

# Análise de Dados e Inferência Estatística com Python



-



Existem muitas ferramentas para Visualização de Dados, e a cada dia novas ferramentas surgem. Nós iremos focar na ferramenta que é provavelmente a mais difundida para este fim no Python, que é o **matplotlib**.

O matplotlib já está incluso na distribuição Anaconda (e em praticamente todas as distribuições de Python focadas em Análise de Dados e programação científica). Então, se você a está utilizando, não há necessidade de qualquer instalação.

Bem, vamos começar com um exemplo bem simples. Em primeiro lugar, é uma convenção usualmente adotada importar a coleção de comandos `pyplot` do `matplotlib`. De acordo com a documentação oficial, esta coleção de comandos faz com que o `matplotlib` funcione em estilo semelhante ao `MATLAB`.

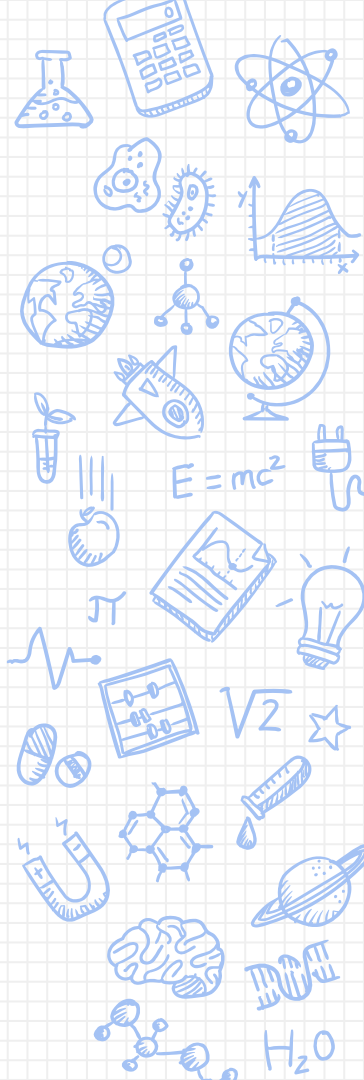
Normalmente, importa-se esta coleção como plt. Feito isso, no matplotlib, cada comando executa uma alteração no gráfico, como criação da área, traçado dos pontos, mudança do label nos eixos, e o comando final exibe o gráfico.

Vamos passar uma lista para o comando `plot()` e depois usar o comando `show()` para exibir o gráfico.

# Traçando os primeiros gráficos

Por padrão, ao receber uma lista, o comando `plot()` traça um gráfico de linha, com o índice da lista no eixo X e os itens no eixo Y.

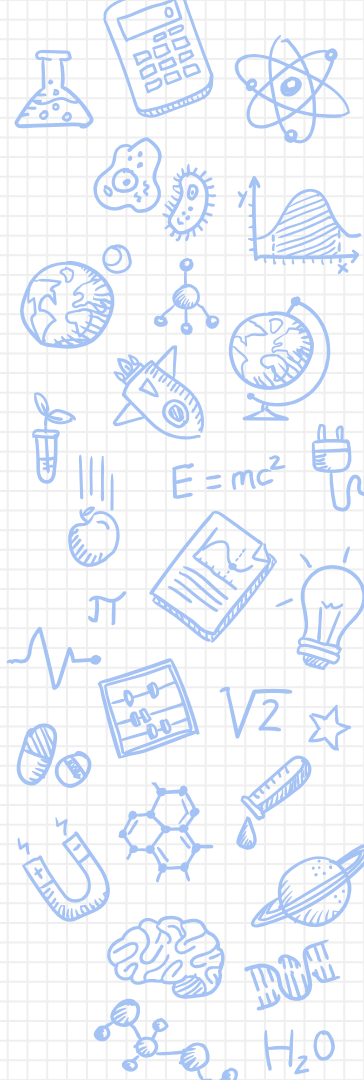
Para o segundo exemplo, vamos definir os pontos do eixo X. Adicionalmente, um terceiro argumento passado à função `plot()` define o formato dos pontos e da linha do gráfico. Vamos fazer o gráfico com quadrados vermelhos e uma linha tracejada, passando como argumentos *color*, *linestyle* e *marker*. Usaremos “r” para a cor *red*, *linestyle* “-” para o tracejado e *marker* “s” para *square*, ou seja, quadrado. Por fim, vamos deixar a linha mais grossa com o parâmetro *linewidth*. Vamos também definir o eixo com a função *axis*, para melhorar um pouco a visualização. No *axis* você passa, dentro de uma lista, os valores para o eixo X e depois para o eixo Y, indicando os pontos iniciais e finais para cada eixo.



