

Observations and Insights

Type *Markdown* and LaTeX: α^2

```
In [609]: 1 # Dependencies and Setup
2 import matplotlib.pyplot as plt
3 import pandas as pd
4 import scipy.stats as st
5 import numpy as np
6
7 # Study data files
8 mouse_metadata_path = "Data/Mouse_metadata.csv"
9 study_results_path = "Data/Study_results.csv"
10
11 # Read the mouse data and the study results
12 mouse_metadata = pd.read_csv(mouse_metadata_path)
13 study_results = pd.read_csv(study_results_path)
14
15 # Combine the data into a single dataset
16 df1 = mouse_metadata.merge(study_results, on="Mouse ID")
17
18 # Display the data table for preview
19 df1
```

Out[609]:

	Mouse ID	Drug Regimen	Sex	Age_months	Weight (g)	Timepoint	Tumor Volume (mm3)	Metastatic Sites
0	k403	Ramicane	Male	21	16	0	45.000000	0
1	k403	Ramicane	Male	21	16	5	38.825898	0
2	k403	Ramicane	Male	21	16	10	35.014271	1
3	k403	Ramicane	Male	21	16	15	34.223992	1
4	k403	Ramicane	Male	21	16	20	32.997729	1
...
1888	z969	Naftisol	Male	9	30	25	63.145652	2
1889	z969	Naftisol	Male	9	30	30	65.841013	3
1890	z969	Naftisol	Male	9	30	35	69.176246	4
1891	z969	Naftisol	Male	9	30	40	70.314904	4
1892	z969	Naftisol	Male	9	30	45	73.867845	4

1893 rows × 8 columns

```
In [384]: 1 # Checking the number of mice.
2 number_of_mice = df1.nunique()["Mouse ID"]
3 number_of_mice
```

Out[384]: 249

```
In [401]: 1 # Getting the duplicate mice by ID number that shows up for Mouse ID and
2 df1.loc[df1.duplicated(subset=['Mouse ID', 'Timepoint']), 'Mouse ID'].u
3 df1.loc[df1['Mouse ID'] == 'g989']
```

Out[401]:

	Mouse ID	Drug Regimen	Sex	Age_months	Weight (g)	Timepoint	Tumor Volume (mm3)	Metastatic Sites
908	g989	Propriva	Female	21	26	0	45.000000	0
910	g989	Propriva	Female	21	26	5	48.786801	0
912	g989	Propriva	Female	21	26	10	51.745156	0
914	g989	Propriva	Female	21	26	15	51.325852	1
916	g989	Propriva	Female	21	26	20	55.326122	1
918	g989	Propriva	Female	21	26	25	56.045564	1
919	g989	Propriva	Female	21	26	30	59.082294	1
920	g989	Propriva	Female	21	26	35	62.570880	2

```
In [402]: 1 # Create a clean DataFrame by dropping the duplicate mouse by its ID.
2 df1.drop_duplicates(subset = ['Mouse ID', 'Timepoint'], inplace = True)
3 df = df1
4 df
```

	ID	Regimen	Sex	Age_months	Weight (g)	Timepoint	(mm3)	Sites
0	k403	Ramicane	Male	21	16	0	45.000000	0
1	k403	Ramicane	Male	21	16	5	38.825898	0
2	k403	Ramicane	Male	21	16	10	35.014271	1
3	k403	Ramicane	Male	21	16	15	34.223992	1
4	k403	Ramicane	Male	21	16	20	32.997729	1
...
1888	z969	Naftisol	Male	9	30	25	63.145652	2
1889	z969	Naftisol	Male	9	30	30	65.841013	3
1890	z969	Naftisol	Male	9	30	35	69.176246	4
1891	z969	Naftisol	Male	9	30	40	70.314904	4
1892	z969	Naftisol	Male	9	30	45	73.867845	4

```
In [403]: 1 # Checking the number of mice in the clean DataFrame.
2 number_of_mice = df.nunique()["Mouse ID"]
3 number_of_mice
```

