Introduzione Il codice del progetto è open source con licenza MIT ed è possibile trovarlo al seguente link; https://github.com/OB-UNISA/Statistica. E' possibile visualizzare una versione interattiva online al seguente link: https://datalore.jetbrains.com/view/notebook/JY0gPFblPn0G4xonfjdSWh. 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 non sono molto precisi in quanto la quantità del cibo non veniva pesata, 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 non sono dinamici, bensì basati su un campione di esso. Poiché i dati sono molti riguardanti il peso corporeo, circa 500, e sulle calorie assunte, quasi 1000, ne verrà preso un campione calcolato nel sequente 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 sul 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 [115... !pip install matplotlib !pip install pandas Requirement already satisfied: matplotlib in c:\python397\lib\site-packages (3.5.2) Requirement already satisfied: numpy>=1.17 in c:\python397\lib\site-packages (from matplotlib) (1.22.4) Requirement already satisfied: cycler>=0.10 in c:\python397\lib\site-packages (from matplotlib) (0.11.0) Requirement already satisfied: python-dateutil>=2.7 in c:\python397\lib\site-packages (from matplotlib) (2.8.2) 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: pyparsing>=2.2.1 in c:\python397\lib\site-packages (from matplotlib) (3.0.9) Requirement already satisfied: fonttools>=4.22.0 in c:\python397\lib\site-packages (from matplotlib) (4.33.3) 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->matplotli b) (1.16.0)Requirement already satisfied: pandas in c:\python397\lib\site-packages (1.4.2) 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: numpy>=1.18.5 in c:\python397\lib\site-packages (from pandas) (1.22.4) Requirement already satisfied: six>=1.5 in c:\python397\lib\site-packages (from python-dateutil>=2.8.1->pandas) (1.16.0)In [116... import matplotlib.pyplot as plt import pandas as pd import csv from math import sqrt 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 [117... with open('peso.csv', 'r') as csvfile: reader = csv.reader(csvfile, delimiter=',') peso = [] reader. next () reader. next () for row in reader: peso.append([row[6], 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.reset index() df peso Out[117]: data peso **0** 2019-12-29 140.0 **1** 2020-01-05 NaN **2** 2020-01-12 140.0 3 2020-01-19 NaN **4** 2020-01-26 139.0 **123** 2022-05-08 105.0 **124** 2022-05-15 104.0 **125** 2022-05-22 104.0 **126** 2022-05-29 103.0 **127** 2022-06-05 102.0 128 rows × 2 columns Sono state eliminate le settimane in cui non vi erano dati. In  $[118... x_not sorted = []$ date valide = [] for i in range(len(df peso)): if not pd.isnull(df peso.iloc[i]['peso']): x not sorted.append(int(df peso.iloc[i]['peso'])) date valide.append(df peso.iloc[i]['data']) x = sorted(x not sorted)n = len(x)print(f'Ampiezza del dato: {n}') for i, x i in enumerate(x): print(f'x\_({i + 1})={x\_i}', end=' ') Ampiezza del dato: 115 x (1) = 102 x (2) = 102 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 $\times$  (12)=104  $\times$  (13)=104  $\times$  (14)=104  $\times$  (15)=104  $\times$  (16)=104  $\times$  (17)=104  $\times$  (18)=104  $\times$  (19)=104  $\times$  (20)=104  $\times$  (21)=105  $\times$  $(22) = 105 \times (23) = 105 \times (24) = 105 \times (25) = 105 \times (26) = 105 \times (27) = 105 \times (28) = 105 \times (29) = 105 \times (30) = 105 \times (31) = 106 \times (31$  $(32) = 106 \times (33) = 106 \times (34) = 106 \times (35) = 106 \times (36) = 107 \times (37) = 107 \times (38) = 108 \times (39) = 108 \times (40) = 108 \times (41) = 109 \times (41$  $(42) = 109 \times (43) = 110 \times (44) = 111 \times (45) = 111 \times (46) = 112 \times (47) = 112 \times (48) = 113 \times (49) = 114 \times (50) = 115 \times (51) = 116 \times (42) = 114 \times (43) = 114 \times (44) = 114 \times (44$  $(52) = 117 \times (53) = 117 \times (54) = 117 \times (55) = 117 \times (56) = 117 \times (57) = 117 \times (58) = 117 \times (59) = 118 \times (60) = 118 \times (61) = 118 \times (61$  $(62) = 118 \times (63) = 118 \times (64) = 118 \times (65) = 118 \times (66) = 119 \times (67) = 120 \times (68) = 121 \times (69) = 121 \times (70) = 121 \times (71) = 121 \times (71$  $(72) = 122 \times (73) = 122 \times (74) = 122 \times (75) = 122 \times (76) = 122 \times (77) = 123 \times (78) = 124 \times (79) = 124 \times (80) = 125 \times (81) = 125 \times (81$  $(82) = 125 \times (83) = 126 \times (84) = 126 \times (85) = 126 \times (86) = 126 \times (87) = 126 \times (88) = 127 \times (89) = 127 \times (90) = 127 \times (91) = 128 \times (91$  $(92) = 128 \times (93) = 129 \times (94) = 129 \times (95) = 129 \times (96) = 130 \times (97) = 130 \times (98) = 131 \times (99) = 131 \times (100) = 132 \times (101) =$  $(102) = 133 \times (103) = 133 \times (104) = 134 \times (105) = 134 \times (106) = 135 \times (107) = 136 \times (108) = 137 \times (109) = 137 \times (110) = 138 \times$ 11)=138 x (112)=139 x (113)=139 x (114)=140 x (115)=140Costruzione della tabella delle frequenze Calcolo delle modalità In [119... v x = [x[0]]for i in range (1, len(x)): **if** v x[-1] != x[i]: v x.append(x[i]) k x = len(v x)for i, v i in enumerate(v x): print(f'v {i + 1}={v i}', end=' ')  $v\_1 = 102 \ v\_2 = 103 \ v\_3 = 104 \ v\_4 = 105 \ v\_5 = 106 \ v\_6 = 107 \ v\_7 = 108 \ v\_8 = 109 \ v\_9 = 110 \ v \ 10 = 111 \ v \ 11 = 112 \ v \ 12 = 113 \ v \ 13 = 114 \ v \ 12 = 113 \ v \ 13 = 114 \ v \ 12 = 113 \ v \ 13 = 114 \ v \ 12 = 113 \ v \ 13 = 114 \ v \ 12 = 113 \ v \ 13 = 114 \ v \ 13 = 114 \ v \ 14 = 112 \ v$ 4=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=1 27 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 Calcolo della frequenza assoluta delle modalità In [120...  $f_x = [1]$ for i in range(1, n): **if** x[i - 1] == x[i]: f x[-1] += 1else: f x.append(1)for i, f i in enumerate(f x): print(f'f {i + 1}={f i}', end=' ') 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=7 f 17=7 f 18=1 f 19=1 f 20=4 f 21=5 f 22=1 f 23=2 f 24=3 f 25=5 f 26=3 f 27=2 f 28=3 f 29=2 f 30=2 f 31=2 f 32=2 f 33=2 f 34=1 f 35=1 f 36=2 f 37=2 f 38=2 f 39=2 Calcolo della frequenza cumulativa assoluta delle modalità In [121... F x = [f x[0]]for i in range(1,  $k_x$ ): # Uso della relazione di ricorrenza  $F_x.append(F_x[-1] + f_x[i])$ 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=58 F\_17=65 F\_18=66 F\_19=67 F\_20=71 F\_21=76 F\_22=77 F\_23=79 F\_24=82 F\_25=87 F\_26=90 F\_27=92 F\_28=95 F\_29=97 F 30=99 F 31=101 F 32=103 F 33=105 F 34=106 F 35=107 F 36=109 F 37=111 F 38=113 F 39=115 Calcolo della frequenza relativa delle modalità In [122...  $p_x = []$ for f i in f x: p x.append(f i / n) for i, p i in enumerate(p x): print(f'p {i + 1}={p i:.3f}', end=' ') p\_1=0.017 p\_2=0.078 p\_3=0.078 p\_4=0.087 p\_5=0.043 p\_6=0.017 p\_7=0.026 p\_8=0.017 p\_9=0.009 p\_10=0.017 p\_11=0.017 p 12=0.009 p 13=0.009 p 14=0.009 p 15=0.009 p 16=0.061 p 17=0.061 p 18=0.009 p 19=0.009 p 20=0.035 p 21=0.043 p 22=0.009 p\_23=0.017 p\_24=0.026 p\_25=0.043 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=0.017 Calcolo della frequenza cumulativa relativa delle modalità In  $[123... P_x = []$ for F\_i in F\_x: P\_x.append(F\_i / n) for i, P i in enumerate(P x): print(f'P\_{i + 1}={P\_i:.3f}', end=' ') P 1=0.017 P 2=0.096 P 3=0.174 P 4=0.261 P 5=0.304 P 6=0.322 P 7=0.348 P 8=0.365 P 9=0.374 P 10=0.391 P 11=0.409 P 12=0.417 P 13=0.426 P 14=0.435 P 15=0.443 P 16=0.504 P 17=0.565 P 18=0.574 P 19=0.583 P 20=0.617 P 21=0.661 P 22=0.670 P 23=0.687 P 24=0.713 P 25=0.757 P 26=0.783 P 27=0.800 P 28=0.826 P 29=0.843 P 30=0.861 P 31=0.878 P 32=0.896 P 33=0.913 P 34=0.922 P 35=0.930 P 36=0.948 P 37=0.965 P 38=0.983 P 39=1.000 Tabella delle frequenze In [124... data table = [] for i in range(k x):  $\label{eq:data_table.append} \texttt{data\_table.append}([\texttt{i} + \texttt{1}, \texttt{v}_\texttt{x}[\texttt{i}], \texttt{f}_\texttt{x}[\texttt{i}], \texttt{p}_\texttt{x}[\texttt{i}], \texttt{F}_\texttt{x}[\texttt{i}], \texttt{P}_\texttt{x}[\texttt{i}]])$ df data table = pd.DataFrame(data table, columns=['i', 'v i', 'f i', 'p i', 'F i', 'P i']) df data table.set index('i', inplace=True) df data table Out[124]: v\_i f\_i p\_i F\_i Ρi 2 0.017391 2 0.017391 **1** 102 2 103 9 0.078261 11 0.095652 **3** 104 20 0.173913 9 0.078261 4 105 10 0.086957 30 0.260870 **5** 106 5 0.043478 35 0.304348 37 0.321739 107 2 0.017391 7 108 40 0.347826 3 0.026087 **8** 109 2 0.017391 42 0.365217 9 110 1 0.008696 43 0.373913 **10** 111 2 0.017391 45 0.391304 **11** 112 47 0.408696 2 0.017391 **12** 113 1 0.008696 48 0.417391 **13** 114 1 0.008696 49 0.426087 14 115 1 0.008696 50 0.434783 0.443478 **15** 116 1 0.008696 **16** 117 7 0.060870 58 0.504348 **17** 118 7 0.060870 65 0.565217 **18** 119 66 0.573913 1 0.008696 19 120 67 0.582609 1 0.008696 **20** 121 71 0.617391 4 0.034783 76 0.660870 5 0.043478 **21** 122 **22** 123 1 0.008696 77 0.669565 79 0.686957 23 124 2 0.017391 **24** 125 3 0.026087 82 0.713043 87 0.756522 **25** 126 5 0.043478 **26** 127 90 0.782609 3 0.026087 92 0.800000 **27** 128 2 0.017391 **28** 129 3 0.026087 95 0.826087 2 0.017391 97 0.843478 29 130 131 2 0.017391 0.860870 **31** 132 2 0.017391 101 0.878261 2 0.017391 103 0.895652 **33** 134 2 0.017391 105 0.913043 135 1 0.008696 106 0.921739 1 0.008696 107 0.930435 **35** 136 **36** 137 2 0.017391 109 0.947826 2 0.017391 111 0.965217 **37** 138 2 0.017391 113 0.982609 38 139 **39** 140 2 0.017391 115 1.000000 Grafici della distribuzione delle frequenze assolute Grafico a linee In [125... fig, ax = plt.subplots()ax.plot(v\_x, f\_x, '-o') ax.set title('Distribuzione delle frequenze assolute') ax.set xlabel(r'\$v i\$') ax.set ylabel(r'\$f i\$') ax.grid() plt.show() Distribuzione delle frequenze assolute 10 8 6 4 110 115 120 140 125 130 Istogramma def plot histogram classe(ampiezza classe): In [126... 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, ampiezza classe)), ax.set xticks(bins) ax.bar label(bars) ax.set title(f'Distribuzione delle frequenze assolute con classe di ampiezza {ampiezza classe}') ax.set xlabel(r'\$v i\$') ax.set ylabel(r'\$f i\$') plt.show() In [127... plot\_histogram\_classe(3) Distribuzione delle frequenze assolute con classe di ampiezza 3 17.5 15.0 12.5 10.0 7.5 5.0 2.5 0.0 102 105 108 111 114 117 120 123 126 129 132 135 138 141 In [128... plot\_histogram classe(5) Distribuzione delle frequenze assolute con classe di ampiezza 5 30 25 20 15 10 10 5 0 107 112 117 122 127 132 137 142 102 Indici di posizione Calcolo della media campionaria In [129... x mean = 0for x i in x: x mean += x i x mean = x mean / nprint(f'Media campionaria: {x mean}') Media campionaria: 117.3391304347826 Calcolo della media pesata In  $[130... x_mean_pesata = 0]$ for i in range(k x):  $x_{mean\_pesata} += (f_x[i] / n) * v_x[i]$ print(f'Media pesata: {x\_mean\_pesata}') Media pesata: 117.33913043478262 Calcolo della mediana campionaria In [131... x mediana = (x[n // 2] + x[n // 2 - 1]) / 2else: x mediana = x[n // 2]print(f'Mediana campionaria: {x mediana}') Mediana campionaria: 117 Calcolo della moda campionaria In  $[132... | max_f x = 0]$ index = 0count = 0for i, f\_i in enumerate(f\_x): if f i > max f x:  $\max f x = f i$ count = 1index = 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}') print(f'Moda campionaria multimodale: v {index + 1} = {x moda}, f {index + 1} = {max f x}') Moda campionaria unimodale: v 4 = 105, f 4 = 10Indici di variabilità Calcolo della varianza campionaria La varianza dei dati rispetto alla media campionaria è un po' alta. In  $[133... s2_x = 0]$ for x i in x:  $s2_x += (x_i - x_mean) ** 2$  $s2_x = s2_x / (n - 1)$ print(f'Varianza campionaria: {s2 x}') Varianza campionaria: 132.59450800915334 Calcolo della deviazione standard campionaria In  $[134... s_x = sqrt(s2_x)]$ print(f'Deviazione standard campionaria: {s x}') Deviazione standard campionaria: 11.51496886705098 Calcolo dello scarto medio assoluto In [135...] sa x = 0for x i in x: sa x += abs(x i - x mean)sa x = sa x / nprint(f'Scarto medio assoluto: {sa x}') Scarto medio assoluto: 9.907296786389397 Calcolo dell'ampiezza del campo di variazione In [136...  $w_x = x[n - 1] - x[0]$ print(f'Ampiezza del campo di variazione: {w x}') Ampiezza del campo di variazione: 38 Calcolo del coefficiente di variazione In [137... cv\_x = s\_x / x mean print(f'Coefficiente di variazione: {cv x}') Coefficiente di variazione: 0.09813409068555379 Indici di forma Calcolo dell'indice di simmetria Dall'istogramma si nota una coda a destra, infatti l'indice è positivo. In  $[138... g_x = 0]$ for x i in x:  $g_x += (x_i - x_{mean}) ** 3$  $g_x = g_x / n$  $g_x = g_x / (s_x ** 3)$ **if** g x < 0: print(f'Indice di simmetria: {g\_x}, asimmetria negativa') print(f'Indice di simmetria: {g\_x}, asimmetria positiva') Indice di simmetria: 0.2760600229825662, asimmetria positiva In [139... plot\_histogram\_classe(1) Distribuzione delle frequenze assolute con classe di ampiezza 1 10 8 6 4 2 101/18/4KNG 7889(01117131419161718190) 777374756778990373734896778940 Calcolo dell'indice di curtosi In [140... curtosi x = 0for x i in x:  $curtosi_x += (x_i - x_mean) ** 4$ curtosi x = curtosi x / n 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:</pre> print(f'Indice di curtosi: {curtosi x}, è presente una carenza di dati nelle classi centrali') print(f'Indice di curtosi: {curtosi x}, è presente una distribuzione di dati come quella di una distribuzione Indice di curtosi: -1.196898991757091, è presente una carenza di dati nelle classi centrali Quartili Calcolo dei quartili def get\_quartile(k): In [141... np = n \* (k / 100)**if** np % 1 == 0: return (x[int(np)] + x[int(np) - 1]) / 2 return x[int(np)] In  $[142... Q1_x = get quartile(25)]$  $Q2 \times = get quartile(50)$ Q3 x = get quartile(75)print(f'Principali quartili: Q1 =  $\{Q1_x\}$ , Q2 =  $\{Q2_x\}$ , Q3 =  $\{Q3_x\}$ ') Principali quartili: Q1 = 105, Q2 = 117, Q3 = 126 Calcolo dello scarto interquartile In [143...] si\_x = Q3\_x - Q1\_x print(f'Scarto interquartile: {si\_x}') Scarto interquartile: 21 **Box plot** Come si evince dal grafico, non vi sono outliers. Se essi fosse presenti, sarebbero indicati nel grafico come dei punti rossi. In [144... fig, ax = plt.subplots() ax.boxplot(x, flierprops=dict(marker='s', markerfacecolor='red')) ax.set\_title('Box plot') ax.yaxis.set