

Estudo Estatístico Referente ao Peso do Cérebro Humano

Benício Ramos Magalhães

Base de dados e motivação de estudo fornecidos por Mariana Cúri - Cemeal - ICMC/USP São Carlos

Os dados do arquivo Brain, referem-se ao peso do cérebro (g), tamanho da cabeça (cm3) de 237 adultos, identificados por sexo e grupo etário. O estudo tem por objetivo verificar se:

- 1) Há diferença no peso cerebral entre os sexos? E entre os grupos etários?
- 2) O tamanho da cabeça é preditor do peso cerebral e, neste caso, há diferença nessa relação entre os sexos e entre os grupos etários?
- 3) Estime o peso médio do cérebro de homens e de mulheres (pontual e intervalar).</>

```
In [277]: #bibliotecas
import math
import numpy as np
import pandas as pd
import seaborn as sns
import scipy.stats as sm
import statmodels.api as sm
import matplotlib.pyplot as plt
import statmodels.stats.api as sms
import statmodels.formula.api as smf

from scipy.stats import f
from scipy.stats import t
from scipy.stats import sem
from scipy.stats import norm
from scipy.stats import kstest
from scipy.stats import levene
from scipy.stats import shapiro
from scipy.stats import ks_2samp
from scipy.stats import anderson
from scipy.stats import probplot
from scipy.stats import ttest_ind
from statmodels.formula.api import ols
from scipy.stats import chi2_contingency
from statmodels.stats.stattools import durbin_watson
from statmodels.stats.outliers_influence import variance_inflation_factor

In [175]: #carregando dados originais
data = pd.read_csv('data/Brain.csv')
data.describe()
```

Out[175]:

	Homem	Acima45	Tamanho	Peso
count	237.000000	237.000000	237.000000	237.000000
mean	0.568401	0.538865	363.991561	1282.873446
std	0.498753	0.498768	365.261422	120.340498
min	0.000000	0.000000	2720.000000	955.000000
25%	0.000000	0.000000	3359.000000	1027.000000
50%	1.000000	1.000000	3614.000000	1280.000000
75%	1.000000	1.000000	3876.000000	1350.000000
max	1.000000	1.000000	4747.000000	1635.000000

```
In [185]: #separando as variáveis
#series de peso
peso_homem = data['Peso'][data['Homem']==1]
peso_mulher = data['Peso'][data['Homem']==0]
peso_idade_maior45 = data['Peso'][data['Acima45']==1]
peso_idade_menor45 = data['Peso'][data['Acima45']==0]
```

1) Há diferença no peso cerebral entre os sexos? E entre os grupos etários?

Como o objetivo é comparar duas populações com relação a uma variável quantitativa (peso cerebral) baseado numa amostra, o ideal é usarmos um teste de hipótese. Portanto, iremos verificar se existe diferença no peso cerebral entre os sexos aplicando um teste de hipótese para igualdade das médias.

Para isso, algumas suposições que devemos verificar:

- Se existe ou não dependência entre as amostras (medições pareadas ou não pareadas)
- Se os dados seguem uma distribuição normal.

OBS: Após análise para os sexos, iremos fazer a mesma análise para faixa etária.

Passos para o teste de hipótese:

a) Especificar as hipóteses  $H_0$  e  $H_1$ :

$H_0: \mu_1 = \mu_2 \rightarrow$  **Peso cerebral do homem é igual ao peso cerebral da mulher (hipótese nula)**

$H_1: \mu_1 \neq \mu_2 \rightarrow$  **Peso cerebral do homem não é igual ao peso cerebral da mulher (hipótese alternativa)**

</>

b) Especificar a estatística do teste e sua distribuição, sob  $H_0$ :

```
In [3]: print('Estatística dos dados para Homem:')
peso_homem.describe()

Out[3]: count    134.000000
mean    1331.858209
std      108.933390
min     1220.000000
25%     1252.750000
50%     1313.500000
75%     1400.000000
max     1635.000000
Name: Peso, dtype: float64

In [4]: print('Estatística dos dados para Mulher:')
peso_mulher.describe()

Out[4]: count    103.000000
mean    1219.145631
std      103.829933
min      955.000000
25%     1146.000000
50%     1220.000000
75%     1290.000000
max     1520.000000
Name: Peso, dtype: float64
```

Vamos verificar a normalidade da distribuição com a realização de alguns testes.

