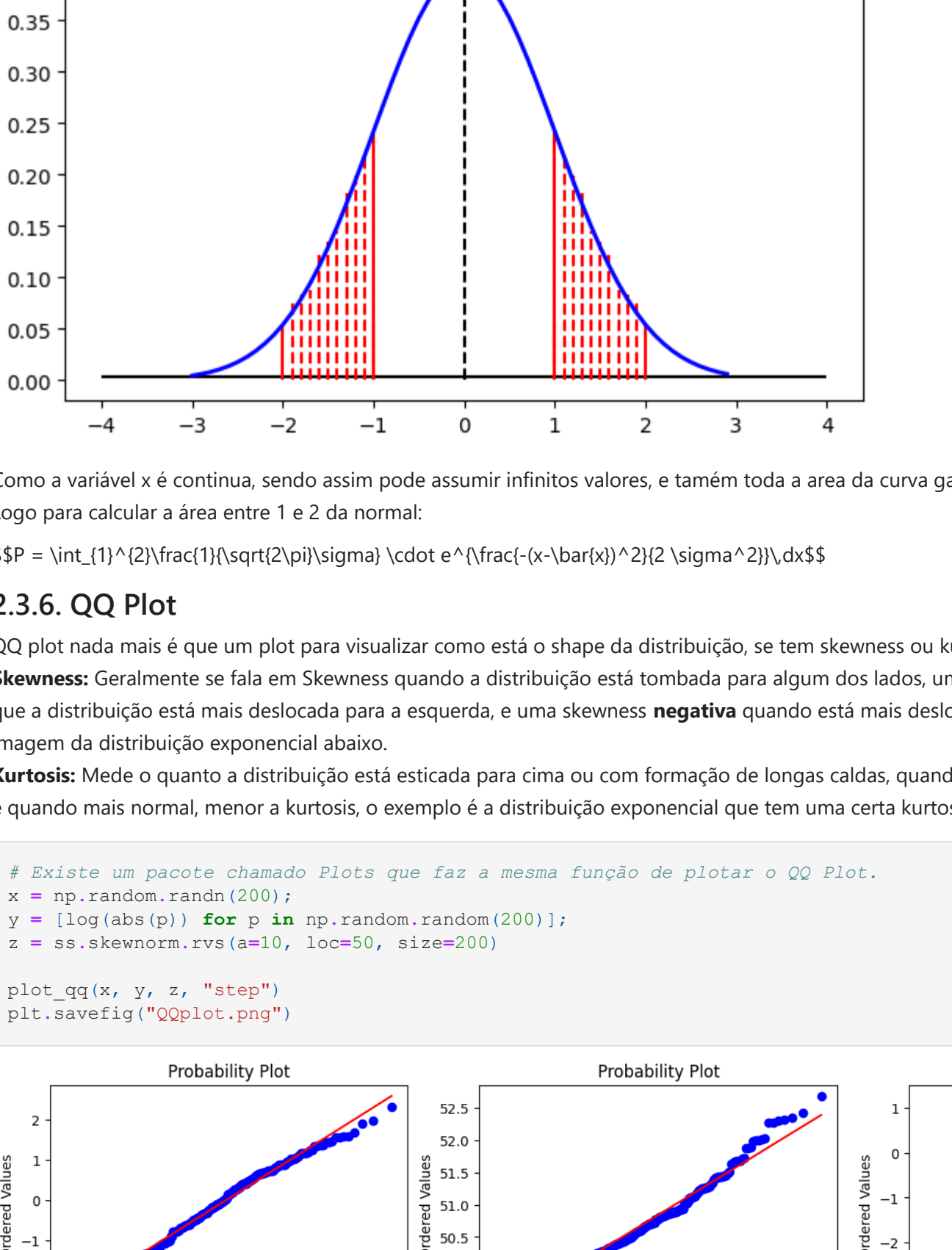


2.3.5 Distribuição Normal

A distribuição Normal é simétrica à média e as outras distribuições são geralmente moldadas de forma normal. Em uma distribuição normal 68% dos dados ficam dentro de um desvio-padrão da média e 90% dos dados em dois desvios-padrões. A diferença entre a distribuição normal das outras distribuições (binomial e poisson) é que na noção de distribuição discreta e contínua, ambas são distribuições discretas pois as possibilidades dos eventos eram discretos, agora x pode assumir uma probabilidade, logo a função é chamada de densidade de probabilidade. Onde para calcular a área em baixo da curva usa-se a ferramenta de Integral. $\int_{-\infty}^x f(x) dx$ Ou utiliza a tabela da normal.

