



FUNDAÇÃO EDSON QUEIROZ
UNIVERSIDADE DE FORTALEZA - UNIFOR
CENTRO DE TECNOLOGIA E ENGENHARIA
CURSO MBA EM ENGENHARIA DE DADOS

TRABALHO FINAL PARA DISCIPLINA: LINGUAGEM DE PROGRAMAÇÃO
PARA ENGENHARIA DE DADOS

FELIPE ALVES DA SILVA

Fortaleza-CE

2024

FELIPE ALVES DA SILVA

**TRABALHO FINAL PARA DISCIPLINA: LINGUAGEM DE PROGRAMAÇÃO
PARA ENGENHARIA DE DADOS**

Trabalho final apresentado ao curso de MBA em engenharia de dados da Universidade de Fortaleza, para obtenção de nota final para a disciplina de linguagem de programação para engenharia de dados lecionado pelo professor Thiago Bluhm.

Fortaleza-CE

2024

SUMÁRIO

SUMÁRIO	3
1 DIFICULDADES E PONTOS A FAVOR.....	4
2 EVIDÊNCIAS DE PERSISTÊNCIA DE DADOS NO POSTGRESQL	5
3 FONTES DE PESQUISA	7

LINK PARA NOTEBOOK NO GITHUB:

[HTTPS://GITHUB.COM/FELIPEALVSS/PROJETO_DADOS_IBGE/BLOB/MAIN/DESAFIO%20IBGE.IPYNB](https://github.com/felipealvss/projeto_dados_ibge/blob/main/desafio%20ibge.ipynb)

1 DIFICULDADES E PONTOS A FAVOR

O primeiro ponto de dificuldade que deparei nesse trabalho foi buscar substitutos para o *Pandas*, pois era a única biblioteca que até o momento conhecia para uso de diversos dados, porém com o uso do *CSV* e do *OS* pude implementar a lógica de carga de vários arquivos de forma ágil. Em seguida, a outra problemática percebida foi entender como manipular o dado de forma competente sem que eu estourasse a memória da minha máquina pessoal. Através de sites e (na maior parte) de diversas tentativas práticas, pude finalmente entender como aplicar uma carga incremental através de lotes de dados de forma que a inserção dos dados funcionou de forma competente.

Como ponto a favor da estratégia que escolhi, percebi que escolher utilizar a função *session.bulk_insert_mappings()* (*sqlalchemy*) juntamente com a divisão dos dados em lote se tornou muito mais rápida do que a inserção de dados utilizando a função *session.add()*, a qual mostrou ótima performance utilizando apenas um arquivo como parâmetro para carga (cerca de 950mil linhas), porém se mostrou extremamente onerosa com grande volume de dados.

Outro ponto positivo foi ter a necessidade de entender como realizar limpezas das variáveis e a manter um controle de memória com o *garbage collector* do python, foi nítido como o uso dessas boas práticas tornam o processamento mais leve para execução em hardwares não tão potentes (caso do meu computador pessoal).

Ao buscar por soluções presentes de mercado, pude perceber que muito se indica o uso do *Pandas* e *NumPy* juntamente com o *Dask*, permitindo aplicação de escalabilidade e utilizando a sintaxe padrão do Python. Outro que é muito comentado é o *Spark*, o qual funciona em cluster e consegue manipular de forma ágil e simples grandes volumes de dados. Há também a presença de ferramentas de ETL que podem realizar estes serviços, tais com o Pentaho Data integration, o Talend Open Studio e até mesmo o SSIS (Microsoft SQL Server Integration Services).

LINK PARA NOTEBOOK NO GITHUB:

[HTTPS://GITHUB.COM/FELIPEALVSS/PROJETO_DADOS_IBGE/BLOB/MAIN/DESAFIO%20IBGE.IPYNB](https://github.com/felipealvss/projeto_dados_ibge/blob/main/desafio%20ibge.ipynb)

2 EVIDÊNCIAS DE PERSISTÊNCIA DE DADOS NO POSTGRESQL

Seguem imagens com evidências que todos os dados foram persistidos no SGBD PostgreSQL:

1. Print da query presente no fim do Notebook informando a quantidade total de dados persistidos no PostgreSQL:

```
1 print('-----')
2 quantidade = session.query(func.count()).select_from(Localidade).scalar()
3 print(f'-- Total de dados registrados no DB: {quantidade}')
4 print('-----')
```

[10]

...

-- Total de dados registrados no DB: 111102875

2. Prints de queries executadas no próximo PostgreSQL exibindo uma amostra dos dados persistidos e executando uma contagem total dos dados persistidos na tabela:

The screenshot shows a PostgreSQL client interface with a sidebar on the left containing a tree view of database objects. The main area displays a query and its results.

Query:

```
1 -- AMOSTRA DE DADOS PERSISTIDOS
2 SELECT *
3 FROM ibge
4 LIMIT 10;
5
6 -- VERIFICAÇÃO DE TOTAL DE DADOS PERSISTIDOS
7 -- Quantidade de dados nos CSVs: 111102875
8 SELECT
9 COUNT (*) AS qtde_dados
10 FROM ibge;
11
```

Data Output:

	id [PK] integer	cod_uf integer	cod_mun integer	cod_especie integer	latitude double precision	longitude double precision	nv_geo_coord integer
1	1	11	1100015	1	-11.929621	-61.999058	1
2	2	11	1100015	7	-11.929365	-61.999312	1
3	3	11	1100015	6	-11.9294	-61.999499	1
4	4	11	1100015	6	-11.929364	-61.999695	1
5	5	11	1100015	1	-11.929381	-61.999262	1
6	6	11	1100015	1	-11.930439	-61.999046	1
7	7	11	1100015	7	-11.930571	-61.999006	1
8	8	11	1100015	1	-11.930311	-61.998982	1
9	9	11	1100015	1	-11.930717	-61.999009	1
10	10	11	1100015	1	-11.930282	-61.999566	1

The interface also shows a sidebar with a tree view of database objects, including Tables (1), Columns (7), Constraints (1), Indexes, RLS Policies, Rules, Triggers, Trigger Functions, Types, and Views. The table 'ibge' is selected, and its columns are listed: id, cod_uf, cod_mun, cod_especie, latitude, longitude, and nv_geo_coord.

LINK PARA NOTEBOOK NO GITHUB:

[HTTPS://GITHUB.COM/FELIPEALVSS/PROJETO_DADOS_IBGE/BLOB/MAIN/DESAFIO%20IBGE.IPYNB](https://github.com/felipealvss/projeto_dados_ibge/blob/main/desafio%20ibge.ipynb)

3 FONTES DE PESQUISA

<https://coderpad.io/blog/development/sqlalchemy-with-postgresql/>

https://docs.sqlalchemy.org/en/20/orm/mapping_styles.html

https://docs-sqlalchemy.readthedocs.io/en/latest/orm/persistence_techniques.html

<https://www.pythonforbeginners.com/basics/generator-comprehension-in-python>

https://www.w3schools.com/python/python_dictionaries.asp

<https://towardsdatascience.com/how-to-perform-bulk-inserts-with-sqlalchemy-efficiently-in-python-23044656b97d>

<https://realpython.com/python-timer/>