Vamos considerar os seguintes testes para distribuição normal:

- Kolmogorov-Smirnov
- Anderson-Darling
- Shapiro-Wilk

```
In [5]: #Teste de Kolmogorov-Smirnov:
ks_stat_homem, ks_p_homem = kstest(peso_homem, 'norm', args=(peso_homem.describe()[1], peso_homem.describe()[2], N=len(peso_homem)))
ks_stat_mulher, ks_p_mulher = kstest(peso_mulher, 'norm', args=(peso_mulher.describe()[1], peso_mulher.describe()[2], N=len(peso_mulher)))

#para número da amostra > 50:
ks_val_homem = 1.36/(np.sqrt(len(peso_homem)))
ks_val_mulher = 1.36/(np.sqrt(len(peso_mulher)))

ks_res_homem = 'Distribuição é normal' if ks_val_homem >= ks_stat_homem else 'Distribuição NÃO é normal'
ks_res_mulher = 'Distribuição é normal' if ks_val_mulher >= ks_stat_mulher else 'Distribuição NÃO é normal'

In [6]: #Teste Anderson-Darling:
ad_stat_homem, ad_p_homem, ad_ic_homem = anderson(peso_homem, 'norm')
ad_stat_mulher, ad_p_mulher, ad_ic_mulher = anderson(peso_mulher, 'norm')

#para ic 5%
ad_res_homem = 'Distribuição é normal' if ad_stat_homem < ad_p_homem[2] else 'Distribuição NÃO é normal'
ad_res_mulher = 'Distribuição é normal' if ad_stat_mulher < ad_p_mulher[2] else 'Distribuição NÃO é normal'

In [7]: #Teste de Shapiro-Wilk:
sh_stat_homem, sh_p_homem = shapiro(peso_homem)
sh_stat_mulher, sh_p_mulher = shapiro(peso_mulher)

#para ic 5%
sh_res_homem = 'Distribuição é normal' if sh_p_homem > 0.05 else 'Distribuição NÃO é normal'
sh_res_mulher = 'Distribuição é normal' if sh_p_mulher > 0.05 else 'Distribuição NÃO é normal'

In [8]: #Resultados dos testes de normalidade HOMEM
print('Resultados dos testes de normalidade para homem:')
print('-----Teste-----|-----Estatística-----|-----P-Valor-----|-----Result')
print('ado (IC 5%)-----|-----|-----|-----|')
print('Kolmogorov-Smirnov |', 'ks_stat_homem,', '|', 'ks_p_homem,', '|', 'ks_res_homem,', '|')
print('-----|-----|-----|-----|')
print('Anderson-Darling |', 'ad_stat_homem,', '|', 'ad_p_homem[2],', '|', 'ad_res_homem,', '|')
print('-----|-----|-----|-----|')
print('Shapiro-Wilk |', 'sh_stat_homem,', '|', 'sh_p_homem,', '|', 'sh_res_homem,', '|')
print('-----|-----|-----|-----|')

Resultados dos testes de normalidade para homem:
-----Teste-----|-----Estatística-----|-----P-Valor-----|-----Result
ado (IC 5%)-----|-----|-----|-----|
Kolmogorov-Smirnov | 0.07319296053954616 | 0.4535070134115867 | Distribuição é normal
Anderson-Darling | 0.7592589134145046 | 0.765 | Distribuição é normal
Shapiro-Wilk | 0.9780169129371643 | 0.02875436283648014 | Distribuição NÃO é normal

In [9]: #Resultados dos testes de normalidade MULHER
print('Resultados dos testes de normalidade para mulher:')
print('-----Teste-----|-----Estatística-----|-----P-Valor-----|-----Result')
print('Resultado (IC 5%)-----|-----|-----|-----|')
print('Kolmogorov-Smirnov |', 'ks_stat_mulher,', '|', 'ks_p_mulher,', '|', 'ks_res_mulher,', '|')
print('-----|-----|-----|-----|')
print('Anderson-Darling |', 'ad_stat_mulher,', '|', 'ad_p_mulher[2],', '|', 'ad_res_mulher,', '|')
print('-----|-----|-----|-----|')
print('Shapiro-Wilk |', 'sh_stat_mulher,', '|', 'sh_p_mulher,', '|', 'sh_res_mulher,', '|')
print('-----|-----|-----|-----|')

Resultados dos testes de normalidade para mulher:
-----Teste-----|-----Estatística-----|-----P-Valor-----|-----Result
ado (IC 5%)-----|-----|-----|-----|
Kolmogorov-Smirnov | 0.04574294927683009 | 0.9823676862451942 | Distribuição é normal
Anderson-Darling | 0.1429274743232656 | 0.759 | Distribuição é normal
Shapiro-Wilk | 0.9959982633590698 | 0.9919323921203613 | Distribuição é normal
```

