

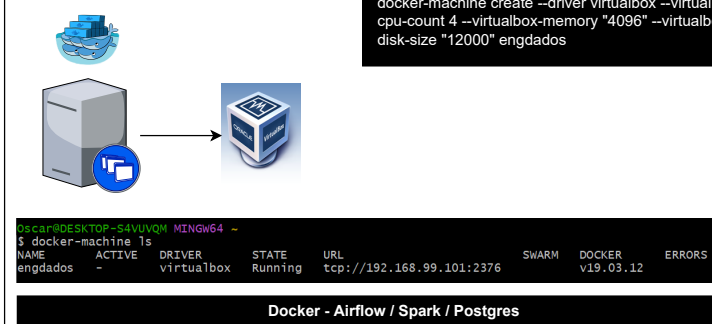
Projeto - Engenharia de Dados

Objetivo do projeto: desenvolver uma estrutura de dados, que possibilite a entrega dos dados com maior velocidade e acessibilidade para a equipe de ciência de dados.
Data: 20/10/2025

1 - Infraestrutura

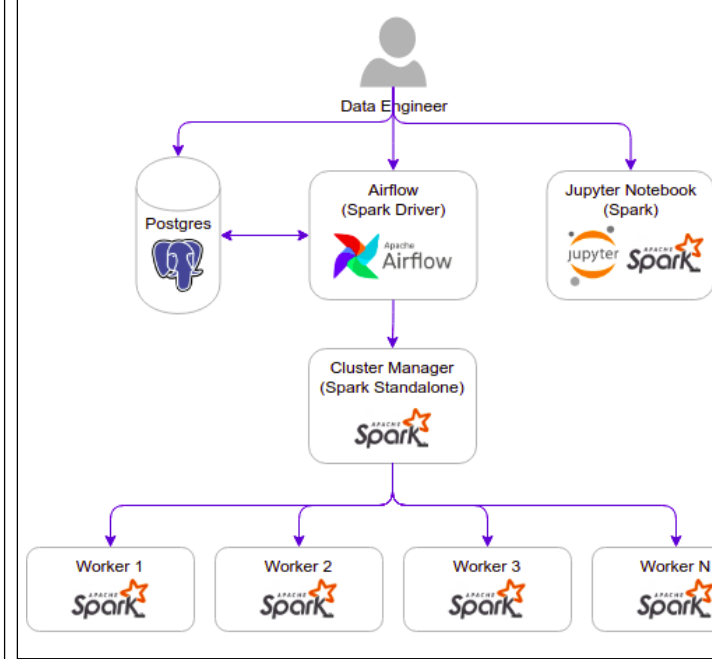
Para o projeto piloto de arquitetura de dados, foi criado um laboratório de dados no postgres, para gerar e desenvolver o laboratório de dados e experimentação no ambiente.

O objetivo da criação do laboratório de dados é gerar a massa de dados, para posteriormente exportar os arquivos csv que serão utilizados no processo de carga.

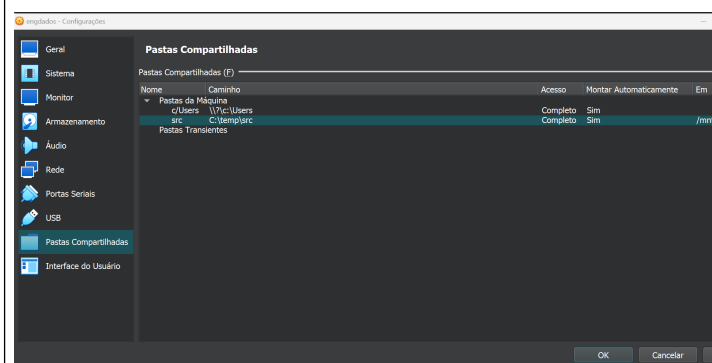


Para aplicar o processo de instalação dos softwares, criei um projeto Docker com o nome que necessitava para o projeto piloto, gostei muito do projeto criado pelo "Thiago Cardan Rodrigues" onde ele traz a arquitetura ideal para o projeto com Kafka, spark e Postgres.

GitHub: <https://github.com/thiagocardan/ThiagoCardanSpark>



Criei uma pasta em meu PC e compatível via VirtualBox, para deixar os dados permanentes nesta pasta.



