

CODING IN BRAZIL

# Ciência de dados

Wellington Silva

Brasil

2021

# Sumário

<b>1</b>	<b>Introdução . . . . .</b>	<b>1</b>
1.1	Ferramentas . . . . .	2
1.2	Material complementar . . . . .	4
<b>2</b>	<b>Python . . . . .</b>	<b>6</b>
2.1	Gerenciamento de pacote . . . . .	6
2.2	Essencial . . . . .	6
2.3	Matemática . . . . .	7
2.4	Visualização de dados . . . . .	7
2.5	Web Scraping . . . . .	7
2.6	Aprendizagem de máquina . . . . .	9
2.7	Processamento de imagens . . . . .	9
	<b>REFERÊNCIAS . . . . .</b>	<b>10</b>

## Lista de ilustrações

Figura 1 – Ilustração das etapas da extração de conhecimento . . . . .	1
Figura 2 – NumPy . . . . .	6
Figura 3 – Pandas . . . . .	7
Figura 4 – Matplotlib . . . . .	8
Figura 5 – Seaborn . . . . .	8

# 1 Introdução

## Origem

A mineração de dados surgiu como área de pesquisa e aplicação independente em meados da década de 1990. Entretanto, as suas origens na matemática, estatística e computação são muito anteriores a esse período.

## Objetivo

Preparação e análise das grandes massas de dados, tendo a finalidade de encontrar o conhecimento. Portanto, para cumprir tal finalidade, reuni áreas distintas, como estatística; matemática; engenharia; inteligência artificial; banco de dados; sistemas de informação; visualização; antropologia; e o especialista do domínio dos dados, que se complementam e formam a área de ciência de dados.

## KDD

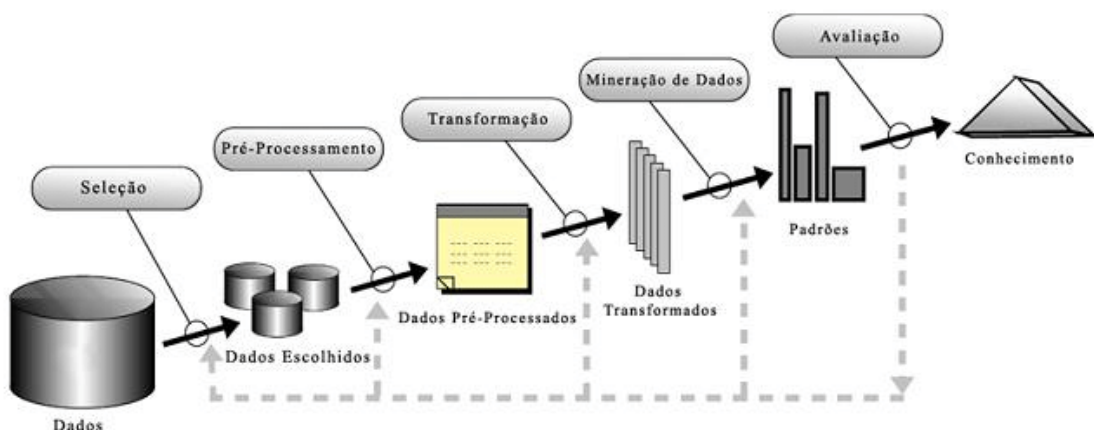


Figura 1 – Ilustração das etapas da extração de conhecimento

1. Dados: conjunto de dados organizados de forma *qualitativa* ou *quantitativa* sobre determinado tema, no qual possibilidade a extração de informação que pode resultar em conhecimento.
2. Pré-processamento dos dados: Selecionar os dados de acordo com a demanda do estudo, descartando assim dados irrelevantes, a fim de tornar a análise dos eficiente e eficaz. As etapas são distribuídas:
  - limpeza: remoção de ruídos de dados inconsistentes e ausentes;
  - integração: combinação dos dados de diferentes fontes;

- seleção: escolha de dados relevantes à análise; e
  - transformação: consolidação dos dados em formato apropriado.
3. Mineração de dados: Utilização de métricas e medidas estatísticas, para representar o conjunto de dados e a sua distribuição. Tais medidas são análise descritiva, agrupamento, predição, associação e detecção de anomalias.
  4. Avaliação: Identificar os padrões obtidos pela representação do conhecimento são válidos, ou seja, representativo.

## 1.1 Ferramentas

Edição auxiliada por histórico de comando:

IPython é um interpretador interativo para várias linguagens de programação, mas especialmente focado em Python.