Na Integral, o primeiro valor de baixo (0) é o primeiro valor da esquerda para direita na distribuição, e o valor de cima (1) é justamente até onde vai a área.

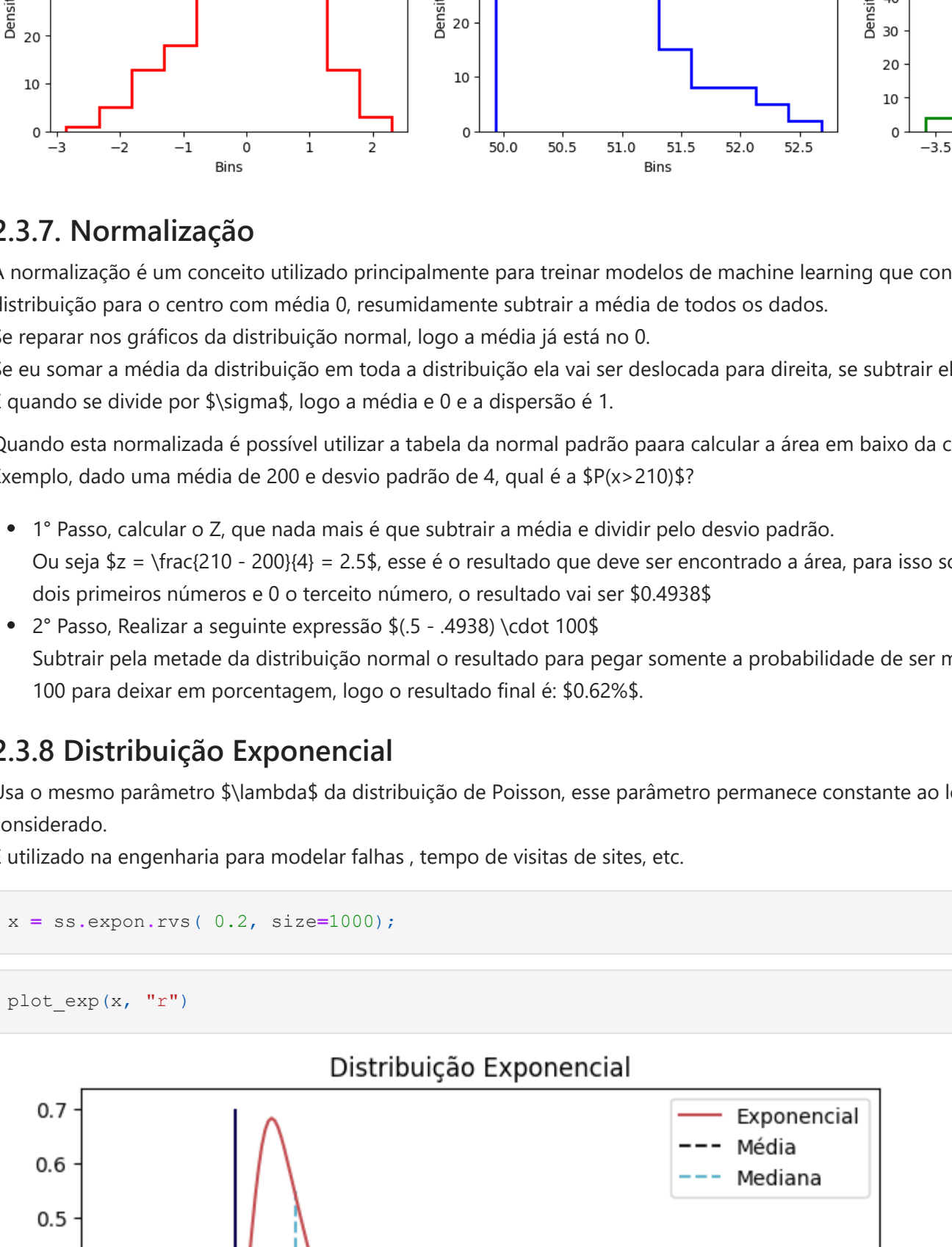
```
In [97]: range = np.arange(-3, 3, 0.1)
x = np.random.randn(200)
plot_normal(range, true)
```



Função densidade de probabilidade

$$f(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2}$$

```
In [290]: plot_normal(range, true)
```