Montei a pasta em /tmp/docker e criei o projeto com o seguinte comando:

```
sudo mount -t vboxsf -o uid=1000,gid=1000 /tmp/docker /tmp/docker
cd /tmp/docker
mkdir projeto
```

Depois disso, criei o arquivo Dockerfile com o seguinte conteúdo:

```
FROM postgres:15.3
COPY ./projeto /projeto
WORKDIR /projeto
RUN apt-get update && apt-get install -y postgresql-client
```

Salvando o arquivo Dockerfile no mesmo diretório do projeto.

Salvando o arquivo Dockerfile no mesmo diretório do projeto.

Salvando o arquivo Dockerfile no mesmo diretório do projeto.

Salvando o arquivo Dockerfile no mesmo diretório do projeto.

Salvando o arquivo Dockerfile no mesmo diretório do projeto.

Salvando o arquivo Dockerfile no mesmo diretório do projeto.

Salvando o arquivo Dockerfile no mesmo diretório do projeto.

Salvando o arquivo Dockerfile no mesmo diretório do projeto.

Salvando o arquivo Dockerfile no mesmo diretório do projeto.

Salvando o arquivo Dockerfile no mesmo diretório do projeto.

Salvando o arquivo Dockerfile no mesmo diretório do projeto.

Salvando o arquivo Dockerfile no mesmo diretório do projeto.

Salvando o arquivo Dockerfile no mesmo diretório do projeto.

Salvando o arquivo Dockerfile no mesmo diretório do projeto.

Salvando o arquivo Dockerfile no mesmo diretório do projeto.

Salvando o arquivo Dockerfile no mesmo diretório do projeto.

Salvando o arquivo Dockerfile no mesmo diretório do projeto.

Salvando o arquivo Dockerfile no mesmo diretório do projeto.

Salvando o arquivo Dockerfile no mesmo diretório do projeto.

Salvando o arquivo Dockerfile no mesmo diretório do projeto.

Salvando o arquivo Dockerfile no mesmo diretório do projeto.

Salvando o arquivo Dockerfile no mesmo diretório do projeto.

Salvando o arquivo Dockerfile no mesmo diretório do projeto.

Salvando o arquivo Dockerfile no mesmo diretório do projeto.

Salvando o arquivo Dockerfile no mesmo diretório do projeto.

Salvando o arquivo Dockerfile no mesmo diretório do projeto.

Salvando o arquivo Dockerfile no mesmo diretório do projeto.

Salvando o arquivo Dockerfile no mesmo diretório do projeto.

Salvando o arquivo Dockerfile no mesmo diretório do projeto.

Salvando o arquivo Dockerfile no mesmo diretório do projeto.

Salvando o arquivo Dockerfile no mesmo diretório do projeto.

Salvando o arquivo Dockerfile no mesmo diretório do projeto.

Salvando o arquivo Dockerfile no mesmo diretório do projeto.

Salvando o arquivo Dockerfile no mesmo diretório do projeto.

Salvando o arquivo Dockerfile no mesmo diretório do projeto.

Salvando o arquivo Dockerfile no mesmo diretório do projeto.

Salvando o arquivo Dockerfile no mesmo diretório do projeto.

Salvando o arquivo Dockerfile no mesmo diretório do projeto.

Salvando o arquivo Dockerfile no mesmo diretório do projeto.

Salvando o arquivo Dockerfile no mesmo diretório do projeto.

Salvando o arquivo Dockerfile no mesmo diretório do projeto.

Salvando o arquivo Dockerfile no mesmo diretório do projeto.

Salvando o arquivo Dockerfile no mesmo diretório do projeto.

Salvando o arquivo Dockerfile no mesmo diretório do projeto.

Salvando o arquivo Dockerfile no mesmo diretório do projeto.

Salvando o arquivo Dockerfile no mesmo diretório do projeto.

Salvando o arquivo Dockerfile no mesmo diretório do projeto.

Salvando o arquivo Dockerfile no mesmo diretório do projeto.

Salvando o arquivo Dockerfile no mesmo diretório do projeto.

Salvando o arquivo Dockerfile no mesmo diretório do projeto.