# Gráfico de Barras

---

Outro gráfico popular presente no matplotlib é o gráfico de barras. Para o gráfico de barras, usamos a função `bar()`, onde definimos a posição das barras no eixo X e sua altura no eixo Y. Adicionalmente, também podemos configurar outras características, como a espessura das barras, cor, entre outros. O eixo X será um *range* com a mesma quantidade de itens do eixo Y. Vejamos um exemplo bem simples, onde vamos guardar as configurações em variáveis e então passá-los para a função.



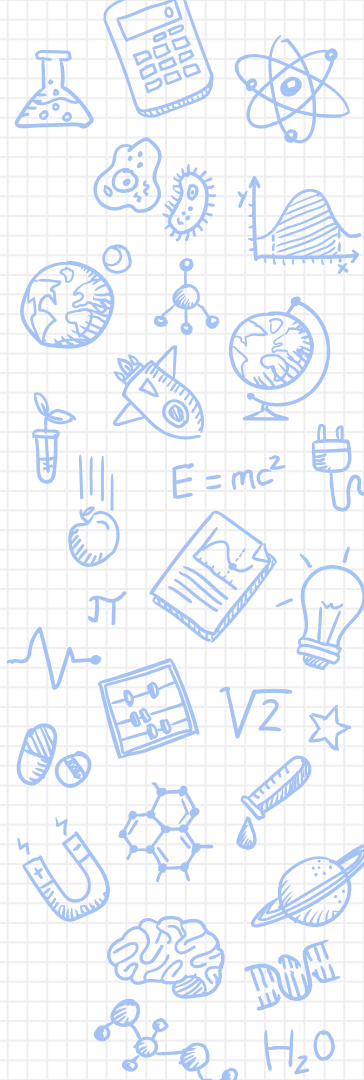


# Gráfico de Barras

---

Agora vamos usar o *Dataset* do Titanic para fazer um gráfico de barras empilhadas, onde poderemos visualizar, para cada sexo, os sobreviventes.

Primeiro usaremos a função `pivot_table()` do Pandas para criar uma tabela que agregue esses dados mencionados, passando qual deve ser a variável das linhas e colunas, a função que deve agregar os valores (pode ser soma, contagem, entre outros) e os valores aos quais esta função será aplicada, dividindo entre os valores já definidos nas linhas e colunas. Para quem já trabalhou com tabela dinâmica no Excel, esta função fica mais fácil de entender.

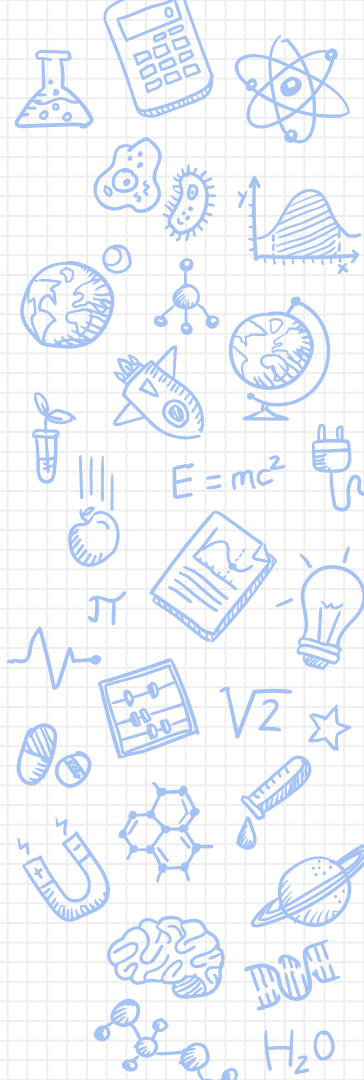


Porém, melhor ainda para fazer esta comparação seriam dois gráficos de pizza (*pie charts*), mostrando percentualmente quem sobreviveu e quem não sobreviveu para cada sexo.