Os resultados de 2 testes foram favoráveis a dizer que os dados tem distribuição normal dado o intervalo de confiança de 5%. Assim, iremos considerar o peso para homem com uma distribuição normal gaussiana.

```
In [9]: #Resultados dos testes de normalidade MULHER
print('Resultados dos testes de normalidade para mulher:')
print('-----Teste-----|-----Estatística-----|-----P-Valor-----|-----Result')
print('Resultado (IC 5%)-----|-----|-----|-----|')
print('Kolmogorov-Smirnov |', 'ks_stat_mulher,', '|', 'ks_p_mulher,', '|', 'ks_res_mulher,', '|')
print('-----|-----|-----|-----|')
print('Anderson-Darling |', 'ad_stat_mulher,', '|', 'ad_p_mulher[2],', '|', 'ad_res_mulher,', '|')
print('-----|-----|-----|-----|')
print('Shapiro-Wilk |', 'sh_stat_mulher,', '|', 'sh_p_mulher,', '|', 'sh_res_mulher,', '|')
print('-----|-----|-----|-----|')

Resultados dos testes de normalidade para mulher:
-----Teste-----|-----Estatística-----|-----P-Valor-----|-----Result
ado (IC 5%)-----|-----|-----|-----|
Kolmogorov-Smirnov | 0.04574294927683009 | 0.9823676862451942 | Distribuição é normal
Anderson-Darling | 0.1429274743232656 | 0.759 | Distribuição é normal
Shapiro-Wilk | 0.9959982633590698 | 0.9919323921203613 | Distribuição é normal
```

Os resultados de todos os testes foram favoráveis a dizer que os dados tem distribuição normal dado o intervalo de confiança de 5%. Assim, iremos considerar o peso para mulher com uma distribuição normal gaussiana.

Considerando a natureza dos dados, as diferenças sugerem que as amostras são independentes, pois a medição do peso do cérebro são feitas em indivíduos diferentes e, portanto, não aparentam ter relação de dependência. Para fundamentarmos melhor essa afirmação, iremos realizar um teste Z para a hipótese nula de que a média das amostras são iguais.

```
In [280]: #teste Z de normalidade
#para este caso vamos considerar peso como uma variável quantitativa e compará-la entre os dois sexos.
stat_peso_homem = [np.mean(peso_homem), np.std(peso_homem), np.std(peso_homem)**2]
stat_peso_mulher = [np.mean(peso_mulher), np.std(peso_mulher), np.std(peso_mulher)**2]

#valor da estatística do teste
z = (stat_peso_homem[0] - stat_peso_mulher[0]) / (math.sqrt((stat_peso_homem[2]) + math.sqrt(len(peso_homem) + len(peso_mulher)))) # estatística do teste Z
z_p = norm.ppf(0.05) #cálculo do ponto crítico (x_corte)
z_pvalor = norm.cdf(z) #cálculo do p-valor

z_res = 'Amostras independentes' if z_pvalor > 0.05 else 'Amostras dependentes'

print('Resultado:', z_pvalor)
print('res:', z_res)

Resultado: 0.999976904283696
Amostras independentes

c) Fixar o nível de significância do teste (α)
```

Neste caso, iremos considerar o nível de significância do teste de 5%, ou seja,  $\alpha = 0.05$

d) Calcular o p-valor (ou região crítica do teste)

Como já concluímos que as amostras são independentes e ambas seguem uma distribuição normal, precisamos avaliar agora se existe uma relação de igualdade das variâncias para decidir qual teste t de Student para médias de duas amostras iremos aplicar. Para isso, iremos aplicar um teste de levene.

```
In [11]: #teste de levene
lv_stat, lv_p = levene(peso_homem, peso_mulher)

lv_res = 'Não existe grande diferença na variância.' if lv_p > 0.05 else 'Existe diferença na variância'

print('Resultado:', lv_p)
print('res:', lv_res)

Resultado: 0.76701602271913
Não existe grande diferença na variância.
```

