Missão Prática – Mundo 05 – Nível 03 Gilvan Pereira de Oliveira – 2023.01.53256-61197

Polo Centro – São Lourenço Da Mata – PE

RPG0033 - TRATANDO A IMENSIDÃO DOS DADOS - 9001 - 2025.1

https://github.com/GilvanPOliveira/FullStack/tree/main/Mundo05/tratandoDados

Contextualização

Para resolução das micro atividades será necessário ter em mãos um conjunto de dados no formato CSV. Tais dados podem ser obtidos a partir de fontes gratuitas, na Web (como, exemplo, dataset disponível disponíveis por 0 https://archive.ics.uci.edu/dataset/352/online+retail). assim como um ambiente contendo o interpretador da linguagem python.

Ainda em relação aos conjuntos de dados, a atividade "pico web" terá como base o seguinte conjunto de dados (que deverá ser copiado e salvo num arquivo "csv",

usando como separados de colunas o";"):

		'		, ,		dados.csv X
ID	Duration	Date	Pulse	Maxpulse	Calories	data
0	60	'2020/12/01'	110	130	4091	<pre>1 ID;Duration;Date;Pulse;Maxpulse;Calories</pre>
1	60	'2020/12/02'	117	145	4790	2 0;60;'2020/12/01'; 110;130; 4091
2	60	'2020/12/03'	103	135	3400	3 1;60;'2020/12/02'; 117;145; 4790
3	45	'2020/12/04'	109	175	2824	4 2;60;'2020/12/03'; 103;135; 3400
4	45	'2020/12/05'	117	148	4060	5 3;45;'2020/12/04'; 109 <i>;</i> 175; 2824
5	60	'2020/12/06'	102	127	3000	6 4;45;'2020/12/05'; 117 <i>;</i> 148; 4060
6	60	'2020/12/07'	110	136	3740	7 5;60;'2020/12/06';102;127;3000
7	450	'2020/12/08'	104	134	2533	8 6;60;'2020/12/07'; 110 <i>;</i> 136; 3740
8	30	'2020/12/09'	109	133	1951	9 7;450; '2020/12/08';104;134;2533
9	60	'2020/12/10'	98	124	2690	10 8;30;'2020/12/09'; 109;133; 1951
10	60	'2020/12/11'	103	147	3293	11 9;60; '2020/12/10';98;124;2690
11	60	'2020/12/12'	100	120	2507	12 10;60; '2020/12/11';103;147;3293 13 11;60; '2020/12/12';100;120;2507
12	60	'2020/12/12'	100	120	2507	14 12;60; '2020/12/12';100;120;2507
13	60	'2020/12/13'	106	128	3453	15 13;60; '2020/12/13';106;128;3453
14	60	'2020/12/14'	104	132	3793	16 14;60; '2020/12/14'; 104;132; 3793
15	60	'2020/12/15'	98	123	2750	17 15;60;'2020/12/15';98;123;2750
16	60	'2020/12/16'	98	120	2152	18 16;60; '2020/12/16'; 98;120; 2152
17	60	'2020/12/17'	100	120	3000	19 17;60;'2020/12/17'; 100;120; 3000
18	45	'2020/12/18'	90	112	NaN	20 18;45;'2020/12/18'; 90 <i>;</i> 112; NaN
19	60	'2020/12/19'	103	123	3230	21 19;60;'2020/12/19'; 103;123; 3230
20	45	'2020/12/20';	97	125	2430	22 20;45;'2020/12/20'; 97 <i>;</i> 125; 2430
21	60	'2020/12/21'	108	131	3642	23 21;60;'2020/12/21'; 108;131; 3642
22	45	NaN	100	119	2820	24 22;45;NaN; 100;119; 2820
23	60	'2020/12/23'	130	101	3000	25 23;60; '2020/12/23';1 30;101; 3000
24	45	'2020/12/24'	105	132	2460	26 24;45; '2020/12/24'; 105;132; 2460
25	60	'2020/12/25'	102	126	3345	27 25;60; '2020/12/25' ; 102; 126; 3345
26	60	'2020/12/26'	100	120	2500	28 26;60;20201226;100;120;2500
27	60	'2020/12/27'	92	118	2410	29 27;60; '2020/12/27'; 92;118; 2410 30 28;60; '2020/12/28'; 103;132; NaN
28	60	'2020/12/28'	103	132	NaN	30 28;60; 2020/12/28 ;103;132; NaN 31 29;60; '2020/12/29' ;100;132; 2800
29	60	'2020/12/29'	100	132	2800	32 30;60; '2020/12/30' ;102;129; 3803
30	60	'2020/12/30'	102	129	3803	33 31;60; '2020/12/31' ;92;115; 2430
31	60	'2020/12/31'	92	115	2430	34

