# Python - Pandas - Parte 2

## 1. Introdução

Nesta aula utilizaremos algumas coisas novas, a saber:

## 1.1. SciPy

Um ecossistema de ferramentas de computação científica para Python, com muitas rotinas numéricas. Mais detalhes em: https://www.scipy.org/. Em nossos código importaremos o SciPy da seguinte forma:

import scipy as sp

# 1.2. Operações de expressão regular

Uma expressão regular (ou RE) especifica um conjunto de strings que corresponde a ela; as funções neste módulo permitem verificar se uma determinada string corresponde a uma determinada expressão regular (ou se uma determinada expressão regular corresponde a uma determinada string, que se resume à mesma coisa). Mais detalhes em: <a href="https://docs.python.org/3/library/re.html">https://docs.python.org/3/library/re.html</a> . Em nossos código importaremos o Re da seguinte forma:

import re

# 1.3. display

É uma API pública para ferramentas de exibição em IPython. Em nossos código importaremos o display da seguinte forma:

from IPython.display import display

# 2. Indexação e seleção de dados

Dataframe é uma estrutura de dados bidimensional, ou seja, os dados são alinhados de forma tabular em linhas e colunas, que pode ser criado usando várias entradas:

- Listas;
- Dicionários;
- Séries;
- NumPy ndarrays;
- outro DataFrame.

Indexar no pandas significa simplesmente selecionar linhas e colunas específicas de dados de um DataFrame:

- selecionar todas as linhas e algumas das colunas;
- selecionar algumas linhas e todas as colunas;
- selecionar algumas das linhas e colunas;
- selecionar subconjuntos.

```
import pandas as pd
games = pd.read_csv(
           "vgsales.csv",
           sep=',
           header=0,
           encoding='latin1',
           names=["Rank","NomeJogo","Plataforma","ano","Gênero",
                   "Fabricante", "Vendas_EUA", "Vendas_Europa",
"Vendas_Japão", "Vendas_Resto_Mundo", "Total_Vendas"]
print(games['Total_Vendas'].head())#primeira forma de fazer
print(games.Total_Vendas.head())#segunda forma de fazer
     82.74
1
     40.24
     35.82
     33.00
     31.37
Name: Total_Vendas, dtype: float64
     82.74
0
1
     40.24
     35.82
     33.00
     31.37
Name: Total Vendas, dtype: float64
```

## 2.1. Selecionando mais de uma coluna

Também podemos selecionar mais de uma coluna para análise.

```
games = pd.read csv(
          "vgsales.csv",
          sep=',',
          header=0,
          encoding='latin1',
          names=["Rank","NomeJogo","Plataforma","ano","Gênero",
                 "Fabricante", "Vendas_EUA", "Vendas_Europa",
                 "Vendas Japão", "Vendas Resto Mundo", "Total Vendas"]
print(games[['Total_Vendas','Vendas_EUA','Fabricante']].head())
   Total_Vendas Vendas_EUA Fabricante
0
         82.74 41.49 Nintendo
1
         40.24
                    29.08 Nintendo
         35.82
33.00
31.37
2
                    15.85 Nintendo
3
                    15.75 Nintendo
4
                    11.27 Nintendo
```

#### 2.2. Selecionando trechos de linhas

Você pode selecionar trechos de linhas específicos.

```
games = pd.read_csv(
          "vgsales.csv",
          sep=',
          header=0,
          encoding='latin1',
          names=["Rank","NomeJogo","Plataforma","ano","Gênero",
                 "Fabricante","Vendas_EUA","Vendas_Europa",
"Vendas_Japão","Vendas_Resto_Mundo","Total_Vendas"]
print(games[10:15])
    Rank
                             NomeJogo Plataforma
                                                                  Gênero \
                                                      ano
                        Nintendogs DS 2005.0
Mario Kart DS DS 2005.0
10
     11
                                                              Simulation
11
                                                                  Racing
                             on Silver
Wii Fit
12
     13 Pokemon Gold/Pokemon Silver
                                               GB 1999.0 Role-Playing
                                             Wii 2007.0
                                                                  Sports
14
                         Wii Fit Plus
                                             Wii 2009.0
   Fabricante Vendas_EUA Vendas_Europa Vendas_Japão Vendas_Resto_Mundo
   Nintendo 9.07 11.00 1.93
Nintendo 9.81 7.57 4.13
10
11
    Nintendo
                     9.00
                                    6.18
                                                   7.20
                                                                        0.71
12
13
    Nintendo
                     8.94
                                     8.03
                                                   3.60
                                                                        2.15
14 Nintendo
                     9.09
                                    8.59
                                                   2.53
                                                                        1.79
    Total_Vendas
10
           24.76
11
           23.42
12
           23.10
13
           22.72
14
           22.00
```