# Gráficos de pizza e gráficos múltiplos

Para gráficos múltiplos usamos a função `subplots()`. Ela cria o que o matplotlib chama de “figure”, que seria o espaço onde os múltiplos gráficos podem ser criados. Os espaços onde são criados os gráficos funcionam como listas, ou *arrays*. Nos casos de uma linha ou coluna, os espaços são *arrays* de uma dimensão, e nos casos de mais de uma linha e coluna, os *arrays* são de duas dimensões.

Aqui criaremos uma *figure* com uma linha e duas colunas, onde colocaremos nossos gráficos de pizza lado a lado usando a função *pie*, e definindo nela os valores que irão compor o gráfico, *labels* para esses valores, vamos incluir a porcentagem de cada valor relativo ao total da pizza, junto com sua formatação e as cores de cada pedaço. Depois disso, incluiremos um título em cada *subplot* e a opção `axis('equal')` garantirá que os gráficos serão redondos (meramente estético).



Também é possível configurar para que o tamanho dos pontos varie de acordo com o valor de uma dada variável, bastando passar a dada variável para o parâmetro 's', de *size*, ao invés do parâmetro 'c', de *color*.

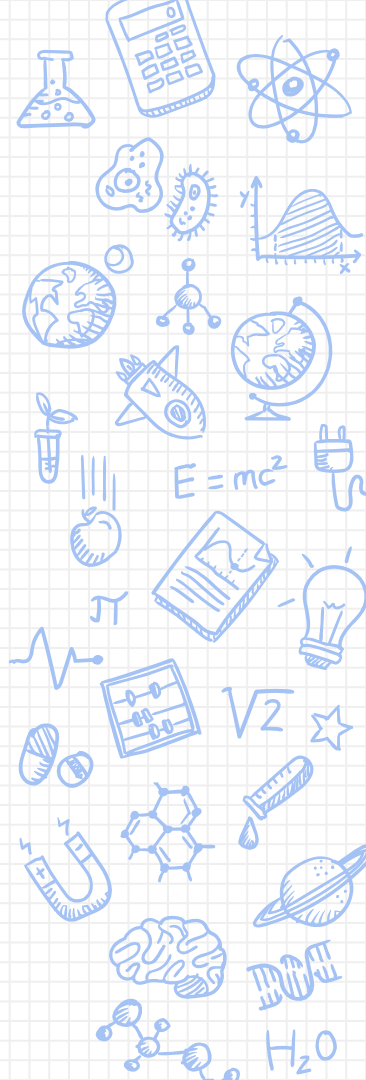


# Outros gráficos no Matplotlib

---

Existe ainda uma série de outros tipos de gráficos no matplotlib, como diagramas de caixa (*boxplot*), gráfico radar e outros.

Os gráficos apresentados nesta aula são alguns dos mais úteis e utilizados na análise de dados.



Probabilidade é a medida numérica que representa a chance de um determinado evento acontecer. Qual é a chance de eu jogar um dado e tirar um número 6? De eu jogar duas moedas para cima e as duas resultarem em Coroa? Essas perguntas (e muitas outras) são respondidas pela Probabilidade.

Para o cálculo básico da probabilidade de um evento, somamos todas as possibilidades que atendem este evento e dividimos pelo número total de possibilidades. Por exemplo, a probabilidade de rolar um 3 em um dado só é atendida pelo número 3. Entretanto, temos 6 possibilidades em um dado, os números que vão de 1 a 6 (o 3 inclusive). Desta forma, temos uma possibilidade que atende ao evento desejado dividido por 6 possibilidades no total, ou  $1/6$ .

Um primeiro conceito importante da Probabilidade é o conceito de eventos dependentes e independentes.

Os eventos independentes são eventos nos quais saber o resultado de um deles não nos dá qualquer informação sobre o possível resultado do segundo. O rolar de um dado não nos dirá nada sobre o resultado do rolar do segundo dado.

