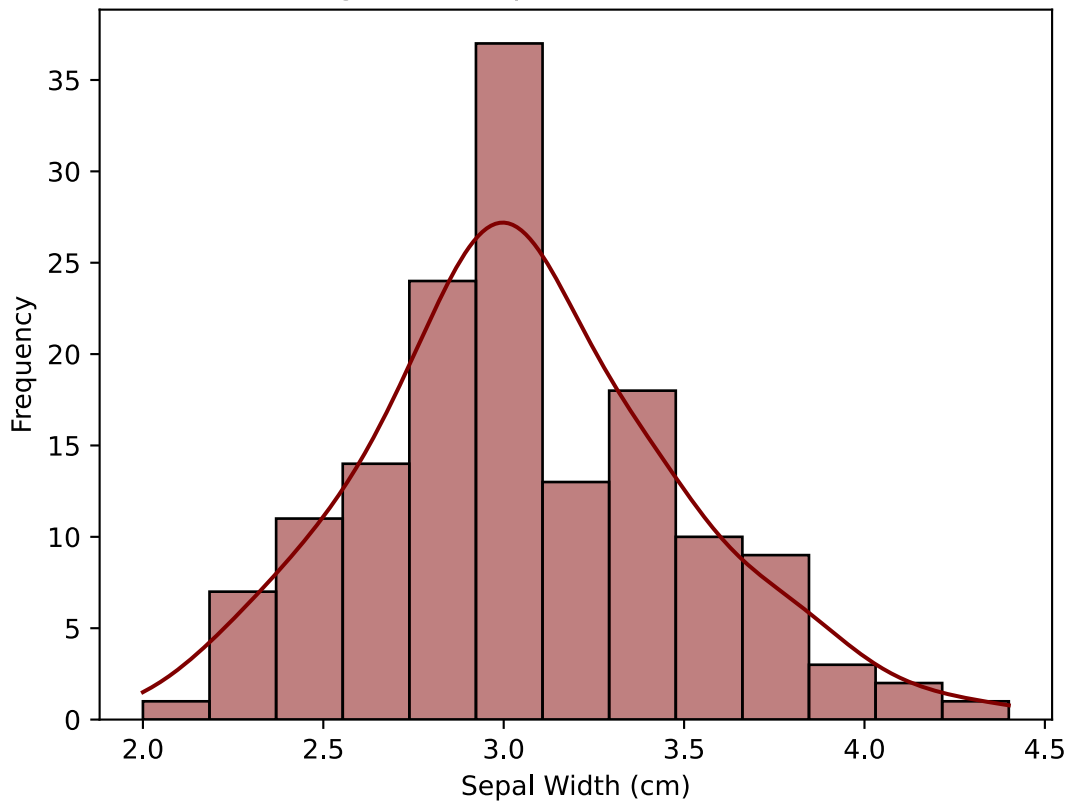


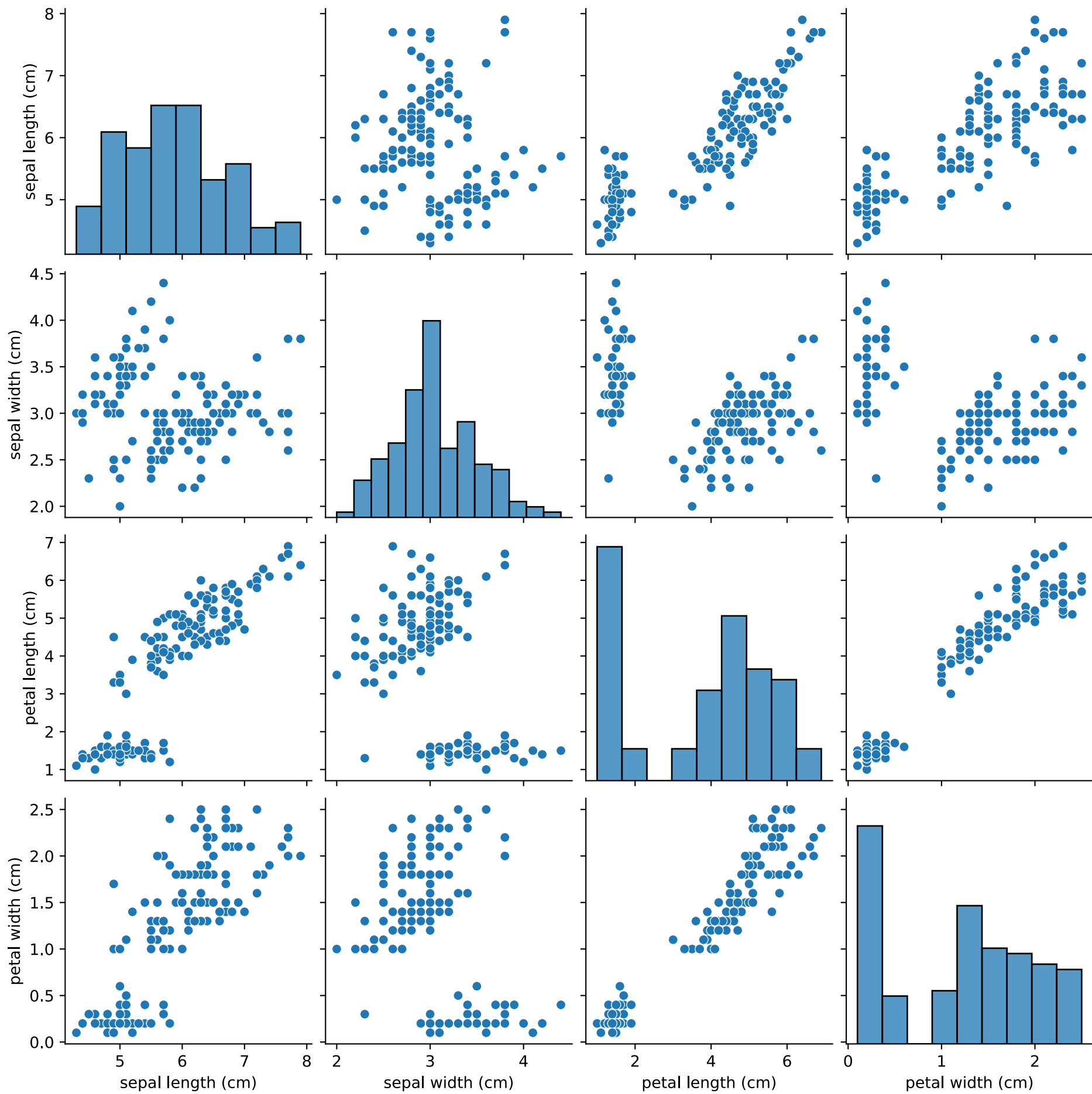
**Week 3 Assignment.py**

```
1 # Week 3 Assignment
2 # Due Date: 5/25/2025
3 # Author: Drake Shaub
4
5 # import required packages
6 from sklearn.datasets import load_iris
7 import pandas as pd
8 import matplotlib.pyplot as plt
9 import numpy as np
10 import seaborn as sns
11
12 # import iris dataset as pandas DataFrame
13 iris = load_iris(as_frame = True)
14 df_iris = iris.frame
15
16 # create PlantGrowth dataset as pandas DataFrame
17 data = { "weight": [4.17, 5.58, 5.18, 6.11, 4.50, 4.61, 5.17, 4.53, 5.33, 5.14,
18 4.81, 4.17,
19 4.41, 3.59, 5.87, 3.83, 6.03, 4.89, 4.32, 4.69, 6.31, 5.12,
20 5.54, 5.50, 5.37, 5.29,
21 4.92, 6.15, 5.80, 5.26], "group": ["ctrl"] * 10 + ["trt1"] *
22 10 + ["trt2"] * 10}