Salvando o arquivo Dockerfile no mesmo diretório do projeto.

Salvando o arquivo Dockerfile no mesmo diretório do projeto.

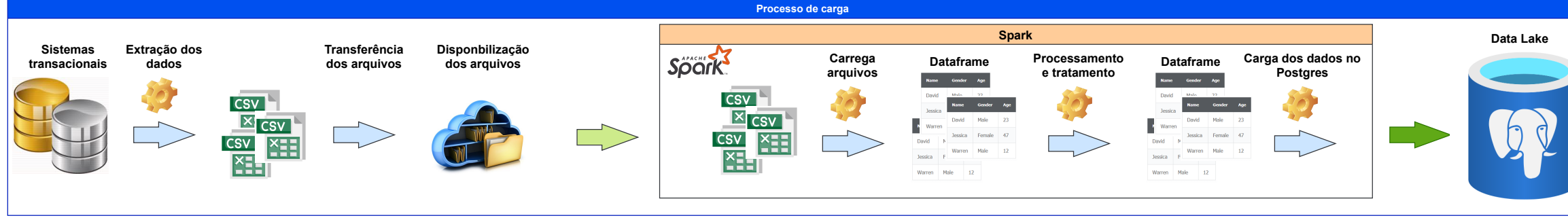
Arquitetura dos Dados

3 - Arquitetura do processo de carga dos dados

A arquitetura do processo de carga dos dados escolhido seria o modelo ETL.

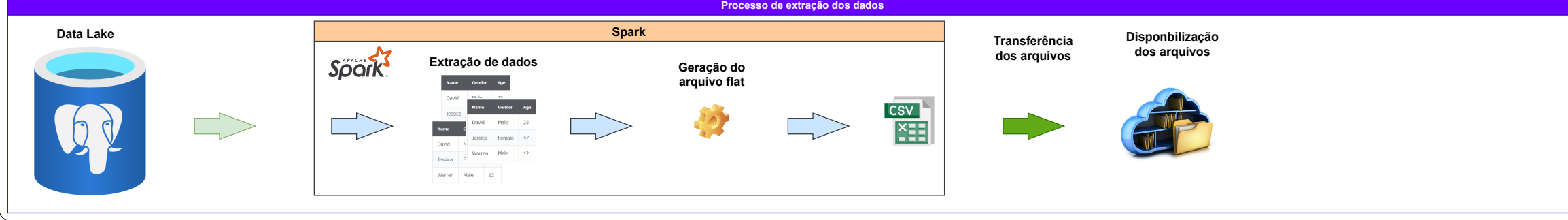
No cenário atual, onde os dados relacionais são armazenados em um banco de dados relacional como o postgres, talvez o modelo ideal seria o ETL, onde extrairíamos os dados em CSV, carregariamos os dados diretamente no postgres utilizando os utilitários do postgres e todo o tratamento seria dentro do postgres.

A escolha do modelo ETL, é para executar o processamento de dados utilizando o motor de processamento do Spark.



O objetivo do exercício, é executar o processamento no Spark. A extração dos dados e relacionamento entre os dados, serão processados no Spark.

Por esse razão, os dados não serão relacionados no postgres.



