

Projeto 11



Um pipeline **ETL (Extract, Transform, and Load)** extrai dados de fontes, transforma-os e carrega-os em um sistema de armazenamento.

Ele ajuda a criar formatos de dados limpos e utilizáveis para análise.

O PySpark é ideal para construir pipelines ETL para processamento de dados em larga escala.

Ele oferece computação distribuída, alto desempenho e lida com dados estruturados e não estruturados de forma eficiente.

Este projeto o levará pela construção de um pipeline ETL usando o PySpark.

Construindo um Pipeline ETL usando PySpark

O conjunto de dados que usaremos para construir um ETL Pipeline contém dados relacionados à temperatura para vários países de 1961 a 2022.

As colunas incluem identificadores como ObjectId, Country, ISO2 e ISO3, juntamente com dados de temperatura anuais, como F1961, F1962, etc., como valores de ponto flutuante.

Algumas colunas contêm valores ausentes.

Você pode baixar este conjunto de dados [aqui](#).

Estudo de Caso

Este estudo de caso utiliza dois conjuntos de dados abrangentes para explorar o impacto global das emissões de carbono.

O primeiro conjunto de dados, **Mudança Anual de Temperatura da Superfície (Annual Surface Temperature Change)**, fornece anomalias de temperatura em vários países em relação a uma linha de base, medida em graus Celsius.

O conjunto de dados abrange de 1960 a 2022 e inclui observações anuais, oferecendo uma perspectiva temporal sobre como as temperaturas globais mudaram ao longo de décadas.

Os principais recursos incluem nomes de países e valores de mudança de temperatura para cada ano.

O segundo conjunto de dados, **Concentrações Atmosféricas Mensais de CO₂ (Monthly Atmospheric CO₂ Concentrations)**, inclui medições globais de CO₂ de 1958 a 2024.

Os dados são representados em partes por milhão (ppm), com observações mensais que permitem uma análise temporal detalhada.

Os principais recursos incluem a data da medição e os valores de concentração de CO₂.

Juntos, esses conjuntos de dados fornecem uma base sólida para analisar tendências climáticas e entender a ligação entre emissões de carbono e aquecimento global.

Desenvolveremos um Pipeline ETL usando PySpark para processar este conjunto de dados para lidar com as seguintes tarefas:

- 1. Extrair:** Carregar o conjunto de dados do arquivo CSV;
- 2. Transformar:** Limpar os dados, lidar com valores ausentes e pivotar dados de temperatura anuais para análise;
- 3. Carregar:** Salvar os dados processados em um novo formato de armazenamento (por exemplo, Parquet ou um banco de dados).

Solução

Etapa 1: Configurando o ambiente e inicializando uma sessão do PySpark

Certifique-se de que o PySpark esteja instalado e configurado.

Execute o seguinte comando para instalar o PySpark se ele ainda não estiver instalado:

```
In [1]: #pip install pyspark
```

Inicialize uma sessão PySpark para permitir a interação com a estrutura Spark:

```
In [2]: # Importa a classe SparkSession do módulo pyspark.sql
from pyspark.sql import SparkSession

# Configura a SparkSession
builder = SparkSession.builder
builder = builder.appName("ETL Pipeline") # Define o nome do aplicativo
spark = builder.getOrCreate() # Cria ou reutiliza uma SparkSession existente

# Exibe uma mensagem indicando que a SparkSession foi inicializada com sucesso
print("SparkSession inicializada com sucesso!")
```

SparkSession inicializada com sucesso!