23 PlantGrowth = pd.DataFrame(data)
24
25 # Question 1a
26 # Make a histogram of the variable Sepal.Width
27 sns.histplot(df_iris['sepal width (cm)'], kde=True, color='maroon')
28 plt.title('Histogram of Sepal Width in iris Dataset')
29 plt.xlabel('Sepal Width (cm)')
30 plt.ylabel('Frequency')
31
32 # Save the figure as .pdf
33 plt.savefig('/Users/drakeshaub/Documents/Future/Education/Purdue University 2025-
34 2027/Summer 2025/GRAD 505 - Foundations in Data Science/Week 3/Sepal Width
35 Histogram.pdf')
36
37 # Show the histogram
38 plt.show()
39
40 # Question 1b
41 # Based on the histogram from 1a, which would you expect to be higher, the mean or
```

Histogram of Sepal Width in iris Dataset



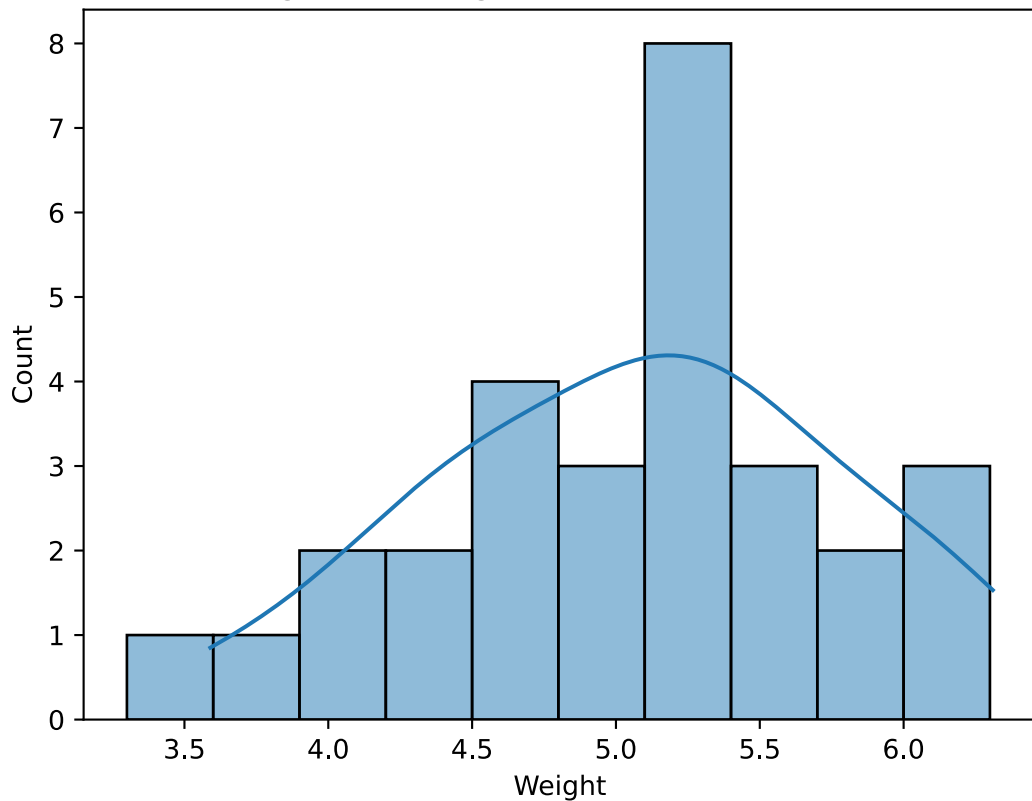
```
the median? why?
39
40 print("See commented code for Question 1b")
41 # I would expect the mean to be slightly higher because it appears that the
dataset is slightly right-skewed.
42 # Because it's right-skewed, I would expect the mean to be higher than the median
because the overall sum would be greater
43 # due to the higher values on the right hand side of the dataset (towards the
maximum).
44
45
46 # Question 1c
47 # Confirm #1b by finding the median and mean of Sepal.Width values.
48 mean = np.mean(df_iris['sepal width (cm)'])
49 median = np.median(df_iris['sepal width (cm)'])
50
51 # making the mean and median print out look cleaner
52 print("Question 1c:")
53 print(f"Mean: {mean}")
54 print(f"Median: {median}")
55
56
57 # Question 1d
58 # Only 27% of flowers have Sepal.Width higher than ____ cm. Fill in the blank.