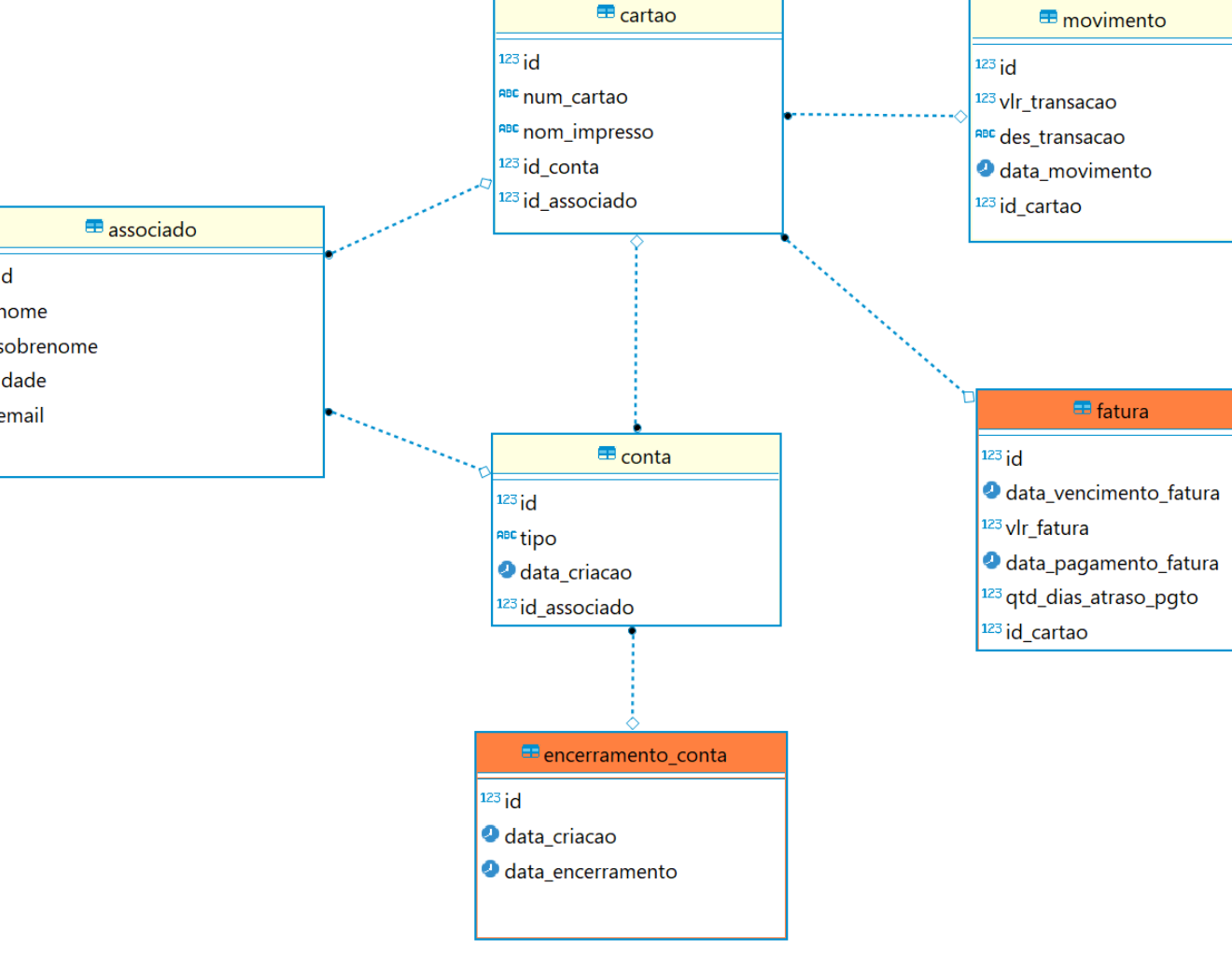
4 - Modelo dos dados no Data Lake (Postgres)

Para o desenvolvimento do projeto, foi desenvolvido o seguinte modelo de dados:

Obs: Foi adicionado duas tabelas auxiliares (fatura e encerramento de conta), para acrescentar 2 cenários para o projeto.

Tabelas em amarelo: tabelas propostas no projeto.

Tabelas em laranja: tabelas para enriquecimento da análise de inadimplência e análise de churn.

