

Introduzione

Il codice del progetto è in open source con [licenza MIT](#) ed è possibile trovarlo al seguente link: <https://github.com/08-UNISA/Statistica>.

È possibile visualizzare una versione interattiva online al seguente link:
<https://datafore.jupyter.com/view/notebook/Hwybys4WA4X8LTa3qNvN>.

I dati unidimensionali esaminati sono quelli del peso corporeo nel periodo di tre anni di un uomo di età 24 attuali.

I dati bidimensionali riguardano il peso corporeo della medesima persona nello stesso periodo e le calorie assunte giornalmente.

Il peso corporeo è stato registrato con due bilance differenti, una per il primo anno e un'altra nei restanti e quindi vi potrebbe essere una differenza di misura.

I dati dei primi mesi sulle calorie assunte sono molto precisi in quanto la quantità del cibo non veniva pestata, ma approssimata.

Allo stato dell'arte, importando i dati da Samsung Health, è possibile visualizzare le proprie statistiche. Questo potrebbe non funzionare in futuro se la forma dei dati venisse cambiata, ossia il parsing dei dati non funzionerebbe più. I commenti sui dati, tuttavia, risulteranno inefficienti in quanto i più significativi non sono dinamici, bensì basati su un campione di esso.

Poiché i dati sono molto riguardanti il peso corporeo, circa 500, e sulle calorie assunte, quasi 13000, non verrà preso un campione calcolato nel seguente modo:

- Per il peso corporeo, il campione sarà il peso medio delle pesate in ogni settimana.
- Quindi x_i = peso medio nella settimana i .
- Per le calorie assunte, il campione sarà la media delle calorie assunte in ogni settimana.
- Quindi y_i = media delle calorie assunte nella settimana i .

Nota: Se nella settimana i non vi sono dati al peso, essa verrà saltata e quindi anche le rispettive calorie assunte.

Si è deciso di usare la media campionaria e non la mediana perché i picchi di valore sono importanti da considerare.

Installazione dipendenze

```
In [275]: pip install matplotlib
Requirement already satisfied: matplotlib in c:\python397\lib\site-packages (3.5.2)
Requirement already satisfied: fonttools>=4.22.0 in c:\python397\lib\site-packages (from matplotlib) (4.33.3)
Requirement already satisfied: numpy>=1.17 in c:\python397\lib\site-packages (from matplotlib) (1.22.4)
Requirement already satisfied: python-dateutil>=2.7 in c:\python397\lib\site-packages (from matplotlib) (2.8.2)
Requirement already satisfied: cycler>=0.10 in c:\python397\lib\site-packages (from matplotlib) (0.11.0)
Requirement already satisfied: kiwisolver>=1.0.1 in c:\python397\lib\site-packages (from matplotlib) (1.4.2)
Requirement already satisfied: packaging>=20.0 in c:\python397\lib\site-packages (from matplotlib) (21.3)
Requirement already satisfied: pillow>=6.2.0 in c:\python397\lib\site-packages (from matplotlib) (9.1.1)
Requirement already satisfied: six>=1.5 in c:\python397\lib\site-packages (from python-dateutil>=2.7->matplotlib) (1.16.0)
Requirement already satisfied: pandas in c:\python397\lib\site-packages (1.4.2)
Requirement already satisfied: numpy>=1.18.5 in c:\python397\lib\site-packages (from pandas) (1.22.4)
Requirement already satisfied: pytz>=2020.1 in c:\python397\lib\site-packages (from pandas) (2022.1)
Requirement already satisfied: python-dateutil>=2.8.1 in c:\python397\lib\site-packages (from pandas) (2.8.2)
Requirement already satisfied: six>=1.5 in c:\python397\lib\site-packages (from python-dateutil>=2.8.1->pandas) (1.16.0)
```

```
In [276]: import pandas as pd
import csv
from math import sqrt

plt.rcParams['figure.figsize'] = [8, 6]
plt.rcParams['figure.dpi'] = 100
```

Parsing dei dati del peso corporeo

È stato necessario arrotondare i valori del peso corporeo perché altrimenti si avrebbe avuto il numero di modalità del carattere quasi uguale all'ampiezza del dato.

```
In [277]: with open('peso.csv', 'r') as csvfile:
    reader = csv.reader(csvfile, delimiter=',')
    peso = []
    reader.__next__()
    reader.__next__()
    for row in reader:
        peso.append([row[1], float(row[4])])

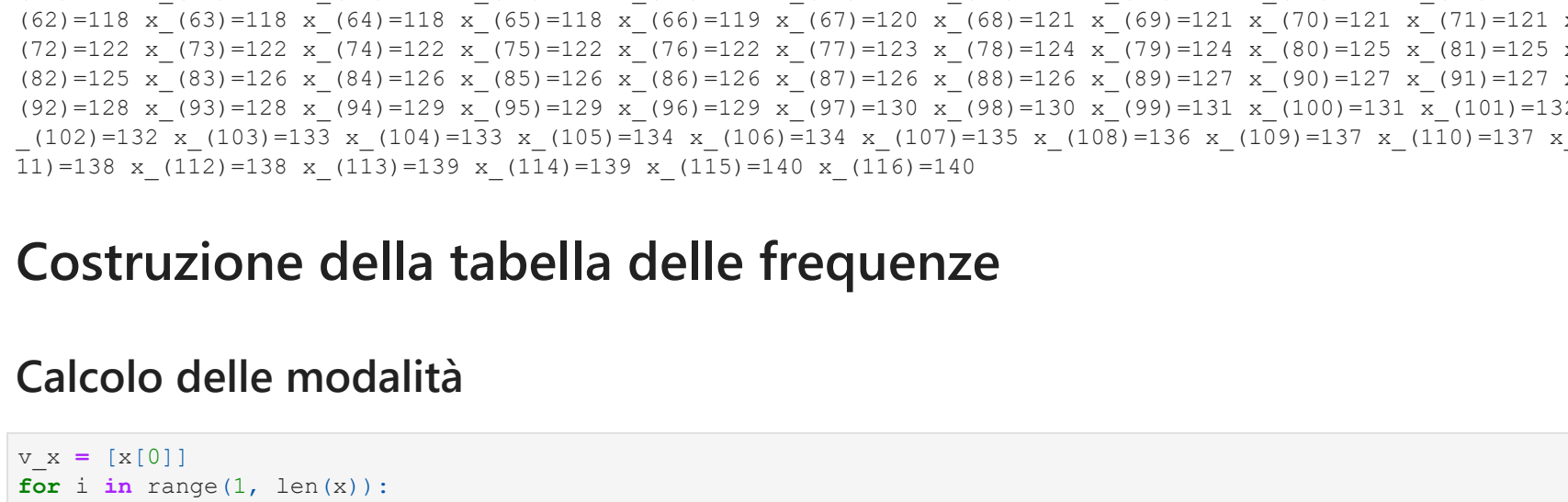
df_peso = pd.DataFrame(peso, columns=['data', 'peso'])
df_peso['data'] = pd.to_datetime(df_peso['data'])
df_peso = df_peso.groupby(pd.Grouper(key='data', freq='W')).mean().round(0)
df_peso = df_peso.dropna()
df_peso = df_peso.reset_index()
df_peso
```

