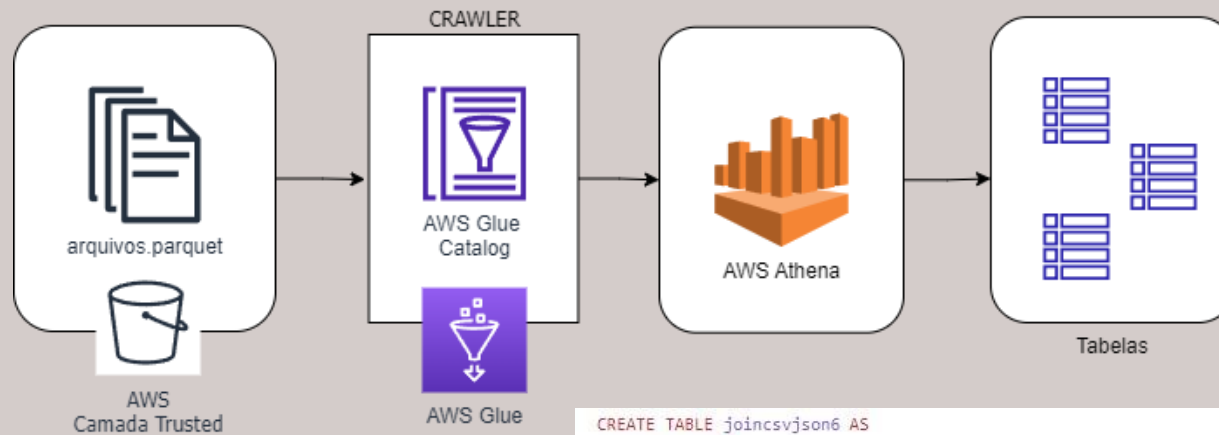




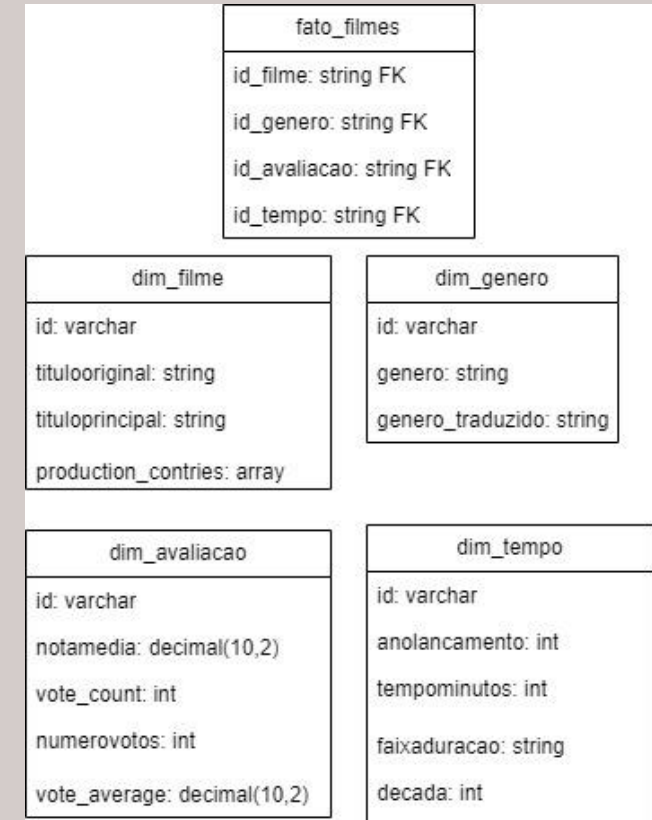
TERCEIRA ETAPA DO DESAFIO – PARTE 2

Nesta etapa crucial do desafio, iremos extrair os arquivos Parquet da camada Trusted utilizando o Glue Data Catalog. Em seguida, criaremos tabelas e aplicaremos manipulações no AWS Athena, visando a modelagem dos dados no formato dimensional. No meu caso específico, realizarei um JOIN entre o CSV e o JSON, selecionando os dados relevantes para a criação do modelo dimensional.

DIAGRAMA E CÓDIGOS DA TERCEIRA ETAPA



```
CREATE TABLE joincsvjson6 AS
SELECT DISTINCT
  CAST(csv.id AS VARCHAR) AS id,
  CAST(csv.tituloPrincipal AS VARCHAR) AS tituloPrincipal,
  CAST(csv.tituloOriginal AS VARCHAR) AS tituloOriginal,
  CAST(csv.tempoMinutos AS INT) AS tempoMinutos,
  CAST(csv.anoLancamento AS INT) AS anoLancamento,
  CAST(csv.genero AS VARCHAR) AS genero,
  CAST(csv.notaMedia AS DECIMAL(10, 2)) AS notaMedia,
  CAST(csv.numeroVotos AS INT) AS numeroVotos,
  json.imdb_id,
  json.production_countries,
  CAST(json.vote_average AS DECIMAL(10, 2)) AS vote_average,
  CAST(json.vote_count AS INT) AS vote_count
FROM csvrefined2 AS csv
JOIN "AwsDataCatalog"."jsonrefined"."12" AS json
ON csv.id = json.imdb_id;
```



