

Grupo:

Andreza Aparecida dos Santos (164213)

Leonardo Marçal (225240)

Ligia Vasconcellos (081938)

Mariana Amaral (262866)

```
In [1]: import sklearn
import pandas as pd
import numpy as np
import seaborn as sns
import matplotlib.pyplot as plt
from scipy.stats import pearsonr, spearmanr

from sklearn.linear_model import LogisticRegressionCV, LinearRegression
from sklearn.model_selection import train_test_split
from sklearn.metrics import classification_report, confusion_matrix

In [2]: def calculate_p_values(df, method):
    df = df.dropna().get_numeric_data()
    df_columns = pd.DataFrame(columns=df.columns)
    p_values = df_columns.transpose().join(df_columns, how='outer')
    for r in df.columns:
        for c in df.columns:
            if method == 'pearson':
                p_values[r][c] = str(pearsonr(df[r], df[c])[1])
            else:
                p_values[r][c] = str(spearmanr(df[r], df[c])[1])
    return p_values
```

Leitura dos dados

```
In [3]: df_zombie_survey = pd.read_csv('data/zombie-survey.csv')
df_zombie_meals = pd.read_csv('data/zombie-meals.csv')
```

Análise 1: Teste de Hipóteses

Tarefa: Considerando a base Zombies Survey (zombie-survey.csv), apresente um teste de hipóteses que realize um estudo sobre a relação entre altura e gênero dos zumbis. Construa suas hipóteses, apresente a análise e as conclusões.

Hipóteses

H0 (Hipótese nula) = Não há diferença de altura (HEIGHT) entre homens (MALE) e mulheres (FEMALE)

