

2.3.2 Distribuição de Poisson

alta concentração de eventos próximos ao eixo y, uma das principais características é que não tem repetições como na distribuição binomial trabalha em um intervalo contínuo.

Ex: Em um estudo de chuva ou cliques em um site, no exemplo da chuva, qual é a chance de uma chuva, só que não existe o evento "não chuva" entre duas chuvas.

Imagine um intervalo, que começa e o até uma variável W por exemplo. E eu divido em n intervalos muito pequenos, onde n tende ao infinito, logo a probabilidade está tendendo a 0 pois existem n intervalos, com essa quantidade de intervalos, virou uma binomial, ou seja, choveu ou não por exemplo.

$$P(r \frac{\lambda}{r}, n) = \lim_{n \rightarrow \infty} (\frac{\lambda}{n})^r \cdot (1 - \frac{\lambda}{n})^{n-r} = \frac{n!}{r!(n-r)!}$$

Distribuição de Poisson, logo λ (Quantidade de Chuva) = $p \cdot n$, então $p = \frac{\lambda}{n}$. No Limite que n tende ao infinito, o produto de n e n não vai mudar pois sempre vai ficar menor e n sempre vai ficando maior.

$$P(r/\lambda) = \frac{e^{-\lambda} \cdot \lambda^r}{r!}$$

A média e a variância da distribuição binomial é dada pelo: λ
E o Desvio Padrão é $\sqrt{\lambda}$

In [11]: `P(x, z) = np.e**(-lambda) * lambda**x / np.math.factorial(x)`

Dado que eu esperava em média 35 carros entrando no shopping, qual a probabilidade de aparecer 20?

In [15]: `#print "A probabilidade de somente uma chuva no mês é de: %.3f%%" P(35, 20)*100`
A probabilidade de somente uma chuva no mês é de: -0.000%

In [140]: `# Julia tem problema com elevar x a o expoente y
function f(x, y)
[x**y for _ in 1:y]
end
f(35, 20)`

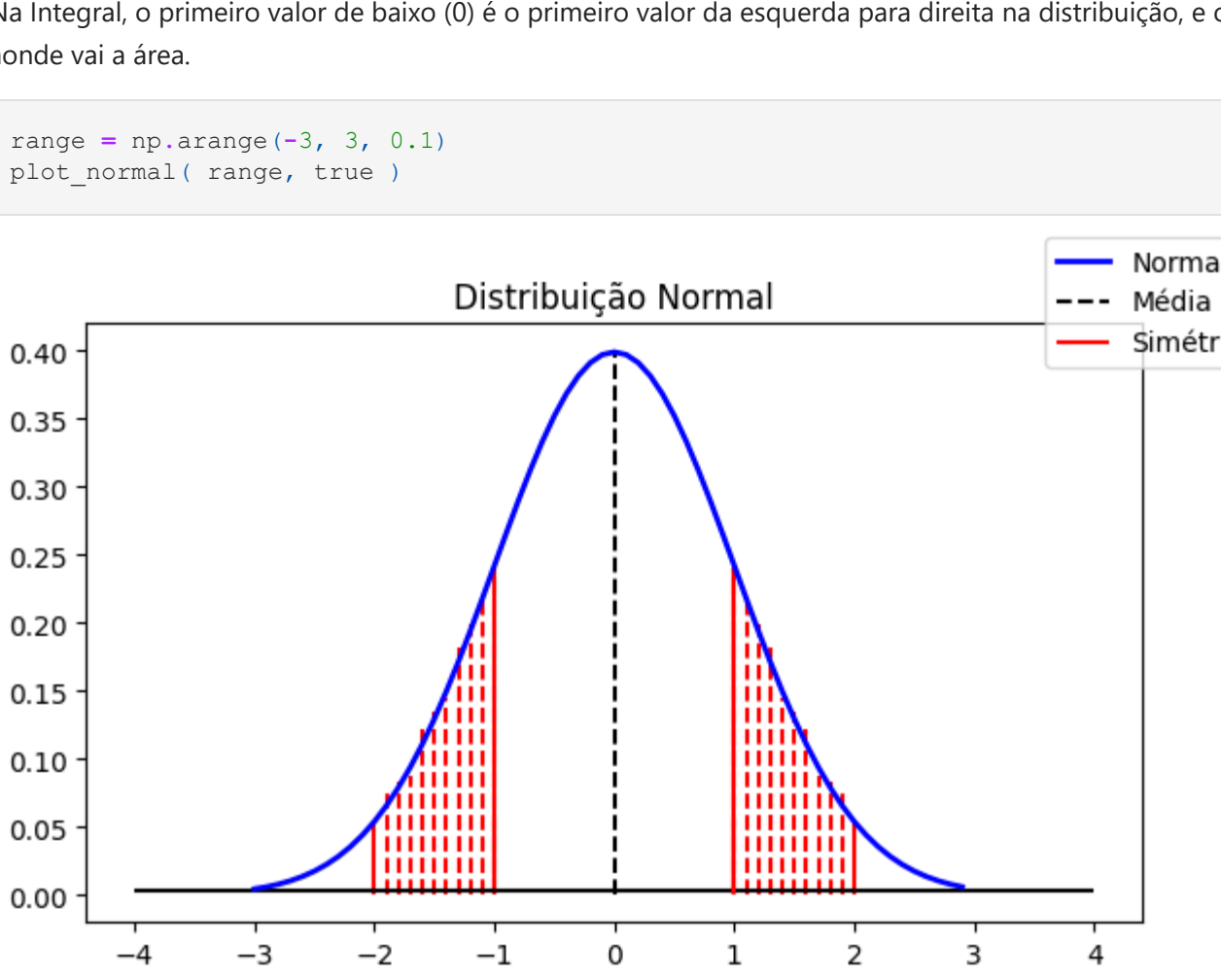
Out[140]: 20-element Vector{Int64}:
1225
1500625
225181330625
661601603543685869
7665930784382691969
48227646868404185
-165024524321314303
450275795304469505
-639265039275616255
7551947002216534017
3654036140188672001
9358444485278177281
4785494631104806913
-2679742982427181055
2893767076708193281
8010019347241369601
-48246705104617471
-578681105722249663
-8394158839848501247
541221425122377279