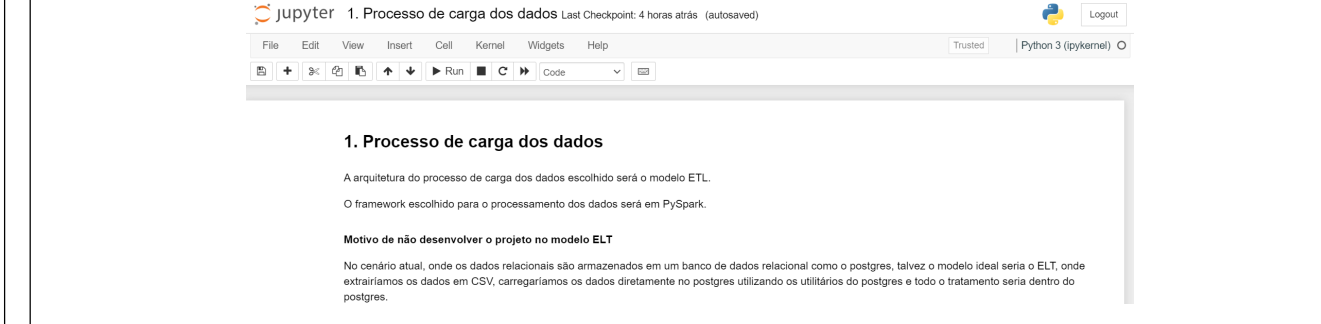


Estrutura das tabelas									
Os scripts de criação do banco está disponibilizado na pasta: DDL - Criação do banco Target									
Tabela Nome: associado									
Objeto ID: 17642									
Tabela Nome: movimento									
Objeto ID: 17644									
Tabela Nome: cartão									
Objeto ID: 17643									
Tabela Nome: fatura									
Objeto ID: 17673									
Tabela Nome: encerramento de conta									
Objeto ID: 17672									
Tabela Nome: conta									
Objeto ID: 17652									

5 - Fluxo do processo de carga

O fluxo do processo de carga, apresenta de forma visual as etapas do processo de carga, facilitando o entendimento do funcionamento do processo de carga.

Obs: Para facilitar o entendimento do código de carga, criei um notebook via Jupyter, que explica cada etapa do processo de carga, por esse razão resolvi não utilizar o arquivo. O arquivo está disponível no GitHub.



O objetivo do exercício, é executar o processamento no Spark. A extração dos dados e relacionamento entre os dados, serão processados no Spark.

Por esse razão, os dados não serão relacionados no postgres.

A escolha do modelo ETL, é para executar o processamento de dados utilizando o motor de processamento do Spark.

Obs: Como o objetivo do exercício não é criar um laboratório de dados, e sim focar no processo de carga dos dados, criei a massa de dados utilizando scripts SQL.

Todos os scripts SQL, que foram utilizados para o desenvolvimento do laboratório de dados, estão disponíveis no GitHub do projeto.

Salvando o arquivo Dockerfile no mesmo diretório do projeto.

Salvando o arquivo Dockerfile no mesmo diretório do projeto.

Salvando o arquivo Dockerfile no mesmo diretório do projeto.

Salvando o arquivo Dockerfile no mesmo diretório do projeto.

Salvando o arquivo Dockerfile no mesmo diretório do projeto.

Salvando o arquivo Dockerfile no mesmo diretório do projeto.

Salvando o arquivo Dockerfile no mesmo diretório do projeto.

Salvando o arquivo Dockerfile no mesmo diretório do projeto.

Salvando o arquivo Dockerfile no mesmo diretório do projeto.

Salvando o arquivo Dockerfile no mesmo diretório do projeto.

Salvando o arquivo Dockerfile no mesmo diretório do projeto.

Salvando o arquivo Dockerfile no mesmo diretório do projeto.

Salvando o arquivo Dockerfile no mesmo diretório do projeto.

Salvando o arquivo Dockerfile no mesmo diretório do projeto.

Salvando o arquivo Dockerfile no mesmo diretório do projeto.

Salvando o arquivo Dockerfile no mesmo diretório do projeto.

Salvando o arquivo Dockerfile no mesmo diretório do projeto.

Salvando o arquivo Dockerfile no mesmo diretório do projeto.

Salvando o arquivo Dockerfile no mesmo diretório do projeto.

Salvando o arquivo Dockerfile no mesmo diretório do projeto.

Salvando o arquivo Dockerfile no mesmo diretório do projeto.

Salvando o arquivo Dockerfile no mesmo diretório do projeto.

Salvando o arquivo Dockerfile no mesmo diretório do projeto.

Salvando o arquivo Dockerfile no mesmo diretório do projeto.

Salvando o arquivo Dockerfile no mesmo diretório do projeto.

Salvando o arquivo Dockerfile no mesmo diretório do projeto.

Salvando o arquivo Dockerfile no mesmo diretório do projeto.

Salvando o arquivo Dockerfile no mesmo diretório do projeto.

Salvando o arquivo Dockerfile no mesmo diretório do projeto.

Salvando o arquivo Dockerfile no mesmo diretório do projeto.

