

FACULDADE ESTÁCIO DE SÁ CURSO: DESENVOLVIMENTO FULL STACK 4º SEMESTRE – MATRÍCULA 202302595341

Repositório GitHub - <u>alaimalmeida/tratantolmensidaoDosDados</u>

ALAIM ALMEIDA DE OLIVEIRA

Tratando a imensidão dos dados

1. Introdução

A manipulação de dados é uma tarefa essencial em diversas áreas, como análise de dados, ciência de dados e gestão de projetos. O formato CSV (Comma-Separated Values) é amplamente utilizado para armazenar e transferir dados devido à sua simplicidade e compatibilidade com diversas ferramentas. Neste relatório, abordaremos como ler arquivos CSV utilizando a linguagem Python e a biblioteca Pandas.

2. Objetivo

Demonstrar o processo de leitura de arquivos CSV utilizando a biblioteca Pandas, configurando corretamente parâmetros como separador de colunas, engine e encoding.

Microatividade 1: Descrever como ler um arquivo CSV usando a biblioteca Pandas (Python)

1. Procedimento

- a. Instalar a biblioteca Pandas, caso ainda não esteja instalada:
 "pip install pandas"
- b. Importar a biblioteca Pandas no código Python:"import pandas as pd"
- c. Definir o nome do arquivo CSV que será lido:"arquivo_csv = "dados.csv"
- d. Ler o arquivo CSV usando a função read_csv, especificando o separador de colunas (sep), a engine (engine='python') e o encoding (encoding='utf-8'):
 - "tabela = pd.read_csv(arquivo_csv, sep=';', engine='python', encoding='utf-8')"
- e. Exibir as primeiras linhas do DataFrame para verificar a leitura dos dados:
 - "print(tabela)"

Ficando dessa forma do código completo:

```
projeto.py >
    projeto.py > ...
    import pandas as pd

    # Definindo o nome do arquivo CSV
    arquivo_csv = "dados.csv"

# Lendo o arquivo CSV
# Usando o separador ';', a engine 'python' e especificando o encoding (se necessário)
# tabela = pd.read_csv(arquivo_csv, sep=';', engine='python', encoding='utf-8')

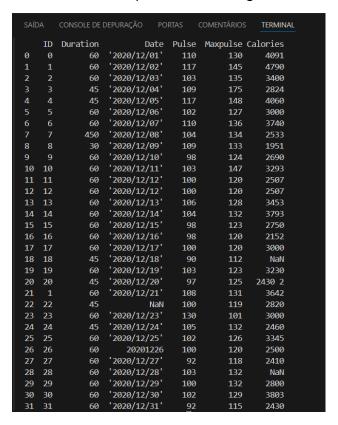
# Exibindo os dados lidos
print(tabela)

print(tabela)
```

A parte do arquivo CSV ficou dessa forma:

```
■ dados.csv > 🗅 data
         ID;Duration;Date;Pulse;Maxpulse;Calories
         0;60;'2020/12/01';110;130;4091
        1;60; '2020/12/02';117;145;4790
2;60; '2020/12/03';103;135;3400
         3;45; '2020/12/04';109;175;2824
        4;45; '2020/12/05';117;148;4060
        5;60; '2020/12/06';102;127;3000
6;60; '2020/12/07';110;136;3740
         7;450; '2020/12/08';104;134;2533
        8;30;'2020/12/09';109;133;1951
        9;60; '2020/12/10';98;124;2690
10;60; '2020/12/11';103;147;3293
        11;60; '2020/12/12';100;120;2507
12;60; '2020/12/12';100;120;2507
        13;60; '2020/12/13';106;128;3453
14;60; '2020/12/14';104;132;3793
        15;60; '2020/12/15';98;123;2750
16;60; '2020/12/16';98;120;2152
        17;60;'2020/12/17';100;120;3000
18;45;'2020/12/18';90;112;NaN
        19;60; '2020/12/19';103;123;3230
20;45; '2020/12/20';97;125;2430 2
         1;60;'2020/12/21';108;131;3642
         22;45;NaN;100;119;2820
         23;60; '2020/12/23';130;101;3000
24;45; '2020/12/24';105;132;2460
        25;60;'2020/12/25';102;126;3345
         26;60;20201226;100;120;2500
         28;60;'2020/12/28';103;132;NaN
        29;60; '2020/12/29';100;132;2800
30;60; '2020/12/30';102;129;3803
B1;60; '2020/12/31';92;115;2430
```

Ao executar o arquivo, fica da seguinte forma:



Microatividade 2: Descrever como criar um subconjunto de dados a partir de um conjunto existente usando a biblioteca Pandas (Python)

Procedimentos:

- 1. No mesmo arquivo/script utilizado na microatividade 1, crie uma nova variável;
- 2. Atribua, a essa nova variável, um subconjunto de dados contendo apenas parte
- colunas (recomenda-se a utilização de 3 colunas) disponíveis no conjunto de dados
- 4. original;
- Salve as alterações realizadas;
- 6. Imprima/exiba em tela os dados da nova variável (que contém o subconjunto de
- 7. dados).