59
60 # If 27% of flowers have sepal width higher than this number, this number would
represent the (100-27) percentile, i.e. 73rd percentile.
61 # Use the np.percentile() function to calculate the 73rd percentile value, which
represents the number at which 73% of the values
62 # fall below, but that 27% of the values fall above.
63 percentile = np.percentile(df_iris['sepal width (cm)'], 73)
64
65 # make percentile print out look cleaner
66 print("Question 1d:")
67 print(f"27% of the flowers have a sepal width greater than {percentile} cm")
68
69
70 #Question 1e
71 # Make scatterplots of each pair of the numerical variables in iris (there should
be 6 plots)
72
73 # can use scatterplot matrix (pairplot) to perform this in one go.
74 iris_num_vars = ['sepal length (cm)', 'sepal width (cm)', 'petal length (cm)',
'petal width (cm)']
75
```



```
76 sns.pairplot(df_iris[iris_num_vars], markers='o')
77
78 # save the figure
79 plt.savefig('/Users/drakeshaub/Documents/Future/Education/Purdue University 2025-
2027/Summer 2025/GRAD 505 - Foundations in Data Science/Week 3/Pairs of Numerical
Variables Scatterplots.pdf')
80
81 # show the scatterplots
82 plt.show()
83
84 # Additional answers for 1e
85 # Can also make all 6 individual scatterplots (so you don't get redundant
scatterplots). That code is shown below.
86
87 # Sepal Length vs Sepal Width
88 # sns.scatterplot(data=df_iris, x='sepal length (cm)', y='sepal width (cm)',
color='blue')
89 # plt.xlabel('Sepal Length (cm)')
90 # plt.ylabel('Sepal Width (cm)')
91 # plt.title('Sepal Length vs Sepal Width')
92 # plt.show()
93
94 # Sepal Length vs Petal Length
95 # sns.scatterplot(data=df_iris, x='sepal length (cm)', y='petal length (cm)',
color='red')
96 # plt.xlabel('Sepal Length (cm)')
97 # plt.ylabel('Petal Length (cm)')
98 # plt.title('Sepal Length vs Petal Length')
99 # plt.show()
100
101 # Sepal Length vs Petal Width
102 # sns.scatterplot(data=df_iris, x='sepal length (cm)', y='petal width (cm)',
color='cyan')
103 # plt.xlabel('Sepal Length (cm)')
104 # plt.ylabel('Petal Width (cm)')
105 # plt.title('Sepal Length vs Petal Width')
106 # plt.show()
107
108 # Sepal Width vs Petal Length
109 # sns.scatterplot(data=df_iris, x='sepal width (cm)', y='petal length (cm)',
color='pink')
110 # plt.xlabel('Sepal Width (cm)')
111 # plt.ylabel('Petal Length (cm)')
112 # plt.title('Sepal Width vs Petal Length')
113 # plt.show()
```

```
114
115 # Sepal Width vs Petal Width
116 # sns.scatterplot(data=df_iris, x='sepal width (cm)', y='petal width (cm)',
117 #                  color='maroon')
118 # plt.xlabel('Sepal Width (cm)')
119 # plt.ylabel('Petal Width (cm)')
120 # plt.title('Sepal Width vs Petal Width')
121 # plt.show()
122
123 # Petal Length vs PEtal Width
124 # sns.scatterplot(data=df_iris, x='petal length (cm)', y='petal width (cm)',
125 #                  color='green')
126 # plt.xlabel('Petal Length (cm)')
127 # plt.ylabel('Petal Width (cm)')
128 # plt.title('Petal Length vs Petal Width')
129 # plt.show()
130
131 # Question 1f
132 # Based on #1e, which two variables appear to have the strongest relationship? And
133 # which two appear to have the weakest relationship?
134 # Petal width and petal length appear to have the strongest relationship. Sepal
135 # legnth and sepal width appear to have the weakest
136 # relationship.