H1 (Hipótese alternativa) = Há diferença de altura entre homens e mulheres

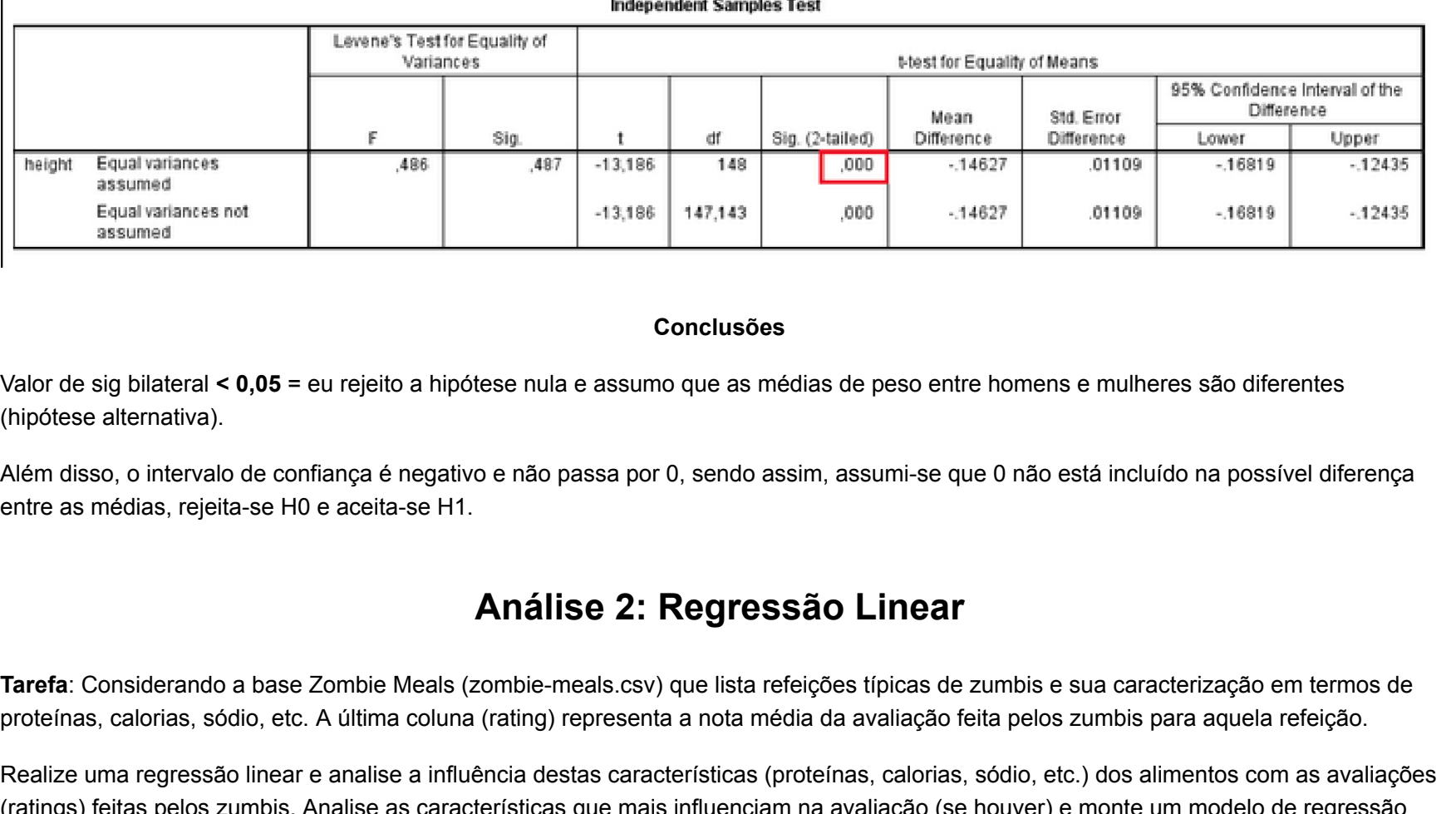
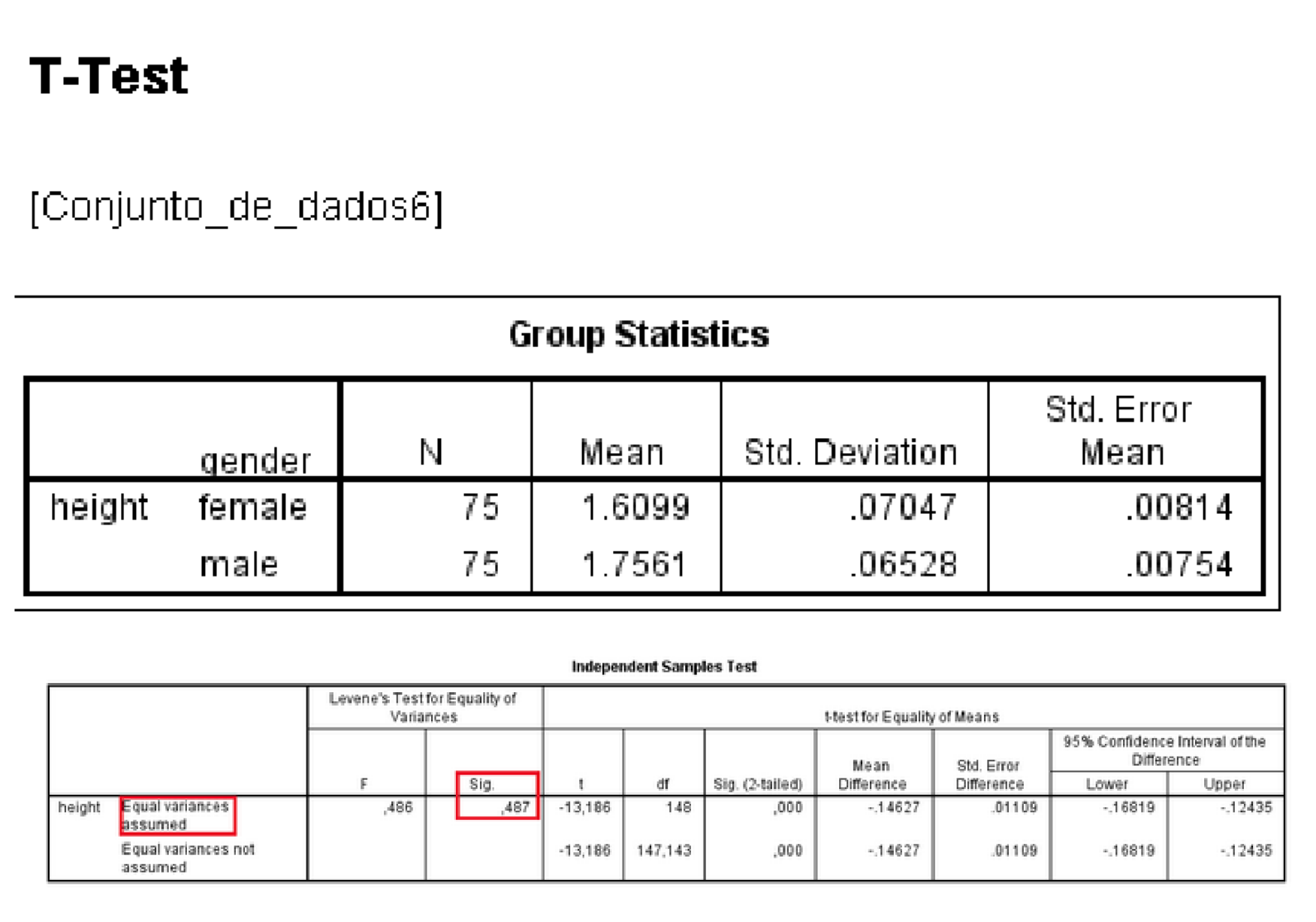