Out[277]:

```
data    peso
0  2019-12-29    140
1  2020-01-12    130
2  2020-01-26    130
3  2020-02-02    130
4  2020-02-23    136.0
...      ...
111 2022-05-18    105.0
112 2022-05-18    104.0
113 2022-05-22    104.0
114 2022-05-29    103.0
115 2022-06-05    102.0

116 rows x 2 columns
```

```
In [278]: fig, ax = plt.subplots(figsize=(10, 8))
ax.plot(df_peso['data'], df_peso['peso'])
ax.set_title('Peso corporeo nel tempo')
ax.set_xlabel('Data')
ax.grid()
plt.show()
```



Sono state eliminate le settimane in cui non vi erano dati.

```
In [279]: f_x = sorted(int(p) for p in df_peso['peso'])
n = len(x)
print(f'Ampiezza dei dati: {n}')
for i, x_i in enumerate(x):
    print(f'x_{i+1}={x_i}', end=' ')

Ampiezza dei dati: 116
x_1=102 x_2=103 x_3=103 x_4=103 x_5=103 x_6=103 x_7=103 x_8=103 x_9=103 x_10=103 x_11=103 x_12=104 x_13=104 x_14=104 x_15=104 x_16=104 x_17=104 x_18=104 x_19=104 x_20=104 x_21=105 x_22=105 x_23=105 x_24=105 x_25=105 x_26=105 x_27=105 x_28=105 x_29=105 x_30=105 x_31=106 x_32=106 x_33=106 x_34=106 x_35=106 x_36=107 x_37=107 x_38=108 x_39=108 x_40=108 x_41=109 x_42=109 x_43=110 x_44=111 x_45=111 x_46=112 x_47=112 x_48=113 x_49=114 x_50=115 x_51=116 x_52=117 x_53=117 x_54=117 x_55=117 x_56=117 x_57=117 x_58=118 x_59=118 x_60=118 x_61=118 x_62=118 x_63=118 x_64=118 x_65=118 x_66=119 x_67=120 x_68=121 x_69=121 x_70=121 x_71=121 x_72=122 x_73=122 x_74=122 x_75=122 x_76=122 x_77=123 x_78=124 x_79=124 x_80=125 x_81=125 x_82=125 x_83=126 x_84=126 x_85=126 x_86=126 x_87=126 x_88=126 x_89=127 x_90=127 x_91=127 x_92=127 x_93=128 x_94=129 x_95=129 x_96=129 x_97=130 x_98=130 x_99=131 x_100=131 x_101=132 x_102=132 x_103=133 x_104=133 x_105=134 x_106=134 x_107=135 x_108=136 x_109=137 x_110=137 x_111=138 x_112=138 x_113=139 x_114=139 x_115=140 x_116=140
```

Costruzione della tabella delle frequenze

Calcolo delle modalità

```
In [280]: f_x = f_x[1:]
for i in range(1, len(x)):
    if x[i-1] != x[i]:
        f_x.append(1)
    else:
        f_x.append(f_x[i] + 1)

p_x = len(f_x)
print(f'Numero di modalità: {p_x}')

Numero di modalità: 39
p_1=102 v_2=103 v_3=104 v_4=105 v_5=106 v_6=107 v_7=108 v_8=109 v_9=110 v_10=111 v_11=112 v_12=113 v_13=114 v_14=115 v_15=116 v_16=117 v_17=118 v_18=119 v_19=120 v_20=121 v_21=122 v_22=123 v_23=124 v_24=125 v_25=126 v_26=127 v_27=128 v_28=129 v_29=130 v_30=131 v_31=132 v_32=133 v_33=134 v_34=135 v_35=136 v_36=137 v_37=138 v_38=139 v_39=140
```

Calcolo della frequenza assoluta delle modalità

Dalle frequenze assolute si può notare quando si è avuto più difficoltà nella perdita del peso o quando si era stabilizzato per un periodo di tempo.

```
In [281]: f_x = f_x[1:]
for i in range(1, n):
    if x[i-1] != x[i]:
        f_x[i-1] += 1
    else:
        f_x.append(1)
    f_x.append(f_x[i])

f_1=2 f_2=9 f_3=9 f_4=10 f_5=5 f_6=2 f_7=3 f_8=2 f_9=1 f_10=2 f_11=2 f_12=1 f_13=1 f_14=1 f_15=1 f_16=6 f_17=8 f_18=1 f_19=1 f_20=4 f_21=5 f_22=1 f_23=9 f_24=2 f_25=6 f_26=4 f_27=2 f_28=3 f_29=2 f_30=2 f_31=2 f_32=2 f_33=10 f_34=10 f_35=12 f_36=10 f_37=10 f_38=14 f_39=116
```