Etapa 2: Extrair - Carregar o Conjunto de Dados

A próxima etapa é carregar o conjunto de dados em um PySpark DataFrame:

```
In [3]: # Define o caminho do arquivo CSV
file_path = "C:/Users/giova/Projetos/Projeto_11/temperature.csv"

# Carrega o arquivo CSV em um DataFrame Spark
# - header=True indica que o arquivo contém cabeçalho
```

```
# - inferSchema=True permite que o Spark infira automaticamente os tipos das colunas
df = spark.read.csv(file_path, header=True, inferSchema=True)
```

```
In [4]: # Exibe o esquema do DataFrame (nomes das colunas e seus tipos de dados)
df.printSchema()
```

```
root
|-- ObjectId: integer (nullable = true)
|-- Country: string (nullable = true)
|-- ISO2: string (nullable = true)
|-- ISO3: string (nullable = true)
|-- F1961: double (nullable = true)
|-- F1962: double (nullable = true)
|-- F1963: double (nullable = true)
|-- F1964: double (nullable = true)
|-- F1965: double (nullable = true)
|-- F1966: double (nullable = true)
|-- F1967: double (nullable = true)
|-- F1968: double (nullable = true)
|-- F1969: double (nullable = true)
|-- F1970: double (nullable = true)
|-- F1971: double (nullable = true)
|-- F1972: double (nullable = true)
|-- F1973: double (nullable = true)
|-- F1974: double (nullable = true)
|-- F1975: double (nullable = true)
|-- F1976: double (nullable = true)
|-- F1977: double (nullable = true)
|-- F1978: double (nullable = true)
|-- F1979: double (nullable = true)
|-- F1980: double (nullable = true)
|-- F1981: double (nullable = true)
|-- F1982: double (nullable = true)
|-- F1983: double (nullable = true)
|-- F1984: double (nullable = true)
|-- F1985: double (nullable = true)
|-- F1986: double (nullable = true)
|-- F1987: double (nullable = true)
|-- F1988: double (nullable = true)
|-- F1989: double (nullable = true)
|-- F1990: double (nullable = true)
|-- F1991: double (nullable = true)
|-- F1992: double (nullable = true)
|-- F1993: double (nullable = true)
|-- F1994: double (nullable = true)
|-- F1995: double (nullable = true)
|-- F1996: double (nullable = true)
|-- F1997: double (nullable = true)
|-- F1998: double (nullable = true)
|-- F1999: double (nullable = true)
|-- F2000: double (nullable = true)
|-- F2001: double (nullable = true)
|-- F2002: double (nullable = true)
|-- F2003: double (nullable = true)
|-- F2004: double (nullable = true)
|-- F2005: double (nullable = true)
|-- F2006: double (nullable = true)
|-- F2007: double (nullable = true)
|-- F2008: double (nullable = true)
|-- F2009: double (nullable = true)
|-- F2010: double (nullable = true)
|-- F2011: double (nullable = true)
|-- F2012: double (nullable = true)
|-- F2013: double (nullable = true)
|-- F2014: double (nullable = true)
|-- F2015: double (nullable = true)
|-- F2016: double (nullable = true)
|-- F2017: double (nullable = true)
|-- F2018: double (nullable = true)
|-- F2019: double (nullable = true)
|-- F2020: double (nullable = true)
|-- F2021: double (nullable = true)
|-- F2022: double (nullable = true)
```



```
In [10]: # Importa a função "expr" que é usada para executar expressões SQL diretamente em uma coluna ou conjunto de colunas em um DataFrame PySpark
from pyspark.sql.functions import expr
```

```
In [11]: # Reorganizar Dados de Temperatura para ter as Colunas 'Year' e 'Temperature'
df_pivot = df.selectExpr(
    # Manter as colunas originais fixas no DataFrame resultante
    "ObjectId",
    "Country",
    "IS03",

    # Usar a função stack para pivotar os dados de temperatura
    # Gerar 62 pares para os anos de 1961 a 2022
    "stack(62, " +
    ", ".join([f"F{1961 + i}", F{1961 + i}" for i in range(62)]) +
    # Nomear as novas colunas como 'Year' e 'Temperature'
    ") as (Year, Temperature)"
)
```

```
In [12]: # Atualizar a Coluna 'Year' Convertendo seu Valor para Inteiro
df_pivot = df_pivot.withColumn(
    "Year", # Nome da coluna a ser substituída
    expr("int(substring(Year, 2, 4))" # Extrair a parte numérica e converter para inteiro
)
```

```
In [13]: # Exibe as primeiras 5 Linhas do DataFrame no Console
df_pivot.show(5)
```

```
+-----+-----+-----+-----+-----+
|ObjectId|      Country|IS03|Year|Temperature|
+-----+-----+-----+-----+-----+
|      1|Afghanistan, Isla...| AFG|1961|      -0.113|
|      1|Afghanistan, Isla...| AFG|1962|      -0.164|
|      1|Afghanistan, Isla...| AFG|1963|       0.847|
|      1|Afghanistan, Isla...| AFG|1964|      -0.764|
|      1|Afghanistan, Isla...| AFG|1965|      -0.244|
+-----+-----+-----+-----+-----+
```

only showing top 5 rows

Etapa 4: Carregar – Salvar os Dados Processados

Após concluir todas as etapas de processamento, salvamos os dados transformados em um arquivo Parquet para armazenamento e consulta eficientes:

```
In [14]: # Define o caminho de saída para o arquivo Parquet
output_path = "/processed_temperature.parquet"
```

```
In [15]: # Salva o DataFrame no formato Parquet
# Sobrescreve o arquivo se ele já existir
# Especifica que o formato de saída será Parquet
df_pivot.write \
    .mode("overwrite") \
    .parquet(output_path)
```

Esta operação salva o DataFrame transformado como um arquivo Parquet, o que o otimiza para armazenamento e consulta em um ambiente distribuído.

Podemos carregar o arquivo Parquet salvo para garantir que os dados foram salvos corretamente:

```
In [16]: # Carrega o Arquivo Parquet Salvo
processed_df = spark.read.parquet(output_path)
```

```
In [17]: # Exibe as primeiras 5 Linhas do Arquivo Salvo
processed_df.show(5)
```

ObjectID	Country	ISO3	Year	Temperature
1	Afghanistan, Isla...	AFG	1961	-0.113
1	Afghanistan, Isla...	AFG	1962	-0.164
1	Afghanistan, Isla...	AFG	1963	0.847
1	Afghanistan, Isla...	AFG	1964	-0.764
1	Afghanistan, Isla...	AFG	1965	-0.244

only showing top 5 rows

Resumo

O PySpark é ideal para construir pipelines ETL (Extract, Transform, and Load) para processamento de dados em larga escala.

Ele oferece computação distribuída, alto desempenho e manipula dados estruturados e não estruturados de forma eficiente.

Espero que este projeto sobre "Como Construir um Pipeline ETL" usando Python tenha sido útil para seu aprendizado.

Referência do Projeto: <https://thedeveloperprogrammer.com/2024/11/28/building-an-etl-pipeline-using-pyspark/>

Base de Dados: <https://statso.io/carbon-emissions-worldwide-case-study/>

Autor: Aman Kharwal - <https://amankharwal.medium.com/>

Projeto Estudado por: Giovani Aloísio da Luz

Disponível em: <https://www.giovani-luz.com.br/>