Out[403]: 249

```
In [404]: 1 #WATCH END OF 1/26 LECTURE FOR HW INSIGHTS
```

Summary Statistics

In [405]: Generate a summary statistics table of mean, median, variance, standard deviation, and SEM of the tumor volume. Use groupby and summary statistical methods to calculate the following properties: mean, median, variance, standard deviation, and SEM of the tumor volume. Assemble the resulting series into a single summary dataframe.

```
5
tumor_mean = df.groupby('Drug Regimen')['Tumor Volume (mm3)'].mean()
tumor_med = df.groupby('Drug Regimen')['Tumor Volume (mm3)'].median()
tumor_var = df.groupby('Drug Regimen')['Tumor Volume (mm3)'].var(ddof=0)
tumor_std = df.groupby('Drug Regimen')['Tumor Volume (mm3)'].std(ddof=0)
tumor_SEM = df.groupby('Drug Regimen')['Tumor Volume (mm3)'].sem(ddof=0)
data = zip(tumor_mean, tumor_med, tumor_var, tumor_std, tumor_SEM)
summary1 = pd.DataFrame(data, columns = ['Drug Regimen', 'Mean of Tumor Volume', 'Median of Tumor Volume', 'Variance of Tumor Volume', 'Standard Deviation of Tumor Volume', 'SEM of Tumor Volume'])
```

Out[405]:

	Drug Regimen	Mean of Tumor Volume Per Drug	Median of Tumor Volume Per Drug	Variance of Tumor Volume Per Drug	Standard Deviation of Tumor Volume Per Drug	SEM of Tumor Volume Per Drug
0	Ramicane	40.675741	41.557809	24.839296	4.983904	0.328629
1	Capomulin	52.591172	51.776157	39.069446	6.250556	0.468499
2	Infubinol	52.884795	51.820584	42.886388	6.548770	0.490851
3	Placebo	55.235638	53.698743	68.188930	8.257659	0.602252
4	Ceftamin	54.331565	52.509285	65.817708	8.112811	0.594860
5	Stelasyn	54.033581	52.288934	60.830138	7.799368	0.579722
6	Zoniferol	52.393463	50.909965	42.862273	6.546928	0.524174
7	Ketapril	40.216745	40.673236	23.383692	4.835669	0.320250
8	Propriva	54.233149	52.431737	59.122106	7.689090	0.571526
9	Naftisol	53.236507	51.818479	48.266689	6.947423	0.514977

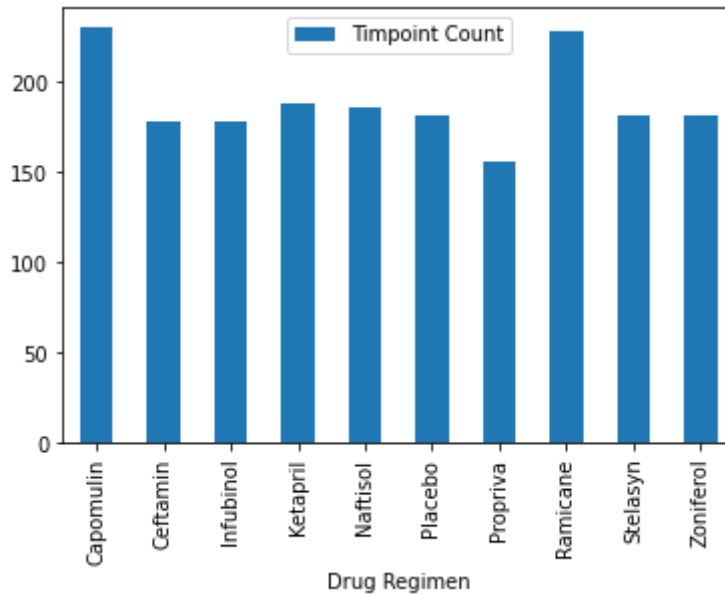
```
In [406]: 1 # Using the aggregation method, produce the same summary statistics in
          2 # .grouby then .agg
          3 summ2 = df.groupby('Drug Regimen')['Tumor Volume (mm3)'].agg(['mean', 'm
          4 summ2
```