Calcolo della frequenza cumulativa assoluta delle modalità

```
In [282]: F_x = F_x[0:]
for i in range(1, k_x):
    # Uso della relazione di ricorrenza
    F_x.append(F_x[i-1] + f_x[i])

for i, F_i in enumerate(F_x):
    print(f'F_{i+1}={F_i}', end=' ')

F_1=2 F_2=11 F_3=20 F_4=30 F_5=35 F_6=37 F_7=40 F_8=42 F_9=43 F_10=45 F_11=47 F_12=48 F_13=49 F_14=50 F_15=51 F_16=57 F_17=65 F_18=66 F_19=70 F_20=75 F_21=80 F_22=85 F_23=94 F_24=96 F_25=101 F_26=106 F_27=111 F_28=114 F_29=119 F_30=129 F_31=140 F_32=154 F_33=166 F_34=177 F_35=187 F_36=197 F_37=211 F_38=225 F_39=341
```

Calcolo della frequenza relativa delle modalità

```
In [283]: f_x = f_x[1:]
for i in range(1, n):
    f_x[i-1] += 1
    else:
        f_x.append(1)
    f_x.append(f_x[i] / n)

p_1=0.017 p_2=0.095 p_3=0.102 p_4=0.086 p_5=0.043 p_6=0.017 p_7=0.026 p_8=0.017 p_9=0.009 p_10=0.009 p_11=0.017 p_12=0.009 p_13=0.009 p_14=0.009 p_15=0.009 p_16=0.052 p_17=0.069 p_18=0.009 p_19=0.009 p_20=0.034 p_21=0.043 p_22=0.009 p_23=0.017 p_24=0.026 p_25=0.052 p_26=0.026 p_27=0.017 p_28=0.026 p_29=0.017 p_30=0.017 p_31=0.017 p_32=0.017 p_33=0.017 p_34=0.009 p_35=0.009 p_36=0.017 p_37=0.017 p_38=0.017 p_39=0.117
```

Calcolo della frequenza cumulativa relativa delle modalità

```
In [284]: f_x = f_x[1:]
for i in range(1, n):
    f_x[i-1] += 1
    else:
        f_x.append(1)
    f_x.append(f_x[i] / n)

p_1=0.017 p_2=0.095 p_3=0.102 p_4=0.086 p_5=0.043 p_6=0.017 p_7=0.026 p_8=0.017 p_9=0.009 p_10=0.009 p_11=0.017 p_12=0.009 p_13=0.009 p_14=0.009 p_15=0.009 p_16=0.052 p_17=0.069 p_18=0.009 p_19=0.009 p_20=0.034 p_21=0.043 p_22=0.009 p_23=0.017 p_24=0.026 p_25=0.052 p_26=0.026 p_27=0.017 p_28=0.026 p_29=0.017 p_30=0.017 p_31=0.017 p_32=0.017 p_33=0.017 p_34=0.009 p_35=0.009 p_36=0.017 p_37=0.017 p_38=0.017 p_39=0.117
```

Tabella delle frequenze

```
In [285]: data_table = []
for k in range(k_x):
    data_table.append([1 + v_x[i], f_x[i], F_x[i], f_x[i], F_x[i]])

df_data_table = pd.DataFrame(data_table, columns=['v_i', 'f_i', 'F_i', 'F_i', 'F_i'])
df_data_table.set_index('v_i', inplace=True)
df_data_table
```

Out[285]:

v_i	f_i	F_i	F_i	F_i
102	2	0.017241	2	0.017241
103	9	0.077586	11	0.094828
104	9	0.077586	20	0.172414
105	10	0.086207	30	0.258621
106	5	0.043103	35	0.301724
107	2	0.017241	37	0.318966
108	3	0.025862	40	0.344828
109	2	0.017241	42	0.362069
110	1	0.008621	43	0.370690
111	2	0.017241	45	0.387931
112	2	0.017241	47	0.405172
113	1	0.008621	48	0.413793
114	1	0.008621	49	0.422414
115	1	0.008621	50	0.431034
116	1	0.008621	51	0.439655
117	6	0.051724	57	0.491379
118	8	0.068966	65	0.560345
119	1	0.008621	66	0.568966
120	1	0.008621	67	0.577586
121	4	0.034483	71	0.612069
122	5	0.043103	76	0.655172
123	1	0.008621	77	0.663793
124	2	0.017241	79	0.681034
125	3	0.025862	82	0.706897
126	6	0.051724	88	0.758621
127	3	0.025862	91	0.784483
128	2	0.017241	93	0.801724
129	3	0.025862	96	0.827586
130	2	0.017241	98	0.844828
131	2	0.017241	100	0.862069
132	2	0.017241	102	0.879310
133	2	0.017241	104	0.896552
134	2	0.017241	106	0.913793
135	1	0.008621	107	0.922414
136	1	0.008621	108	0.931034
137	2	0.017241	110	0.948276
138	2	0.017241	112	0.965517
139	2	0.017241	114	0.982759
140	2	0.017241	116	1.000000