O uso do dataframe acima é imprescindível, uma vez que ele contém dados não válidos que deverão ser tratados posteriormente. Vide as linhas 18 e 28 (coluna Calories); 22 e 26 (coluna Date).

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

Criação do Script de leitura local

```
leitura_csv.py X

...

import pandas as pd

dadosInformados = './dados.csv'

dados = pd.read_csv(dadosInformados, sep=';', engine='python', encoding='utf-8')

print(dados)

print(dados)
```

- Importar a biblioteca pandas, necessária para a manipulação de dados;
- Definir o arquivo CSV que conterá os dados informados e a sua localização;
- Ler o arquivo CSV utilizando os seguintes parâmetros:
 - o sep: define o separador de colunas, que é o ponto e vírgula;
 - o engine: define a linguagem utilizada, que é o Python;
 - o encoding: define a codificação do arquivo, como 'uft-8', que é a mais comum.
- Exibe o conteúdo lido do arquivo CSV, no terminal.

[Ru	ınnin	g] python	-u "(\microAtividade01\leitura_csv.py"			
	ID	Duration	Date	Pulse	Maxpulse	Calories			
0	0	60	'2020/12/01'	110	130	4091.0			
1	1	60	'2020/12/02'	117	145	4790.0			
2	2	60	'2020/12/03'	103	135	3400.0			
3	3	45	'2020/12/04'	109	175	2824.0			
4	4	45	'2020/12/05'	117	148	4060.0			
5	5	60	'2020/12/06'	102	127	3000.0			
6	6	60	'2020/12/07'	110	136	3740.0			
7	7	450	'2020/12/08'	104	134	2533.0			
8	8	30	'2020/12/09'	109	133	1951.0			
9	9	60	'2020/12/10'	98	124	2690.0			
10	10	60	'2020/12/11'	103	147	3293.0			
11	11	60	'2020/12/12'	100	120	2507.0			
12	12	60	'2020/12/12'	100	120	2507.0			
13	13	60	'2020/12/13'	106	128	3453.0			
14	14	60	'2020/12/14'	104	132	3793.0			
15	15	60	'2020/12/15'	98	123	2750.0			
16	16	60	'2020/12/16'	98	120	2152.0			
17	17	60	'2020/12/17'	100	120	3000.0			
18	18	45	'2020/12/18'	90	112	NaN			
19	19	60	'2020/12/19'	103	123	3230.0			
20	20	45	'2020/12/20'	97	125	2430.0			
21	21	60	'2020/12/21'	108	131	3642.0			
22	22	45	NaN	100	119	2820.0			
23	23	60	'2020/12/23'	130	101	3000.0			
24	24	45	'2020/12/24'	105	132	2460.0			
25	25	60	'2020/12/25'	102	126	3345.0			
26	26	60	20201226	100	120	2500.0			
27	27	60	'2020/12/27'	92	118	2410.0			
28	28	60	'2020/12/28'	103	132	NaN			
29	29	60	'2020/12/29'	100	132	2800.0			
30	30	60	'2020/12/30'	102	129	3803.0			
31	31	60	'2020/12/31'	92	115	2430.0			
[Do	[Done] exited with code=0 in 8.576 seconds								