Iremos agora fazer o teste de hipótese das médias da variável média do peso cerebral entre os sexos serem iguais. Para isso, será realizado o teste t de Student (bicudal) para média de duas populações Normais com variâncias iguais.

```
In [12]: #teste t de Student
t_stat, t_p = ttest_ind(peso_homem, peso_mulher)

t_res = 'Aceita hipótese H0' if t_p > 0.05 else 'Rejeita hipótese H0'

print('Resultado:', t_p)
print('res:', t_res)

Resultado: 3.919241152559185e-14
Rejeita hipótese H0

e) Decidir entre  $H_0$  e  $H_1$ , comparando com o p-valor com α
```

Considerando os resultados obtidos no teste t de Student, precisamos rejeitar a hipótese H0, ou seja, com  $\alpha = 0.05$  podemos afirmar que o peso cerebral do homem não é igual ao peso cerebral da mulher.

Resposta:

Sim. Existe diferença no peso cerebral entre os sexos.

Entre os grupos etários?

Realizando mesma análise anterior para os grupos etários.

$H_0: \mu_1 = \mu_2 \rightarrow$  **Peso cerebral de pessoas acima de 45 anos é igual ao de pessoas abaixo de 45 anos. (hipótese nula)**

$H_1: \mu_1 \neq \mu_2 \rightarrow$  **Peso cerebral de pessoas acima de 45 anos não é igual ao de pessoas abaixo de 45 anos. (hipótese alternativa)**

```
In [13]: #estatística e normalidade
print('Idade maior que 45 anos:')
print(peso_idade_maior45.describe())
print('Idade menor que 45 anos:')
print(peso_idade_menor45.describe())

#Teste de Kolmogorov-Smirnov:
ks_stat_maior45, ks_p_maior45 = kstest(peso_idade_maior45, 'norm', args=(peso_idade_maior45.describe()[1], peso_idade_maior45.describe()[2], N=len(peso_idade_maior45)))
ks_stat_menor45, ks_p_menor45 = kstest(peso_idade_menor45, 'norm', args=(peso_idade_menor45.describe()[1], peso_idade_menor45.describe()[2], N=len(peso_idade_menor45)))

#para número da amostra > 50:
ks_val_maior45 = 1.36/(np.sqrt(len(peso_idade_maior45)))
ks_val_menor45 = 1.36/(np.sqrt(len(peso_idade_menor45)))

ks_res_maior45 = 'Distribuição é normal' if ks_val_maior45 >= ks_stat_maior45 else 'Distribuição NÃO é normal'
ks_res_menor45 = 'Distribuição é normal' if ks_val_menor45 >= ks_stat_menor45 else 'Distribuição NÃO é normal'

#Teste Anderson-Darling:
ad_stat_maior45, ad_p_maior45, ad_ic_maior45 = anderson(peso_idade_maior45, 'norm')
ad_stat_menor45, ad_p_menor45, ad_ic_menor45 = anderson(peso_idade_menor45, 'norm')

#para ic 5%
ad_res_maior45 = 'Distribuição é normal' if ad_stat_maior45 < ad_p_maior45[2] else 'Distribuição NÃO é normal'
ad_res_menor45 = 'Distribuição é normal' if ad_stat_menor45 < ad_p_menor45[2] else 'Distribuição NÃO é normal'

#Teste de Shapiro-Wilk:
sh_stat_maior45, sh_p_maior45 = shapiro(peso_idade_maior45)
sh_stat_menor45, sh_p_menor45 = shapiro(peso_idade_menor45)

#para ic 5%
sh_res_maior45 = 'Distribuição é normal' if sh_p_maior45 > 0.05 else 'Distribuição NÃO é normal'
sh_res_menor45 = 'Distribuição é normal' if sh_p_menor45 > 0.05 else 'Distribuição NÃO é normal'

print()
print('Resultados dos testes de normalidade para idade maior que 45 anos:')
print('-----Teste-----|-----Estatística-----|-----P-Valor-----|-----|')
print('Kolmogorov-Smirnov |', 'ks_stat_maior45,', '|', 'ks_p_maior45,', '|', 'ks_res_maior45,', '|')
print('-----|-----|-----|-----|')
print('Anderson-Darling |', 'ad_stat_maior45,', '|', 'ad_p_maior45[2],', '|', 'ad_res_maior45,', '|')
print('-----|-----|-----|-----|')
print('Shapiro-Wilk |', 'sh_stat_maior45,', '|', 'sh_p_maior45,', '|', 'sh_res_maior45,', '|')