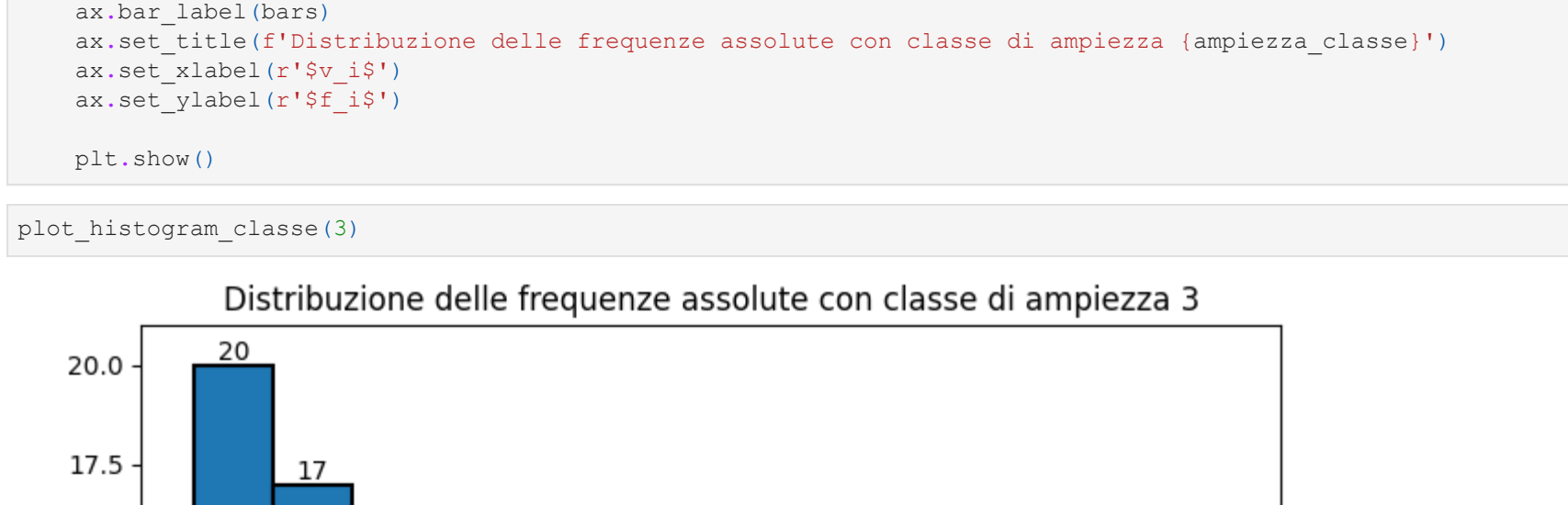
Grafici distribuzione delle frequenze assolute

Grafico a linee

```
In [286]: fig, ax = plt.subplots()
ax.plot(v_x, f_x, 'o')
ax.set_title('Distribuzione delle frequenze assolute')
ax.set_xlabel('v_i')
ax.set_ylabel('f_i')
ax.grid()

for i, v_i in enumerate(v_x):
    ax.text(v_i, f_x[i], f'f_{i+1}={f_x[i]:.01f}', ha='left', va='bottom')

plt.show()
```