Out[406]:

	mean	median	var	std	sem
Drug Regimen					
Capomulin	40.675741	41.557809	24.947764	4.994774	0.329346
Ceftamin	52.591172	51.776157	39.290177	6.268188	0.469821
Infubinol	52.884795	51.820584	43.128684	6.567243	0.492236
Ketapril	55.235638	53.698743	68.553577	8.279709	0.603860
Naftisol	54.331565	52.509285	66.173479	8.134708	0.596466
Placebo	54.033581	52.288934	61.168083	7.821003	0.581331
Propriva	52.393463	50.909965	43.138803	6.568014	0.525862
Ramicane	40.216745	40.673236	23.486704	4.846308	0.320955
Stelasyn	54.233149	52.431737	59.450562	7.710419	0.573111
Zoniferol	53.236507	51.818479	48.533355	6.966589	0.516398

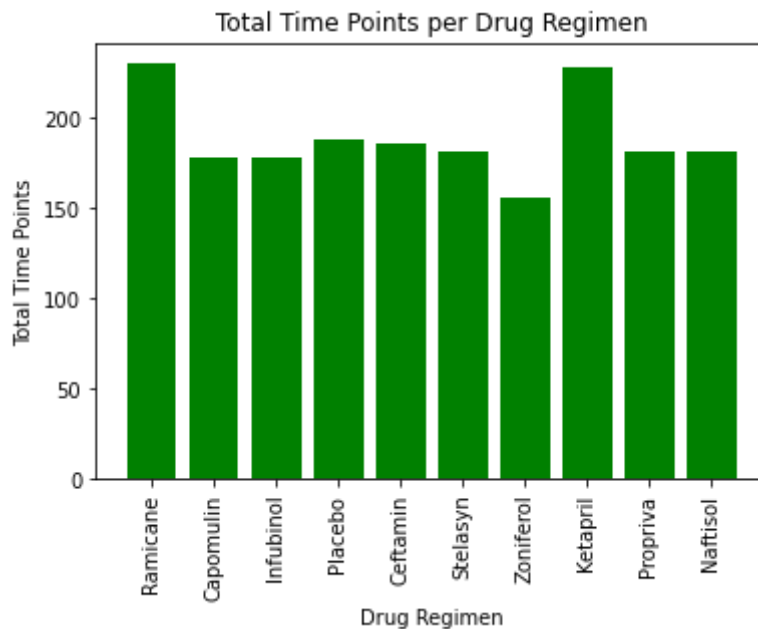
Bar and Pie Charts

```
In [407]: 1 # Generate a bar plot showing the total number of timepoints for all mic
2 x = drugs
3 yax = df.groupby('Drug Regimen').count()['Timepoint']
4 df1 = pd.DataFrame({'Drugs': x, 'Timepoint Count': yax})
5 ax = df1.plot.bar(rot=90)
```



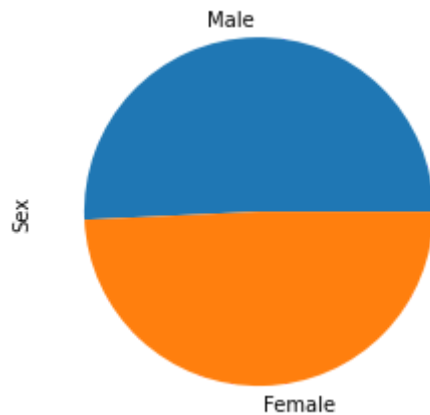
```
In [408]: 1 # Generate a bar plot showing the total number of timepoints for all mic
2 plt.bar(drugs, yax, color = 'g')
3 plt.xticks(rotation = 90)
4 plt.xlabel("Drug Regimen")
5 plt.ylabel("Total Time Points")
6 plt.title("Total Time Points per Drug Regimen")
```

Out[408]: Text(0.5, 1.0, 'Total Time Points per Drug Regimen')