137
138 #Question 2a
139 # Make a histogram of the variable weight with breakpoints (bin edges) at every
140 # 0.3 units, starting at 3.3
141
142 min_edge = 3.3 # defined per question statement
143 max_edge = np.max(PlantGrowth['weight']) # maximum value from PlantGrowth weight
144 # column
145 breakpoints = 0.3 # defined per question statement
146
147 # develop an array with values starting at 3.3, going to the max value, with
148 # interval of 0.3
149 bin_array = np.arange(min_edge, max_edge, breakpoints)
150
151 # pass this bin_array into seaborn histplot function
152 sns.histplot(PlantGrowth['weight'], bins=bin_array, kde=True)
153
154 # give the graph a title and label axes
155 plt.title('Histogram of Weights in PlantGrowth Dataset')
```

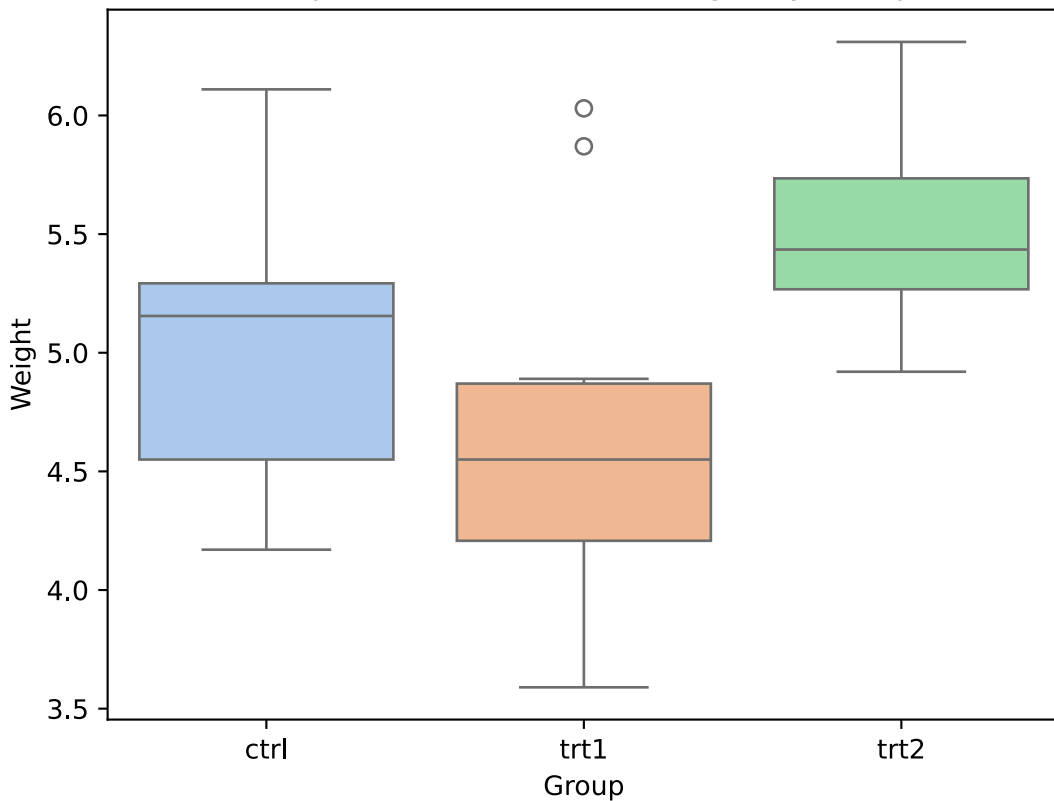
Histogram of Weights in PlantGrowth Dataset



```
152 plt.xlabel('Weight')
153 plt.ylabel('Count')
154
155 # save figure
156 plt.savefig('/Users/drakeshaub/Documents/Future/Education/Purdue University 2025-
2027/Summer 2025/GRAD 505 - Foundations in Data Science/Week 3/PlantGrowth Weight
Histogram with Breakpoints at Every 0.3 Units.pdf')
157
158 # show the plot
159 plt.show()
160
161
162 # Question 2b
163 # Make boxplots of weight separated by group in a single graph
164
165 sns.boxplot(x='group', y='weight', data=PlantGrowth, hue='group', legend=False,
palette='pastel')
166
167 # give the graph a title and label axes
168 plt.title('Boxplots of PlantGrowth Weight by Group')
169 plt.xlabel('Group')
170 plt.ylabel('Weight')
171
172 # save figure
173 plt.savefig('/Users/drakeshaub/Documents/Future/Education/Purdue University 2025-
2027/Summer 2025/GRAD 505 - Foundations in Data Science/Week 3/Boxplots of Weight
Separated by Group.pdf')
174
175 # show the plot
176 plt.show()
177
178 # Question 2c
179 # Based on the boxplots in #2b, approximately what percentage of the "trt1"
weights are below the minimum "trt2" weight?