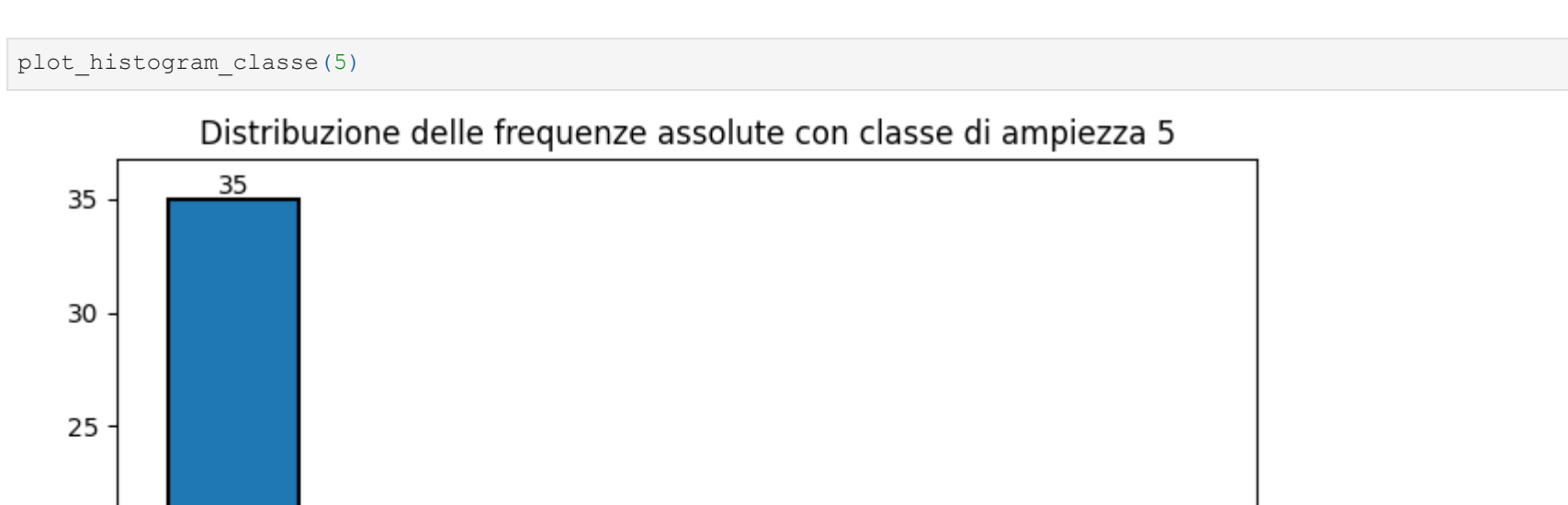


Istogramma

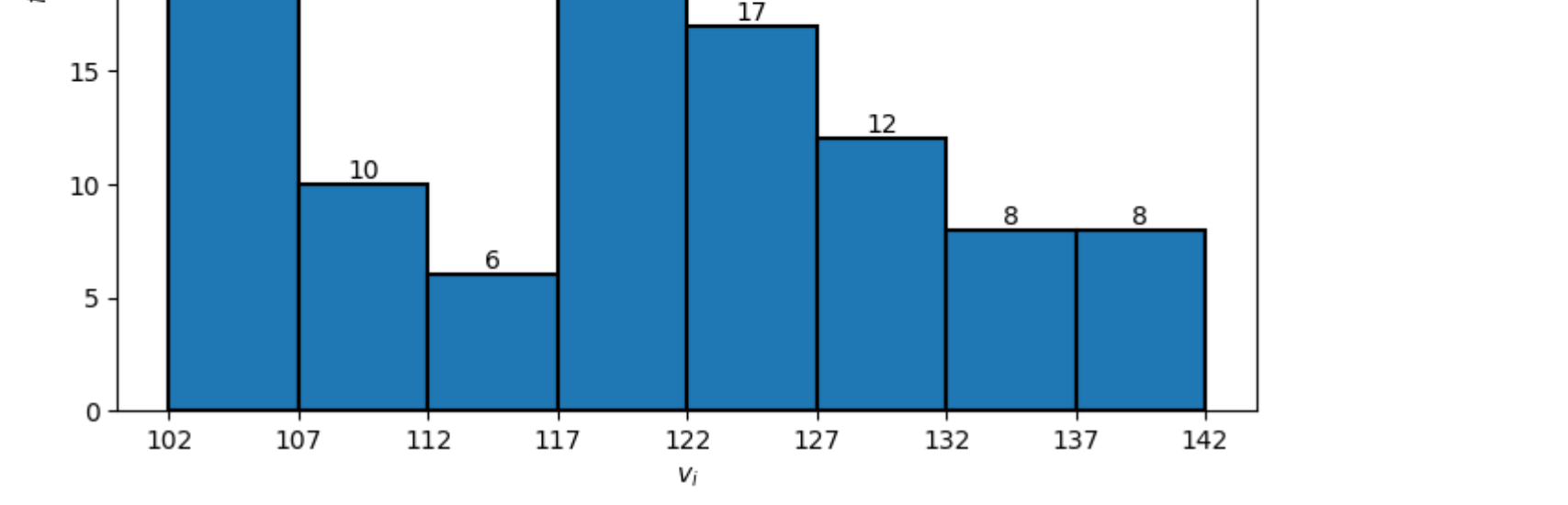
```
In [287]: def plot_histogram_classe(ampiezza_classe):
    min_value = x(0)
    max_value = x(n-1)

    fig, ax = plt.subplots()
    values, bins, bars = ax.hist(x, bins=list(range(min_value, max_value + ampiezza_classe)),
                                ax_title='Distribuzione delle frequenze assolute con classe di ampiezza (ampiezza_classe)')
    ax.set_xlabel('v_i')
    ax.set_ylabel('f_i')
    plt.show()
```

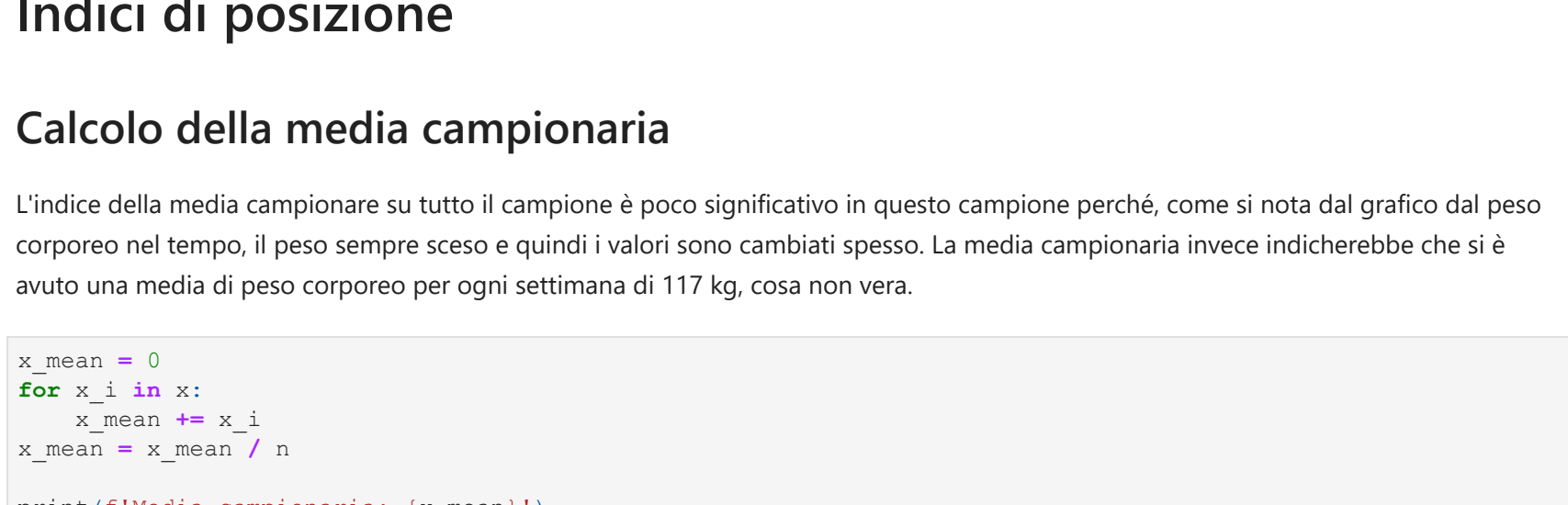
Out[287]:



In [288]: plot_histogram_classe(3)



In [289]: plot_histogram_classe(5)



Indici di posizione

Calcolo della media campionaria

L'indice della media campionaria su tutto il campione è poco significativo in questo campione perché, come si nota dal grafico dal peso corporeo nel tempo, il peso sempre scende e quindi i valori sono cambiati spesso. La media campionaria invece indicherebbe che si è avuto una media di peso corporeo per ogni settimana di 117 kg, cosa non vera.

```
In [290]: x_mean = 0
for x_i in x:
    x_mean += x_i
x_mean = x_mean / n

print(f'Media campionaria: {x_mean}')

Media campionaria: 117.42241379310344
```

Calcolo della media pesata

```
In [291]: x_mean_pesata = 0
for i in range(k_x):
    x_mean_pesata += (f_x[i] / n) * v_x[i]
x_mean_pesata = x_mean_pesata / n

print(f'Media pesata: {x_mean_pesata}')

Media pesata: 117.42241379310344
```

Calcolo della mediana campionaria

Anche la mediana campionaria per la tipologia del dato, è poco significativa.

```
In [292]: if n % 2 == 0:
    x_mediana = (x[n // 2] + x[n // 2 - 1]) / 2
else:
    x_mediana = x[n // 2]

print(f'Mediana campionaria: {x_mediana}')

Mediana campionaria: 118.0
```