Logo em Python, aplicando a mesma função irá retornar a probabilidade de 0.0019

In [14]: `#print "A probabilidade de somente uma chuva no mês é de: %.3f%%" P(5, 1)*100`
A probabilidade de somente uma chuva no mês é de: 3.369%

In [89]: `x = ss.poisson.rvs(2, size=500)`

In [89]: `plot_poisson(x, "x")`



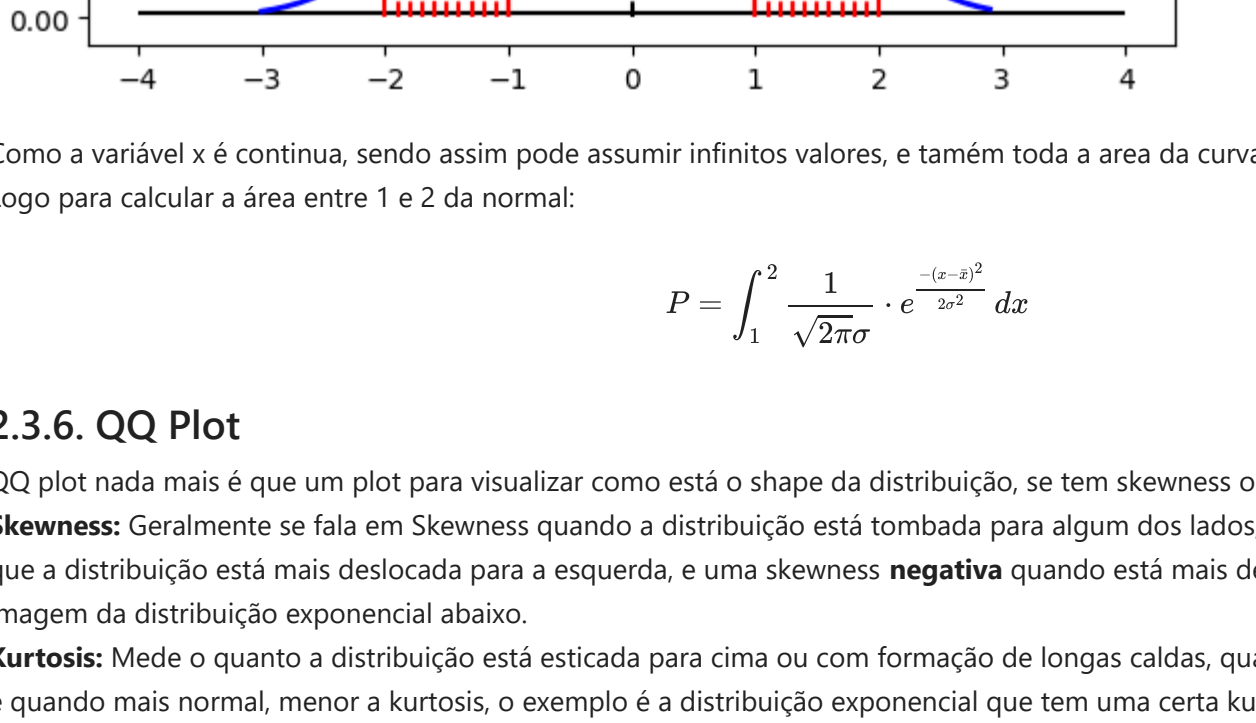
2.3.5 Distribuição Normal

A distribuição Normal é simétrica a 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: $\int_0^1 f(x) dx$
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.

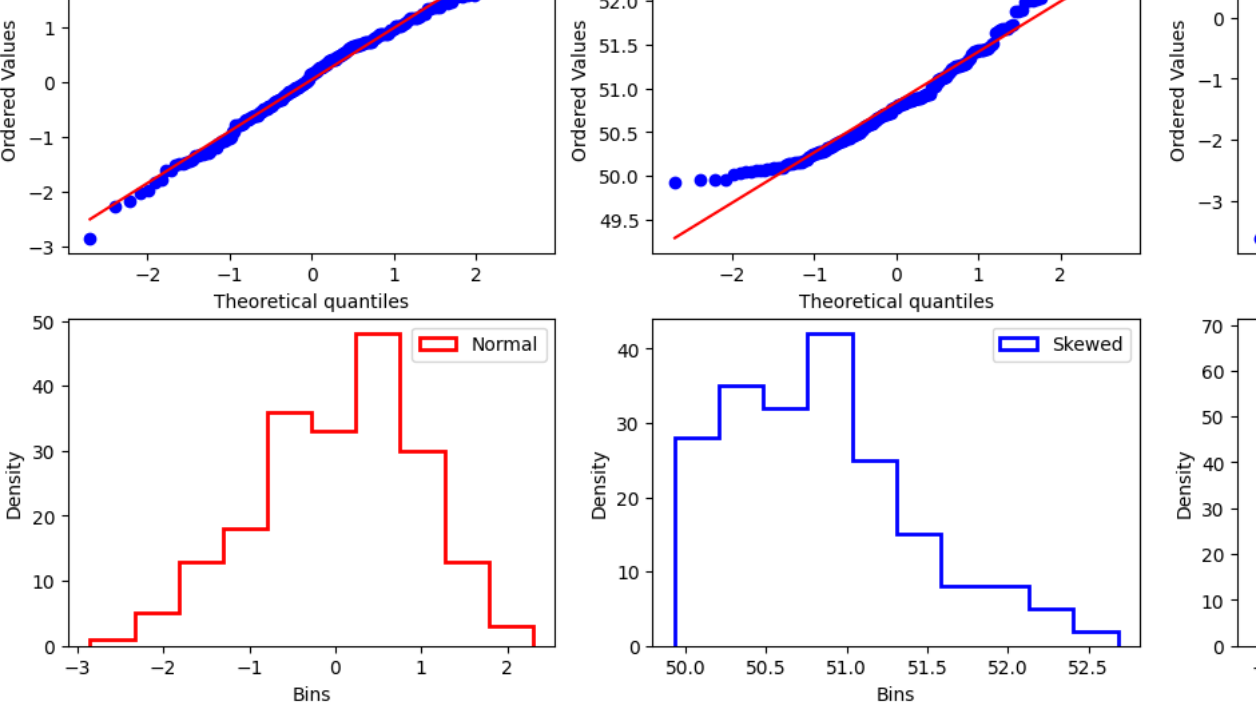
In [97]: `range = np.arange(-3, 3, 0.1)
plot_normal(range, true)`



