panda@ufscar:~\$ intro

# Introdução ao pandas em Python

Por Joãozinho



#### O que é pandas?



- Pacote do Python, criado para manipular dados tabulares
- Muitas similaridades com SQL
  - o **SELECT** vira indexação com [ ]
  - o WHERE vira indexação booleana
  - O GROUP BY, JOIN, UNION...
- Análise exploratória, limpeza dos dados, interface com outras
  - bibliotecas de data science
- Não confundir com o nome do grupo (foi mal)

#### Mexer com dados não é uma tarefa fácil

- Dados faltantes
- Formatos e encodings inutilizáveis
  - o Exemplo: dados advindos de scraping
- Dados duplicados
- Outliers
- Falta de documentação
- ...E outras tarefas de pré-processamento

Ninguém: a

O mundo se todos os dados fossem limpos automaticamente:



#### You better work, data engineer!



Work B\*\*ch - Britney Spears

### Antes, um pouco sobre arrays em Python

```
import numpy as np
minha matriz = np.arange(0, 20).reshape(4,5)
print(minha matriz)
[[0 1 2 3 4]
[56789]
[10 11 12 13 14]
 [15 16 17 18 19]
 [20 21 22 23 24]]
```

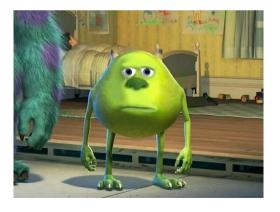
0	1	2	3	4
5	6	7	8	9
10	11	12	13	14
15	16	17	18	19
20	21	22	23	24

```
minha_matriz[linha, coluna]
print(minha_matriz[1,2])
7
```

0	1	2	3	4
5	6	7	8	9
10	11	12	13	14
15	16	17	18	19
20	21	22	23	24

minha\_matriz[linha, coluna]

print(minha\_matriz[-1,-1])
24



0	1	2	3	4
5	6	7	8	9
10	11	12	13	14
15	16	17	18	19
20	21	22	23	24

```
minha_matriz[linha, coluna]
print(minha_matriz[1,:])
[5 6 7 8 9]
```

0	1	2 3		4
5	6	7	8	9
10	11	12	13	14
15	16	17	18	19
20	21	22	23	24

```
minha_matriz[linha, coluna]
print(minha_matriz[:4,2:4])
[[ 2 3]
 [ 7 8]
 [12 13]
 [17 18]]
```

0	1	2	3	4
5	6	7	8	9
10	11	12	13	14
15	16	17	18	19
20	21	22	23	24

#### Outra maneira!

minha\_matriz[linha, coluna]

print(minha\_matriz[:-1,2:4])

[[ 2 3]

[ 7 8]

[12 13]

[17 18]]



#### Exemplo: Livros do John Green

- Vamos criar um dataframe simples contendo os livros do John Green
- Título, data de publicação, se eu já li, minha avaliação e a média de avaliação no Goodreads



#### Series pandas.Series

- Estrutura mais básica do pandas
- Array nomeado ou não, de 1 dimensão e que possui um índice
- Pode ser uma observação ou uma característica

Usada para formar <u>dataframes ou é derivad</u>a de dataframes



```
livros_jg_series

0 Quem é Você, Alasca?
1 O Teorema Katherine
2 Deixe a Neve Cair
3 Cidades de Papel
4 Will e Will
5 A Culpa É das Estrelas
dtype: object
```

```
datas_series

0 2005-03-03
1 2006-09-21
2 2008-10-02
3 2008-10-16
4 2010-04-06
5 2012-01-10
dtype: datetime64[ns]
```

```
lidos = pd.Series([True, False, False, True, True, False])
lidos
      True
     False
     False
     True
     True
     False
dtype: bool
minha_nota = pd.Series([5, np.nan, np.nan, 4.5, 4, np.nan])
minha_nota
     5.0
    NaN
    NaN
     4.5
     4.0
    NaN
dtype: float64
nota_gr = pd.Series([4.02, 3.57, 3.76, 3.80, 3.77, 4.21])
nota_gr
     4.02
     3.57
     3.76
     3.80
     3.77
     4.21
dtype: float64
```

	livros	ano_publicacao	lido	minha_nota	nota_goodreads
0	Quem é Você, Alasca?	2005-03-03	True	5.0	4.02
1	O Teorema Katherine	2006-09-21	False	NaN	3.57
2	Deixe a Neve Cair	2008-10-02	False	NaN	3.76
3	Cidades de Papel	2008-10-16	True	4.5	3.80
4	Will e Will	2010-04-06	True	4.0	3.77
5	A Culpa É das Estrelas	2012-01-10	False	NaN	4.21
6	Tartarugas Até Lá Embaixo	2017-10-10	True	4.5	3.95



#### Esqueci de algo... 😌

dtype: object

### Agora sim! 🧐

	livro	ano_publicacao	lido	minha_nota	nota_goodreads
0	Quem é Você, Alasca?	2005-03-03	True	5.0	4.02
1	O Teorema Katherine	2006-09-21	False	NaN	3.57
2	Deixe a Neve Cair	2008-10-02	False	NaN	3.76
3	Cidades de Papel	2008-10-16	True	4.5	3.80
4	Will e Will	2010-04-06	True	4.0	3.77
5	A Culpa É das Estrelas	2012-01-10	False	NaN	4.21
6	Tartarugas Até Lá Embaixo	2017-10-10	True	4.5	3.95

#### Análise exploratória

- Descobrir padrões nos dados
- Encontrar problemas
- Testar hipóteses
- Confirmar intuições
- Se exibir pros amigos



	livro	data_publicacao	lido	minha_nota	nota_goodreads	
0	Quem é Você, Alasca?	2005-03-03	True	5.0	4.02	
1	O Teorema Katherine	2006-09-21	False	NaN	3.57	
2	Deixe a Neve Cair	2008-10-02	False	NaN	3.76	
3	Cidades de Papel	2008-10-16	True	4.5	3.80	
4	Will e Will	2010-04-06	True	4.0	3.77	
5	A Culpa É das Estrelas	2012-01-10	False	NaN	4.21	
6	Tartarugas Até Lá Embaixo	2017-10-10	True	4.5	3.95	
			df[[	'data_publ	icacao','livro'	]]
df['livro']			da	ta_publicacao	I	ivro
0 Quemé	é Você, Alasca?		0	2005-03-03	Quem é Você, Ala	sca
1 0 Tec	orema Katherine		1	2006-09-21	O Teorema Kathe	rine
	ixe a Neve Cair idades de Papel		2	2008-10-02	Deixe a Neve	Cai
4	. Will e Will		3	2008-10-16	Cidades de P	ape
	É das Estrelas Até Lá Embaixo		4	2010-04-06	Will e	Wil
Name: livro, dty	/pe: object		5	2012-01-10	A Culpa É das Estr	elas
			6	2017-10-10	Tartarugas Até Lá Emb	aixc

	atherine leve Cair de Papel l e Will Estrelas Embaixo	6 Tart		ma Ka a Ne des de Will das E é Lá	therine ve Cair e Papel e Will strelas Embaixo	1 2 3 4 5 A 6 Tarta	Quem é Você, Alasca? O Teorema Katherine Deixe a Neve Cair Cidades de Papel Will e Will Culpa É das Estrelas arugas Até Lá Embaixo ro, dtype: object
df['livro']		df.loc[:,	'livro']			df.iloc[:,	0]
		s Até Lá Embaixo	2012-01-10		NaN 4.5	4.21 3.95	
	4 F. A.C.	Will e Will pa É das Estrelas	2010-04-06	True	4.0	3.77	
	3	Cidades de Papel	2008-10-16	True	4.5	3.80	
	2	eixe a Neve Cair	2008-10-02	False	NaN	3.76	
	1 O Te	orema Katherine	2006-09-21	False	NaN	3.57	
	<b>U</b> Quen	i e voce, Alasca:	2005-03-03	irue	5.0	4.02	

Quem é Você Alasca?

0

livro data\_publicacao lido minha\_nota nota\_goodreads

5.0

4 02

2005-03-03 True

#### df loc[5:6]

5

uı	٠	COC	[5.0	J	

A Culpa É das Estrelas

Tartarugas Até Lá Embaixo

2012-01-10 False

2017-10-10 True

NaN

4.5

livro data\_publicacao lido minha\_nota nota\_goodreads

4.21

3.95

Quen	n é Você, Alasc	a? 2005-03-03	True	5.0	4.02
О Те	orema Katherir	ne 2006-09-21	False	NaN	3.57
0	eixe a Neve Ca	air 2008-10-02	False	NaN	3.76
(	Cidades de Pap	el 2008-10-16	True	4.5	3.80
	Will e W	ill 2010-04-06	True	4.0	3.77
A Cul	pa É das Estrel	as 2012-01-10	False	NaN	4.21
rugas	s Até Lá Embaix	co 2017-10-10	True	4.5	3.95
	df.des	cribe()			
	ı	minha_nota	nota_	goodreads	
	count	4.000000		7.000000	
	mean	4.500000		3.868571	
	std	0.408248		0.208761	
	min	4.000000		3.570000	
	25%	4.375000		3.765000	
	50%	4.500000		3.800000	
	75%	4.625000		3.985000	
	max	5.000000		4.210000	

6 Tartar

livro ano\_publicacao lido minha\_nota nota\_goodreads

# df['nota\_goodreads'].describe() count 6.000000 mean 3.855000 std 0.225278

25% 3.762500 50% 3.785000 75% 3.965000 max 4.210000 Name: nota\_goodreads, dtype: float64

3.570000

min

#### df.grouphy(by=df[']ido'])['nota\_goodreads'].describe()

arigioups) (	by-ui[	cruo j	/[ "	ora_g	joour	cuus	] rucser isc()
count	mean	std	min	25%	50%	75%	max

## lido

False

True

- 3.0 3.846667 0.328684 3.57 3.665 3.76 3.985 4.21
- 3.0 3.863333 0.136504 3.77 3.785 3.80 3.910 4.02

# Brigado 🖭 Agora vamos ao Jupyter!!!!