Calcolo della moda campionaria

La moda campionaria è molto utile, in quanto come detto precedentemente, si riesce a determinare in quali fasi del peso corporeo si è avuto più difficoltà. Calcolando più mode campionarie, escludendo ogni volta quella precedente, si riesce a capire quali sono le fasi del peso corporeo in cui si è avuto più difficoltà.

```
In [293]: max_f_x = 0
index = 0
count = 0

for i, f_i in enumerate(f_x):
    if f_i > max_f_x:
        max_f_x = f_i
        count = 1
        index = i
    elif f_i == max_f_x:
        count += 1

x_moda = v_x[index]
if count == 1:
    print(f'Moda campionaria unimodale: v_{index + 1} = {x_moda}, f_{index + 1} = {max_f_x}')
elif count == 2:
    print(f'Moda campionaria bimodale: v_{index + 1} = {x_moda}, f_{index + 1} = {max_f_x}')
else:
    print(f'Moda campionaria multimodale: v_{index + 1} = {x_moda}, f_{index + 1} = {max_f_x}')

Moda campionaria unimodale: v_4 = 105, f_4 = 10
Moda campionaria bimodale: v_105 = 105, f_105 = 10
Moda campionaria multimodale: v_105 = 105, f_105 = 10
```

Indici di variabilità

Calcolo della varianza campionaria

La varianza dei dati rispetto alla media campionaria è un po' alta.

```
In [294]: s2_x = 0
for x_i in x:
    s2_x += (x_i - x_mean) ** 2

s2_x = s2_x / (n - 1)
s_x = s2_x ** 0.5

print(f'Varianza campionaria: {s2_x}')

Varianza campionaria: 132.08958020969505
```

Calcolo della deviazione standard campionaria

```
In [295]: s_x = sqrt(s2_x)
print(f'Deviazione standard campionaria: {s_x}')

Deviazione standard campionaria: 11.493023110126206
```

Calcolo dello scarto medio assoluto

```
In [296]: s_a_x = 0
for x_i in x:
    s_a_x += abs(x_i - x_mean)

s_a_x = s_a_x / n
print(f'Scarto medio assoluto: {s_a_x}')

Scarto medio assoluto: 9.89788941736267
```

Calcolo dell'ampiezza del campo di variazione

Questo dato indica anche il massimo di chili persi in tutto il periodo.

```
In [297]: w_x = x[n-1] - x[0]
print(f'Ampiezza del campo di variazione: {w_x}')

Ampiezza del campo di variazione: 38
```

Calcolo del coefficiente di variazione

```
In [298]: cv_x = s_x / x_mean
print(f'Coefficiente di variazione: {cv_x}')

Coefficiente di variazione: 0.0978759201047205
```

Indici di forma

Calcolo dell'indice di simmetria

Dall'istogramma si nota una coda a destra, infatti l'indice è positivo.

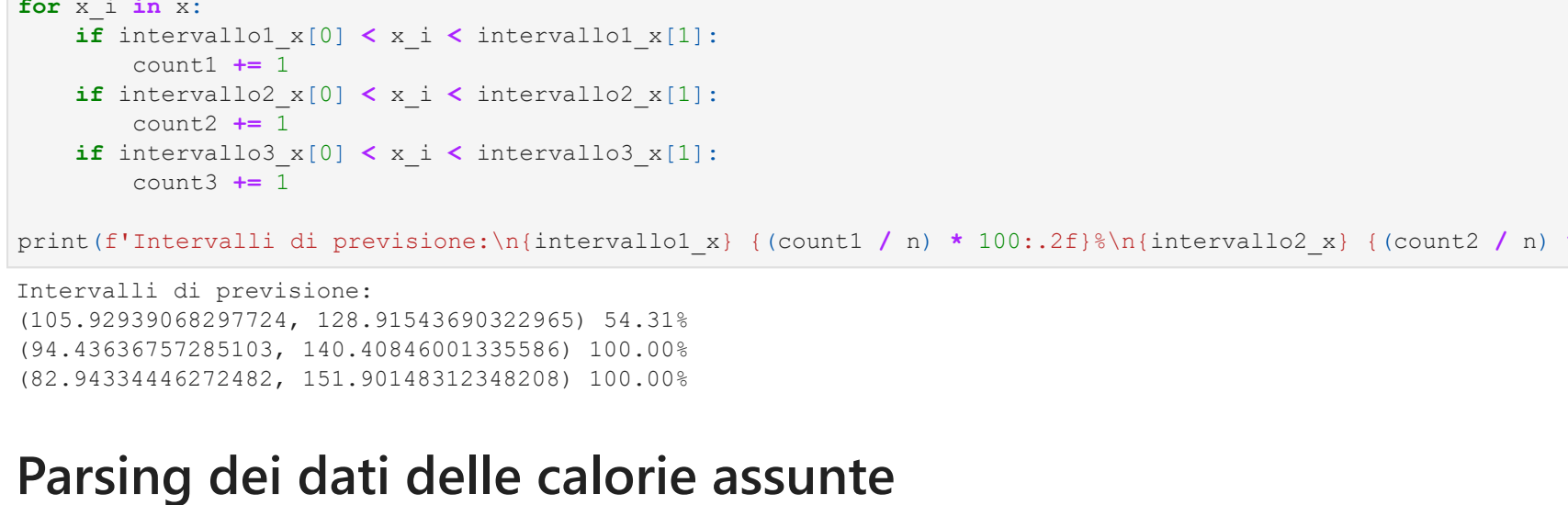
```
In [299]: g_x = 0
for x_i in x:
    g_x += (x_i - x_mean) ** 3

g_x = g_x / (n ** 3)
g_x = g_x / (s_x ** 3)

if g_x < 0:
    print(f'Indice di simmetria: {g_x}, asimmetria negativa')
else:
    print(f'Indice di simmetria: {g_x}, asimmetria positiva')

Indice di simmetria: 0.2573896882376946, asimmetria positiva
```

Out[299]:



Calcolo dell'indice di curtosi

```
In [300]: curtosi_x = 0
for x_i in x:
    curtosi_x += (x_i - x_mean) ** 4

curtosi_x = curtosi_x / (n ** 4)
curtosi_x = curtosi_x / (s_x ** 4)
curtosi_x = curtosi_x - 3

if curtosi_x > 0:
    print(f'Indice di curtosi: {curtosi_x}, è presente un eccesso di dati nelle classi centrali')
elif curtosi_x < 0:
    print(f'Indice di curtosi: {curtosi_x}, è presente una carenza di dati nelle classi centrali')
else:
    print(f'Indice di curtosi: {curtosi_x}, è presente una distribuzione di dati come quella di una distribuzione normale')

Indice di curtosi: -1.2037434028051623, è presente una carenza di dati nelle classi centrali
```