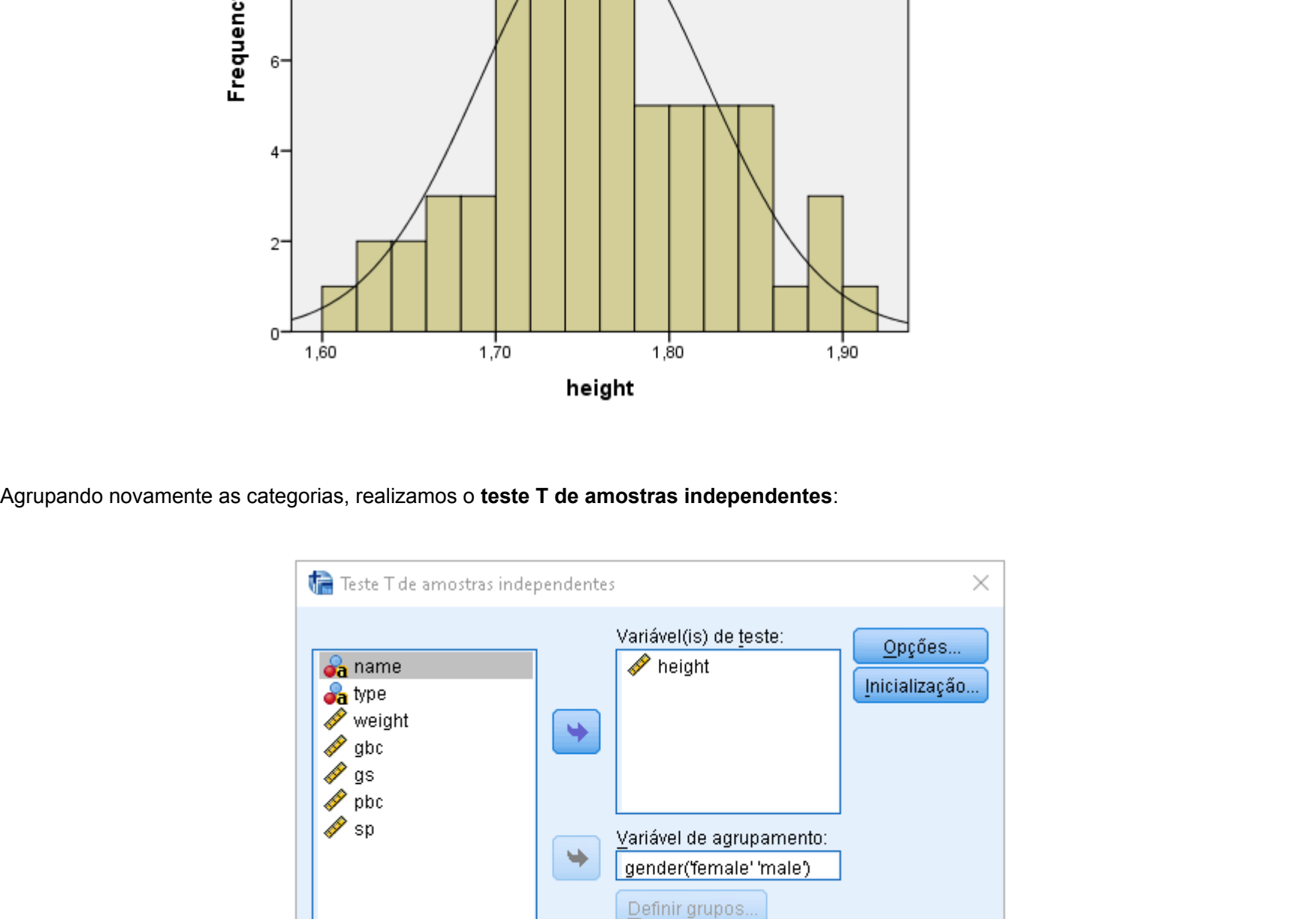