Como a variável x é contínua, sendo assim pode assumir infinitos valores, e também toda a área da curva gaussiana é 1*.

$$\int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2} dx = 1$$

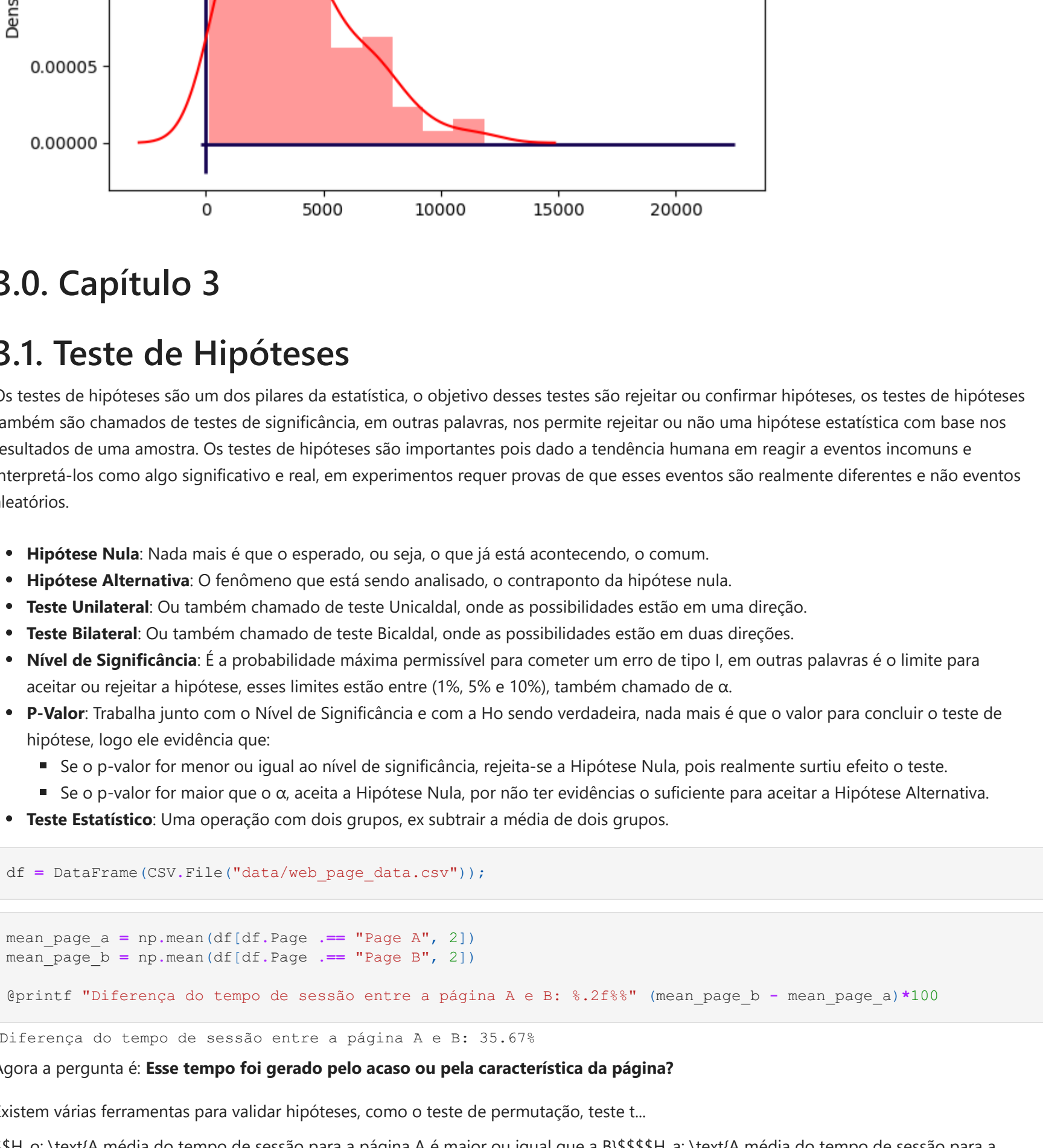
2.3.6. QQ Plot

QQ plot nada mais é que uma forma para visualizar como está o shape da distribuição, se tem skewness ou kurtosis.