print('-----|-----|-----|-----|')

print()
print('Resultados dos testes de normalidade para idade menor que 45 anos:')
print('-----Teste-----|-----Estatística-----|-----P-Valor-----|-----|')
print('Kolmogorov-Smirnov |', 'ks_stat_menor45,', '|', 'ks_p_menor45,', '|', 'ks_res_menor45,', '|')
print('-----|-----|-----|-----|')
print('Anderson-Darling |', 'ad_stat_menor45,', '|', 'ad_p_menor45[2],', '|', 'ad_res_menor45,', '|')
print('-----|-----|-----|-----|')
print('Shapiro-Wilk |', 'sh_stat_menor45,', '|', 'sh_p_menor45,', '|', 'sh_res_menor45,', '|')
print('-----|-----|-----|-----|')

Idade maior que 45 anos:
count    127.000000
mean    1263.937008
std      120.925712
min     1027.000000
25%     1180.000000
50%     1250.000000
75%     1332.000000
max     1620.000000
Name: Peso, dtype: float64

Idade menor que 45 anos:
count    110.000000
mean    1304.736364
std      116.409366
min     1027.000000
25%     1227.500000
50%     1301.000000
75%     1370.750000
max     1635.000000
Name: Peso, dtype: float64

Resultados dos testes de normalidade para idade maior que 45 anos:
-----Teste-----|-----Estatística-----|-----P-Valor-----|-----Result
ado (IC 5%)-----|-----|-----|-----|
Kolmogorov-Smirnov | 0.054175986782391106 | 0.8500369044822889 | Distribuição é normal
Anderson-Darling | 0.32573942719172445 | 0.764 | Distribuição é normal
Shapiro-Wilk | 0.9917106302371521 | 0.65547625170135498 | Distribuição é normal

Resultados dos testes de normalidade para idade menor que 45 anos:
-----Teste-----|-----Estatística-----|-----P-Valor-----|-----Result
ado (IC 5%)-----|-----|-----|-----|
Kolmogorov-Smirnov | 0.06000072085469621 | 0.8232879168286493 | Distribuição é normal
Anderson-Darling | 0.41369054165920716 | 0.761 | Distribuição é normal
Shapiro-Wilk | 0.987789169598389 | 0.42216619849205017 | Distribuição é normal
```

```
In [281]: #teste Z de dependência
#para este caso vamos considerar peso como uma variável quantitativa e compará-la entre os dois sexos.
stat_peso_idade_maior45 = [np.mean(peso_idade_maior45), np.std(peso_idade_maior45), np.std(peso_idade_maior45)**2]
stat_peso_idade_menor45 = [np.mean(peso_idade_menor45), np.std(peso_idade_menor45), np.std(peso_idade_menor45)**2]

#valor da estatística do teste
z = (stat_peso_idade_maior45[0] - stat_peso_idade_menor45[0]) / (math.sqrt((stat_peso_idade_maior45[2]) + math.sqrt(len(peso_idade_maior45) + len(peso_idade_menor45)))) # estatística do teste Z
z_p = norm.ppf(0.05) #cálculo do ponto crítico (x_corte)
z_pvalor = norm.cdf(z) #cálculo do p-valor

z_res = 'Amostras independentes' if z_pvalor > 0.05 else 'Amostras dependentes'

print('Resultado:', z_pvalor)
print('res:', z_res)

Resultado: 0.0919172475388757
Amostras independentes
```

Iremos considerar o nível de significância do teste de 5%, ou seja,  $\alpha = 0.05$

```
In [15]: #teste t de Student (bicudal, independente e distribuição normal)
t_stat_idade, t_p = ttest_ind(peso_idade_maior45, peso_idade_menor45)

t_res = 'Aceita hipótese H0' if t_p > 0.05 else 'Rejeita hipótese H0'

print('Resultado:', t_p_idade)
print('res:', t_res_idade)

Resultado: 0.00895602315452554
Rejeita hipótese H0
```

Considerando os resultados obtidos no teste t de Student, precisamos rejeitar a hipótese H0, ou seja, com  $\alpha = 0.05$  podemos afirmar que o peso cerebral de pessoas acima de 45 anos não é igual ao de pessoas abaixo de 45 anos.