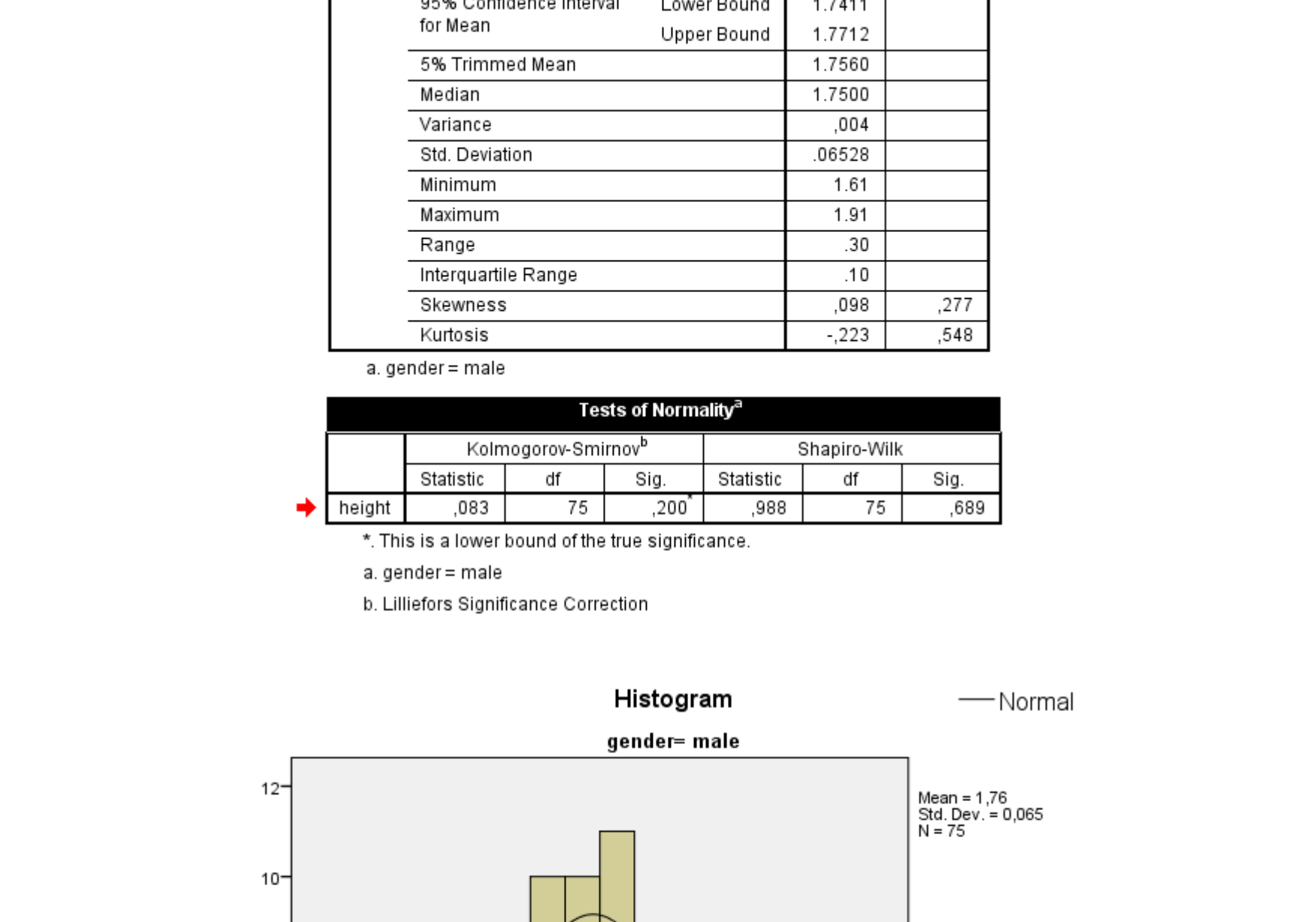
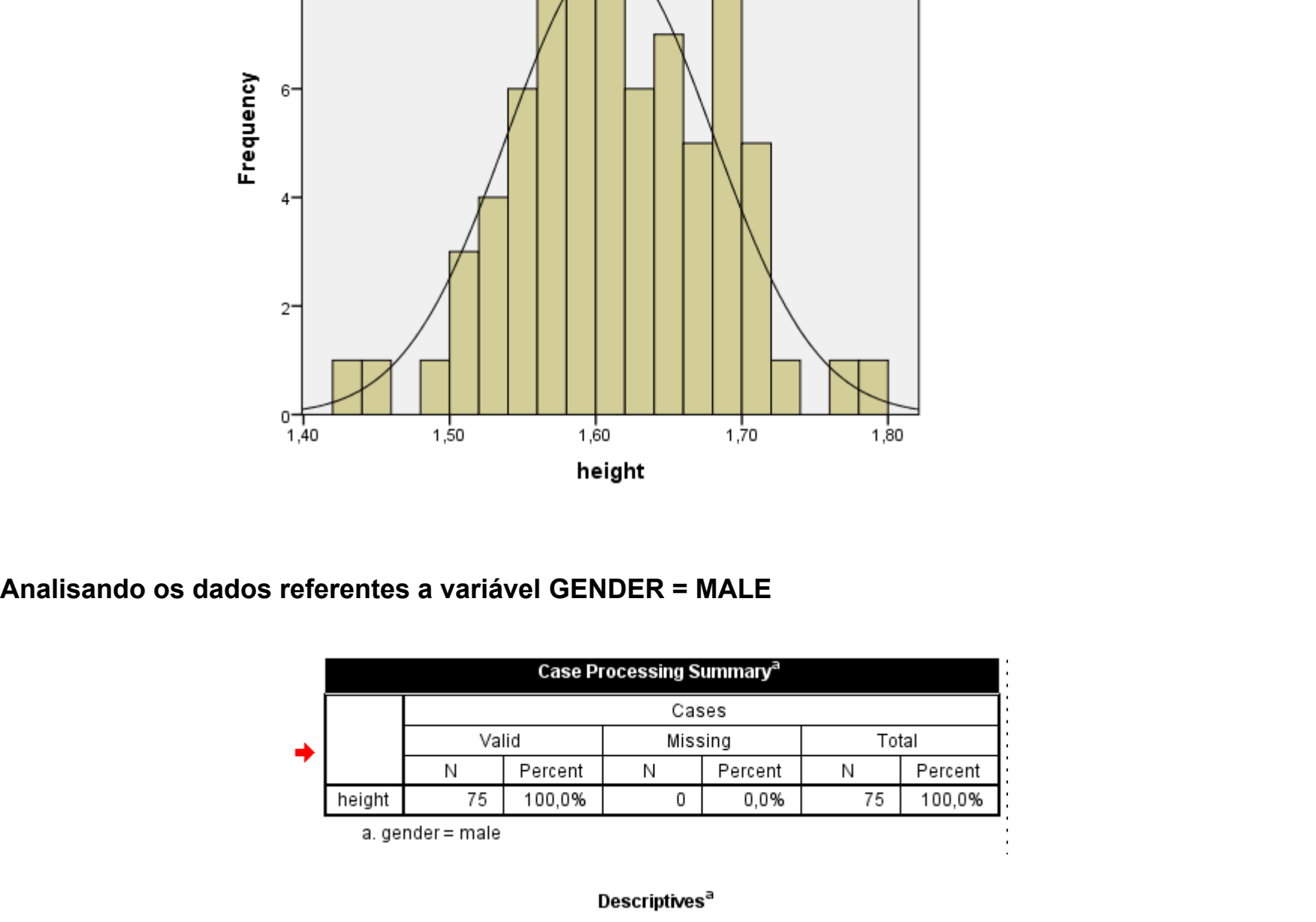
Análises