Salvando o arquivo Dockerfile no mesmo diretório do projeto.

Salvando o arquivo Dockerfile no mesmo diretório do projeto.

Salvando o arquivo Dockerfile no mesmo diretório do projeto.

Salvando o arquivo Dockerfile no mesmo diretório do projeto.

Salvando o arquivo Dockerfile no mesmo diretório do projeto.

Salvando o arquivo Dockerfile no mesmo diretório do projeto.

Salvando o arquivo Dockerfile no mesmo diretório do projeto.

Salvando o arquivo Dockerfile no mesmo diretório do projeto.

Salvando o arquivo Dockerfile no mesmo diretório do projeto.

Salvando o arquivo Dockerfile no mesmo diretório do projeto.

Salvando o arquivo Dockerfile no mesmo diretório do projeto.

Salvando o arquivo Dockerfile no mesmo diretório do projeto.

Salvando o arquivo Dockerfile no mesmo diretório do projeto.

Salvando o arquivo Dockerfile no mesmo diretório do projeto.

Salvando o arquivo Dockerfile no mesmo diretório do projeto.

Salvando o arquivo Dockerfile no mesmo diretório do projeto.

Salvando o arquivo Dockerfile no mesmo diretório do projeto.

Salvando o arquivo Dockerfile no mesmo diretório do projeto.

Salvando o arquivo Dockerfile no mesmo diretório do projeto.

Salvando o arquivo Dockerfile no mesmo diretório do projeto.

Salvando o arquivo Dockerfile no mesmo diretório do projeto.

Salvando o arquivo Dockerfile no mesmo diretório do projeto.

Salvando o arquivo Dockerfile no mesmo diretório do projeto.

Salvando o arquivo Dockerfile no mesmo diretório do projeto.

Salvando o arquivo Dockerfile no mesmo diretório do projeto.

Salvando o arquivo Dockerfile no mesmo diretório do projeto.

Salvando o arquivo Dockerfile no mesmo diretório do projeto.

Salvando o arquivo Dockerfile no mesmo diretório do projeto.

Salvando o arquivo Dockerfile no mesmo diretório do projeto.

Salvando o arquivo Dockerfile no mesmo diretório do projeto.

Salvando o arquivo Dockerfile no mesmo diretório do projeto.

Salvando o arquivo Dockerfile no mesmo diretório do projeto.

Salvando o arquivo Dockerfile no mesmo diretório do projeto.

Salvando o arquivo Dockerfile no mesmo diretório do projeto.

Salvando o arquivo Dockerfile no mesmo diretório do projeto.

Salvando o arquivo Dockerfile no mesmo diretório do projeto.

Salvando o arquivo Dockerfile no mesmo diretório do projeto.

Salvando o arquivo Dockerfile no mesmo diretório do projeto.

Salvando o arquivo Dockerfile no mesmo diretório do projeto.

Salvando o arquivo Dockerfile no mesmo diretório do projeto.

Salvando o arquivo Dockerfile no mesmo diretório do projeto.

Salvando o arquivo Dockerfile no mesmo diretório do projeto.

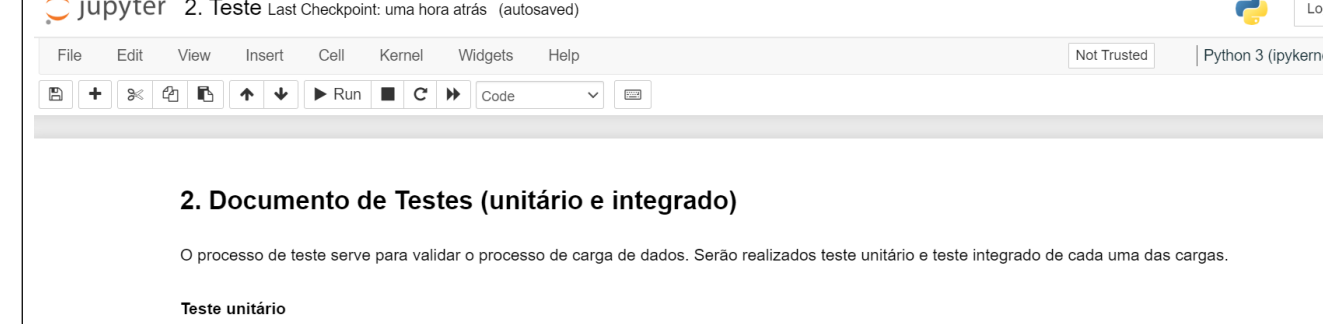
Salvando o arquivo Dockerfile no mesmo diretório do projeto.