Resposta:

Sim. Existe diferença no peso cerebral entre os grupos etários.

2) O tamanho da cabeça é preditor do peso cerebral e, neste caso, há diferença nessa relação entre os sexos e entre os grupos etários?

Inicialmente iremos verificar se o tamanho da cabeça é realmente preditor do peso cerebral, utilizando um modelo de regressão linear simples sendo X composto de uma única variável explicativa (tamanho da cabeça). A interpretação para este modelo nos ajudará a entender a resposta para essa primeira parte da pergunta.

```
In [168]: #regressão linear simples
#Modelo 1: Tamanho
mod1 = ols('Peso ~ Tamanho', data=data)
res1 = mod1.fit()
print(res1.summary())

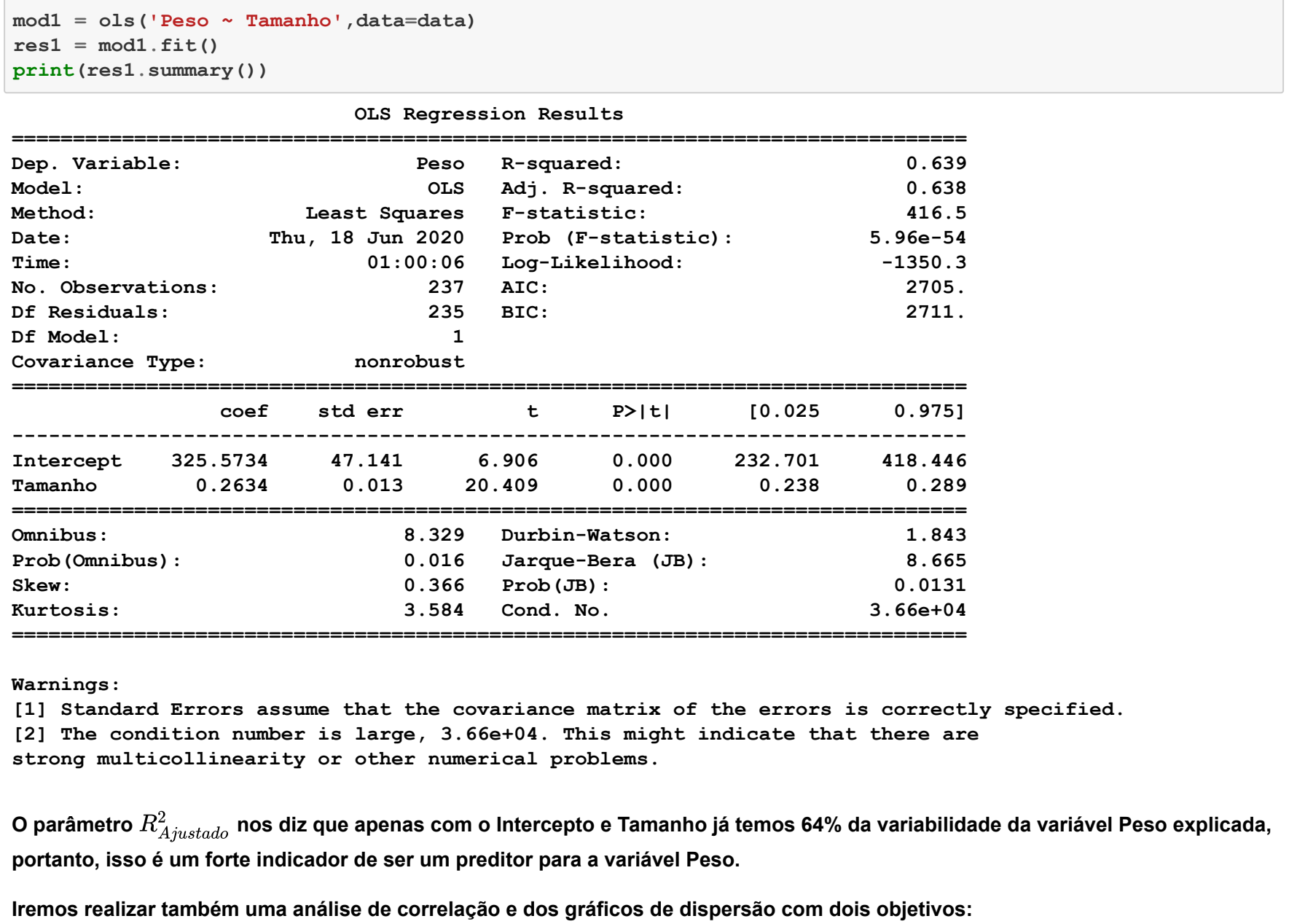
OLS Regression Results
=====
Dep. Variable:      Peso      R-squared:      0.639
Model:              OLS      Adj. R-squared:    0.638
Method:              Least Squares      P-statistic:    416.5
Date:                Thu, 18 Jun 2020      Prob (F-statistic):    5.96e-54
Time:                10:00:06      Log-Likelihood:    -1350.3
No. Observations:    237      AIC:                2705.
DF Residuals:        235      BIC:                2721.
DF Model:            1
Covariance Type:     nonrobust

=====
coef    std err      t      P>|t|      [0.025      0.975]
-----
Intercept    325.8990    102.157    3.184    0.002    152.369    647.399
Tamanho      2.2619    0.036    62.006    0.000    207.515    610.083
Homem         -0.2434    0.046    -0.947    0.345    -0.134    0.047
Acima45       -180.9501    167.602    -1.076    0.283    -510.598    149.880
Tamanho:Acima45    -196.4052    104.039    -1.888    0.060    -401.398    8.585
Homem:Acima45    -65.0907    222.855    -0.292    0.770    -504.200    374.018
Tamanho:Homem:Acima45    0.0076    0.062    0.122    0.900    -0.115    0.130
=====
Omnibus:            8.421      Durbin-Watson:      1.922
Prob(Omnibus):      0.015      Jarque-Bera (JB):    8.927
Skew:               0.359      Prob(JB):          0.015
Kurtosis:           3.624      Cond. No.          3.43e+05
=====
Warnings:
[1] Standard Errors assume that the covariance matrix of the errors is correctly specified.
[2] The condition number is large, 3.43e+05. This might indicate that there are strong multicollinearity or other numerical problems.
```