Realizada no software estatístico SPSS

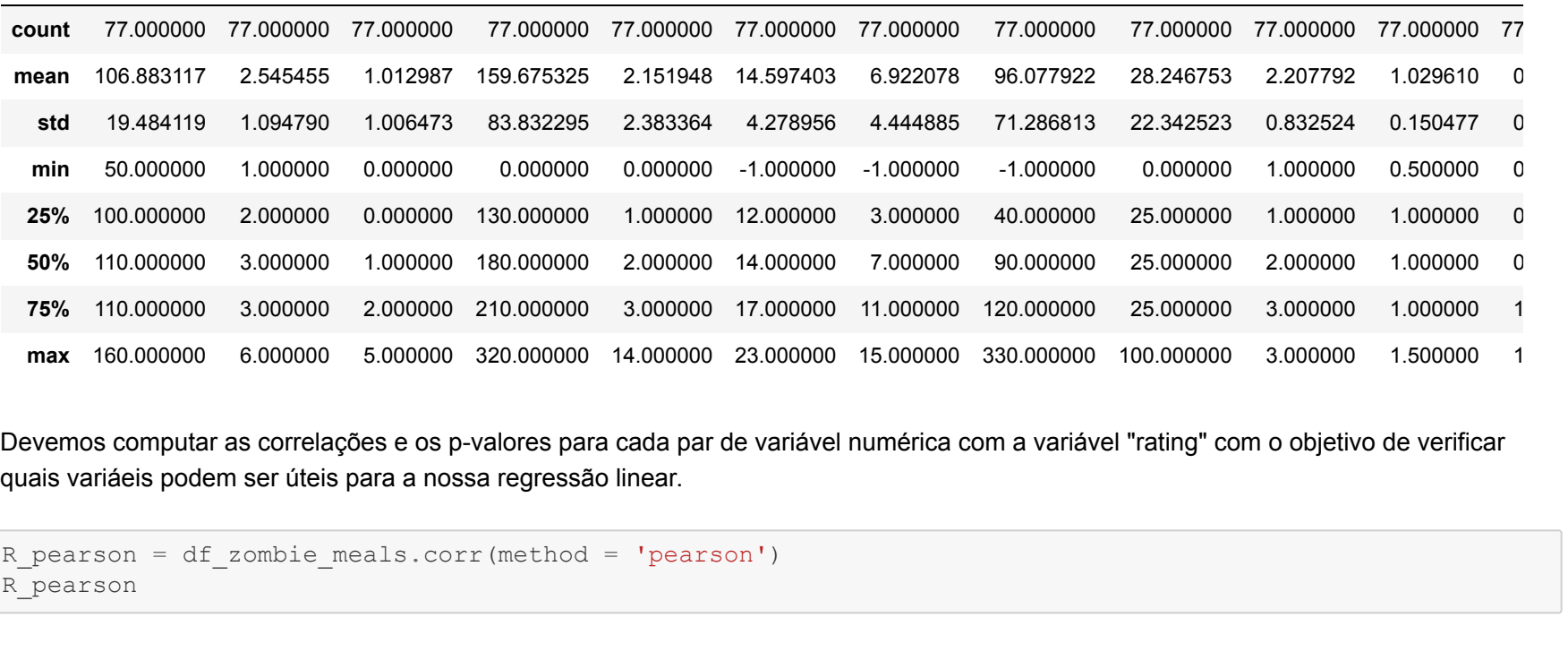
O conjunto de dados foi analisado dividindo-se a variável categórica independente **GENDER**

Em seguida foi analisado a relação entre as categorias da variável independente (**FEMALE e MALE**) e a variável dependente **HEIGHT**

Analisando os dados referentes a variável GENDER = FEMALE



Agrupando novamente as categorias, realizamos o teste T de amostras independentes:



T-Test

[Conjunto de dados6]

Group Statistics					
gender		N	Mean	Std. Deviation	Std. Error Mean
height	female	75	1.6099	.07047	.00814
	male	75	1.7561	.06528	.00754

Independent Samples Test											
Levene's Test for Equality of Variances						t-test for Equality of Means					
height	Equal variances assumed	F	Sig.	t	df	Sig. (2-tailed)	Mean Difference	Std. Error Difference	95% Confidence Interval of the Difference		
									Lower	Upper	
height	Equal variances not assumed	.486	.487	-13.186	148	.000	-0.14627	.01109	-0.16819	-0.12435	

O valor de sig > 0,05 mostra que as variâncias são iguais

Independent Samples Test											
Levene's Test for Equality of Variances						t-test for Equality of Means					
height	Equal variances assumed	F	Sig.	t	df	Sig. (2-tailed)	Mean Difference	Std. Error Difference	95% Confidence Interval of the Difference		
									Lower	Upper	
height	Equal variances not assumed	.486	.487	-13.186	148	.000	-0.14627	.01109	-0.16819	-0.12435	

Valor de sig bilateral < 0,05 = eu rejeito a hipótese nula e assumo que as médias de peso entre homens e mulheres são diferentes (hipótese alternativa).