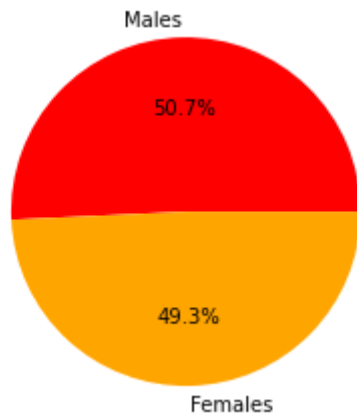
```
In [409]: 1 # Generate a pie plot showing the distribution of female versus male mi
2 genderdist = df['Sex'].value_counts()
3 print(genderdist)
4 malecount = genderdist[0]
5 femalecount = genderdist[1]
6 dist_df = pd.DataFrame({'count': [malecount, femalecount]}, index = ['M
7 plot = genderdist.plot.pie(y = 'count')
```

```
Male      958
Female    930
Name: Sex, dtype: int64
```



```
In [410]: 1 # Generate a pie plot showing the distribution of female versus male mi
2 labels = ['Males', 'Females']
3 sizes = [int(malecount), int(femalecount)]
4 plt.pie(sizes, labels=labels, colors = ['red', 'orange'], autopct="%1.1f
5
```

```
Out[410]: ([<matplotlib.patches.Wedge at 0x7fbe2c728e50>,
<matplotlib.patches.Wedge at 0x7fbe2c735340>],
[Text(-0.025622895044835736, 1.0997015355311284, 'Males'),
Text(0.025622792083347525, -1.099701537930112, 'Females')],
[Text(-0.013976124569910401, 0.5998372011987972, '50.7%'),
Text(0.01397606840909865, -0.5998372025073339, '49.3%')])
```



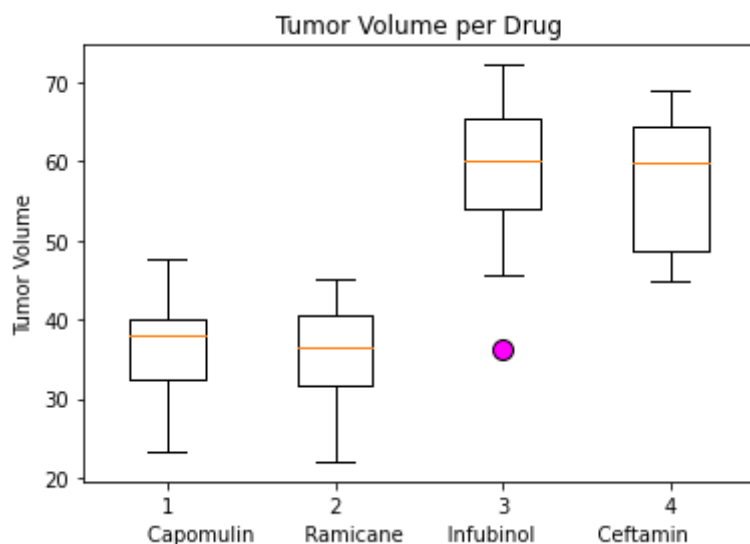
Quartiles, Outliers and Boxplots

```
In [711]: 1 # Calculate the final tumor volume of each mouse across four of the tre
2
3 grouped = df1.drop_duplicates('Mouse ID', keep='last')
4 cap = grouped.loc[grouped['Drug Regimen'] == 'Capomulin', 'Tumor Volume
5 ram = grouped.loc[grouped['Drug Regimen'] == 'Ramicane', 'Tumor Volume
6 inf = grouped.loc[grouped['Drug Regimen'] == 'Infubinol', 'Tumor Volume
7 ceft = grouped.loc[grouped['Drug Regimen'] == 'Ceftamin', 'Tumor Volume
8 d = {'Capomulin': cap, 'Ramicane': ram, 'Infubinol': inf, 'Ceftamin': c
9 ser = pd.Series(data = d, index = ['Capomulin', 'Ramicane', 'Infubinol'
10
```

```

In [726]: 1 # Generate a box plot of the final tumor volume of each mouse across fo
2 fig1, ax1 = plt.subplots()
3 ax1.set_title('Tumor Volume per Drug')
4 ax1.set_ylabel('Tumor Volume')
5 ax1.set_xlabel('Capomulin      Ramicane      Infubinol      Ceft
6 ax1.boxplot(ser, flierprops={'marker': 'o', 'markersize': 10, 'markerfa
7 plt.show()
8 outliers = inf[(inf <= lower_bound) | (inf >= upper_bound)]
9 print('The following are the outliers in the boxplot:{}'.format(outlier