180
181 print("See commented code for Question 2c")
182 # Minimum "trt2" weight = ~ 4.9
183 # ~ 75% of the trt1 weights are below the minimum weight for trt2. The 75th
percentile (Q3) is less than
184 # the minimum value for trt2. Therefore, at least 75% of the values fall below the
minimum value for trt2.
185
186
187 # Question 2d
188 # Find the exact percentage of the "trt1" weights that are below the minimum
"trt2" weight.
```



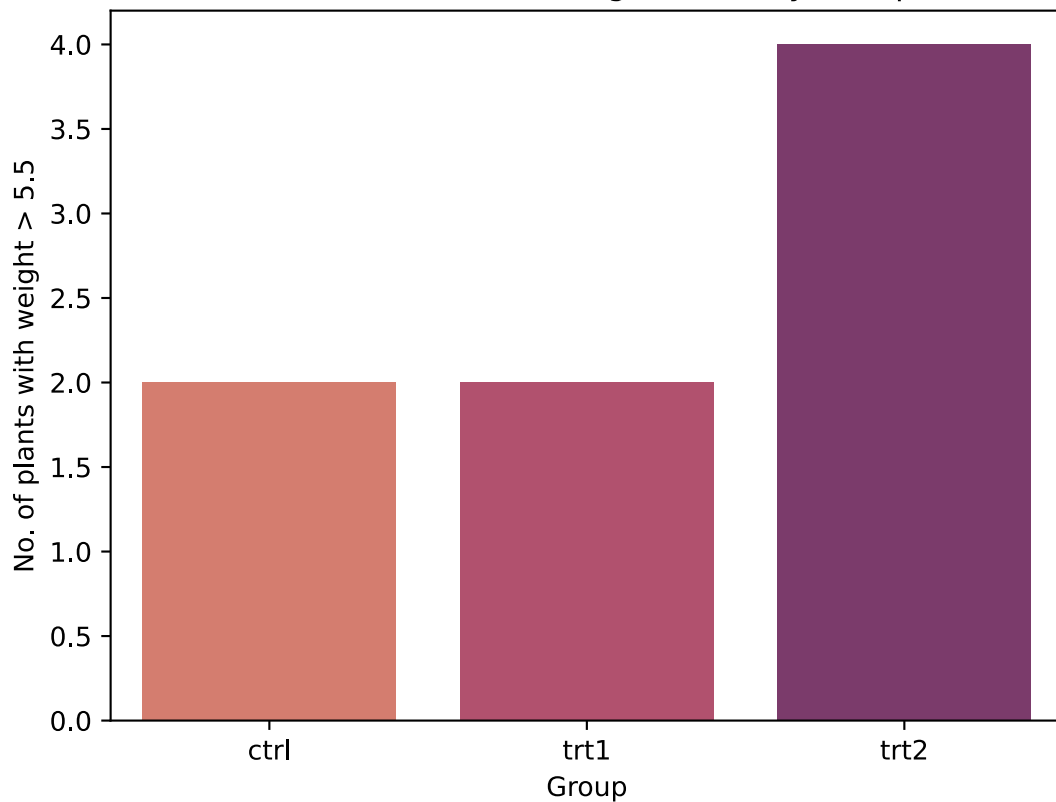
Boxplots of PlantGrowth Weight by Group



```
189
190 # filter datasets to create new dataframes grouped by column
191 ctrl_df = PlantGrowth[PlantGrowth['group'] == 'ctrl']
192 trt2_df = PlantGrowth[PlantGrowth['group'] == 'trt2']
193
194 # find minimum value from trt2_df dataframe and assign to variable
195 min_trt2_df = np.min(trt2_df['weight'])
196
197 # filter PlantGrowth to only include trt1 group and those values less than the
198 # minium value of trt2
199 trt1_df = PlantGrowth[(PlantGrowth['group'] == 'trt1') & (PlantGrowth['weight'] <
200 min_trt2_df)]
201
202 # determine percentile by dividing the filtered count by the total count of trt1
203 # group in PlantGrowth
204 percentage = (trt1_df.count()[0] / PlantGrowth[PlantGrowth['group'] ==
205 'trt1'].count()[0]) * 100
206
207 # print out answer
208 print("Question 2d:")
209 print(f"{percentage}% of trt1 weights are below the minimum value of trt2
210 weights.")