Função densidade de probabilidade

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

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*.
Logo para calcular a área entre 1 e 2 da normal.

$$P = \int_1^2 \frac{1}{\sqrt{2\pi\sigma}} \cdot e^{-\frac{(x-\mu)^2}{2\sigma^2}} dx$$

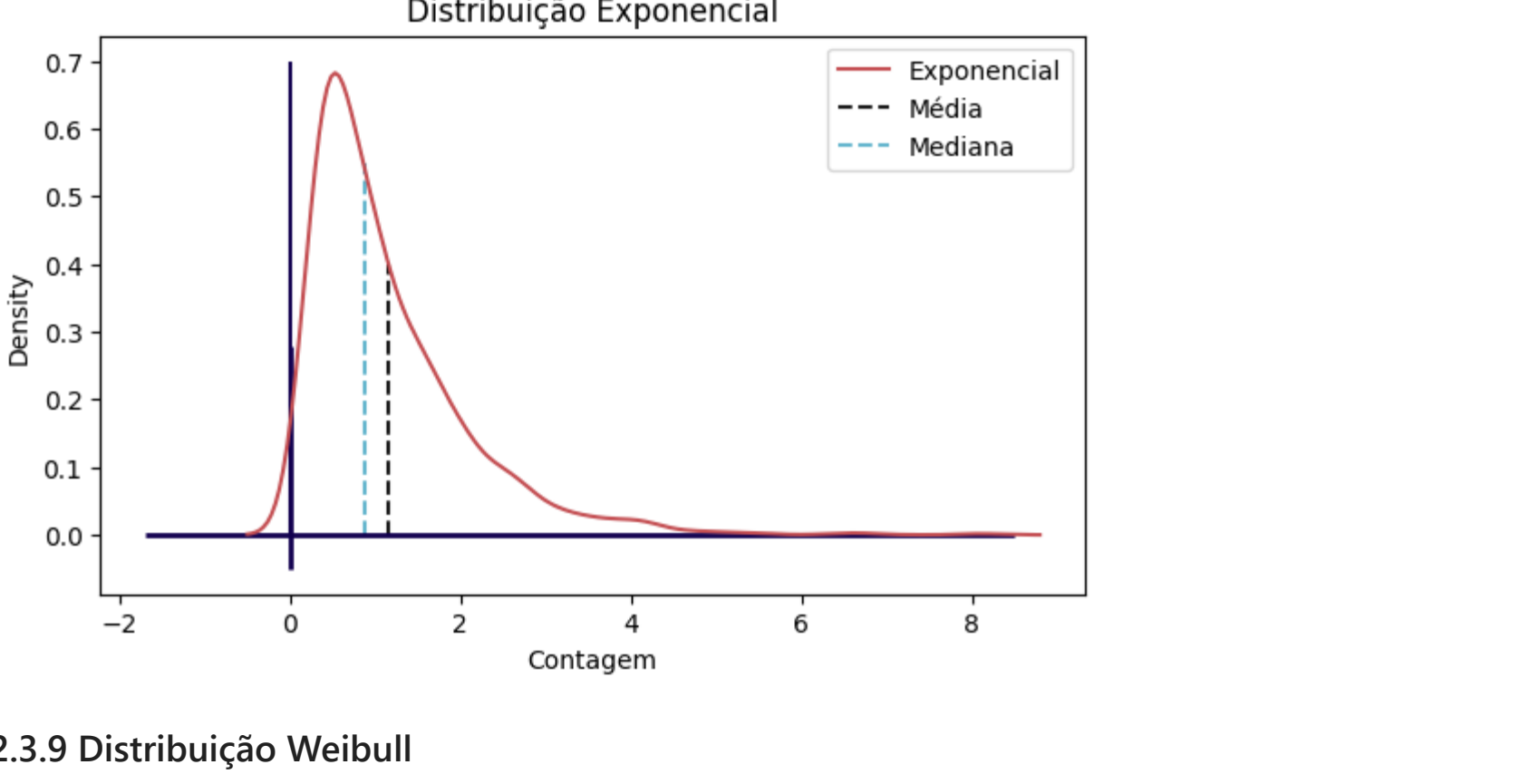
2.3.6 QQ Plot

QQ plot nada mais é que um plot 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.
np.random.randn(200)
p = [log(abs(p)) for p in np.random.randn(200)]
z = ss.skewnorm.rvs(n=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 separar 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 σ , 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 $P(x > 210)$?

- 1º Passo, calcular o z, que nada mais é que subtrair a média e dividir pelo desvio padrão.
Ou seja $z = \frac{210-200}{4} = 2.5$, esse é o resultado que deve ser encontrado a área, para isso se checar a tabela, onde são 2.5 os primeiros números e o 0 terceiro número, o resultado vai ser 0.4938
- 2º Passo, Realizar a seguinte expressão $(.5 - .4938) \cdot 100$
Subtrair pela metade da distribuição normal o resultado para pegar 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, deve ser permanente 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, "x")`



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 o α da distribuição de Weibull é 1, retoma a distribuição exponencial.

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

Cumulativa:

$$f(x, \alpha, \beta) = 1 - e^{-(\frac{x}{\alpha})^\beta}$$

Densidade de Probabilidade:

$$f(x, \alpha, \beta) = \frac{\alpha}{\beta} x^{\alpha-1} \cdot e^{-(\frac{x}{\alpha})^\beta}$$