Salvando o arquivo Dockerfile no mesmo diretório do projeto.

6 - Testes

Foi realizado testes unitários e integrados para validar o processo de carga dos dados.

Obs: Para facilitar o entendimento do código de carga, criei um notebook via Jupyter, que explica cada etapa do processo de teste. O arquivo está disponível no GitHub.



O objetivo do exercício, é executar o processamento no Spark. A extração dos dados e relacionamento entre os dados, serão processados no Spark.

Por esse razão, os dados não serão relacionados no postgres.

A escolha do modelo ETL, é para executar o processamento de dados utilizando o motor de processamento do Spark.

Obs: Como o objetivo do exercício não é criar um laboratório de dados, e sim focar no processo de carga dos dados, criei a massa de dados utilizando scripts SQL.

Todos os scripts SQL, que foram utilizados para o desenvolvimento do laboratório de dados, estão disponíveis no GitHub do projeto.

Salvando o arquivo Dockerfile no mesmo diretório do projeto.

Salvando o arquivo Dockerfile no mesmo diretório do projeto.

Salvando o arquivo Dockerfile no mesmo diretório do projeto.

Salvando o arquivo Dockerfile no mesmo diretório do projeto.

Salvando o arquivo Dockerfile no mesmo diretório do projeto.

Salvando o arquivo Dockerfile no mesmo diretório do projeto.

Salvando o arquivo Dockerfile no mesmo diretório do projeto.

Salvando o arquivo Dockerfile no mesmo diretório do projeto.

Salvando o arquivo Dockerfile no mesmo diretório do projeto.

Salvando o arquivo Dockerfile no mesmo diretório do projeto.

Salvando o arquivo Dockerfile no mesmo diretório do projeto.

Salvando o arquivo Dockerfile no mesmo diretório do projeto.

Salvando o arquivo Dockerfile no mesmo diretório do projeto.

Salvando o arquivo Dockerfile no mesmo diretório do projeto.

Salvando o arquivo Dockerfile no mesmo diretório do projeto.

Salvando o arquivo Dockerfile no mesmo diretório do projeto.

Salvando o arquivo Dockerfile no mesmo diretório do projeto.

Salvando o arquivo Dockerfile no mesmo diretório do projeto.

Salvando o arquivo Dockerfile no mesmo diretório do projeto.

Salvando o arquivo Dockerfile no mesmo diretório do projeto.

Salvando o arquivo Dockerfile no mesmo diretório do projeto.

Salvando o arquivo Dockerfile no mesmo diretório do projeto.

Salvando o arquivo Dockerfile no mesmo diretório do projeto.

Salvando o arquivo Dockerfile no mesmo diretório do projeto.

Salvando o arquivo Dockerfile no mesmo diretório do projeto.

Salvando o arquivo Dockerfile no mesmo diretório do projeto.

Salvando o arquivo Dockerfile no mesmo diretório do projeto.

Salvando o arquivo Dockerfile no mesmo diretório do projeto.

Salvando o arquivo Dockerfile no mesmo diretório do projeto.

Salvando o arquivo Dockerfile no mesmo diretório do projeto.

Salvando o arquivo Dockerfile no mesmo diretório do projeto.

Salvando o arquivo Dockerfile no mesmo diretório do projeto.

Salvando o arquivo Dockerfile no mesmo diretório do projeto.

Salvando o arquivo Dockerfile no mesmo diretório do projeto.

Salvando o arquivo Dockerfile no mesmo diretório do projeto.

Salvando o arquivo Dockerfile no mesmo diretório do projeto.

Salvando o arquivo Dockerfile no mesmo diretório do projeto.

Salvando o arquivo Dockerfile no mesmo diretório do projeto.

Salvando o arquivo Dockerfile no mesmo diretório do projeto.

Salvando o arquivo Dockerfile no mesmo diretório do projeto.