Além disso, o intervalo de confiança é negativo e não passa por 0, sendo assim, assumi-se que 0 não está incluído na possível diferença entre as médias, rejeita-se H0 e aceita-se H1.

Análise 2: Regressão Linear

Tarefa: Considerando a base Zombie Meals (zombie-meals.csv) que lista refeições típicas de zumbis e sua caracterização em termos de proteínas, calorias, sódio, etc. A última coluna (rating) representa a nota média da avaliação feita pelos zumbis para aquela refeição.

Realizei uma regressão linear e analise a influência destas características (proteínas, calorias, sódio, etc.) dos alimentos com as avaliações (ratings) feitas pelos zumbis. Analise as características que mais influenciaram na avaliação (se houver) e monte um modelo de regressão (com uma ou múltiplas variáveis independentes) que você acha que melhor prediz o rating. Considere aspectos como influência dos parâmetros, correlação, R2 e overfitting.

```
In [4]: df_zombie_meals.describe()

Out[4]:
```

	calories	protein	fat	sodium	fiber	carbo	sugars	potass	vitamins	shelf	weight
count	77.00000	77.00000	77.00000	77.00000	77.00000	77.00000	77.00000	77.00000	77.00000	77.00000	77.00000
mean	106.883117	2.545455	0.102987	159.675325	2.151948	14.597403	6.922078	96.077922	28.246753	2.207792	1.029610
std	19.484119	1.004790	0.1006473	83.832295	2.383364	4.278956	0.270819	19.193279	0.031156	0.214625	0.150477
min	50.00000	0.000000	0.000000	0.000000	0.000000	-1.000000	-1.000000	-1.000000	0.000000	1.000000	0.500000
25%	100.00000	2.000000	0.000000	130.000000	1.000000	12.000000	3.000000	40.000000	25.000000	1.000000	1.000000
50%	110.000000	3.000000	1.000000	180.000000	2.000000	14.000000	7.000000	90.000000	25.000000	2.000000	1.000000
75%	110.000000	3.000000	2.000000	210.000000	3.000000	17.000000	11.000000	120.000000	25.000000	3.000000	1.000000
max	160.00000	6.000000	5.000000	320.000000	14.000000	23.000000	15.000000	330.000000	100.000000	3.000000	1.500000

Quais valores computar as correlações e os p-valores para cada par de variável numérica com a variável "rating" com o objetivo de verificar quais variáveis podem ser úteis para a nossa regressão linear.

```
In [5]: R_pearson = df_zombie_meals.corr(method='pearson')
R_pearson

Out[5]:
```