Analogamente, eventos dependentes são eventos nos quais o resultado de um nos dá mais informações sobre a probabilidade do outro. Neste caso, o rolar de um dado nos dá informações se queremos saber a probabilidade, por exemplo, de rolar 6 em dois dados.



# Eventos dependentes e independentes

---

Para eventos independentes, a probabilidade de dois eventos acontecerem é igual a probabilidade de um multiplicada pela probabilidade do outro. Desta forma, se a probabilidade de rolar um determinado valor no dado é de  $1/6$ , a probabilidade de rolar dois 6 seguidos é calculada da seguinte forma:

Sendo  $P(A)$  a probabilidade do evento A,  $P(B)$  a probabilidade do evento B e  $P(A,B)$  a probabilidade de ocorrerem os eventos A e B:

$$P(A,B) = P(A) \times P(B)$$



$$E = mc^2$$



Vejamos por exemplo, o seguinte caso. Qual a probabilidade de jogar dois dados e os dois números serem diferentes? Neste caso, a lista de resultados que atendem as condições entre os 36 possíveis são muitos. Mas o evento complementar deste é quando jogamos dois dados e os dois números rolados são iguais. Desta forma, os resultados que atendem a condição limitam-se a (1,1), (2,2), (3,3), (4,4), (5,5) e (6,6). Assim, podemos facilmente calcular a probabilidade do evento complementar em  $6/36$ , e a probabilidade do evento desejada em  $1 - 6/36$ , equivalente a  $30/36$  ou  $5/6$ .

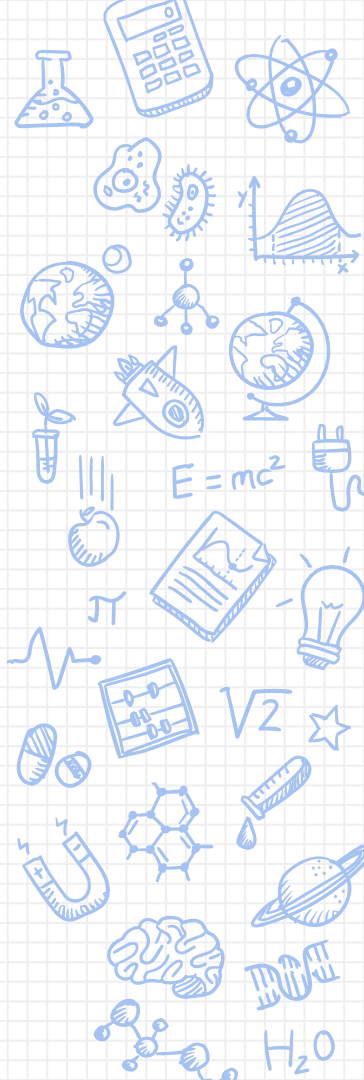
Eventos mutuamente exclusivos ocorrem quando, entre dois eventos, apenas um pode acontecer. Ao jogar uma moeda, ela tem que resultar em cara ou coroa. Os dois não podem acontecer simultaneamente. Em um dado, o número rolado tem que ser par ou ímpar. Nenhum número atende às duas condições simultaneamente.

# Eventos Mutuamente Exclusivos

Para dois eventos mutuamente exclusivos A e B,  $P(A, B)$ , ou seja, a probabilidade que os dois ocorram, é igual a zero. Entretanto, podemos calcular a probabilidade de  $P(A \text{ ou } B)$  que, para estes tipos de eventos, será igual a  $P(A) + P(B)$ .

Para o caso dos dados, por exemplo,  $P(A \text{ ou } B)$  será igual a 1, pois o número rolado necessariamente será par ou ímpar.