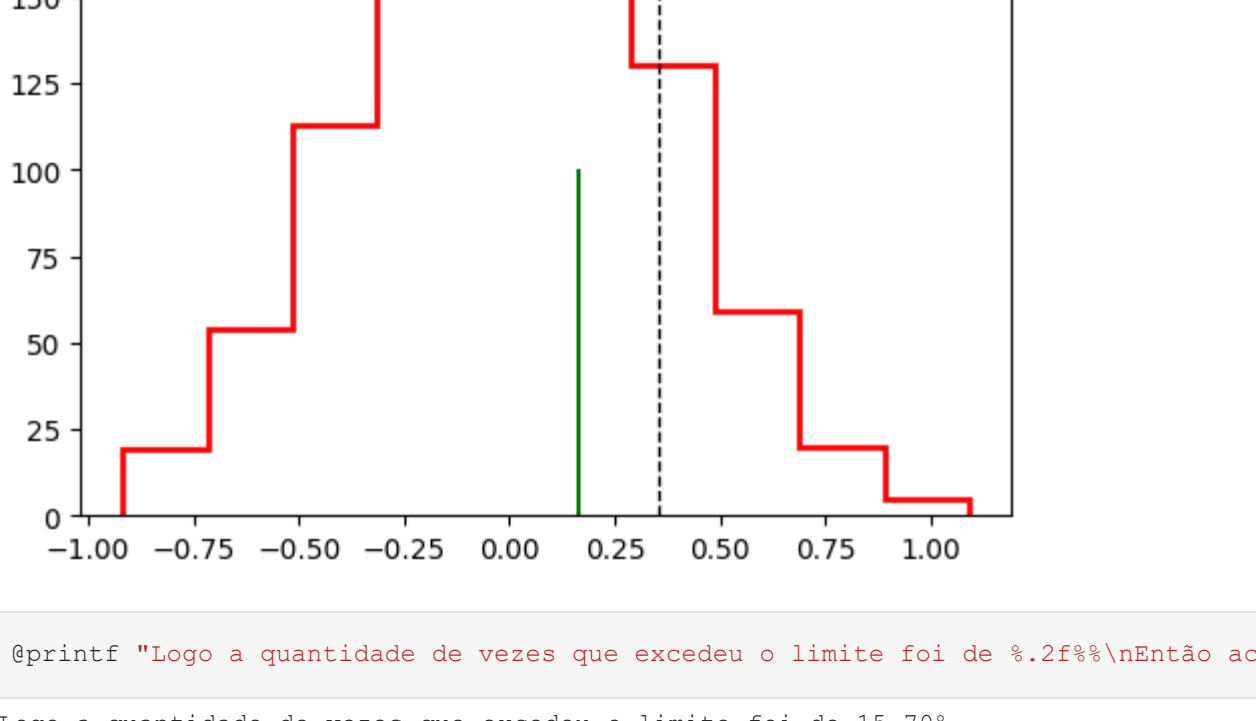
In [276]: `f(x, a, b) = (a / (b**a)) * (x**(a-1)) * a**(-(x/b)**a)`
`f(3, 3, 3)`

Out[276]: 0.3678794417144233

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, "x")`



Capítulo 3

3.1. Teste de Hipóteses

Os testes de hipóteses são de dois 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 incertos e interpreta-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:** Ou também chamado de teste Unicaidal, onde as possibilidades estão em uma direção.
- **Teste Bilateral:** Ou também chamado de teste Bicaidal, onde as possibilidades estão em duas direções.
- **Nível de Significância:** É a probabilidade máxima permissível 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 .
- **P-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 excede que o nível de significância, rejeita-se a Hipótese Nula, pois realmente surtiu efeito o teste.
 - Se o p-valor for menor 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 [301]: `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...

$$H_0: \text{A média do tempo de sessão para a página A é maior ou igual que a B}$$

$$H_a: \text{A média do tempo de sessão para a página A é menor que a B}$$

$$\alpha = .05$$

3.1. Reamostragem

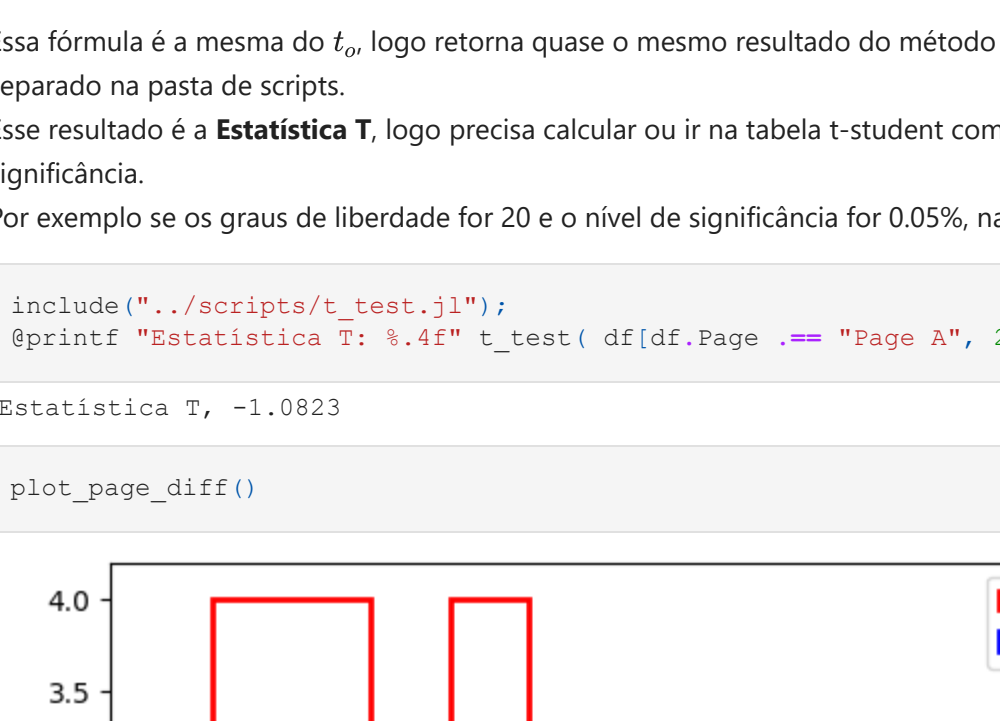
3.1.1. Teste de Permutação em tempos de sessão.

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.
2. Depois do teste, calcular alguma estatística, exemplo a média da diferença dos dois grupos.
3. Juntar em uma base de dados o grupo de controle e o grupo de tratamento.
4. Amostragem aleatória de diferentes indivíduos dessa base de dados e calcular a mesma estatística e armazenar o resultado.
5. Montar uma distribuição com os resultados.
6. 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]: `diff = [permutation(df.Time, 21, 15) for _ in 1:1000]
p_valor = py"test"(diff, diff_mean)`

Out[246]: `plot_permutation(diff, diff_mean, "x", "y")`



In [302]: `#print "Logo a quantidade de vezes que excedeu o limite foi de %.2f%%\nEntão aceita a Ho." p_valor*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
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: %.3f%%" obs_diff`