Para saber mais sobre: <https://github.com/jademarinho/Sprint1/blob/main/Sprint9/readme.md>

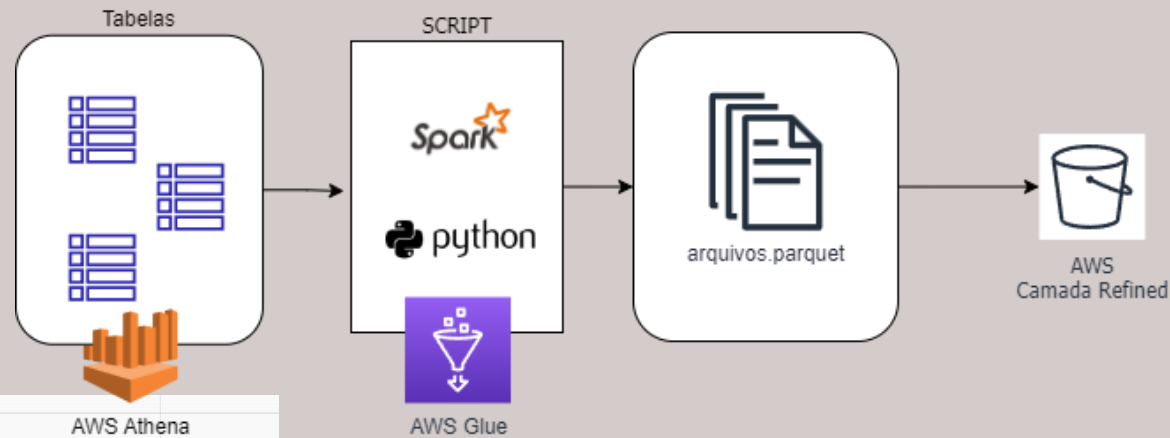


TERCEIRA ETAPA DO DESAFIO – PARTE 3

Persistir dados da AWS Athena na camada Refined do data lake no formato PARQUET através de Jobs em PySpark no AWS Glue.

Observação importante: Durante o criação do meu dash eu tive que adicionar mais campos para atender minhas necessidades, tal como faixaduracao e genero_traducao e tirar o production_countries.

DIAGRAMA E CÓDIGOS DA TERCEIRA ETAPA



```
1 CREATE OR REPLACE VIEW "traducao_genero" AS
2 SELECT
3     id,
4     CASE
5         WHEN genero = 'Action,Adventure,Western' THEN 'Ação, Aventura, Western'
6         WHEN genero = 'Action,Adventure,Thriller' THEN 'Ação, Aventura, Suspense'
7         WHEN genero = 'Action,Adventure,History' THEN 'Ação, Aventura, História'
8         WHEN genero = 'Action,Adventure,Horror' THEN 'Ação, Aventura, Terror'
9         WHEN genero = 'Action,Adventure,Romance' THEN 'Ação, Aventura, Romance'
10        WHEN genero = 'Action,Adventure,Fantasy' THEN 'Ação, Aventura, Fantasia'
11        WHEN genero = 'Action,Adventure,Crime' THEN 'Ação, Aventura, Crime'
12        WHEN genero = 'Action,Adventure,War' THEN 'Ação, Aventura, Guerra'
13        WHEN genero = 'Action,Adventure,Comedy' THEN 'Ação, Aventura, Comédia'
14        WHEN genero = 'Action,Adventure,Mystery' THEN 'Ação, Aventura, Mistério'
15        WHEN genero = 'Action,Adventure,Sci-Fi' THEN 'Ação, Aventura, Ficção Científica'
16        WHEN genero = 'Action,Adventure,Animation' THEN 'Ação, Aventura, Animação'
17        WHEN genero = 'Action,Adventure,Drama' THEN 'Ação, Aventura, Drama'
18        WHEN genero = 'Action,Adventure,Family' THEN 'Ação, Aventura, Família'
19        WHEN genero = 'Action,Adventure,Biography' THEN 'Ação, Aventura, Biografia'
20        WHEN genero = 'Action,Adventure' THEN 'Ação, Aventura'
21        ELSE genero
22    END AS generos
23 FROM "csvrefined"."joinscsvjson6"
```

AWS Athena

```
1 CREATE OR REPLACE VIEW "filmes_por_duracao_em_horas" AS
2 SELECT
3     id
4     , faixa_duracao
5     , "count"(*) quantidade_filmes
6 FROM
7     (
8         SELECT
9             id
10            , (CASE WHEN (tempominutos < 90) THEN 'Menos de 1h e 30min' WHEN (tempominutos BETWEEN 90 AND 119) THEN 'Entre 1h e 30min e 2h' WHEN (tempominutos BETWEEN 120 AND 149) THEN 'Entre 2h e 2h e 30min' WHEN (tempominutos BETWEEN 150 AND 179) THEN 'Entre 2h e 30min e 3h' WHEN (tempominutos BETWEEN 180 AND 209) THEN 'Entre 3h e 3h e 30min' ELSE 'Mais de 3h e 30min' END) faixa_duracao
11        FROM
12            "csvrefined"."joinscsvjson6"
13        ) subquery
14 GROUP BY id, faixa_duracao
15
```