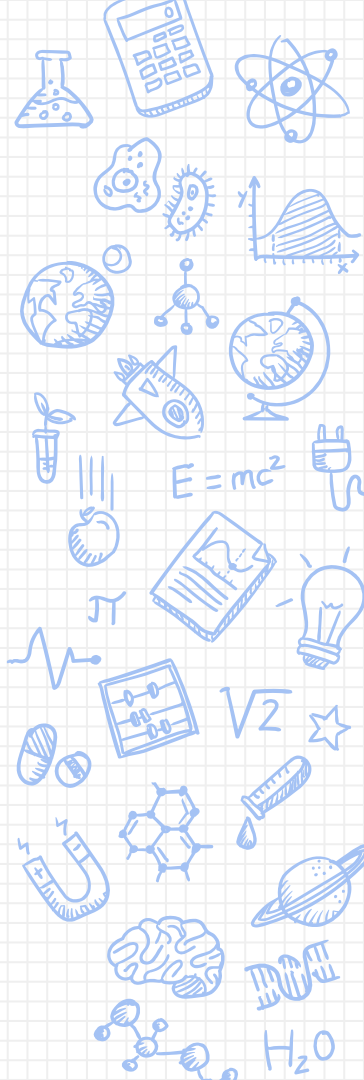
Para outro exemplo, em um baralho, consideremos a probabilidade de uma carta ser um valete ( $P(A)$ ), igual a  $4/52$ , ou  $1/13$ , e a probabilidade dela ser um Ás ( $P(B)$ ), também  $1/13$ . Estas duas probabilidades dão-se por termos 4 de cada uma destas cartas em um baralho de 52 cartas. São eventos mutuamente exclusivos, pois uma carta não pode ser os dois. Desta forma,  $P(A, B)$  é igual a zero. Entretanto, a probabilidade dela ser um Valete ou um Ás ( $P(A \text{ ou } B)$ ) será igual à soma de  $P(A)$  e  $P(B)$ , totalizando  $1/13 + 1/13 = 2/13$ .



# Probabilidade Condicional

Para eventos dependentes, o cálculo da Probabilidade muda. Vamos estabelecer que:

- $P(A|B)$  -> Probabilidade condicional de A dado B, ou seja, probabilidade do evento A ocorrer, dado que ocorreu o evento B.
- $P(A,B)$  -> Como já vimos, é a probabilidade dos dois eventos ocorrerem.
- $P(A)$  e  $P(B)$  -> Também, como já vimos, é a probabilidade de cada evento acontecer.



# Probabilidade Condicional

Para eventos dependentes, o cálculo é o seguinte:

$$P(A|B) = P(A,B)/P(B)$$

E algumas vezes, passamos  $P(B)$  para o outro lado da igualdade, e a equação fica assim:

$$P(A,B) = P(A|B) \times P(B)$$

Vejamos um exemplo e um pouco de código. Consideremos um dado. Seja o primeiro evento tirar um número ímpar e o segundo tirar um 5 ou 6 no dado. Vamos calcular a probabilidade de A dado B.

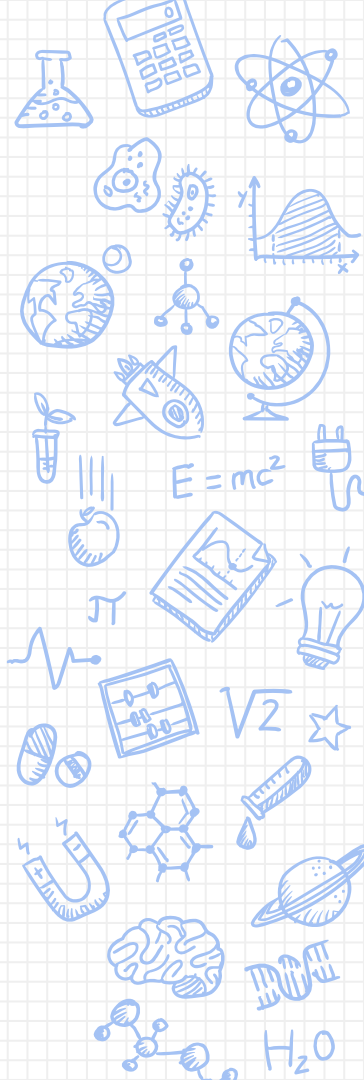
$$P(A|B) = ??$$

$$P(A) = 3/6$$

$$P(B) = 2/6$$

$P(A,B) = 2/6 \times 3/6 = 0,16666$  ou 16,66%, ou ainda  $1/6$  – Este resultado equivale a nossa única possibilidade, que é o número 5.

$$P(A|B) = (1/6)/(2/6)$$





$$E = mc^2$$

$$E = mc^2$$






# Obrigado!