- JupyterLab
- Jupyter Notebook
- Colab Notebooks

Kaggle (<https://www.kaggle.com/>)

Tradicional:

- \* Vim
- \* Atom
- \* Sublime-text
- \* VSCode
- \* Spyder3
- \* PyDev
- \* PyCharm da JetBrains
- \* Komodi IDE
- \* [Kite](<https://www.kite.com/>)

IDLE

- Weka ([www.cs.waikato.ac.nz/ml/weka](http://www.cs.waikato.ac.nz/ml/weka))
- Matlab
- R Studio ([www.r-project.org](http://www.r-project.org))
  - Bioconductor ([www.bioconductor.org](http://www.bioconductor.org))
- Wolfram Mathematica ([www.wolfram.com/mathematica](http://www.wolfram.com/mathematica))

- RapidMiner ([rapidminer.com](http://rapidminer.com))
- SAS ([sas.com](http://sas.com))
- SSPS by IBM ([www-01.ibm.com/software/analytics/spss](http://www-01.ibm.com/software/analytics/spss))
- Orange ([orange.biolab.si](http://orange.biolab.si))
- Mahout by Apache ([mahout.apache.org](http://mahout.apache.org))
- ELKI ([elki.dbs.ifi.lmu.de](http://elki.dbs.ifi.lmu.de)): aprendizagem não supervisionado
- LIBSVM ([www.csie.ntu.edu.tw/~cjlin/libsvm](http://www.csie.ntu.edu.tw/~cjlin/libsvm))

#### Visualização de dados

- \* Tableau: visualização dinâmica dos dados e dashboard personalizado
- \* PowerBI: integração com os serviços Microsoft
- \* Data Studio: integração com os serviços Google
- \* Qlik: licitação

#### SGBD

- \* [PostgreSQL]([www.postgresql.org](http://www.postgresql.org))
- \* [ORACLE]([www.oracle.com](http://www.oracle.com))
- \* [SQL Server]([www.microsoft.com/en-us/server-cloud/products/sql-server](http://www.microsoft.com/en-us/server-cloud/products/sql-server))
- \* [DB2]([www-01.ibm.com/software/data/db2](http://www-01.ibm.com/software/data/db2))

#### Apache

- \* [Airflow](<https://airflow.apache.org/>): plataforma de gerenciamento de fluxo de trabalho de código aberto

#### Integração dos dados

- \* [Kafka](<https://kafka.apache.org/>)
- \* [Spark](<https://spark.apache.org/>)
- \* [NiFi](<https://nifi.apache.org/>)

#### Microserviços

- \* Docker
- \* Docker Compose
- \* Kubernetes

## 1.2 Material complementar

Livros:

- Introdução a mineração de dados por Ferrari e Silva (2017)
- Data Science para Negócios por Fawcett e Provost (2018)
- Python para análise de dados por McKinney (2019)
- Introdução à Ciência de Dados Fundamentos e Aplicações <sup>1</sup>

Cursos:

- ML4all - UFPR <sup>2</sup>

Blog:

- DIKW by Towards Data Science <sup>3</sup>
- Curso R <sup>4</sup>
- Tests as linear by Lindeloev <sup>5</sup>
- JTemporal <sup>6</sup>

Base de dados:

- UCI Machine Learning Repository <sup>7</sup>
- KDnuggets <sup>8</sup>
- Governo Brasileiro <sup>9</sup>
  - Brasil IO <sup>10</sup>
  - Gasto de parlamentar <sup>11</sup>
- Governo Americano <sup>12</sup>

---

<sup>1</sup> <https://www.ime.usp.br/~jmsinger/MAE5755/cdados2019ago06.pdf>

<sup>2</sup> <http://cursos.leg.ufpr.br/ML4all/1parte/>

<sup>3</sup> <https://towardsdatascience.com/rootstrap-dikw-model-32cef9ae6dfb>

<sup>4</sup> <https://blog.curso-r.com/>

<sup>5</sup> <https://lindeloev.github.io/tests-as-linear/>

<sup>6</sup> <https://jtemporal.com/>

<sup>7</sup> <http://archive.ics.uci.edu/ml/index.php>

<sup>8</sup> <https://www.kdnuggets.com/datasets/index.html>

<sup>9</sup> <https://dados.gov.br/>

<sup>10</sup> <https://brasil.io/>

<sup>11</sup> <https://serenata.ai/>

<sup>12</sup> <https://www.data.gov/>

- Governo do Inglês <sup>13</sup>
- PyData Book <sup>14</sup>

---

<sup>13</sup> <https://data.gov.uk/>

<sup>14</sup> <https://github.com/wesm/pydata-book>



## 2 Python

### 2.1 Gerencimaneto de pacote

Python Package Index - PIP <sup>15</sup> Sistema de gerenciamento de pacotes padrão de facto usado para instalar e gerenciar pacotes de software escritos em Python. Muitos pacotes podem ser encontrados na fonte padrão para pacotes e suas dependências.