Criação do Script de leitura remoto

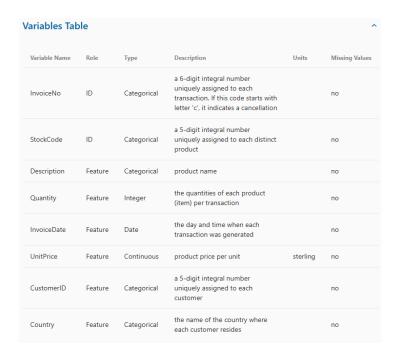
```
role
                                                                                   units missing values
                                                             name
🥏 csv_externo.py 🗙
                                                   0
                                                        InvoiceNo
                                                                        TD
                                                                                    None
                                                                                                     no
                                                   1
                                                        StockCode
                                                                        ID
                                                                                    None
                                                                                                     no
        from ucimlrepo import fetch ucirepo
                                                   2
                                                      Description Feature
                                                                                    None
                                                                                                     no
                                                   3
                                                         Quantity
                                                                  Feature
                                                                                    None
                                                                                                     no
                                                   4
                                                      InvoiceDate
                                                                   Feature
                                                                                    None
       online_retail = fetch_ucirepo(id=352)
                                                   5
                                                        UnitPrice Feature
                                                                                sterling
                                                                                                     no
                                                       CustomerID Feature
                                                                                    None
                                                                                                     no
                                                          Country Feature ...
       # data (as pandas dataframes)
                                                                                    None
                                                                                                     no
       X = online_retail.data.features
                                                   [8 rows x 7 columns]
       y = online_retail.data.targets
                                                   [Done] exited with code=0 in 5.484 seconds
       # metadata
       print(online_retail.metadata)
       print(online retail.variables)
```

- Necessário efetuar a instalação do pacote: ucimlrepo;
- Criar um novo arquivo: csv externo.py;
- O código do exemplo foi fornecido pelo link, fornecido no sway da missao;
 - o https://archive.ics.uci.edu/dataset/352/online+retail

No arquivo criado:

- Importar o pacote instalado para localizar os dados desejados;
- Selecionar o id dos dados a serem exibidos;
- Separar os dados de forma ampla e de forma específica;
- · Exibir ambas as formas no terminal.

Em cima como foi exibido no VsCode, abaixo como consta no site:

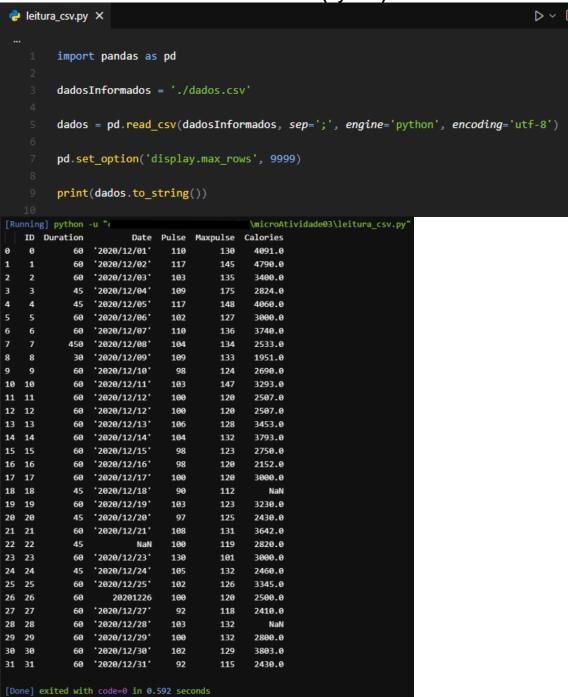


Micro atividade 2: Descrever como criar um subconjunto de dados a partir de um conjunto existente usando a biblioteca Pandas (Python)

```
🍦 leitura_csv.py 🗙
                                                                                       >
        import pandas as pd
       dadosInformados = './dados.csv'
       dados = pd.read_csv(dadosInformados, sep=';', engine='python', encoding='utf-8')
       subConjunto = dados[['ID', 'Date', 'Calories']]
       print(subConjunto)
[Running] python -u "
                                           \microAtividade02\leitura_csv.py"
           Date Calories
   0 '2020/12/01'
                   4091.0
   1 '2020/12/02'
                    4790.0
   2 '2020/12/03'
                    3400.0
                   2824.0
   3 '2020/12/04'
4
   4 '2020/12/05'
                    4060.0
   5
      '2020/12/06'
                    3000.0
6
   6
      '2020/12/07'
                    3740.0
       '2020/12/08'
      '2020/12/09'
8
   8
                    1951.0
      '2020/12/10'
9
   9
                    2690.0
10 10 '2020/12/11'
                    3293.0
11 11 '2020/12/12'
                  2507.0
12 12 '2020/12/12'
                  2507.0
13 13 '2020/12/13'
                   3453.0
14 14 '2020/12/14'
                   3793.0
15 15 '2020/12/15' 2750.0
16 16 '2020/12/16' 2152.0
17 17 '2020/12/17'
                   3000.0
18 18 '2020/12/18'
                     NaN
19 19 '2020/12/19'
                    3230.0
                  2430.0
20 20 '2020/12/20'
21 21 '2020/12/21'
                    3642.0
22 22
            NaN
                    2820.0
   23 '2020/12/23'
23
                    3000.0
   24 '2020/12/24'
24
                     2460.0
  25 '2020/12/25'
                  3345.0
25
        20201226 2500.0
26 26
27 27 '2020/12/27'
                    2410.0
28 28 '2020/12/28'
                     NaN
29 29 '2020/12/29'
                    2800.0
30 30 '2020/12/30'
                   3803.0
31 31 '2020/12/31'
                    2430.0
[Done] exited with code=0 in 1.287 seconds
```

- Utilizando-se o mesmo arquivo da micro atividade anterior;
- Foi criado um subconjunto contendo apenas as colunas 'ID', 'Date' e 'Calories';
- Foi exibido esse subconjunto.

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



- Utilizando-se o mesmo arquivo da micro atividade 01;
- Foi solicitado o uso de uma propriedade de exibição do Pandas:
 o display.max_rows definindo seu valor para 9999;
- Porém, para exibir os valores de 'dados', foi necessário utilizar o método to_string() para converter o DataFrame em uma string completa.

OBS: Foram exibidos apenas 31 linhas devido a quantidade de linhas que o arquivo 'dados' possui, porém o código foi ajustado para exibir até 9999 linhas.