Skewness: Geralmente se fala em Skewness quando a distribuição está tombada para algum dos lados, uma Skewness **positiva** significa que a distribuição está mais deslocada para a esquerda, e uma skewness **negativa** quando está mais deslocada a direita, o exemplo na imagem da distribuição exponencial abaixo.

Kurtosis: Mede o quanto a distribuição está esticada para cima ou com formação de longas caudas, quando mais pontuda, maior a kurtosis e quando mais normal, menor a kurtosis, o exemplo é a distribuição exponencial que tem uma certa kurtosis positiva.

```
In [161]: # Existe um pacote chamado Plots que faz a mesma função de plotar o QQ Plot.
x = np.random.randn(200)
y = [log(abs(p)) for p in np.random.randn(200)]
z = ss.skewnorm.rvs(a=10, loc=50, size=200)
plot_qq(x, y, z, "step")
plt.savefig("QQplot.png")
```



2.3.7. Normalização

A normalização é um conceito utilizado principalmente para treinar modelos de machine learning que consiste em movimentar a distribuição para o centro com média 0, resumidamente subtrair a média de todos os dados.

Se reparar nos gráficos da distribuição normal, logo a média já está no 0.

Se eu somar a média da distribuição em toda a distribuição ela vai ser deslocada para direita, se subtrair ela é deslocada a esquerda. E quando se divide por 2, logo a média é 0 e a dispersão é 1.

Quando esta normalizada é possível utilizar a tabela da normal padrão para calcular a área em baixo da curva.

Exemplo, dado uma média de 200 e desvio padrão de 4, qual é a $P(X > 210)$?

- 1º Passo, calcular o z, que nada mais é que subtrair a média e dividir pelo desvio padrão.
Ous seja $z = \frac{210 - 200}{4} = 2.5$, esse é o resultado que deve ser encontrado a área, para isso só checar a tabela, onde são 2.5 dos primeiros números e 0 o terceiro número, o resultado vai ser 0.49385
- 2º Passo, Realizar a seguinte expressão $1 - 0.4938$ \times 100

Subtrair pela metade da distribuição normal o resultado para obter somente a probabilidade de ser maior que 210, e multiplicar por 100 para deixar em porcentagem, logo o resultado final é: 0.62%.