Anaconda <sup>16</sup> Distribuição gratuita e de código aberto das linguagens de programação Python e R para computação científica, que visa simplificar o gerenciamento e a implantação de pacote.

Conda <sup>17</sup> Gerenciador de pacotes e sistema de gerenciamento de ambiente de código aberto, plataforma cruzada e independente de linguagem.

### 2.2 Essencial

NumPy é um pacote para a linguagem Python que suporta arrays e matrizes multidimensionais, possuindo uma larga coleção de funções matemáticas para trabalhar com estas estruturas.



Figura 2 – NumPy

Pandas é uma biblioteca de software criada para a linguagem Python para manipulação e análise de dados. Em particular, oferece estruturas e operações para manipular tabelas numéricas e séries temporais. O nome é derivado de painel data.

<sup>15</sup> <https://pypi.org/project/pip>

<sup>16</sup> <https://www.anaconda.com>

<sup>17</sup> <https://docs.conda.io/en/latest>

IMPORTANT METHODS IN PANDAS PACKAGE		
DATA IMPORTING	DATA CLEANING	DATA STATISTICS
<ul style="list-style-type: none"> <li>• <code>pd.read_csv()</code></li> <li>• <code>pd.read_table()</code></li> <li>• <code>pd.read_excel()</code></li> <li>• <code>pd.read_sql()</code></li> <li>• <code>pd.read_json()</code></li> <li>• <code>pd.read_html()</code></li> <li>• <code>pd.read_clipboard()</code></li> <li>• <code>pd.DataFrame()</code></li> <li>• <code>pd.concat()</code></li> <li>• <code>pd.Series()</code></li> <li>• <code>pd.date_range()</code></li> </ul>	<ul style="list-style-type: none"> <li>• <code>df.dropna()</code></li> <li>• <code>df.fillna()</code></li> <li>• <code>df.describe()</code></li> <li>• <code>df.sort_values()</code></li> <li>• <code>df.groupby()</code></li> <li>• <code>df.apply()</code></li> <li>• <code>df.append()</code></li> <li>• <code>df.join()</code></li> <li>• <code>df.rename()</code></li> <li>• <code>df.set_index()</code></li> <li>• <code>df.to_csv()</code></li> </ul>	<ul style="list-style-type: none"> <li>• <code>df.head()</code></li> <li>• <code>df.tail()</code></li> <li>• <code>df.info()</code></li> <li>• <code>df.describe()</code></li> <li>• <code>df.mean()</code></li> <li>• <code>df.median()</code></li> <li>• <code>df.std()</code></li> <li>• <code>df.corr()</code></li> <li>• <code>df.count()</code></li> <li>• <code>df.max()</code></li> <li>• <code>df.min()</code></li> </ul>

Figura 3 – Pandas

## 2.3 Matemática

SymPy é uma biblioteca Python para computação simbólica. Ela fornece ferramentas de álgebra computacional tanto como uma aplicação independente como, também, uma biblioteca para outras aplicações.

SciPy é uma biblioteca Open Source em linguagem Python que foi feita para matemáticos, cientistas e engenheiros. Também tem o nome de uma popular conferência de programação científica com Python.

Statsmodels é um pacote Python que permite aos usuários explorar dados, estimar modelos estatísticos e executar testes estatísticos

## 2.4 Visualização de dados

Matplotlib<sup>18</sup> é uma biblioteca para geração de gráficos e visualizações de dados em geral, feita para e da linguagem de programação Python e sua extensão de matemática NumPy.