```



The following are the outliers in the boxplot:669 36.321346
 Name: Tumor Volume (mm3), dtype: float64


```
In [727]: 1 # Only showing statistics for Infubinol as the box plot showed it's the
2
3 quartiles = inf.quantile([.25,.5,.75])
4 lowerq = quartiles[0.25]
5 upperq = quartiles[0.75]
6 iqr = upperq-lowerq
7
8 print(f"The lower quartile of Infubinol is: {lowerq}")
9 print(f"The upper quartile of Infubinol is: {upperq}")
10 print(f"The interquartile range of Infubinol is: {iqr}")
11
12 lower_bound = lowerq - (1.5*iqr)
13 upper_bound = upperq + (1.5*iqr)
14 print(f"Values below {lower_bound} could be outliers.")
15 print(f"Values above {upper_bound} could be outliers.")
```

The lower quartile of Infubinol is: 54.04860769

The upper quartile of Infubinol is: 65.52574285

The interquartile range of Infubinol is: 11.477135160000003

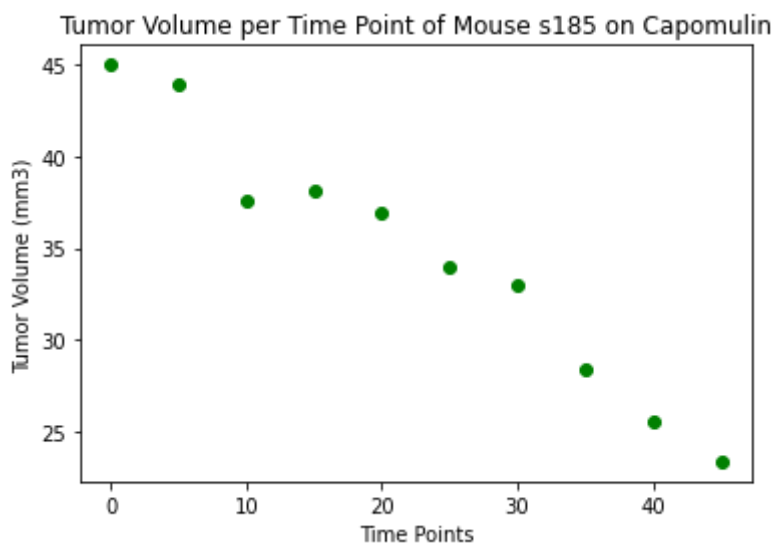
Values below 36.832904949999999 could be outliers.

Values above 82.741445590000001 could be outliers.

Line and Scatter Plots

```
In [515]: 1 # Generate a line plot of tumor volume vs. time point for a mouse treat
2 df.loc[df['Drug Regimen'] == 'Capomulin']
3 df2 = df.loc[df['Mouse ID'] == 's185']
4 tumorvol = list(df2['Tumor Volume (mm3)'])
5 timepoint = list(df2['Timepoint'])
6 plt.title('Tumor Volume per Time Point of Mouse s185 on Capomulin')
7 plt.xlabel('Time Points')
8 plt.ylabel('Tumor Volume (mm3)')
9 plt.scatter(timepoint, tumorvol, color='green')
```