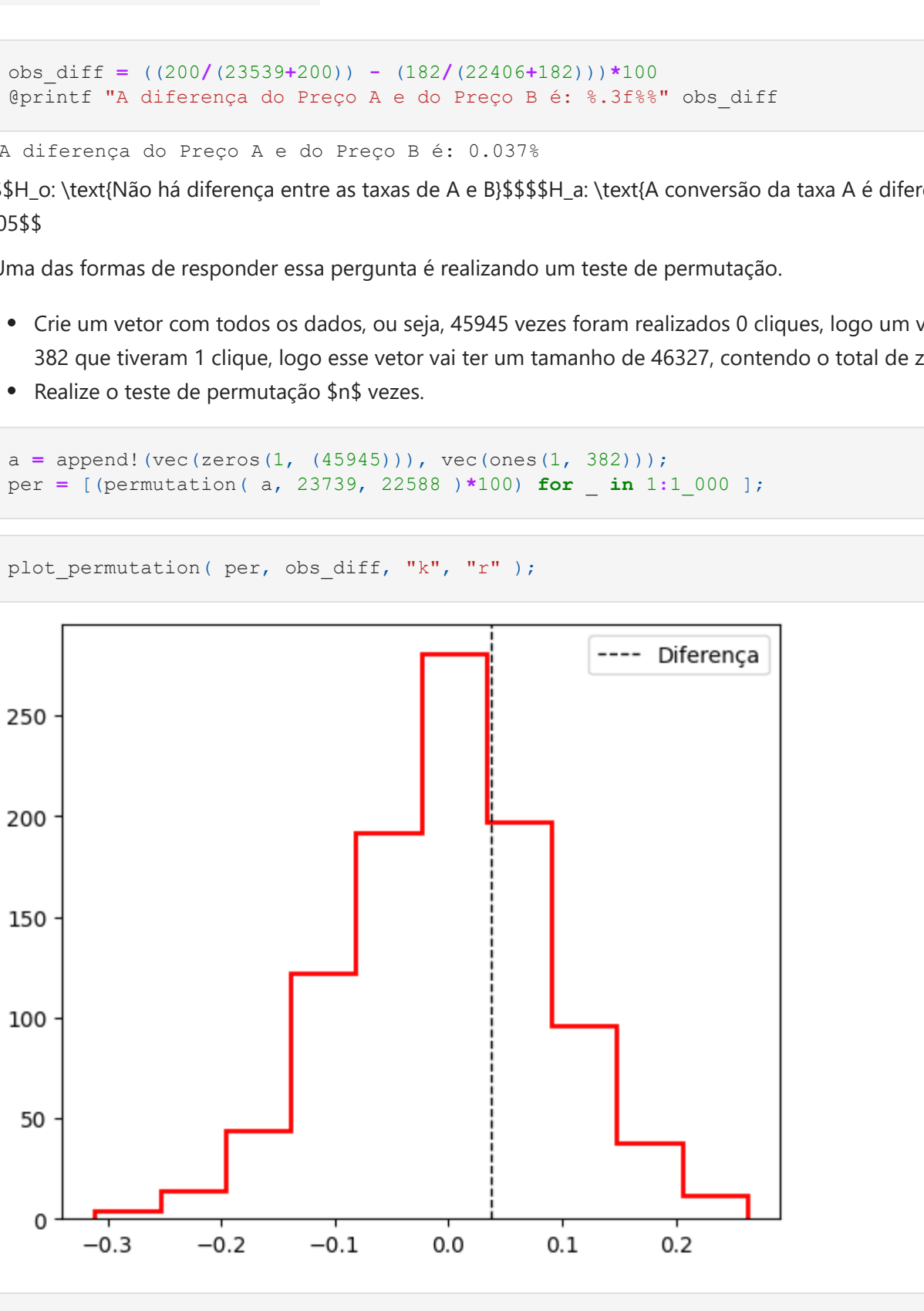
2.3.8 Distribuição Exponencial

Usa o mesmo parâmetro λ da distribuição de Poisson, esse parâmetro permanece constante ao longo do período sendo considerado.

É utilizado na engenharia para modelar falhas, tempo de visitas de sites, etc.

```
In [870]: x = ss.expon.rvs( 0.2, size=1000)
```

```
In [871]: plot_exp(x, "r")
```



2.3.9 Distribuição Weibull

É uma extensão da distribuição Exponencial, na qual a taxa de evento pode mudar de acordo com um "parâmetro de forma" β .

Se $\beta > 1$, a probabilidade de um evento aumenta com o tempo.

Se $\beta < 1$, a probabilidade de um evento diminui com o tempo.

Quando $\beta = 1$, a distribuição de Weibull é 1, retorna a distribuição exponencial.

Sendo assim, pode ser utilizada na análise de sobrevivência e confiabilidade, e sua função é:

$$Cumulativa: F(x, \lambda, \beta) = 1 - e^{-\left(\frac{x}{\lambda}\right)^\beta}$$

$$Densidade de Probabilidade: f(x, \lambda, \beta) = \frac{\beta}{\lambda} \left(\frac{x}{\lambda}\right)^{\beta-1} e^{-\left(\frac{x}{\lambda}\right)^\beta}$$