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

$$H_0: \text{Não há diferença entre as taxas de A e B}$$

$$H_a: \text{A conversão da taxa A é diferente em relação a taxa B}$$

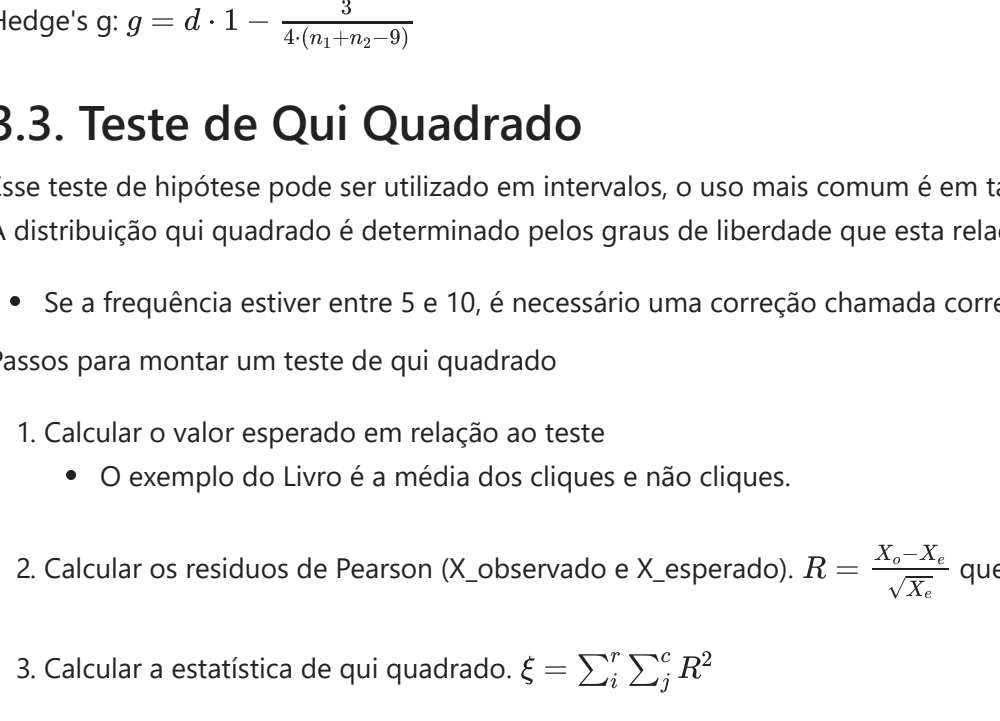
$$\alpha = .05$$

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

- 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.
- Realize o teste de permutação n vezes.

In [228]: `a = append!(vec(zeros(1, (45945))), vec(ones(1, 382)))
per = [permutation(a, 23739, 22588 *100) for _ in 1:1_000]`

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



In [304]: `p_valor = py"test"(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_valor`
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.

Esse teste também é uma aproximação para o teste de permutação, devido que o teste de permutação na época de sua implementação era muito custoso. Para utilizar o teste t, é necessário alguns passos:

- 1. Calcular a média de ambos os grupos.
- A diferença entre as duas médias é significativa?
- 1. Elevar os dados de ambos os grupos ao quadrado o grupo X e o grupo Y e somar $\sum X^2, \sum Y^2$ os grupos ao quadrado.
- 1. Calcular a variância de ambos os grupos (ou o std).
- 1. Colocar na fórmula do teste t, com o resultado da fórmula, localizada na tabela da distribuição t.

Obs: Existem várias fórmulas para a variância conhecida, não conhecida, amostras dependentes e assim vai.

Para aplicar o teste t ou Welch's t-test, quando a variância populacional é desconhecida em amostras independentes (Com ou sem o N-1):

$$t_0 = \frac{\bar{x}_0 - \bar{x}_1}{\sqrt{\frac{(n_1-1)s_1^2 + (n_2-1)s_2^2}{n_1+n_2-2}} \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}}$$

Welch's T-Test:

$$t' = \frac{\bar{x}_0 - \bar{x}_1}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}}$$

Essa fórmula é a mesma do t_0 , logo retorna quase o mesmo resultado do método `test_ind` do scipy que esta armazenada em um script separado na pasta de scripts.

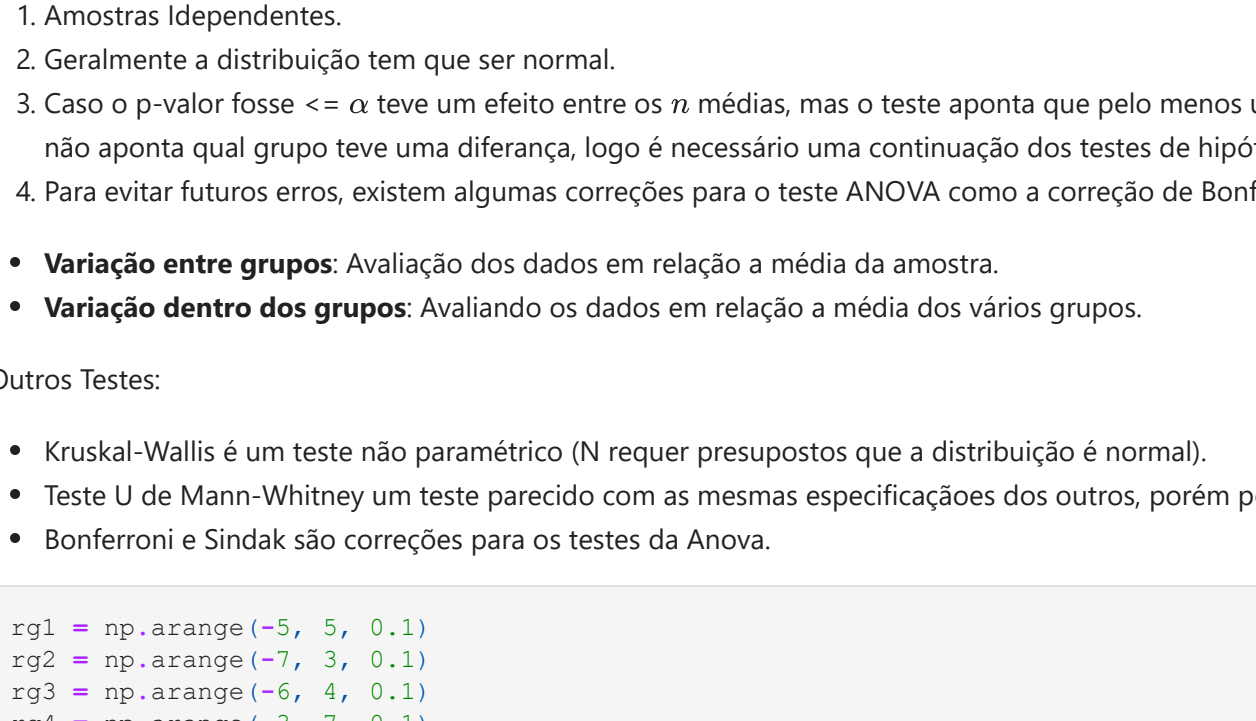
Esse resultado é a **Estatística T**, logo precisa calcular ou ir na tabela t-student com os graus de liberdade ($n_1 + n_2 - 2$) e o nível de significância.