211
212 # Question 2e
213 # Only including plants with a weight above 5.5, make a barplot of the variable
214 # group.
215 # Make the barplot colorful using some color palette.
216
217 # filter PlantGrowth dataset to only include those with weight > 5.5
218 barplot_df = PlantGrowth[PlantGrowth['weight'] > 5.5]
219
220 # use .value_counts() function of dataframes to pull value counts for each label
221 frequency_table = barplot_df['group'].value_counts()
222
223 # create labels (groups) and their values
224 labels = sorted(frequency_table.index)
225 values = sorted(frequency_table.values)
226
227 # create bar plot
228 sns.barplot(x=labels, y=values, hue=labels, legend=False, palette='flare')
229
230 # create title and label axes
231 plt.title('No. of Plants with Weight > 5.5 by Group')
232 plt.xlabel('Group')
```

```
228 plt.ylabel('No. of plants with weight > 5.5')
229
230 # save figure
231 plt.savefig('/Users/drakeshaub/Documents/Future/Education/Purdue University 2025-
2027/Summer 2025/GRAD 505 - Foundations in Data Science/Week 3/Barplot of Plants
with Weight Above 5.5.pdf')
232
233 # show the plot
234 plt.show()
```

No. of Plants with Weight > 5.5 by Group



**extension-output-formulahendry.code-runner-#1-Code**

```
1 [Running] python3 -u "/Users/drakeshaub/Documents/Future/Education/Purdue
  University 2025-2027/GitHub/Assignment---2/Week 3 Assignment.py"
2 See commented code for Question 1b
3 Question 1c:
4 Mean: 3.057333333333337
5 Median: 3.0
6 Question 1d:
7 27% of the flowers have a sepal width greater than 3.3 cm
8 See commented code for Question 2c
9 /Users/drakeshaub/Documents/Future/Education/Purdue University
  2025-2027/GitHub/Assignment---2/Week 3 Assignment.py:201: FutureWarning:
  Series.__getitem__ treating keys as positions is deprecated. In a future version,
  integer keys will always be treated as labels (consistent with DataFrame behavior).
  To access a value by position, use `ser.iloc[pos]`
10 percentage = (trt1_df.count()[0] / PlantGrowth[PlantGrowth['group'] ==
  'trt1'].count()[0]) * 100
11 /Users/drakeshaub/Documents/Future/Education/Purdue University
  2025-2027/GitHub/Assignment---2/Week 3 Assignment.py:201: FutureWarning:
  Series.__getitem__ treating keys as positions is deprecated. In a future version,
  integer keys will always be treated as labels (consistent with DataFrame behavior).
  To access a value by position, use `ser.iloc[pos]`
12 percentage = (trt1_df.count()[0] / PlantGrowth[PlantGrowth['group'] ==
  'trt1'].count()[0]) * 100
13 Question 2d:
14 80.0% of trt1 weights are below the minimum value of trt2 weights.
15
16 [Done] exited with code=0 in 6.886 seconds
17
18
```