#### 2.3. Usando loc e iloc

A seleção de subconjunto é uma das tarefas mais frequentemente executadas ao manipular dados. O Pandas fornece maneiras diferentes de selecionar eficientemente subconjuntos de dados do seu DataFrame. Para selecionar um número específico de linhas e colunas, você pode usar o .loc.

```
games = pd.read_csv(
                                                                games = pd.read csv(
         "vgsales.csv",
                                                                        "vgsales.csv'
         sep=',
         header=0.
                                                                        header=0,
         encoding='latin1',
                                                                         encoding='latin1',
         print(games.loc[10:15,['Vendas_EUA','ano']])
                                                                linhas = [10,12,13,14]
colunas = ["Vendas_EUA","Vendas_Japão"]
print(games.loc[linhas,colunas])
   Vendas EUA
                ano
        9.07
             2005.0
10
                                                                   Vendas_EUA Vendas_Japão
11
         9.81
              2005.0
                                                                        9.07
         9.00
             1999.0
                                                               12
                                                                        9.00
                                                                                     7.20
13
        8.94
             2007.0
                                                               13
                                                                        8.94
                                                                                     3.60
        9.09
              2009.0
14
15
              2010.0
```

Para selecionar um número específico de linhas e colunas, você pode fazer o seguinte usando .iloc.

#### 2.4. Usando máscaras booleanas

Podemos selecionar os dados com mascaras booleanas como em NumPy

```
import pandas as pd
from IPython.display import display
games = pd.read_csv(
          "vgsales.csv",
          sep=',',
          header=0,
          encoding='latin1',
          names=["Rank","NomeJogo","Plataforma","ano","Gênero",
                  "Fabricante", "Vendas_EUA", "Vendas_Europa",
                  "Vendas Japão", "Vendas Resto Mundo", "Total Vendas"]
amostra = games[(games.ano > 2009) & (games.Total_Vendas > 20)]
display(amostra.mean())
Rank
                        16.500
ano
                      2011.500
Vendas EUA
                        10.990
Vendas Europa
                         7.105
Vendas Japão
                         0.605
Vendas Resto Mundo
                         2.905
Total Vendas
                        21.610
dtype: float64
```

#### 3. Introdução a estatísticas

Como python você pode manipular dados, no melhor sentido da palavra, e calcular os resultados de várias operações estatísticas usando o arquivo "statistics". Principais funções:

mean(): esta função retorna média dos dados.

- mode(): retorna o número com o número máximo de ocorrências.
- mediana(): retorna a mediana, ou seja, o elemento do meio dos dados.
- median\_low(): retorna a mediana dos dados em caso de número ímpar de elementos, mas em caso de número par de elementos, retorna o menor dos dois elementos do meio.
- median\_high(): retorna a mediana dos dados no caso de número ímpar de elementos, mas no caso de número par de elementos, retorna o maior dos dois elementos.

```
#usando funções estatísticas básicas
import statistics
dados = [3,3,3,5,5,7,9,9,11,11]
print(f"Média = {statistics.mean(dados)}")
print(f"Moda = {statistics.mode(dados)}")
print(f"Mediana = {statistics.median(dados)}")
print(f"Mediana Acima= {statistics.median_high(dados)}")

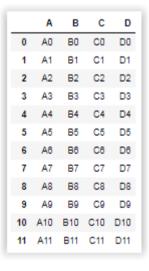
Média = 6.6
Moda = 3
Mediana = 6.0
Mediana Acima= 7
```

## 4. Usando concat, merge e joins

Um quadro de dados é uma estrutura de dados bidimensional com várias linhas e colunas. Em um quadro de dados, os dados são alinhados apenas na forma de linhas e colunas. Um quadro de dados pode executar operações aritméticas e condicionais. Tem tamanho mutável.

# 4.1. concat()

Realiza a concatenação de DataFrames.



## 4.2. merge()

Pandas fornece uma única função, merge (), como o ponto de entrada para todas as operações de junção de banco de dados padrão entre objetos DataFrame, que pode utilizar as seguintes propriedades:

- left: um DataFrame nomeado.
- right: outro DataFrame nomeado.
- on: Nomes de nível de coluna ou índice nos quais participar.
- left on: colunas ou níveis de índice do DataFrame esquerdo para usar como chaves.
- right on: colunas ou níveis de índice do DataFrame direita para usar como chaves.

- left\_index: se True, use o índice (rótulos de linha) do DataFrame esquerdo como sua(s) chave(s) de junção.
- right index: o mesmo uso left indexdo DataFrame direito.
- how: Um dos 'left', 'right', 'outer', 'inner'. O padrão é inner. Veja abaixo uma descrição mais detalhada de cada método.
- sort: classifica o DataFrame do resultado pelas chaves de junção em ordem lexicográfica.
- suffixes: uma tupla de sufixos de sequência a ser aplicada a colunas sobrepostas.
- indicator: adiciona uma coluna ao DataFrame de saída chamado \_merge com informações na origem de cada linha.

# 4.3. Join

O join() usa merge internamente para a junção de índice em índice (por padrão) e coluna(s) em índice.

## 5. Exercícios - base de dados personagens.csv

- Mostre os 5 primeiros registros da base de dados
- Repita o item anterior, mas mostrando apenas nome e poder
- Mostre os 5 primeiros registros da base de dados, que tem poder maior que 90
- Mostre a média do poder

•	Mostre a mediana de combate
•	