Por exemplo se os graus de liberdade for 20 e o nível de significância for 0.05%, na tabela o p-valor vai ser aprox: **0.80 *****

In [316]: `include("scripts/t_test.jl")
println("Estatística T: %.4f") t_test(df[df.Page == "Page A", 2], df[df.Page == "Page B", 2], true)`

Estatística T: -1.0823

In [5]: `plot_page_diff()`



In [155]: `info = pages()
libt = ["less", "two-sided"]
for i in list
t_test_s8 = ss.ttest_ind(df[df.Page == "Page A", 2], df[df.Page == "Page B", 2], equal_var=false, alternative=libt[i])
if i == "less"
println("Two-Sided T-Statistic: %.3f\n" t_test_s8[2])
println("T-Statistic: %.3f\n" t_test_s8[1])
else
println("Two-Sided T-Statistic: %.3f\n" t_test_s8[2])
println("T-Statistic: %.3f\n" t_test_s8[1])
end
end`

F-Value One Sided 0.141
T-Statistic: -1.096
F-Value Two Sided 0.282
T-Statistic: -1.096

In [172]: `q1 = class{del{["mean_a"], ["mean_b"], ["std_a"], ["std_b"]}}
r = pearson_r(t_test_s8[1], 34)
d = cohen_d(info["mean_a"], info["mean_b"], info["std_a"], info["std_b"])
g = hedge_g(d, info["n_a"], info["n_b"])`

#print "Glass A: %.4f" g
#print "Cohen's d: %.4f" d
#print "Hedge's g: %.4f" g
#print "Pearson's r: %.4f" r

Glass A: -0.4131
Pearson's r: 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's r: 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.

$$\text{Glass } g = \frac{\bar{x}_1 - \bar{x}_2}{s_0}$$

$$\text{Pearson } r = \sqrt{\frac{r^2}{r^2 + d}}$$

$$\text{Cohen's d } d = \frac{\bar{x}_1 - \bar{x}_2}{\sqrt{\frac{s_1^2 + s_2^2}{2}}}$$

$$\text{Hedge's g } g = d \cdot \sqrt{\frac{2}{n_1 + n_2 - 2}}$$

3.3. Teste de Qui Quadrado

Esse teste de hipótese pode ser utilizado em intervalos, o uso mais comum é em tabelas de contingência como cita no livro.

A distribuição qui quadrado é determinado pelos graus de liberdade que esta relacionado ao número de dados.

- Se a frequência estiver entre 5 e 10, é necessário uma correção chamada correção de Yates.

Passos para montar um teste de qui quadrado

1. Calcular o valor esperado em relação ao teste
 - O exemplo do Livro é a média dos cliques e não cliques.
2. Calcular os resíduos de Pearson ($X_{\text{observado}}$ e X_{esperado}). $R = \frac{X_{\text{observado}} - X_{\text{esperado}}}{\sqrt{X_{\text{esperado}}}}$ que retorna as contagens diferentes.
3. Calcular a estatística de qui quadrado. $\chi^2 = \sum \frac{R^2}{R^2}$
4. Calcular o grau de liberdade que é $k - 1$ onde k é o n° de categorias. Exemplo do Livro em três Headlines (Teste A, B, C)

In [431]: `df = DataFrame{Dict{"Teste" => ["Headline A", "Headline B", "Headline C"], "Click" => [14, 8, 12], "NoClick" => [986, 992, 988]}}`

Out[431]: 3 rows x 3 columns

	Click	NoClick	Teste
1	14	986	Headline A
2	8	992	Headline B
3	12	988	Headline C

In [283]: `#print "Valores Esperados: N Média de Cliques: %.4f\n Média de Não Cliques C: %.4f" ((14+8+12)/3) ((986+992+988)/3)`

Valores Esperados:
Média de Cliques: 11.3333
Média de Não Cliques C: 988.6667

In [664]: `df = DataFrame{Dict{"Teste" => ["Headline A", "Headline B", "Headline C"], "Click" => [0.792, -0.990, 0.198], "NoClick" => [-0.084, 0.106, -0.021]}}`