Seaborn<sup>19</sup>

## 2.5 Web Scraping

LXML<sup>20</sup>

<sup>18</sup> <https://matplotlib.org/>

<sup>19</sup> <https://seaborn.pydata.org/>

<sup>20</sup> <https://lxml.de>

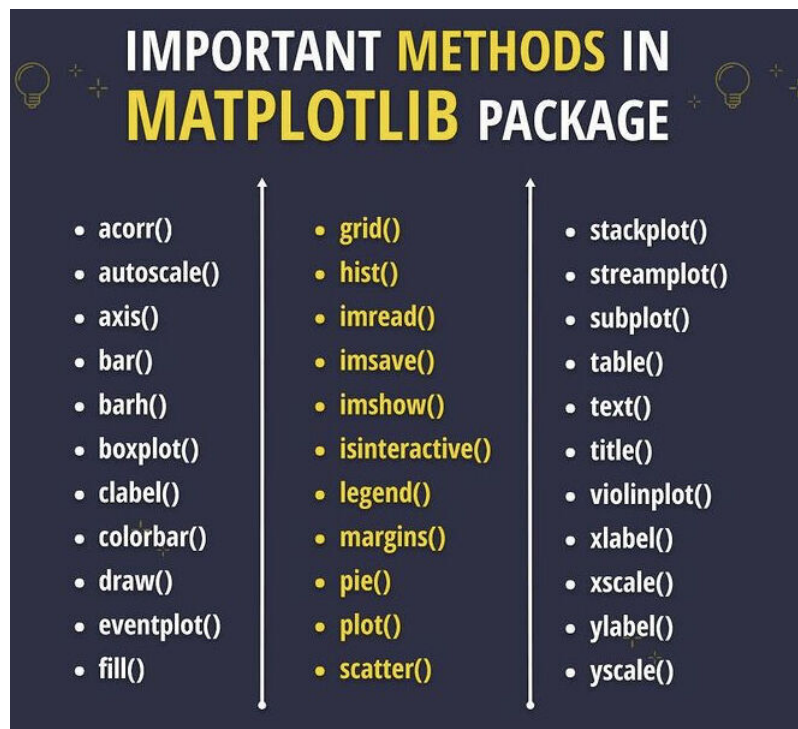


Figura 4 – Matplotlib



Figura 5 – Seaborn

HTLM5LIB <sup>21</sup>

Beautiful Soup <sup>22</sup>

## 2.6 Aprendizagem de máquina

Scikit-learn <sup>23</sup> Biblioteca de aprendizado de máquina de código aberto para a linguagem de programação Python

TensorFlow <sup>24</sup> Biblioteca de código aberto para aprendizado de máquina aplicável a uma ampla variedade de tarefas. É um sistema para criação e treinamento de redes neurais para detectar e decifrar padrões e correlações

Keras <sup>25</sup> Biblioteca de rede neural de código aberto escrita em Python. Ele é capaz de rodar em cima de TensorFlow, Microsoft Cognitive Toolkit, R, Theano, ou PlaidML. Projetado para permitir experimentação rápida com redes neurais profundas, ele se concentra em ser fácil de usar, modular e extensível

## 2.7 Processamento de imagens

Python Imaging Library Biblioteca da linguagem de programação Python que adiciona suporte à abertura e gravação de muitos formatos de imagem diferentes.

OpenCV Biblioteca multiplataforma, totalmente livre ao uso acadêmico e comercial, para o desenvolvimento de aplicativos na área de Visão computacional

Scikit-image Biblioteca de processamento de imagens de código aberto para a linguagem de programação Python. Inclui algoritmos para segmentação, transformações geométricas, manipulação do espaço de cores, análise, filtragem, morfologia, detecção de recursos e muito mais.

PyTorch Biblioteca de aprendizado de máquina de código aberto baseada na biblioteca Torch, usada para aplicativos como visão computacional e processamento de linguagem natural

---

<sup>21</sup> <https://html5lib.readthedocs.io/en/latest>

<sup>22</sup> <https://www.crummy.com/software/BeautifulSoup/bs4/doc>

<sup>23</sup> <https://scikit-learn.org>

<sup>24</sup> <https://www.tensorflow.org>

<sup>25</sup> <https://keras.io>

## Referências

FAWCETT, T.; PROVOST, F. *Data Science para Negócios: O que você precisa saber sobre mineração de dados e pensamento analítico de dados*. [S.l.]: Alta Books Editora, 2018.

FERRARI, D. G.; SILVA, L. N. D. C. *Introdução a mineração de dados*. [S.l.]: Saraiva Educação SA, 2017.

MCKINNEY, W. *Python para análise de dados: Tratamento de dados com Pandas, NumPy e IPython*. [S.l.]: Novatec Editora, 2019.