O parâmetro  $\beta_1$  estimado, nos diz que apenas com o intercepto e Tamanho já temos 64% da variabilidade da variável Peso explicada, portanto, isso é um forte indicador de ser um preditor para a variável Peso.

Iremos realizar também uma análise de correlação e dos gráficos de dispersão com dois objetivos:

- Verificar a correlação entre Tamanho e Peso, para embasar melhor o fato de que Tamanho é preditor de Peso.
- Verificar resultados da análise de multicolinearidade e simplificar o modelo final. Se as variáveis preditoras apresentarem correlações altas entre si, resulta na possibilidade de termos fatores redundantes no modelo e isso poderá aumentar a variância dos coeficientes da regressão, tornando-os instáveis.



É possível observar uma forte relação entre a variável preditora Tamanho e a variável resposta Peso (80%), o que sugere que de fato isso é um bom preditor da variável Peso.

Entre as preditoras, destacamos a relação entre Tamanho e Homem (51%), o que poderia gerar um problema de multicolinearidade.

Vamos iniciar a criação do modelo com todas as variáveis e interações possíveis e através da técnica de stepwise, iremos adequar e otimizar o modelo final, eliminando as variáveis ou interações não significativas. Após ajustarmos o modelo, vamos aplicar mais algumas técnicas para analisar com maiores detalhes as questões de multicolinearidade e verificar se há diferença na relação Tamanho e Peso quando consideramos grupo etário e gênero.

```
In [53]: #regressão linear múltipla
#Modelo 2: Tamanho, Homem e Acima45 com todas as interações possíveis
mod2 = ols('Peso ~ Tamanho * Homem * Acima45', data=data)
res2 = mod2.fit()
print(res2.summary())

OLS Regression Results
=====
Dep. Variable:      Peso      R-squared:      0.661
Model:              OLS      Adj. R-squared:    0.650
Method:              Least Squares      P-statistic:    63.68
Date:                Wed, 17 Jun 2020      Prob (F-statistic):    2.90e-13
Time:                17:35:25      Log-Likelihood:    -1343.0
No. Observations:    237      AIC:                2702.
DF Residuals:        229      BIC:                2730.
DF Model:            4
Covariance Type:     nonrobust

=====
coef    std err      t      P>|t|      [0.025      0.975]
-----
Intercept    399.8990    125.610    3.184    0.002    152.369    647.399
Tamanho      2.2619    0.036    62.006    0.000    207.515    610.083
Homem         -0.2434    0.046    -0.947    0.345    -0.134    0.047
Acima45       -180.9501    167.602    -1.076    0.283    -510.598    149.880
Tamanho:Acima45    -196.4052    104.039    -1.888    0.060    -401.398    8.585
Homem:Acima45    -65.0907    222.855    -0.292    0.770    -504.200    374.018
Tamanho:Homem:Acima45    0.0076    0.062    0.122    0.900    -0.115    0.130
=====
Omnibus:            8.421      Durbin-Watson:      1.922
Prob(Omnibus):      0.015      Jarque-Bera (JB):    8.927
Skew:               0.359      Prob(JB):          0.015
Kurtosis:           3.624      Cond. No.          3.43e+05