Quartili

Calcolo dei quartili

```
In [302]: def get_quartile(k):
    np = k / 100
    if np < 0.25:
        return x(int(np) * n)
    else:
        return x(int(np) * n)
```

```
In [303]: Q1_x = get_quartile(25)
Q2_x = get_quartile(50)
Q3_x = get_quartile(75)

print(f'Principali quartili: Q1 = {Q1_x}, Q2 = {Q2_x}, Q3 = {Q3_x}')

Principali quartili: Q1 = 105.0, Q2 = 118.0, Q3 = 126.0
```

Calcolo dello scarto interquartile

Indica la lunghezza del rettangolo nel box plot

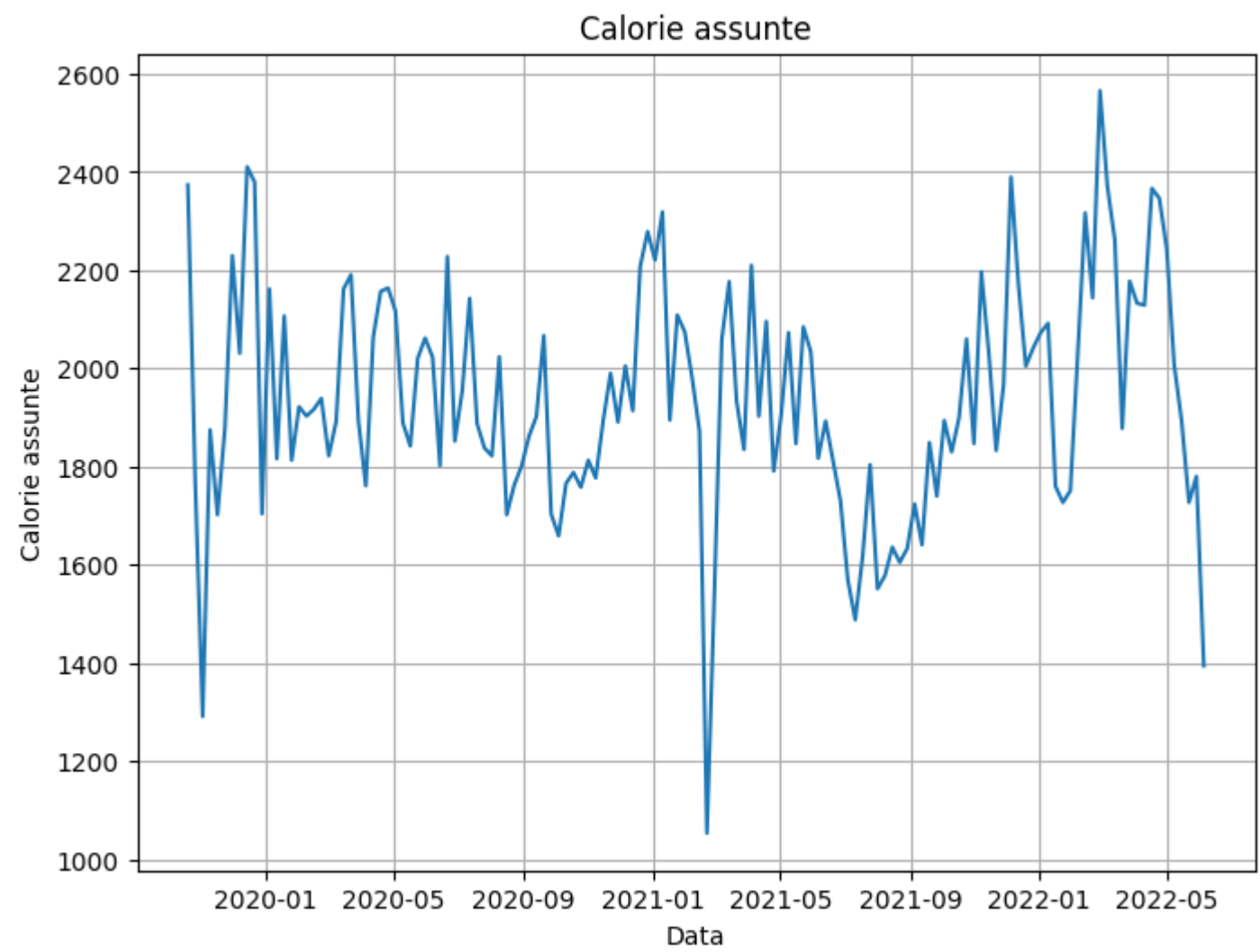
```
In [304]: si_x = Q3_x - Q1_x
print(f'Scarto interquartile: {si_x}')

Scarto interquartile: 21.0
```

Box plot

Come si evince dal grafico, non vi sono outliers. Se essi fosse presenti, sarebbero indicati nel grafico come dei punti rossi. Il segmento arancione rappresenta la media campionaria.


```
In [308... fig, ax = plt.subplots()
ax.plot(df_calorie['data'], df_calorie['calorie'])
ax.set_title('Calorie assunte')
ax.set_xlabel('Data')
ax.set_ylabel('Calorie assunte')
ax.grid()
plt.show()
```



```
In [309... df_merge = pd.merge(df_peso, df_calorie, on='data', how='inner')
df_merge = df_merge.dropna()
df_merge = df_merge.reset_index()

print(f'Numero di ampiezza dei dati dopo l\'eliminazione delle settimane in cui non vi sono state pesate corpor')
df_merge
```

Numero di ampiezza dei dati dopo l'eliminazione delle settimane in cui non vi sono state pesate corporee: 116