ticks position('none') ax.grid(axis='y', linestyle='-', linewidth=0.5, color='lightgrey') ax.spines['top'].set\_visible(False) ax.spines['right'].set visible(False) ax.spines['left'].set\_visible(False) ax.get\_xaxis().set\_visible(False) plt.show() Box plot 140 135 130 125 120 115 110 105 Intervalli di previsione In [145...] intervallol\_x =  $(x_mean - s_x, x_mean + s_x)$  $intervallo2 x = (x_mean - 2 * s_x, x_mean + 2 * s_x)$ intervallo3  $x = (x mean - 3 * s_x, x_mean + 3 * s_x)$ count1 = 0count2 = 0count3 = 0for x i in x: if intervallo1\_x[0] < x\_i < intervallo1\_x[1]:</pre> count1 += 1 if intervallo2\_x[0] < x\_i < intervallo2\_x[1]:</pre> count2 += 1 if intervallo3 x[0] < x i < intervallo3 <math>x[1]: count3 += 1 print(f'Intervalli di previsione:\n{intervallo1\_x} {(count1 / n) \* 100:.2f}%\n{intervallo2\_x} {(count2 / n) \* 1 Intervalli di previsione: (105.82416156773162, 128.85409930183357) 53.91% (94.30919270068064, 140.36906816888455) 100.00% (82.79422383362967, 151.88403703593553) 100.00% Parsing dei dati delle calorie assunte with open('calorie.csv', 'r') as csvfile: In [146... reader = csv.reader(csvfile, delimiter=',') calorie = [] reader.\_\_next\_\_() reader. next () for row in reader: calorie.append([row[3], float(row[8])]) df calorie = pd.DataFrame(calorie, columns=['data', 'calorie']) df calorie['data'] = pd.to datetime(df calorie['data']) df calorie = df calorie.groupby(pd.Grouper(key='data', freq='D')).sum().round(0) df calorie = df calorie.reset index() df calorie = df calorie.groupby(pd.Grouper(key='data', freq='W')).mean().round(0) df calorie = df calorie.reset index() print(f'Numero di ampiezza dei dati prima della eliminazione delle settimane in cui non vi sono state pesate co Numero di ampiezza dei dati prima della eliminazione delle settimane in cui non vi sono state pesate corporee: Out[146]: data calorie **0** 2019-10-20 2374.0 **1** 2019-10-27 1768.0 **2** 2019-11-03 1291.0 **3** 2019-11-10 1875.0 **4** 2019-11-17 1702.0 2022-05-08 2008.0 2022-05-15 1894.0 **135** 2022-05-22 1727.0 **136** 2022-05-29 1780.0 **137** 2022-06-05 1394.0 138 rows × 2 columns In [147... df merge = pd.merge(df peso, df calorie, on='data', how='outer') 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 Numero di ampiezza dei dati dopo l'eliminazione delle settimane in cui non vi sono state pesate corporee: 115 Out[147]: data peso calorie 0 2019-12-29 140.0 1798.0 2 2020-01-12 140.0 1816.0 2 4 2020-01-26 139.0 1813.0 5 2020-02-02 139.0 1922.0 8 2020-02-23 138.0 1538.0 123 2022-05-08 105.0 2008.0 110 124 2022-05-15 104.0 111 1894.0 125 2022-05-22 104.0 1727.0 112 113 126 2022-05-29 103.0 114 127 2022-06-05 102.0 1394.0 115 rows × 4 columns Dati bidimensionali Calcolo dei coefficienti di correlazione campionario In [148... 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() 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') print(f'Coefficiente di correlazione campionario: {r}, è presente una correlazione nulla') Coefficiente di correlazione campionario: -0.06346445052282339, è presente una correlazione negativa Diagramma a dispersione (scatter plot) In [149... fig, ax = plt.subplots() ax.scatter(df\_merge['peso'], df\_merge['calorie'], s=1) ax.set\_xlabel('Peso') ax.set\_ylabel('Calorie assunte') plt.show() 2600

2400

2000 1800

cose interessanti.

115

120

Peso

125

130

135

come un evento negativo, poiché l'aumento del consumo delle calorie fu dato dalla maggiore attività fisica.

140

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

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

Calorie assunte