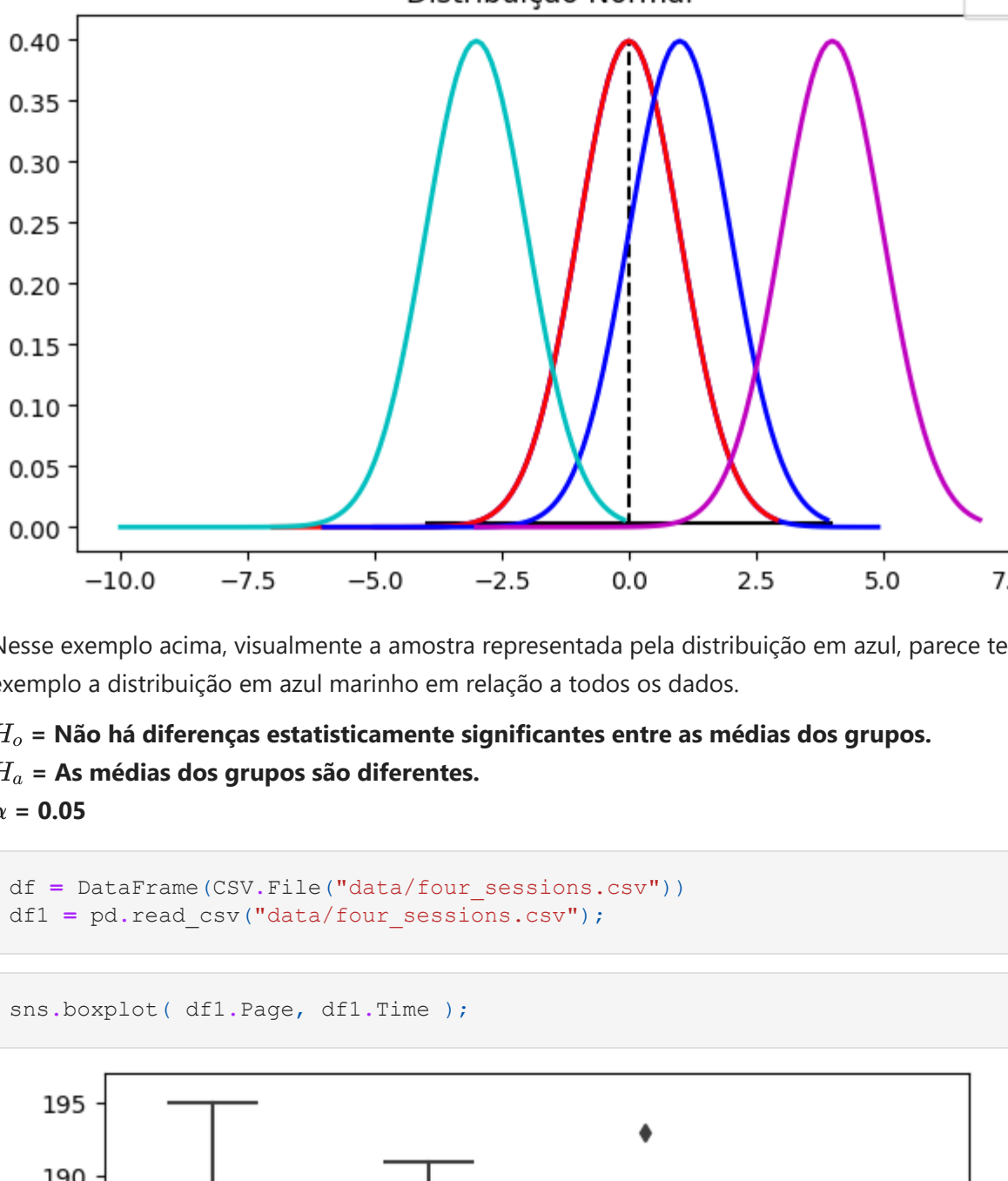
Out[664]: 3 rows x 3 columns

	Click	NoClick	Teste
1	0.792	-0.084	Headline A
2	-0.99	0.106	Headline B
3	0.198	-0.021	Headline C

Mesma lógica do teste T, tem que ir na tabela da distribuição qui quadrado com os graus de liberdade que são 2 e localizar o valor mais próximo de 1.66, que está entre 0.50 na tabela, logo o p-valor vai ser aprox.

In [3

Distribuição Normal



Nesse exemplo acima, visualmente a amostra representada pela distribuição em azul, parece ter dados não tão importantes quanto por exemplo a distribuição em azul marinho em relação a todos os dados.

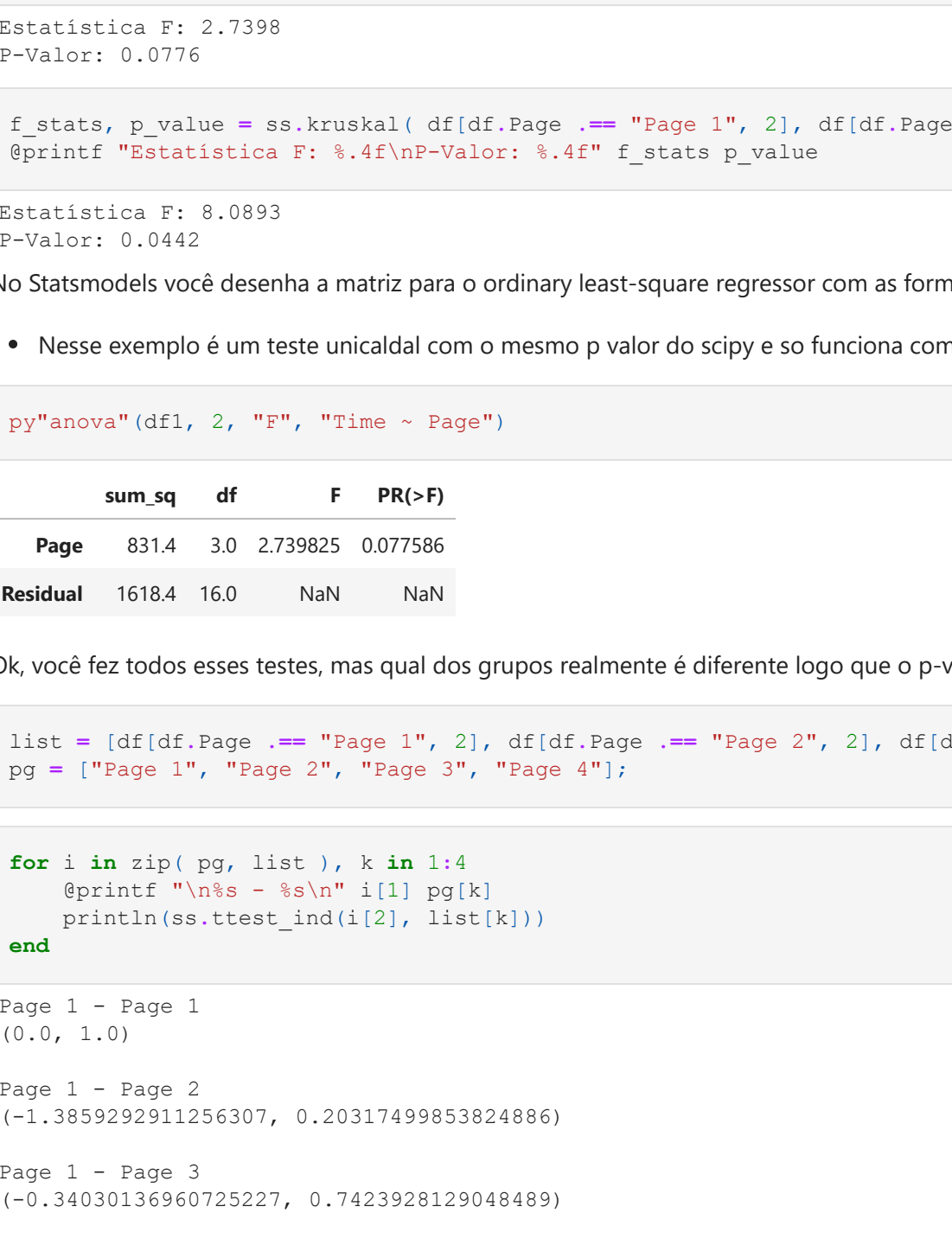
H_0 = Não há diferenças estatisticamente significantes entre as médias dos grupos.