Out[515]: <matplotlib.collections.PathCollection at 0x7fbe2c957d30>

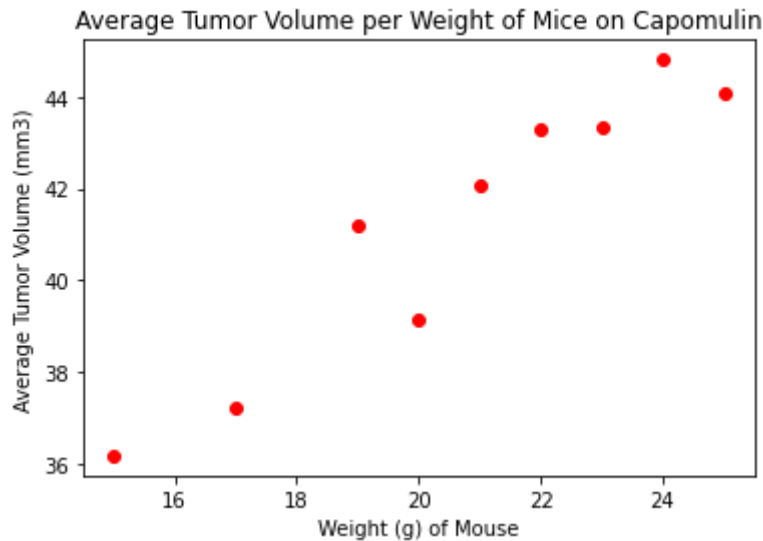


```

In [528]: 1 # Generate a scatter plot of average tumor volume vs. mouse weight for
          2 df3 = df.groupby(['Drug Regimen', 'Weight (g)'])['Tumor Volume (mm3)'].m
          3 new = pd.DataFrame(df3)
          4 df4 = pd.DataFrame(new.groupby(['Drug Regimen']).get_group('Capomulin'))
          5 x = df4['Weight (g)'].tolist()
          6 y = df4['Tumor Volume (mm3)'].tolist()
          7 plt.title('Average Tumor Volume per Weight of Mice on Capomulin')
          8 plt.xlabel('Weight (g) of Mouse')
          9 plt.ylabel('Average Tumor Volume (mm3)')
         10 plt.scatter(x, y, color='red')

```

Out[528]: <matplotlib.collections.PathCollection at 0x7f8e2caa4550>



Correlation and Regression

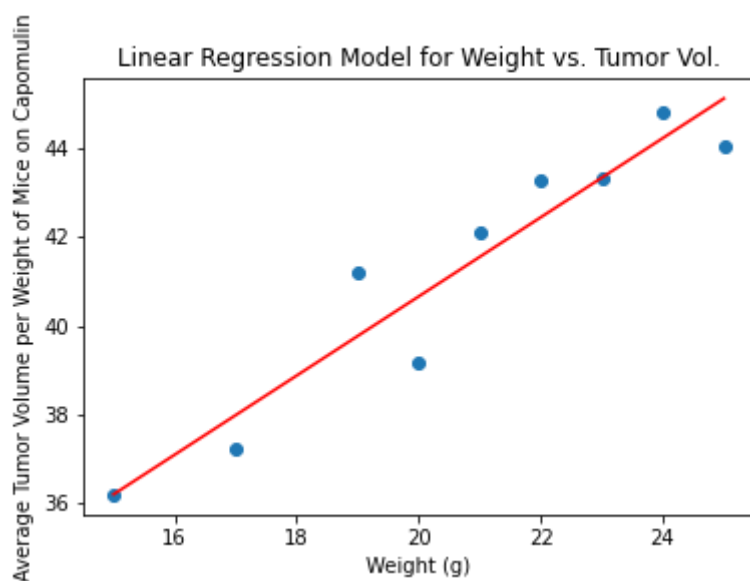
```

In [553]: Calculate the correlation coefficient and linear regression model for mouse
          correlation = st.pearsonr(x, y)
          print(f'The correlation coefficient for mouse weight and average tumor volume f

```

The correlation coefficient for mouse weight and average tumor volume for the Capomulin regimen is: 0.95

```
In [569]: 1 (slope, intercept, rvalue, pvalue, stderr) = st.linregress(x, y)
2 regress_values = np.array(x) * slope + intercept
3 line_eq = "y = " + str(round(slope,2)) + "x + " + str(round(intercept,2)
4 plt.scatter(x,y)
5 plt.plot(x,regress_values, "r-")
6 plt.annotate(line_eq, (6,10), fontsize=15,color="red",)
7 plt.title('Linear Regression Model for Weight vs. Tumor Vol.')
8 plt.xlabel('Weight (g)')
9 plt.ylabel('Average Tumor Volume per Weight of Mice on Capomulin')
10 plt.show()
11 print(line_eq)
```



$$y = 0.89x + 22.76$$

```
In [ ]: 1
```