	calories	protein	fat	sodium	fiber	carbo	sugars	potass	vitamins	shelf	weight	cut
calories	1.000000	0.019086	0.498610	0.300649	-0.239431	0.250681	0.562340	-0.066809	0.265356	0.097234	0.696091	0.08726
protein	0.019086	1.000000	-0.054074	0.005039	-0.130864	-0.320142	0.549407	0.007335	0.133865	0.216158	-0.24446	
fat	0.498610	-0.054074	1.000000	-0.005407	0.016719	-0.318043	0.270819	0.193279	-0.031156	0.214625	-0.17586	
sodium	0.300649	-0.054074	-0.005407	1.000000	-0.070675	0.356983	0.101451	-0.032603	0.361477	-0.069719	0.308576	0.11966
fiber	-0.239431	0.050339	0.016719	-0.070675	1.000000	-0.356983	-0.141265	0.903374	-0.032243	0.297539	0.247226	-0.51306
carbo	0.250681	-0.130864	0.005039	0.356983	-0.356983	1.000000	-0.331665	-0.349685	0.258148	-0.101790	0.135136	-0.36393
sugars	0.562340	-0.320142	0.270819	0.101451	-0.141265	-0.331665	1.000000	0.021696	0.125137	0.100438	0.450648	-0.03325
potass	-0.066809	0.549407	-0.031156	0.361477	-0.032243	0.297539	-0.101790	1.000000	0.020969	0.299262	0.320324	0.12840
vitamins	0.265356	0.007335	-0.031156	0.361477	-0.032243	0.297539	-0.101790	1.000000	0.020969	0.299262	0.320324	0.12840
shelf	0.097234	0.133865	0.216158	-0.069719	0.297539	-0.101790	0.100438	0.360663	1.000000	0.190762	0.190762	0.33526
weight	0.696091	0.216158	-0.214625	0.308576	0.247226	0.135136	0.450648	0.416303	0.320324	1.000000	0.190762	-0.19968
cuts	0.087260	-0.244460	-0.175862	0.119660	-0.513061	-0.363932	-0.759675	-0.480195	0.128405	-0.335269	-0.199683	1.00000
rating	-0.689376	0.470616	-0.409284	-0.401295	0.584160	0.052055	-0.235688	-0.395165	-0.240544	0.025159	-0.298124	-0.20316

```
In [6]: pvalores_pearson = calculate_p_values(df_zombie_meals, method='pearson')
pvalores_pearson

Out[6]:
```

	calories	protein	fat	sodium	fiber
calories	0.0	0.8692725966354424	0.390548311924413e-06	0.00788516337690929	0.0096041896675990717
protein	0.8692725966354424	0.0	0.068890515340739	0.6367328372849335	3.95958478623688773e-06
fat	3.95958478623688773e-06	0.068890515340739	0.0	0.962772369351892	0.8852469047311262
sodium	0.00788516337690929	0.6367328372849335	0.962772369351892	0.0	0.541337645410273
fiber	0.0096041896675990717	3.95958478623688773e-06	0.8852469047311262	0.541337645410273	0.0
carbo	0.0096041896675990717	3.95958478623688773e-06	0.8852469047311262	0.541337645410273	0.0
sugars	1.024972788732055e-07	0.0034581380775225877	0.01720667114054707	0.3799765969627058	0.22059252147567893
potass	0.5624905859478581	2.2854573505193052e-07	0.0921443087173648	0.7783387946982803	2.685747928026386e-09
vitamins	0.019680328024492837	0.9495151700111357	0.7879384064155898	0.0012377837186298345	0.7807260000096288
shelf	0.4001976143518716	0.2457272362611833	0.02049241598123358	0.5468352415315789	0.00858926850758344
weight	2.097529400153122e-12	0.05900322041875147	0.0608712609111044	0.006325115283939827	0.03018283108442263
cuts	0.4507918561687942	0.03213332170354314	0.1259714655961732	0.29991851835580985	1.8328178019207975e-06
rating	4.14027740006343e-12	1.5663097849166376e-06	0.0002190278375102041	0.0002978663125537866	2.44525044336989976e-08

```
In [7]: R_spearman = df_zombie_meals.corr(method='spearman')
R_spearman

Out[7]:
```

	calories	protein	fat	sodium	fiber	carbo	sugars	potass	vitamins	shelf	weight	cut
calories	1.000000	-0.066970	0.562796	0.290765	-0.142801	0.080772	0.596230	-0.013812	0.314042	0.152254	0.622532	0.04741
protein	-0.066970	1.000000	0.230563	-0.114397	0.680429	-0.001905	-0.285281	0.711829	-0.006491	0.182255	-0.292305	-0.36251
fat	0.562796	0.230563	1.000000	0.034785	-0.110626	-0.297343	0.322380	0.116609	0.236715	0.287783	-0.24575	
sodium	0.290765	-0.114397	0.034785	1.000000	-0.168594	0.376945	-0.114589	0.439967	-0.145355	0.263954	0.14561	
fiber	-0.142801	0.680429	-0.110626	-0.168594	1.000000	-0.145859	-0.109703	0.852266	-0.042424	0.316692	-0.348492	-0.50765
carbo	0.080772	-0.001905	-0.297343	0.376945	-0.145859	1.000000	-0.197597	0.101142	-0.110154	0.119288	0.31935	
sugars	0.596230	-0.285281	0.322380	-0.114589	-0.109703	-0.197597	1.000000	0.007495	0.304366	0.07		