H_a = As médias dos grupos são diferentes.

$\alpha = 0.05$

```
In [289]: df = DataFrame(CSV_File("data/four_sessions.csv"))
df1 = pd.read_csv("data/four_sessions.csv");
```

```
In [210]: sns.boxplot( df1.Page, df1.time );
```



```
In [290]: f_stats, p_value = ss.f_owenway( df[df.Page == "Page 1", 2], df[df.Page == "Page 2", 2], df[df.Page == "Page 3", 2], df[df.Page == "Page 4", 2] )
f_stats p_value
Estatística F: 2.7398
p-Valor: 0.0776
```

```
In [296]: f_stats, p_value = ss.kruskal( df[df.Page == "Page 1", 2], df[df.Page == "Page 2", 2], df[df.Page == "Page 3", 2], df[df.Page == "Page 4", 2] )
f_stats p_value
Estatística F: 8.0893
p-Valor: 0.0442
```

No Statsmodels você desenha a matriz para o ordinary least-square regressor com as formulas.

- Nesse exemplo é um teste unidacal com o mesmo p valor do scipy e so funciona com um dataframe do pandas.

```
In [291]: py"anova"(df1, 2, "P", "Time ~ Page")
```

	sum_sq	df	F	PR(>F)
Page	831.4	3.0	2.739825	0.077586
Residual	1618.4	16.0	NaN	NaN

Ok, você fez todos esses testes, mas qual dos grupos realmente é diferente logo que o p-valor foi pequeno?

```
In [362]: list = [df[df.Page == "Page 1", 2], df[df.Page == "Page 2", 2], df[df.Page == "Page 3", 2], df[df.Page == "Page 4", 2]]
pg = ["Page 1", "Page 2", "Page 3", "Page 4"];
```

```
In [369]: for i in zip( pg, list ), k in 1:4
0:printf "\n%o - %\n" i[1] pg[k]
printlin(as.ttest_ind(i[2], list[k]))
end
```

Page 1 ~ Page 1
(0.0, 1.0)

Page 1 ~ Page 2
(-1.3859292931256307, 0.20317499853824886)

Page 1 ~ Page 3
(-0.34030136960725227, 0.7423928129048489)

Page 1 ~ Page 4
(1.136191733238663, 0.28096106091142226)

Page 2 ~ Page 1
(1.3859292931256307, 0.20317499853824886)

Page 2 ~ Page 2
(0.0, 1.0)

Page 2 ~ Page 3
(1.2650143347989609, 0.24146881279653873)

Page 2 ~ Page 4
(4.950587917224828, 0.00112005643441126)

Page 3 ~ Page 1
(0.34030136960725227, 0.7423928129048489)

Page 3 ~ Page 2
(-1.2650143347989609, 0.24146881279653873)

Page 3 ~ Page 3
(0.0, 1.0)

Page 3 ~ Page 4
(1.9782125034804172, 0.08327876770509358)

Page 4 ~ Page 1
(-1.136191733238663, 0.28096106091142226)

Page 4 ~ Page 2
(-4.950587917224828, 0.00112005643441126)

Page 4 ~ Page 3
(-1.9782125034804172, 0.08327876770509358)

Page 4 ~ Page 4
(0.0, 1.0)

Somente em relação a página 4 com a página 2 o p-value ficou pequeno.

3.4.1. Bonferroni & Sidak Correction

Quando se trabalha com a Anova ou múltiplos testes t, geralmente podem superestimar a significância, ou seja, maior a chance de encontrar p valores pequenos, uma forma de corrigir esse problema é com a correção de bonferroni, que nada mais é que dividi o nível de significância pelo número unicos de comparações feitas, por exemplo $0.05 \div 6 = 0.0083$

Com a correção de Sidak, a formula é $1 - (1 - 0.05)^{1/6} = 0.0085$.

Outro comumente utilizado também é o Tukey's test

```
In [427]: py"tukey"(df1.Time, df1.Page, 0.05)
```

Multiple Comparisons Between All Pairs (Tukey)



group1	group2	meandiff	p-adj	lower	upper	reject	
1	Page 1	Page 2	9.8	0.4393	-8.3999	27.9999	False
2	Page 1	Page 3	2.8	0.9	-15.3999	20.9999	False
3	Page 1	Page 4	-8.2	0.5758	-26.3999	9.9999	False
4	Page 2	Page 3	-7.0	0.6774	-25.1999	11.1999	False
5	Page 2	Page 4	-18.0	0.0531	-36.1999	0.1999	False
6	Page 3	Page 4	-11.0	0.3419	-29.1999	7.1999	False

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.cienciaedados.com/8-conceitos-estatisticos-fundamentais-para-data-science/>

IGOR SOARES **Correlação não implica em Causalidade**.

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

annahaensch **Número de Casos de Divórcio em Maine**

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

Wikipédia **Cramer's V**

Link: https://en.wikipedia.org/wiki/Cram%C3%A9r%27s_V

BURKEYACADEMY **What are Skewness and Kurtosis?**

Link: <https://www.youtube.com/watch?v=IK7hLzxiAQQ>

(Discourse) **qgnorm & qqplot**

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

Professor Guru **Tabela Normal Padrão**

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

Univesp **Probabilidade e Estatística**

Link: <https://www.youtube.com/watch?v=7VQE278hXc&list=PLxI8Can9YAHeeWqe3m9HZFibHt33Mfxew>

Todd Grande **Bonferroni**

Link: <https://www.youtube.com/watch?v=pbWkwEz-XBY>

Zach **Kruskal Wallis**

Link: <https://www.statology.org/kruskal-wallis-test-python/#~:text=A%20kruskal%2Dwallis%20Test%20is,o%20the%20One%2DWay%20ANOVA.>

USP **Teste Exato de Fisher**

Link: <http://wiki.icmc.usp.br/images/7/73/Chisq-fisher.pdf>