De acordo com o procedimento, fica dessa forma:

```
projeto.py > ...

import pandas as pd

# Definindo o nome do arquivo CSV

arquivo_csv = "dados.csv"

# Lendo o arquivo CSV

# Usando o separador ';', a engine 'python' e especificando o encoding (se necessário)

# tabela = pd.read_csv(arquivo_csv, sep=';', engine='python', encoding='utf-8')

subconjunto = tabela[['ID', 'Date', 'Calories']]

# Exibindo os dados lidos

print(subconjunto)

# Exibindo os dados lidos
```

Resultado ficando da seguinte forma:

```
CONSOLE DE DEPURAÇÃO PORTAS
                                                TERMINAL
               Date Calories
    0 '2020/12/01'
0
                      4091
       '2020/12/02'
       '2020/12/03'
                       3400
       '2020/12/04'
                       2824
    4 '2020/12/05'
                       4060
       '2020/12/06'
                       3000
    6 '2020/12/07'
                       3740
       '2020/12/08'
    8 '2020/12/09'
                       1951
       '2020/12/10'
                       2690
10 10 '2020/12/11'
                       3293
11 11 '2020/12/12'
                       2507
12 12 '2020/12/12'
                       2507
       '2020/12/13'
                       3453
   14 '2020/12/14'
                       3793
15 15 '2020/12/15'
                       2750
16 16 '2020/12/16'
17 17 '2020/12/17'
                       3000
       '2020/12/18'
                        NaN
   19 '2020/12/19'
19
                       3230
20 20 '2020/12/20'
                     2430 2
   1 '2020/12/21'
21
                       3642
22
   22
               NaN
                       2820
   23 '2020/12/23'
                       3000
24 24 '2020/12/24'
                       2460
25 25 '2020/12/25'
                       3345
         20201226
                       2500
   27 '2020/12/27'
                       2410
28 28 '2020/12/28'
                        NaN
29 29 '2020/12/29'
                       2800
       '2020/12/30'
30 30
                       3803
       '2020/12/31'
                       2430
```

Microatividade 3: Descrever como configurar o número máximo de linhas a serem exibidas na visualização de um conjunto de dados usando a biblioteca Pandas (Python)

Ficando dessa forma o código completo:

Resultado ficando da seguinte forma:

```
projeto.py > ...
6  # Lendo o arquivo CSV
SAÍDA CONSOLE DE DEPURAÇÃO PORTAS COMENTÁRIOS TERMINAL
                           Date Pulse Maxpulse Calories
           60 '2020/12/01'
60 '2020/12/02'
                                                        4790
             60 2020/12/02
60 '2020/12/03'
45 '2020/12/04'
45 '2020/12/06'
60 '2020/12/06'
                                             135
175
                                    103
                                                        3400
                                     109
                                                        2824
                                               148
                                                        4060
                                     102
                                                        3000
              60 '2020/12/07'
                                     110
                                                        3740
             450 '2020/12/08'
                                     104
                                                        2533
               30 '2020/12/09'
                                     109
              60 '2020/12/10'
                                     98
                                                        2690
              60 '2020/12/11'
                                     103
                                               147
                                                        3293
              60 '2020/12/12'
                                     100
                                               120
                                                        2507
              60 '2020/12/12'
                                     100
                                               120
                                                        2507
              60 '2020/12/13'
                                                        3453
                                     106
                                               128
              60 '2020/12/14'
                                     104
                                                        3793
              60 '2020/12/15'
                                     98
                                                        2750
              60 '2020/12/16'
60 '2020/12/17'
45 '2020/12/18'
                                                        3000
              60 '2020/12/19'
45 '2020/12/20'
              60 '2020/12/21'
                                                        3642
               45
                                                        2820
              60 '2020/12/23'
45 '2020/12/24'
                                               101
                                                        3000
                                                        2460
               60 '2020/12/25'
                                     102
                                               126
                                                        3345
               60
                      20201226
                                     100
                                                        2500
               60 '2020/12/27'
                                     92
                                                        2410
              60 '2020/12/28'
                                     103
                                                         NaN
              60 '2020/12/29'
                                     100
                                                        2800
               60 '2020/12/30'
    30
                                     102
                                                        3803
               60 '2020/12/31'
```

Microatividade 4: Descrever como exibir as primeiras e últimas "N" linhas de um conjunto de dados usando a biblioteca Pandas (Python)

Procedimentos

- 1. Abra o arquivo/script utilizado nas microatividades anteriores;
- Imprima na tela as apenas as primeiras 10 linhas do conjunto de dados original
- 3. (criado na microatividade 1);
- Imprima na tela as apenas as últimas 10 linhas do conjunto de dados original
- 5. (criado na microatividade 1).