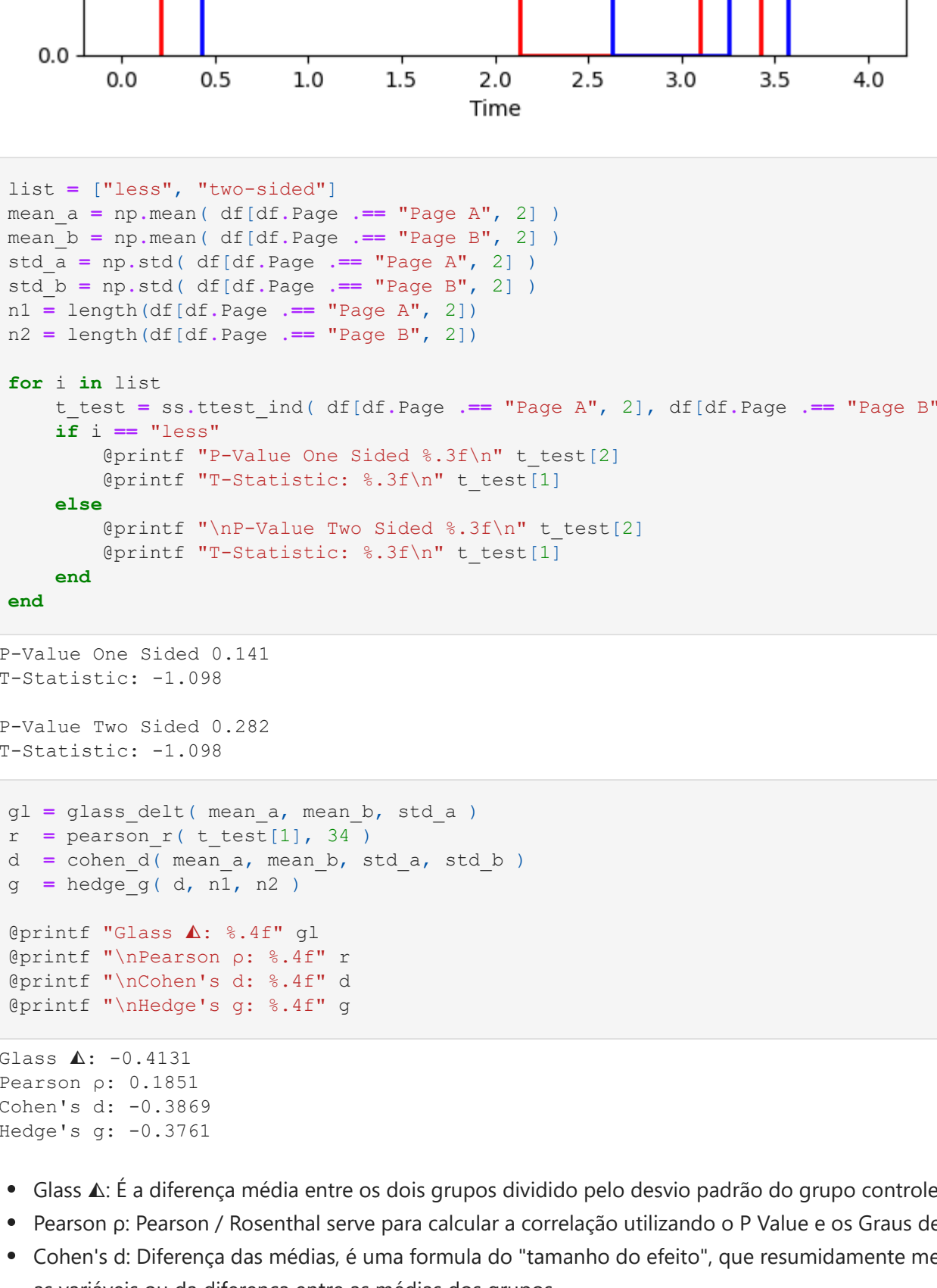
```
In [276]: f(x, a, b) = (a / (b**a)) * (x**(a-1)) * e**(-(x/b)**a)
f(3, 3, 3)
```

```
Out[276]: 0.36787944117144233
```

```
In [275]: 20*30 # não é possível utilizar números grandes pois o mesmo problema.
```

```
Out[275]: -8070450532247928832
```

```
In [86]: x = ss.weibull_min.rvs(1.5, scale=5000, size=100)
plot_weib(x, "r")
```



3.0. Capítulo 3

3.1. Teste de Hipóteses

Os testes de hipóteses são um dos pilares da estatística, o objetivo desses testes são rejeitar ou confirmar hipóteses, os testes de hipóteses também são chamados de testes de significância, em outras palavras, nos permite rejeitar ou não uma hipótese estatística com base nos resultados de uma amostra. Os testes de hipóteses são importantes pois dado a tendência humana em reagir a eventos incógnitos e interpretá-los como algo significativo e real, em experimentos requer provas de que esses eventos são realmente diferentes e não eventos aleatórios.

- Hipótese Nula:** Nada mais é que o esperado, ou seja, o que já está acontecendo, o comum.
- Hipótese Alternativa:** O fenômeno que está sendo analisado, o contraponto da hipótese nula.
- Teste Unilateral:** O fenômeno que está sendo analisado, onde as possibilidades estão em uma direção.
- Teste Bilateral:** O fenômeno que está sendo analisado, onde as possibilidades estão em duas direções.
- Nível de Significância:** A probabilidade máxima permitida para cometer um erro de tipo I, em outras palavras é o limite para aceitar ou rejeitar a hipótese, esses limites estão entre 1%, 5% e 10%, também chamado de α .
- Valor:** Trabalha junto com o Nível de Significância e com a Ho sendo verdadeira, nada mais é que o valor para concluir o teste de hipótese, logo ele evidencia que:
 - Se o p-valor for menor ou igual ao nível de significância, rejeita-se a Hipótese Nula, pois realmente surtiu efeito o teste.
 - Se o p-valor for maior que o α , aceita a Hipótese Nula, por não ter evidências o suficiente para aceitar a Hipótese Alternativa.
- Teste Estatístico:** Uma operação com dois grupos, ex subtrair a média de dois grupos.

```
In [31]: df = DataFrame(CSV.File("data/web_page_data.csv"));
```

```
In [50]: mean_page_a = np.mean(df[df.Page == "Page A", 2])
mean_page_b = np.mean(df[df.Page == "Page B", 2])
print("Diferença do tempo de sessão entre a página A e B: %.2f" % (mean_page_b - mean_page_a)*100)
```