Para saber mais sobre: <https://github.com/jademarinho/Sprint1/blob/main/Sprint9/readme.md>



QUARTA ETAPA DO DESAFIO

A partir do modelo de dados criado na etapa anterior, criei um dashboard interativo para visualização dos resultados. Nesse dashboard, apresentamos gráficos e métricas que nos permitiram obter insights. Por exemplo, as tendências de duração de filme no passar das décadas e padrões interessantes que nos ajudaram a compreender melhor a evolução desses gêneros cinematográficos.



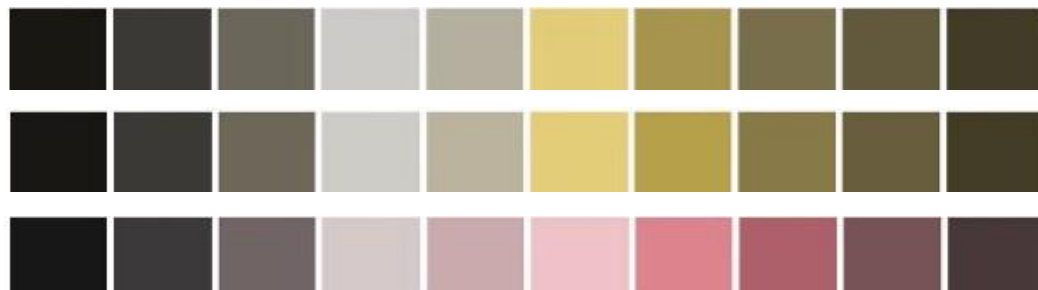
A DURAÇÃO DE UM FILME DEVE
ESTAR DIRETAMENTE RELACIONADA
À RESISTÊNCIA DA BEXIGA HUMANA.

Alfred Hitchcock

VISUALIZAÇÃO DA DASHBOARD

A inspiração para a paleta de cores da dashboard foi extraída do filme que conquistou o Primeiro lugar em termos de Maior Nota e também foi testado o contraste para daltônicos.

Protan(red), Deutan(green) e Tritan(blue), respectivamente.



Senhor dos Anéis: O Retorno do Rei. Fonte: @colorpalette.cinema



ÁREAS DE CONHECIMENTO

EFICIÊNCIA NO PROCESSAMENTO DE DADOS

Aprender sobre o processo de coleta, armazenamento e processamento de dados desde a camada RAW até a camada Trusted permite otimizar e automatizar essas etapas. Isso resulta em uma maior eficiência na manipulação e transformação dos dados, acelerando a disponibilidade de informações relevantes para os clientes da Compass.

ESCALABILIDADE E FLEXIBILIDADE

O aprendizado sobre o uso de serviços como AWS S3, AWS Lambda e AWS Glue proporciona uma base sólida para lidar com grandes volumes de dados e ambientes de análise em constante evolução.

GERAÇÃO DE VALOR

ANÁLISE E VISUALIZAÇÃO DE DADOS

Com o conhecimento adquirido no uso de ferramentas como AWS Glue, AWS Athena e AWS QuickSight, é possível realizar análises avançadas e criar visualizações impactantes dos dados. Isso permite aos clientes da Compass obter insights valiosos de maneira mais eficiente e concisa.

TOMADA DE DECISÕES EMBASADAS EM DADOS

Os conhecimentos adquiridos no Programa de Bolsas permitem uma compreensão mais profunda sobre como coletar, transformar e analisar dados relevantes para o negócio da Compass. Isso capacita os clientes a tomar decisões embasadas em dados sólidos, aumentando a precisão e melhorando os resultados de suas estratégias e ações.

FERRAMENTAS VALIOSAS

Além dos aprendizados da AWS, ao longo do Programa de Bolsas aprendemos sobre Segurança da Informação, Metodologias Ágeis, versionamento com Git, uso de terminal no Linux, SQL para Análise de Dados, Fundamentos de Big Data, Estatística e Docker.



RESUMO

Toda a trilha de conhecimentos fornecidos pelo Programa de Bolsas – da coleta, processamento, refinamento e análise – capacitam os estagiários da Compass a obter insights valiosos e repassar para os clientes tomarem decisões informadas e desenvolver estratégias eficazes com base em análises de dados confiáveis e visualizações intuitivas. Essa capacidade resulta em um maior valor agregado para os clientes, impulsionando o crescimento e o sucesso de seus negócios.



OBRIGADA

Jade Marinho Torres

jademarinho@hotmail.com

github.com/jademarinho

<https://www.linkedin.com/in/jade-marinho/>