Ficando dessa forma o código completo:

```
projeto.py X
projeto.py X
projeto.py X
import pandas as pd

# Definindo o nome do arquivo CSV
arquivo_csv = "dados.csv"

# Lendo o arquivo CSV
tabela = pd.read_csv(arquivo_csv, sep=';', engine='python', encoding='utf-8')

print("Primeiras 10 linhas:")
print(tabela.head(10))

print(tabela.head(10))

# Exibindo os dados lidos
print(tabela.tail(10))
```

O resultado ficando da seguinte forma:

```
TERMINAL
Primeiras 10 linhas:
  ID Duration Date 0 60 '2020/12/01'
                           Date Pulse Maxpulse Calories
                                  110
                                               130
             60 '2020/12/02'
                                                145
                                                         4790
             60 '2020/12/03'
45 '2020/12/04'
                                    103
                                                         3400
                                    109
                                                         2824
             45 '2020/12/05'
60 '2020/12/06'
                                                148
                                                         4060
                                     102
                                                127
                                                         3000
             60 '2020/12/07'
                                    110
             450 '2020/12/08'
                                    104
                                                         2533
             30 '2020/12/09'
                                    109
                                                         1951
             60 '2020/12/10'
                                     98
                                                124
Últimas 10 linhas:
   ID Duration Date 0 60 '2020/12/01'
                           Date Pulse Maxpulse Calories
                                    110
            60 '2020/12/01'
60 '2020/12/03'
60 '2020/12/04'
                                    103
                                                         3400
                                     109
                                                175
                                                         2824
             45 '2020/12/05'
60 '2020/12/06'
                                                148
                                                         4060
                                                         3000
              60 '2020/12/07'
                  '2020/12/08'
                                     104
                  '2020/12/09'
             30
                                     109
                                                133
                                                         1951
             60 '2020/12/10'
                                     98
                                                124
                                                         2690
```

Microatividade 5: Descrever como exibir informações gerais sobre as colunas, linhas e dados de um conjunto de dados usando a biblioteca Pandas (Python)

Procedimentos

- 1. Abra o arquivo/script utilizado nas microatividades anteriores;
- 2. Tendo como base o conjunto de dados original:
 - Imprima as informações gerais sobre o conjunto suas colunas, linhas e dados;
 - 2. Descubra a partir do comando acima:
 - 1. O total de linhas;
 - 2. O total de colunas:
 - 3. A quantidade de dados nulos, caso existam;
 - 4. O tipo de dado de cada coluna;
 - 5. A quantidade de memória utilizada pelo conjunto de dados.

Ficando dessa forma o código completo:

```
projeto.py X
             dados.csv
projeto.py > ...
  import pandas as pd
      # Definindo o nome do arquivo CSV
      arquivo_csv = "dados.csv"
      # Lendo o arquivo CSV
      tabela = pd.read_csv(arquivo_csv, sep=';', engine='python', encoding='utf-8')
      print("Informações gerais do conjunto:")
      print(tabela.info())
      totalLinhas, totalColunas = tabela.shape
      print(f"\nQuantidade de linhas: {totalLinhas}")
 14
      print(f"Quantidade de colunas: {totalColunas}")
      print("\nQuantidade de valores nulos por coluna:")
      print(tabela.isnull().sum())
      print("\nTipo de dado de cada coluna:")
      print(tabela.dtypes)
      print("\nMemória utilizada pelo conjunto de dados:")
      print(tabela.memory_usage(deep=True))
```

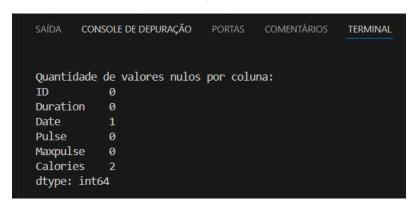
Resultado das informações gerais do conjunto:

SAÍDA	CONSOLE	DE DEPURAÇÃO	PORTAS	COMENTÁRIOS	TERMINAL
• Informações gerais do conjunto: <class 'pandas.core.frame.dataframe'=""> RangeIndex: 32 entries, 0 to 31 Data columns (total 6 columns):</class>					
#	Column	Non-Null Coun	t Dtype		
0	ID	32 non-null	int64		
1	Duration	32 non-null	int64		
2	Date	31 non-null	object	t	
3	Pulse	32 non-null	int64		
4	Maxpulse	32 non-null	int64		
5	Calories	30 non-null	object	t	

Resultado de quantidade de linhas e colunas:



Resultado de valores nulos, caso existam:



Resultado de dado de cada coluna:

```
SAÍDA CONSOLE DE DEPURAÇÃO PORTAS COMENTÁRIOS TERMINAL

Tipo de dado de cada coluna:
ID int64
Duration int64
Date object
Pulse int64
Maxpulse int64
Calories object
dtype: object
```

Resultado da quantidade de memória utilizada pelo conjunto de dados:

