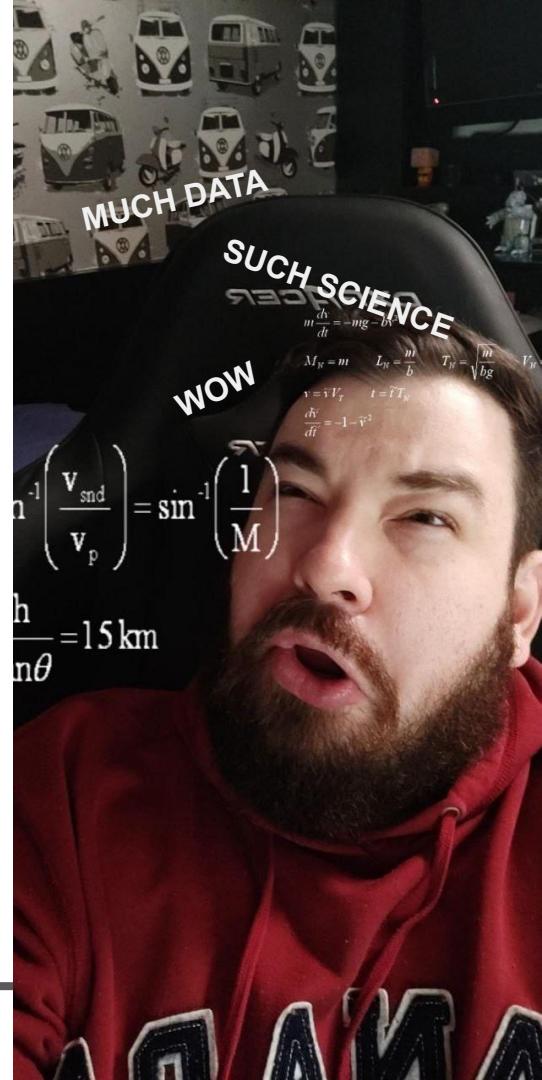


Ferramentas de Softwares para uso em ciência de Dados II

OBJETIVO

Nasser Santiago Boan

- > arquivologia - UNB
- > tecnologia do Petróleo e Gás - UFRJ
- > MBA Processos de Negócio - IBMEC
- > IBM Data Science Professional Certified
- > líder de Ciência de Dados - Certsys
- > Cientista de Dados - Stefanini
- > docente Pós-Graduação Ciência de Dados



Ementa

- Básico
 - pip
 - conda
 - miniconda
 - jupyter notebook
- Nivelamento 2
 - Variáveis e operações matemáticas
 - Operações com strings
 - Listas
 - Tuplas
 - Dicionários
 - Condições lógicas
 - Loops (for + compreensão de lista + while)
 - Funções (def + lambda)
 - Funções built-in (zip, enumerate, map e filter)



Ementa

- Files
 - Lendo arquivos com OPEN
 - Escrevendo arquivos com OPEN
- Numpy
 - Numpy array
 - Operações com arrays
 - Indexing e slicing
 - Stacking
 - Funções estatísticas (mean, var, std, median, quantiles)

Ementa

- Pandas
 - Dataframes
 - read_ / to_
 - Indexing e slicing (.loc , .iloc , filter , head, sample)
 - Tipos de dados e datas (to_datetime)
 - Funções estatísticas (describe, corr, skew, mode, unique, nunique)
 - Transformações (apply + operações com séries)
 - Reshaping e sorting (pivot + transpose + nlargest + nsmallest + value_counts + sort_values)
 - Concat e merge
 - Valores nulos (isna, isnull, fillna-- imputação --, drop, dropna, replace)
 - Method Chaining

Ementa

- matplotlib
 - Arquitetura (pyplot e axes)
 - Anatomia dos gráficos em matplotlib
 - Plotando (figure, plt.bar, plt.scatter, plt.plot, title, labels, legend, colors, axvline, axhline)
 - Subplots
 - Abrindo imagens
 - Exportando gráficos (plt.savefig)

Ementa

- Scikit-Learn
 - Pre-processamento (standartscaler, minmaxscaler, onehotencoding)
 - Conceitos de Regressão / Classificação / Clusterização
 - Dataset de treino e teste (train_test_split)
 - Underfitting / overfitting
 - Regressão linear simples (na mão + coeficientes explícitos)
 - Métricas de regressão e função de custo
 - Regressão linear múltipla e polynomial features
 - DecisionTree / RandomForest
 - Métricas de classificação
 - KMeans
 - métricas de clusterização
 - Kfold e GridsearchCV

Avaliação

- Lista 01 (python)
- Lista 02 (pandas + matplotlib)
- Projeto 01 (análise de dados)
- Lista 03 (conceitos de machine learning 1)
- Lista 04 (conceitos de machine learning 2)
- Projeto 02 (projeto ML)

Módulo Básico

pip

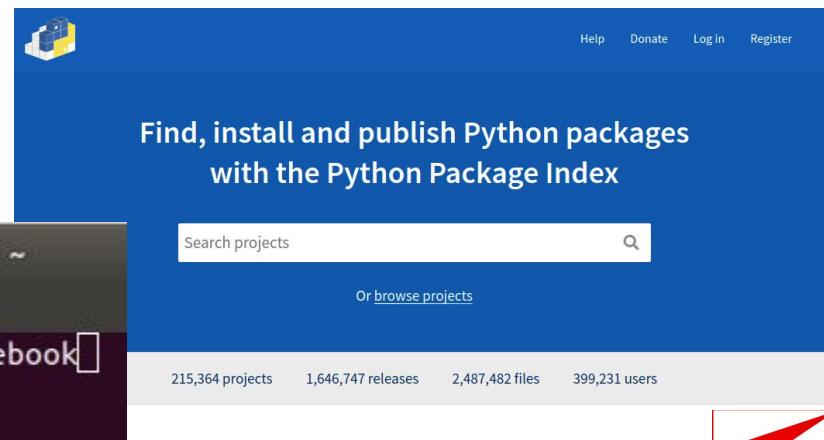
O pip é o gerenciador de pacotes do python.

```
$ pip install <nome do pacote>
```

```
$ pip --version
```

```
$ pip list
```

```
$ pip install -r requirements.txt
```



conda e miniconda

O conda vai além e gerênciambientes e dependências

```
$ conda create env -n <nome do ambiente>
```

```
$ conda create env -n <nome do ambiente> python=3.6
```

```
$ conda env create -f environment.yml
```

```
$ conda list
```

```
$ conda install <nome do pacote>
```

```
$ conda activate <nome do ambiente>
```

```
$ conda deactivate <nome do ambiente>
```



conda e miniconda

O conda vai além e gerênciambientes e dependências

```
(base) nzboan@ubutop:~$ conda list
# packages in environment at /home/nzboan/miniconda3:
#
# Name           Version    Build  Channel
_libgcc_mutex   0.1        main
_tflow_select   2.1.0      gpu
absl-py         0.8.1      py37_0
altair          3.2.0      py37_0
asnincrypto     1.2.0      py37_0
astor           0.8.0      py37_0
astroid          2.3.3      py37_0
attrs            19.3.0     py_0
backcall         0.1.0      py37_0
beautifulsoup4  4.8.1      py37_0
biopython        1.74       pypi_0  pypi
blas             1.0        mkl
bleach           3.1.0      py37_0
blosc            1.16.3     hd408876_0
bokeh            1.4.0      py37_0
boto3            1.9.232    pypi_0  pypi
botocore         1.12.232   pypi_0  pypi
branca           0.3.1      py_0
bzzip2           1.0.8      h7b6447c_0
c-ares            1.15.0     h7b6447c_1001
ca-certificates  2019.11.27  0
cachetools       3.1.1      pypi_0  pypi
cairo             1.14.12    h8948797_3
certifi           2019.11.28  py37_0
cffi              1.13.2     py37h2e261b9_0
chardet          3.0.4      py37_1003
click             7.0        py37_0
click-plugins     1.1.1      pypi_0  pypi
cligj             0.5.0      pypi_0  pypi
cloudpickle      1.2.2      py_0
cogroo-interface 0.3        pypi_0  pypi
colorcet          2.0.2      py_0
conda            4.8.1      py37_0
```



conda e miniconda

O conda vai além e gerênciambientes e dependências

```
(base) nzboan@ubutop:~$ conda activate new
(base) nzboan@ubutop:~$ conda list
# packages in environment at /home/nzboan/miniconda3/envs/new:
#
# Name           Version      Build Channel
(base) nzboan@ubutop:~$
```

```
In [5]: from sklearn.datasets import load_iris
```

```
-----
ModuleNotFoundError          Traceback (most recent call
last)
<ipython-input-5-56d11ecab3a7> in <module>
----> 1 from sklearn.datasets import load_iris

ModuleNotFoundError: No module named 'sklearn'
```

conda e miniconda

```
(new) nzboan@ubutop:~$ conda install scikit-learn
Collecting package metadata (current_repodata.json): done
Solving environment: done

## Package Plan ##

environment location: /home/nzboan/miniconda3/envs/new

added / updated specs:
- scikit-learn

The following packages will be downloaded:

      package          build
-----
ca-certificates-2020.1.1           0
certifi-2019.11.28                  py38_0
ld_impl_linux-64-2.33.1            h53a641e_7
mkl-service-2.3.0                   py38he904b0f_0
mkl_fft-1.0.15                     py38ha843d7b_0
mkl_random-1.1.0                   py38h962f231_0
numpy-1.18.1                       py38hf9e942_0
numpy-base-1.18.1                  py38hde5b4d6_1
openssl-1.1.1d                     h7b6447c_3
pip-20.0.2                          py38_1
python-3.8.1                        h0371630_1
scikit-learn-0.22.1                 py38hd81dba3_0
scipy-1.3.2                         py38h7c811a0_0
setuptools-45.1.0                   py38_0
six-1.14.0                          py38_0
sqlite-3.30.1                      h7b6447c_0
wheel-0.34.1                        py38_0
-----
                                         Total: 8
```

```
In [2]: from sklearn.datasets import load_iris
```

```
In [9]: data = load_iris()
```

```
In [10]: data.target
```



conda e miniconda

```
! environment.yml
1  name: stats
2  dependencies:
3    - jupyterlab
4    - pandas
5    - numpy
6    - scikit-learn
7
```

```
(base) nzboan@ubutop:~$ conda env create -f environment.yml
(base) nzboan@ubutop:~$ conda env create -f environment.yml
Collecting package metadata (repodata.json): done
Solving environment: done

Downloading and Extracting Packages
pandas-1.0.0      | 8.8 MB      | #####| 100%
jupyterlab-1.2.5  | 2.8 MB      | #####| 100%
Preparing transaction: done
Verifying transaction: done
Executing transaction: done
#
# To activate this environment, use
#
#     $ conda activate stats
#
# To deactivate an active environment, use
#
#     $ conda deactivate
(base) nzboan@ubutop:~$
```

jupyter notebook

O jupyter notebook é documento que contém tanto código quanto elementos de texto. Ele executa código célula a célula, pode ser lido por humanos (human-readable) ou pela máquina (machine-readable).

google colab

Welcome To Colaboratory

File Edit View Insert Runtime Tools Help

Share Connect Editing

+ Code + Text Copy to Drive

What is Colaboratory?

Colaboratory, or "Colab" for short, allows you to write and execute Python in your browser, with

- Zero configuration required
- Free access to GPUs
- Easy sharing

Whether you're a **student**, a **data scientist** or an **AI researcher**, Colab can make your work easier. Watch [Introduction to Colab](#) to learn more, or just get started below!

Getting started

The document you are reading is not a static web page, but an interactive environment called a **Colab notebook** that lets you write and execute code.

For example, here is a **code cell** with a short Python script that computes a value, stores it in a variable, and prints the result:

```
[ ] 1 seconds_in_a_day = 24 * 60 * 60
2 seconds_in_a_day
```

86400

To execute the code in the above cell, select it with a click and then either press the play button to the left of the code, or use the keyboard shortcut "Command/Ctrl+Enter". To edit the code, just click the cell and start editing.

Variables that you define in one cell can later be used in other cells:

```
[ ] 1 seconds_in_a_week = 7 * seconds_in_a_day
2 seconds_in_a_week
```

604800

referências importantes

instalação do conda

<https://docs.conda.io/projects/conda/en/latest/user-guide/install/index.html#regular-installation>

instalação do pip

<https://pip.pypa.io/en/stable/installing/>

instalação do jupyter notebook

<https://jupyter.org/install.html>

google colab

<https://colab.research.google.com/>

Nivelamento (again...)

<https://github.com/aulas-iesb/CastoresIndomaveis>

Files e Numpy

File Handling

Em python não há necessidade de importar bibliotecas para ler e escrever arquivos. As funções built-in open e write resolvem bem a grande maioria das necessidades.

The screenshot shows a Jupyter Notebook environment. On the left is a file browser with a folder named 'files-numpy /'. It contains two items: 'mod-files-numpy-file-handling.ipynb' (modified 'seconds ago') and 'primeiro_arquivo.txt' (modified 'a minute ago'). A red arrow points from the 'primeiro_arquivo.txt' entry to the code cell on the right. The code cell at the top right contains the line: [1]: `f = open('primeiro_arquivo.txt', 'w+')`. Below it is an empty output cell: []:.

open()

A função open() abre um arquivo e retorna um objeto. A função possui vários e o objeto pode ou não ser alterado de acordo com o modo passado. Os modos que a função pode receber:

Modo	Descrição
'r'	Modo padrão da função. Abre um arquivo para leitura.
'w'	Esse modo abre um arquivo para escrita. Se o arquivo não existir esse modo cria o arquivo. Se o arquivo já existe ele substitui o arquivo.
'x'	Cria um arquivo novo. Se o arquivo já existir a função retorna um erro
'a'	Abre um arquivo no modo 'append'. Se o arquivo não existe ele cria um novo.
'+'	Esse modo abre um arquivo para leitura e escrita.

read()

```
[33]: f = open('lorem.txt','r')
conteudo = f.read()
f.close()
```

```
[35]: print(conteudo)
```

 Lorem ipsum dolor sit amet, consectetur adipiscing elit. Aliquam sagittis quis massa sed euismod. Aenean at mattis nisl. Vestibulum a imperdiet lacus. Duis hendrerit, justo maximus convallis placerat, purus nibh vulputate nibh, ac varius libero est eu orci. Fusce ac dolor id ex dictum mattis. Donec f inibus pellentesque ligula, et semper tellus varius maximus. In rutrum vitae quam nec suscipit. Orci varius natoque penatibus et magnis dis parturient montes, nascetur ridiculus mus. Nullam laoreet dolor non felis tincidunt dignissim. Suspendisse lacinia accumsan risus vitae rhoncus. Nulla luctus con vallis enim, ut blandit arcu dignissim at. Pellentesque pharetra ornare efficitur. In molestie, quam in tristique ullamcorper, purus nulla accumsan m i, et lobortis purus magna quis risus. Cras magna ipsum, accumsan in purus eget, tincidunt feugiat est. Proin laoreet velit a nunc ullamcorper varius.

 Cras vel pharetra lacus. In imperdiet at erat ac suscipit. Donec eget urna vel lectus fermentum maximus. Quisque eu tempus nisl. Sed ultricies lacus q uis condimentum efficitur. Quisque non metus at eros eleifend eleifend. Fusce et semper nisi, lobortis consequat erat. Praesent condimentum pulvinar n unc in scelerisque. Pellentesque habitant morbi tristique senectus et netus et malesuada fames ac turpis egestas.

 Duis in maximus lacus, vel eleifend ligula. Integer in elementum leo. Nam interdum pharetra neque, a lacinia arcu. Donec scelerisque velit sit amet ju sto sagittis, eu commodo quam viverra. Vivamus non est viverra mi vehicula scelerisque. Donec id suscipit ligula. Nullam condimentum semper elit, non accumsan justo consectetur a. Quisque non nisl elit. Ut ultrices risus ut tortor tincidunt accumsan. Curabitur aliquam mi nec mi aliquam lacinia. Maec enas eget felis non mi ultrices suscipit.

readlines()

```
[37]: f = open('lines.txt','r')
conteudo = f.readlines()
for i in conteudo:
    print(i)
```

linha 1

linha 2

linha 3

linha 4

linha 5

linha 6

write()

```
[51]: with open('primeiro_arquivo.txt', 'w') as f:  
    for i in range(10):  
        f.write(f'{i}\n')
```

The screenshot shows a Jupyter Notebook interface. The top bar has two tabs: 'mod-files-numpy-file-handlir X' and 'primeiro_arquivo.txt'. The main area displays a code cell with the following content:

```
1 | 0  
2 | 1  
3 | 2  
4 | 3  
5 | 4  
6 | 5  
7 | 6  
8 | 7  
9 | 8  
10| 9  
11|
```

The code cell contains a Python script that uses a 'with' statement to open a file named 'primeiro_arquivo.txt' in write mode ('w'). It then enters a 'for' loop that iterates 10 times, writing the value of 'i' followed by a new line character ('\n') to the file.

csv.reader()

```
[38]: import csv

[40]: with open('avocado.csv') as csv_file:
        csv_reader = csv.reader(csv_file, delimiter=',')
        for row in csv_reader:
            print(row)

[', 'Date', 'AveragePrice', 'Total Volume', '4046', '4225', '4770', 'Total Bags', 'Small Bags', 'Large Bags', 'XLarge Bags', 'type', 'year', 'region']
['0', '2015-12-27', '1.33', '64236.62', '1036.74', '54454.85', '48.16', '8696.87', '8603.62', '93.25', '0.0', 'conventional', '2015', 'Albany']
['1', '2015-12-20', '1.35', '54876.98', '674.28', '44638.81', '58.33', '9505.56', '9408.07', '97.49', '0.0', 'conventional', '2015', 'Albany']
['2', '2015-12-13', '0.93', '118220.22', '794.7', '109149.67', '130.5', '8145.35', '8042.21', '103.14', '0.0', 'conventional', '2015', 'Albany']
['3', '2015-12-06', '1.08', '78992.15', '1132.0', '71976.41', '72.58', '5811.16', '5677.4', '133.76', '0.0', 'conventional', '2015', 'Albany']
['4', '2015-11-29', '1.28', '51039.6', '941.48', '43838.39', '75.78', '6183.95', '5986.26', '197.69', '0.0', 'conventional', '2015', 'Albany']
```

csv.reader()

```
[44]: avg_price = []

    with open('avocado.csv') as csv_file:
        csv_reader = csv.reader(csv_file,delimiter=',')
        next(csv_reader)
        for row in csv_reader:
            avg_price.append(row[2])
```

```
[47]: avg_price[:10]
```

```
[47]: ['1.33',
       '1.35',
       '0.93',
       '1.08',
       '1.28',
       '1.26',
       '0.99',
       '0.98',
       '1.02',
       '1.07']
```

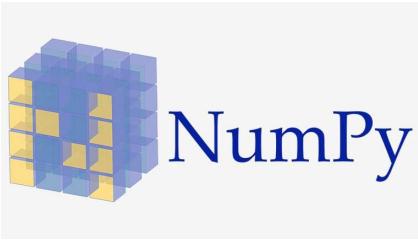
csv.writer()

```
[76]: import csv

with open('empregados.csv', 'w') as f:
    empregados_writer = csv.writer(f, delimiter=',', quotechar='''')

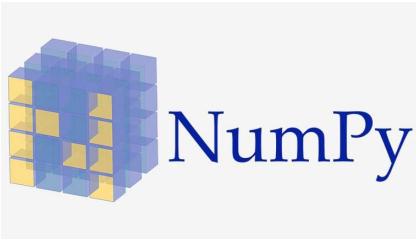
    empregados_writer.writerow(['nome', 'idade', 'salario'])
    empregados_writer.writerow(['nasser', 29, 1000])
    Delimiter: ,
```

	nome	idade	salario
1	nasser	29	1000
2	joao	20	1500



```
[2]: import numpy as np  
  
a = np.array([[1,2,2],[2,2,2]])  
a
```

```
[2]: array([[1, 2, 2],  
           [2, 2, 2]])
```



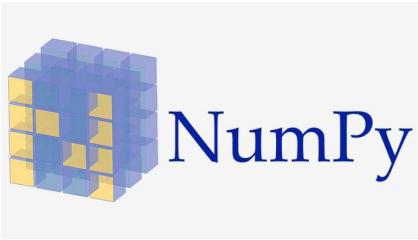
NumPy

```
[3]: a.ndim
```

```
[3]: 2
```

```
[31]: np.zeros((40,50,50,50)).ndim
```

```
[31]: 4
```



NumPy

```
[38]: a.shape
```

```
[38]: (2, 3)
```

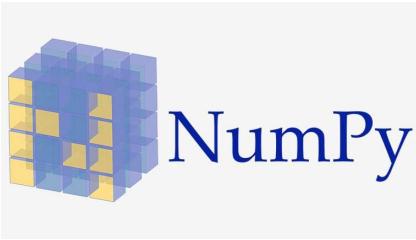
```
[13]: b = np.array([[1,2,3,4,5],[6,7,8,9,10],[11,12,13,14,15]])
```

```
b
```

```
[13]: array([[ 1,  2,  3,  4,  5],  
           [ 6,  7,  8,  9, 10],  
           [11, 12, 13, 14, 15]])
```

```
[39]: b.shape
```

```
[39]: (3, 5)
```



NumPy

```
[43]: a.size
```

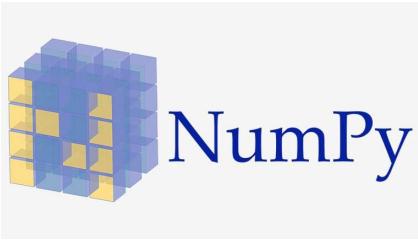
```
[43]: 6
```

```
[44]: b.size
```

```
[44]: 15
```

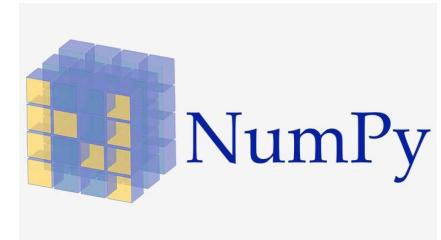
```
[42]: np.zeros((40,50,50,50)).size
```

```
[42]: 5000000
```



```
[61]: a + 1  
  
[61]: array([[2, 3, 3],  
           [3, 3, 3]])  
  
[62]: a ** 2  
  
[62]: array([[1, 4, 4],  
           [4, 4, 4]])  
  
[63]: a - 10  
  
[63]: array([[-9, -8, -8],  
           [-8, -8, -8]])  
  
[66]: (a*3) > 5  
  
[66]: array([[False,  True,  True],  
           [ True,  True,  True]])
```

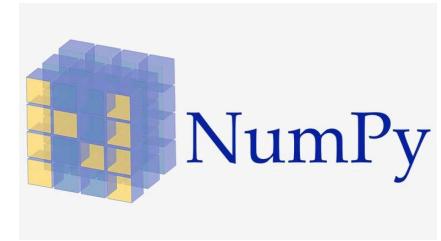
operações



```
[52]: a + b
-----
ValueError                                Traceback (most recent call last)
<ipython-input-52-bd58363a63fc> in <module>
      1 a + b
ValueError: operands could not be broadcast together with shapes (2,3) (3,5)
```

```
[57]: c = np.array([[2,1,1],[1,1,1]])
c
[57]: array([[2, 1, 1],
           [1, 1, 1]])
```

operações



```
[57]: c = np.array([[2,1,1],[1,1,1]])
c
```

```
[57]: array([[2, 1, 1],
           [1, 1, 1]])
```

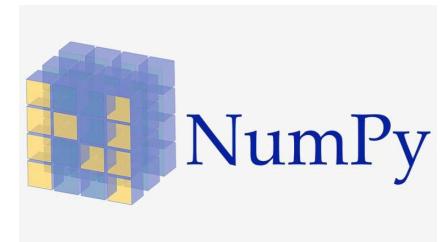
```
[56]: a + c
```

```
[56]: array([[3, 3, 3],
           [3, 3, 3]])
```

```
[67]: a - c
```

```
[67]: array([[-1,  1,  1],
           [ 1,  1,  1]])
```

operações



```
[69]: array1 = np.array([[1,2,3],[4,5,6],[7,8,9]])  
array1
```

```
[69]: array([[1, 2, 3],  
           [4, 5, 6],  
           [7, 8, 9]])
```

```
[71]: array2 = np.array([[  
array2
```

```
[71]: array([[5, 5, 5, 5,  
           [5, 5, 5, 5,  
           [5, 5, 5, 5,
```

```
[73]: array1@array2
```

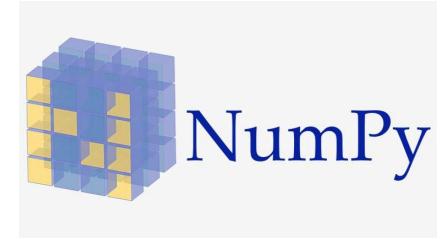
```
[73]: array([[ 30,  30,  30,  30,  30,  30],  
           [ 75,  75,  75,  75,  75,  75],  
           [120, 120, 120, 120, 120, 120]])
```

```
[79]: a = np.arange(12).reshape(3,4)  
a
```

```
[79]: array([[ 0,  1,  2,  3],  
           [ 4,  5,  6,  7],  
           [ 8,  9, 10, 11]])
```

<http://matrixmultiplication.xyz/>

indexing & slicing



```
[79]: a = np.arange(12).reshape(3,4)  
a
```

```
[79]: array([[ 0,  1,  2,  3],  
           [ 4,  5,  6,  7],  
           [ 8,  9, 10, 11]])
```

```
[82]: a[1][2]
```

```
[82]: 6
```

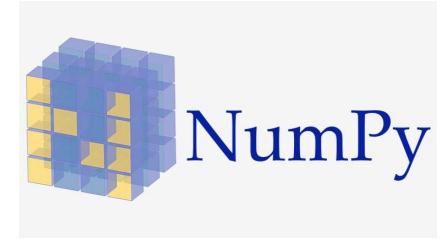
```
[83]: a[:2]
```

```
[83]: array([[0, 1, 2, 3],  
           [4, 5, 6, 7]])
```

```
[86]: a[::-2]
```

```
[86]: array([[ 0,  1,  2,  3],  
           [ 8,  9, 10, 11]])
```

indexing & slicing



```
[87]: a[:::-1]
```

```
[87]: array([[ 8,  9, 10, 11],
           [ 4,  5,  6,  7],
           [ 0,  1,  2,  3]])
```

```
[89]: index_bol = a > 5
```

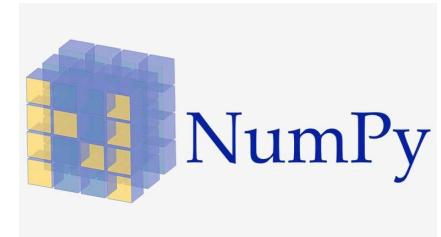
```
[90]: a[index_bol]
```

```
[90]: array([ 6,  7,  8,  9, 10, 11])
```

```
[91]: a
```

```
[91]: array([[ 0,  1,  2,  3],
           [ 4,  5,  6,  7],
           [ 8,  9, 10, 11]])
```

stacking



```
a
```

```
[79]: array([[ 0,  1,  2,  3],  
           [ 4,  5,  6,  7],  
           [ 8,  9, 10, 11]])
```

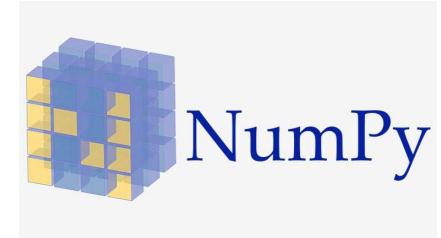
```
[96]: b
```

```
[96]: array([[ 0,  1,  2,  3],  
           [ 4,  5,  6,  7],  
           [ 8,  9, 10, 11],  
           [12, 13, 14, 15]])
```

```
[99]: np.vstack_(a,b)
```

```
[99]: array([[ 0,  1,  2,  3],  
           [ 4,  5,  6,  7],  
           [ 8,  9, 10, 11],  
           [ 0,  1,  2,  3],  
           [ 4,  5,  6,  7],  
           [ 8,  9, 10, 11],  
           [12, 13, 14, 15]])
```

funções estatísticas

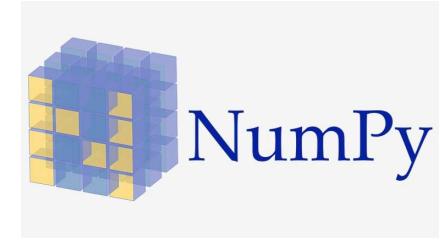


```
[101]: a = np.arange(24).reshape(8,3)

[102]: a

[102]: array([[ 0,  1,  2],
              [ 3,  4,  5],
              [ 6,  7,  8],
              [ 9, 10, 11],
              [12, 13, 14],
              [15, 16, 17],
              [18, 19, 20],
              [21, 22, 23]])
```

funções estatísticas



```
[103]: a.mean()  
[103]: 11.5  
  
[104]: a.var()  
[104]: 47.916666666666664  
  
[105]: a.std()  
[105]: 6.922186552431729  
  
[106]: np.median(a)  
[106]: 11.5  
  
[109]: np.quantile(a,0.25)  
[109]: 5.75
```