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

```
🍦 leitura_csv.py 🗙
                                                                                < <</p>
       import pandas as pd
       dadosInformados = './dados.csv'
       dados = pd.read_csv(dadosInformados, sep=';', engine='python', encoding='utf-8')
       print("Primeiras 10 linhas do DataFrame:")
       print(dados.head(10))
       print("\nUltimas 10 linhas do DataFrame:")
       print(dados.tail(10))
[Running] python -u "
                                                  \microAtividade04\leitura_csv.py"
Primeiras 10 linhas do DataFrame:
       Duration
                         Date Pulse Maxpulse Calories
                 '2020/12/01'
0
    0
             60
                                 110
                                            130
                                                   4091.0
    1
             60 '2020/12/02'
                                                   4790.0
1
                                 117
                                            145
    2
                 '2020/12/03'
2
             60
                                 103
                                            135
                                                   3400.0
    3
3
             45 '2020/12/04'
                                                   2824.0
                                 109
                                            175
             45 '2020/12/05'
4
   4
                                 117
                                            148
                                                   4060.0
5
    5
                 '2020/12/06'
                                 102
                                                   3000.0
             60
                                            127
6
    6
             60 '2020/12/07'
                                 110
                                            136
                                                   3740.0
    7
                 '2020/12/08'
7
            450
                                 104
                                            134
                                                   2533.0
             30
                 '2020/12/09'
                                 109
                                            133
                                                   1951.0
8
    9
             60
                 '2020/12/10'
                                  98
                                            124
                                                   2690.0
Ultimas 10 linhas do DataFrame:
    ID Duration
                          Date Pulse Maxpulse Calories
22 22
              45
                                             119
                                                    2820.0
                           NaN
                                  100
23 23
              60 '2020/12/23'
                                  130
                                             101
                                                    3000.0
24 24
              45 '2020/12/24'
                                  105
                                                    2460.0
                                            132
25 25
              60
                  '2020/12/25'
                                  102
                                            126
                                                    3345.0
26 26
                      20201226
                                  100
              60
                                            120
                                                    2500.0
   27
              60
                  '2020/12/27'
                                   92
                                                    2410.0
27
                                             118
28 28
              60
                 '2020/12/28'
                                  103
                                             132
                                                       NaN
29
                  '2020/12/29'
   29
              60
                                  100
                                             132
                                                    2800.0
30 30
              60
                 '2020/12/30'
                                  102
                                             129
                                                    3803.0
                 '2020/12/31'
              60
                                   92
31 31
                                             115
                                                    2430.0
[Done] exited with code=0 in 0.684 seconds
```

- Utilizando-se o mesmo arquivo da micro atividade 01;
- Foi solicitado exibir as 10 primeiras e as 10 últimas linhas do arquivo DataFrame: 'dados';

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

Utilizando-se o mesmo arquivo da micro atividade 01, foi solicitado:

```
eitura_csv.py X
                                                                                    D ~
       import pandas as pd
       dadosInformados = './dados.csv'
       dados = pd.read_csv(dadosInformados, sep=';', engine='python', encoding='utf-8')
       print("Informacoes gerais sobre o conjunto de dados:")
       dados.info()
       total linhas, total colunas = dados.shape
       print(f"\nTotal de linhas: {total_linhas}")
       print(f"Total de colunas: {total_colunas}")
       print("\nQuantidade de dados nulos por coluna:")
       print(dados.isnull().sum())
       print("\nTipo de dado de cada coluna:")
       print(dados.dtypes)
       print("\nMemoria utilizada pelo conjunto de dados:")
       print(dados.memory_usage(deep=True))
```

Exibir informações gerais sobre o conjunto de dados: suas colunas, linhas e dados

```
[Running] python -u "c
                                              \microAtividade05\leitura_csv.py"
Informacoes gerais sobre o conjunto de dados:
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 32 entries, 0 to 31
Data columns (total 6 columns):
             Non-Null Count Dtype
0
   ID
             32 non-null
                            int64
    Duration 32 non-null
                            int64
2
   Date
            31 non-null
                           object
3
                            int64
    Pulse
             32 non-null
   Maxpulse 32 non-null
                             int64
    Calories 30 non-null float64
dtypes: float64(1), int64(4), object(1)
memory usage: 1.6+ KB
```

Total de Linhas e Colunas

```
Total de linhas: 32
Total de colunas: 6
```

Verificar a quantidade de dados nulos, caso existam

```
Quantidade de dados nulos por coluna:

ID 0

Duration 0

Date 1

Pulse 0

Maxpulse 0

Calories 2

dtype: int64
```

O tipo de dado de cada coluna

```
Tipo de dado de cada coluna:
ID int64
Duration int64
Date object
Pulse int64
Maxpulse int64
Calories float64
dtype: object
```

A quantidade de memória utilizada pelo conjunto de dados

```
Memoria utilizada pelo conjunto de dados:
Index
             132
ID
             256
Duration
             256
Date
            1919
Pulse
             256
Maxpulse
            256
Calories
             256
dtype: int64
[Done] exited with code=0 in 0.701 seconds
```