Diferença do tempo de sessão entre a página A e B: 35.67%

Agora a pergunta é: **Esse tempo foi gerado pelo acaso ou pela característica da página?**

Existem várias ferramentas para validar hipóteses, como o teste de permutação, teste t...

H_0 : A média do tempo de sessão para a página A é maior ou igual que a B H_1 : A média do tempo de sessão para a página A é menor que a B $\alpha = 0.05$

3.1. Reamostragem

3.1.1. Teste de Permutação em tempos de sessão.

O teste de permutação nada mais é que um teste para verificar se tem realmente um significado estatístico.

- 1. Separar o grupo de controle e o grupo de tratamento, o de tratamento que vai ser o grupo que irá receber o teste.
- 1. Depois do teste, calcular alguma estatística, exemplo a média da diferença dos dois grupos.
- 1. Juntar em uma base de dados o grupo de controle e o grupo de tratamento.
- 1. Amostragem aleatória de diferentes indivíduos dessa base de dados e calcular a mesma estatística e armazenar o resultado.
- 1. Montar uma distribuição com os resultados.
- 1. Calcular a quantidade de vezes que os valores maior que a média saíram durante o processo de amostragem e divide pela quantidade de vezes que foi realizado a amostragem, logo esse é o p-valor, o valor da aleatoriedade na escolha dos grupos.

```
In [219]: obs_diff = [permutation(df.Time, 21, 15) for _ in range(1000)]
p_value = pytest(df, obs_diff)
```

```
In [246]: plot_permutation(df, obs_diff, "x", "y")
```



```
In [302]: print("Logo a quantidade de vezes que excedeu o limite foi de %.2f%%\nEntão aceita a Ho." % p_value*100)
```

Logo a quantidade de vezes que excedeu o limite foi de 15.70%. Então aceita a Ho.

3.1.2. Teste de permutação em taxas de conversão.

Nesse exemplo existe 20000 visualizações de um determinado preço e foi mensurado a quantidade de cliques em ambos os preços.

```
In [118]: DataFrame(Dict({"Resultado" => ["Preço A", "Preço B"], "Cliques" => [200, 182], "No Cliques" => [23539, 22406]}))
```

```
Out[118]: 2 rows x 3 columns
```

	Cliques	No Cliques	Resultado
int64	int64	String	
1	200	23539	Preço A
2	182	22406	Preço B

```
In [122]: obs_diff = ((200/(23539+200)) - (182/(22406+182)))*100
print("A diferença do Preço A e do Preço B é: %.03f" % obs_diff)
```

A diferença do Preço A e do Preço B é: 0.037%

H_0 : Não há diferença entre as taxas de A e B H_1 : A conversão da taxa A é diferente em relação a taxa B $\alpha = 0.05$

Uma das formas de responder essa pergunta é realizando um teste de permutação.

- 1. Crie um vetor com todos os dados, ou seja, 45945 vezes foram realizados 0 cliques, logo um vetor com 45945 zeros, e um vetor com 382 que tiveram 1 clique, logo esse vetor vai ter um tamanho de 46327, contendo o total de zeros e o total de 1 cliques.
- 1. Realize o teste de permutação $\alpha = 0.05$

```
In [228]: a = append(vec(zeros(1, 45945))), vec(ones(1, 382))
per = [permutation(a, 23739, 22588)*100 for _ in range(1000)]
```

```
In [293]: plot_permutation(per, obs_diff, "x", "y")
```



```
In [304]: p_value = pytest(per, obs_diff)
```

```
In [308]: print("Em relação aos preços, com uma significância de 0.05, o p-valor foi: %.2f%%\nEntão aceita a Ho." % p_value*100)
Em relação aos preços, com uma significância de 0.05, o p-valor foi: 32.60%
Então aceita a Ho.
```

```
In [314]: # Outras formas
chi2, p_value, df, _ = ss.chi2_contingency([200, 23739 + 200], [182, 22588 + 182])
print("Chi2: %.2f\nP-Value: %.2f" % (chi2, p_value))
Chi2: 0.15
P-Value: 0.70
```

3.2. Teste T de Student

O teste T de Student nada mais é que um teste de comparação de dois grupos em relação a sua média.