Out[309]:

	index	data	peso	calorie
	0	2019-12-29	140.0	1704.0
1	1	2020-01-12	140.0	1816.0
2	2	2020-01-26	139.0	1813.0
3	3	2020-02-02	139.0	1922.0
4	4	2020-02-23	138.0	1939.0
...
111	111	2022-05-08	105.0	2008.0
112	112	2022-05-15	104.0	1894.0
113	113	2022-05-22	104.0	1727.0
114	114	2022-05-29	103.0	1780.0
115	115	2022-06-05	102.0	1394.0

116 rows × 4 columns

Dati bidimensionali

Calcolo dei coefficienti di correlazione campionario

```
In [310... n = len(df_merge)
peso_mean = df_merge['peso'].mean()
calorie_mean = df_merge['calorie'].mean()
s_peso = df_merge['peso'].std()
s_calorie = df_merge['calorie'].std()

r = 0
for i in range(n):
    r += (df_merge['peso'][i] - peso_mean) * (df_merge['calorie'][i] - calorie_mean)
r = r / ((n - 1) * s_peso * s_calorie)

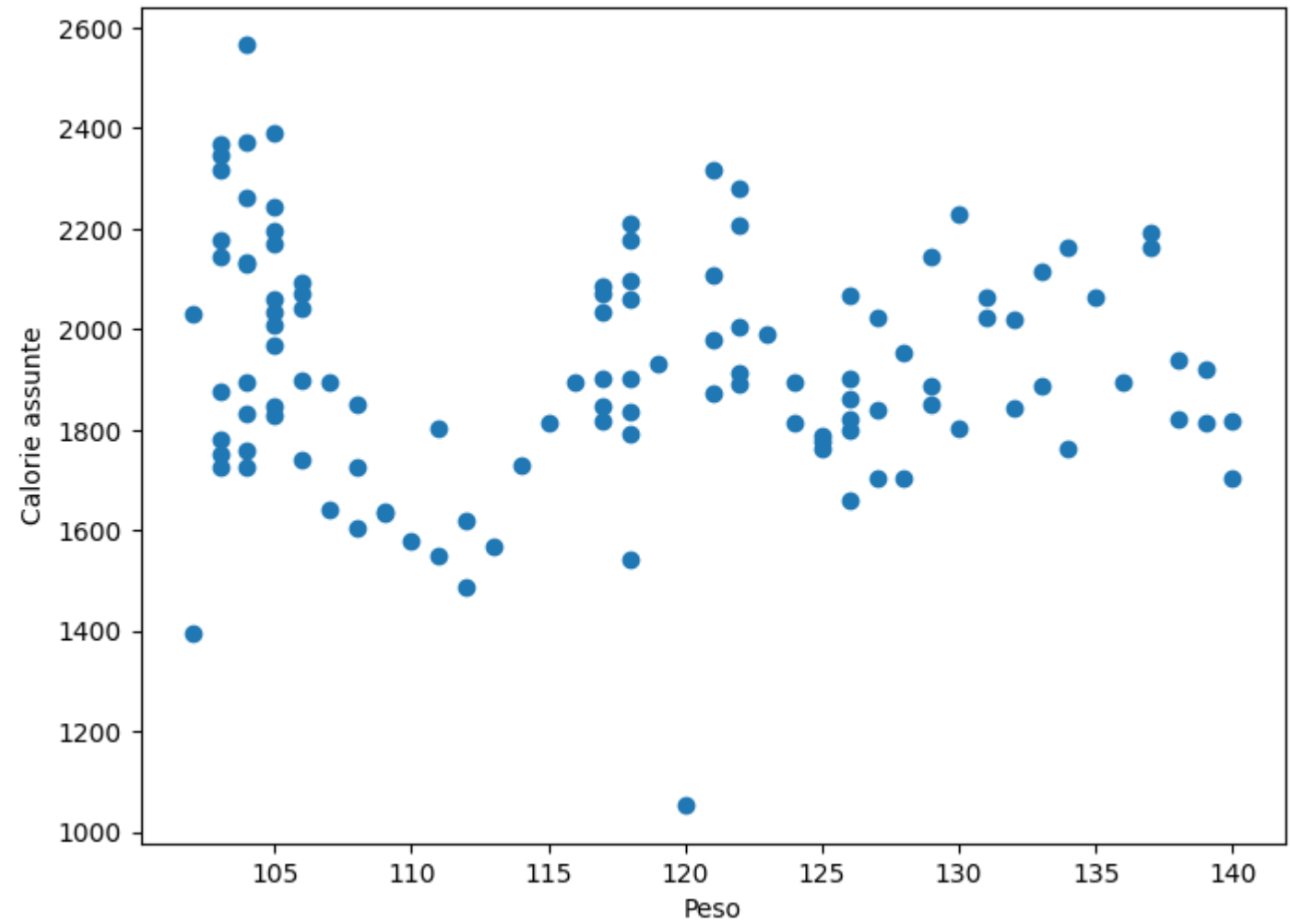
if r > 0:
    print(f'Coefficiente di correlazione campionario: {r}, è presente una correlazione positiva')
elif r < 0:
    print(f'Coefficiente di correlazione campionario: {r}, è presente una correlazione negativa')
else:
    print(f'Coefficiente di correlazione campionario: {r}, è presente una correlazione nulla')
```

Coefficiente di correlazione campionario: -0.07087210023488409, è presente una correlazione negativa

Diagramma a dispersione (scatter plot)

```
In [311... fig, ax = plt.subplots()
ax.scatter(df_merge['peso'], df_merge['calorie'])
ax.set_xlabel('Peso')
ax.set_ylabel('Calorie assunte')

plt.show()
```



Quello che emerge dalla correlazione negativa tra il peso e le calorie è che quando il peso era alto, vi era un minore consumo di calorie. Al diminuire del peso, le calorie assunte sono aumentate. In effetti nella realtà questo accadde. Tuttavia, questo non è da interpretare come un evento negativo, poiché l'aumento del consumo delle calorie fu dato dalla maggiore attività fisica.

Con i dati da Samsung Health si potrebbero fare molti altri studi, come ad esempio quanto ha influito l'attività fisica sul peso, come è cambiato il battito cardiaco e la pressione sanguigna (la pressione scese di molto), ecc. Con l'incrocio dei dati si potrebbero fare tante cose interessanti.