=====
Warnings:
[1] Standard Errors assume that the covariance matrix of the errors is correctly specified.
[2] The condition number is large, 3.43e+05. This might indicate that there are strong multicollinearity or other numerical problems.
```

Vamos analisar a interação Tamanho:Homem:Acima45.

Neste caso, temos um teste de hipótese que considera:

$H_0: \beta_9 = 0$

$H_1: \beta_9 \neq 0$

Com um p-valor de 90%, coeficiente de regressão 0.0076 e intervalo de confiança entre -0.115 e 0.130, nós aceitamos a hipótese que  $\beta_9 = 0$  no modelo, portanto, esta interação torna-se insignificante e será retirada.

Vamos então aplicando a técnica de stepwise para a seleção das demais variáveis para melhorarmos o modelo final.

```
In [55]: #regressão linear múltipla
#Modelo 3: Tamanho, Homem, Acima45 e respectivas interações
mod3 = ols('Peso ~ Tamanho * Homem * Tamanho * Acima45', data=data)
res3 = mod3.fit()
print(res3.summary())

OLS Regression Results
=====
Dep. Variable:      Peso      R-squared:      0.661
Model:              OLS      Adj. R-squared:    0.652
Method:              Least Squares      P-statistic:    63.68
Date:                Wed, 17 Jun 2020      Prob (F-statistic):    2.90e-13
Time:                17:46:27      Log-Likelihood:    -1343.1
No. Observations:    237      AIC:                2702.
DF Residuals:        232      BIC:                2717.
DF Model:            5
Covariance Type:     nonrobust

=====
coef    std err      t      P>|t|      [0.025      0.975]
-----
Intercept    453.0924    73.463    6.168    0.000    308.353    597.831
Tamanho      0.2285    0.020    11.558    0.000    0.188    0.269
Tamanho:Tamanho    0.0034    0.030    0.112    0.912    -0.054    0.061
Acima45       -129.5507    93.796    -1.381    0.169    -314.332    55.270
Tamanho:Acima45    0.0290    0.026    1.131    0.259    -0.022    0.079
=====
Omnibus:            7.248      Durbin-Watson:      1.918
Prob(Omnibus):      0.027      Jarque-Bera (JB):    7.449
Skew:               0.331      Prob(JB):          0.021
Kurtosis:           3.662      Cond. No.          1.04e+05
=====
Warnings:
[1] Standard Errors assume that the covariance matrix of the errors is correctly specified.
[2] The condition number is large, 1.04e+05. This might indicate that there are strong multicollinearity or other numerical problems.
```

Considerando a mesma análise anterior (p-valor > 0.05 e IC inclui o zero), iremos retirar a interação Tamanho:Acima45.

```
In [59]: #regressão linear múltipla
#Modelo 5: Tamanho, Homem, Acima45, Interação: Tamanho - Acima45.
mod5 = ols('Peso ~ Tamanho * Homem * Tamanho * Acima45', data=data)
res5 = mod5.fit()
print(res5.summary())

OLS Regression Results
=====
Dep. Variable:      Peso      R-squared:      0.665
Model:              OLS      Adj. R-squared:    0.651
Method:              Least Squares      P-statistic:    110.0
Date:                Wed, 17 Jun 2020      Prob (F-statistic):    2.05e
```