Nas quais os dados são numéricos, mas para que seja utilizado é necessário usar uma forma padronizada de estatística de teste.

```
In [5]: plot_page_diff()
```



```
In [43]: list = ["less", "two-sided"]
mean_a = np.mean(df[df.Page == "Page A", 2])
mean_b = np.mean(df[df.Page == "Page B", 2])
std_a = np.std(df[df.Page == "Page A", 2])
std_b = np.std(df[df.Page == "Page B", 2])
n1 = length(df[df.Page == "Page A", 2])
n2 = length(df[df.Page == "Page B", 2])

for i in list:
    t_test = st.ttest_ind(df[df.Page == "Page A", 2], df[df.Page == "Page B", 2], equal_var=False, alternative=i)
    if i == "less":
        print("P-Value One Sided %.3f\n" % t_test[2])
        print("T-Statistic: %.3f\n" % t_test[1])
    else:
        print("P-Value Two Sided %.3f\n" % t_test[2])
        print("T-Statistic: %.3f\n" % t_test[1])
    end
end
```

P-Value One Sided 0.141
T-Statistic: -1.098
P-Value Two Sided 0.282
T-Statistic: -1.098

```
In [189]: g1 = glass.del1(mean_a, mean_b, std_a)
g2 = pearson_r(t_test[1], 34)
d = cohen_d(mean_a, mean_b, std_a, std_b)
g = hedge_g(d, n1, n2)
```

Glass A: 0.41
Pearson p: 0.1851
Cohen's d: -0.3869
Hedge's g: -0.3761

- Glass A: É a diferença média entre os dois grupos dividido pelo desvio padrão do grupo controle.
- Pearson p: Pearson / Rosenthal serve para calcular a correlação utilizando o P Value e os Graus de Liberdade.
- Cohen's d: Diferença das médias, é uma fórmula do "tamanho do efeito", que resumidamente mede o tamanho das associações entre as variáveis ou da diferença entre as médias dos grupos.
- Hedge's g: Correção do D de Cohen.

$$Glass: S_g = \frac{\bar{x}_1 - \bar{x}_2}{s_d}$$

$$Pearson: r = \sqrt{\frac{t^2}{t^2 + df}}$$

$$Cohen's d: d = \frac{\bar{x}_1 - \bar{x}_2}{s_d} \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}$$

$$Hedge's g: g = d \sqrt{\frac{n_1 n_2}{n_1 + n_2}}$$

```
In [ ]:
```

x.0. Referências

PETER BRUCE & ANDREW BRUCE Estatística prática para cientistas de dados: 50 conceitos essenciais.

Link: <https://www.amazon.com.br/Estat%C3%AAdica-Pr%C3%A1tica-Para-Cientistas-Dados/dp/855080603X>

DAVID MATOS 8 Conceitos Estatísticos Fundamentais Para Data Science.

Link: <https://www.cientiaedados.com/8-conceitos-estatisticos-fundamentais-para-data-science/>

IGOR SOARES Correlação não implica em Causalidade.

Link: <https://medium.com/@eliemapaolo/correla%C3%A7%C3%A3o-n%C3%A3o-implica-em-causalidade-8459179ad1bc>.

annaheusch Número de Casos de Divórcio em Maine

Link: <https://blogs.ams.org/blog/mathblogs/2017/04/10/divorce-and-margarine/>

Wikipédia Cramer's V

Link: https://en.wikipedia.org/wiki/Cramer%27s_V

BURKEYACADEMY What are Skewness and Kurtosis?

Link: <https://www.youtube.com/watch?v=iK7nLzxiAQQ>

(Discourse) qqnorm & qqplot

Link: <https://discourse.julialang.org/t/qqnorm-and-qqplot/6118/8>

Professor Guru Tabela Normal Padrão

Link: <https://professorguru.com.br/tabela-normal.html>