

▼ Week 4: Mini Project

This notebook will guide you through smaller portions of your final project. For this notebook, we will use the Abalone dataset from the [UCI Machine Learning Repository](https://archive.ics.uci.edu/ml/datasets/Abalone) (originating from the Marine Research Laboratory at the University of Hawaii at Hilo, Hawaii). This dataset should already be in your folder (under `abalone.csv`) or you can download it from the [link](https://archive.ics.uci.edu/ml/datasets/Abalone).



A Brief History of Abalones

An abalone is a sea snail belonging to one of a range of 30 to 130 species (depending on which species). It is commonly prized for its mother-of-pearl shell, pearls, and delicious flesh by a variety of cultures and has been a valuable source of food in its native environments. Sadly, wild populations of abalone have been overfished and poached to the point where commercial farming supplies most of abalone flesh nowadays. It is on the list of current animals threatened by extinction.

Source: <https://en.wikipedia.org/wiki/Abalone>

▼ Part 1: Familiarize Yourself With the Dataset

The purpose of this dataset is to predict the age of an abalone through physical characteristics, determined by cutting the shell through the cone, staining it, and counting the number of rings through a microscope, a tedious and time-consuming task. Good thing it's already been done for us!

Below is the dataset description from the UCI Machine Learning Repository.

Name	Data Type	Measure	Description
Sex	nominal		M, F, and I (infant)
Length	continuous	mm	Longest shell measurement
Diameter	continuous	mm	perpendicular to length
Height	continuous	mm	with meat in shell
Whole weight	continuous	grams	whole abalone
Shucked weight	continuous	grams	weight of meat
Viscera weight	continuous	grams	gut weight (after bleeding)
Shell weight	continuous	grams	after being dried
Rings	integer		+1.5 gives the age in years

Run the cells below to examine the dataset.

```
1 # Load Abalone dataset
2 # Remember to change the file location if needed
```

```

3 import csv
4 f = open("./abalone.csv")
5 all_lines = csv.reader(f, delimiter = ',')
6
7 # We define a header ourselves since the dataset contains only the raw numbers.
8 dataset = []
9 header = ['Sex', 'Length', 'Diameter', 'Height', 'Whole Weight', 'Shucked Weight', 'Viscer
10           'Shell Weight', 'Rings']
11 for line in all_lines:
12     d = dict(zip(header, line))
13     d['Length'] = float(d['Length'])
14     d['Diameter'] = float(d['Diameter'])
15     d['Height'] = float(d['Height'])
16     d['Whole Weight'] = float(d['Whole Weight'])
17     d['Shucked Weight'] = float(d['Shucked Weight'])
18     d['Viscera Weight'] = float(d['Viscera Weight'])
19     d['Shell Weight'] = float(d['Shell Weight'])
20     d['Rings'] = int(d['Rings'])
21     dataset.append(d)

```

```

1 # See first line of dataset
2 dataset[0]

```

```

☞ {'Diameter': 0.365,
   'Height': 0.095,
   'Length': 0.455,
   'Rings': 15,
   'Sex': 'M',
   'Shell Weight': 0.15,
   'Shucked Weight': 0.2245,
   'Viscera Weight': 0.101,
   'Whole Weight': 0.514}

```

▼ Part 2: Simple Statistics

This dataset is already cleaned for us and relatively straightforward, without strings or time data. In your project, you will have to take care of missing or tricky values yourself.

Fill in the following cells with the requested information about the dataset. The answers are given so you can check the output of your own code. For floating numbers, don't worry too much about the exact number as they are quite close – different systems may have different rounding protocols.

Feel free to `import numpy` if you want more practice with it, or just use Python's native structures to play with the numbers.

```

1 # Q: What is the total number of entries in the dataset?
2 # A: 4177
3 len(dataset)

```

↩ 4177

```
1 # Q: What is the average length of an abalone?
2 # A: 0.5239920995930099 or 0.524
3 tot_length=[d['Length'] for d in dataset]
4 sum1=0.0
5 for k in range(len(tot_length)):
6     sum1=sum1+tot_length[k]
7 avg=sum1/len(tot_length)
8 avg
```

↩ 0.5239920995930099

```
1 # Q: What is the widest abalone in the dataset (diameter)?
2 # A: 0.65
3 wide_data=sorted([d['Diameter'] for d in dataset])
4 wide_data[-1]
```

↩ 0.65

```
1 # Q: What is the average number of rings of smaller abalones compared to that of larger at
2 # is, do smaller abalones tend to be younger or older than larger abalones?
3 # We will count small abalones as abalones with lengths less than or equal to the average
4 # an abalone. The average length of an abalone is 0.524.
5 # A: Small Abalones have on average 8.315645514223196 rings.
6 # Large Abalones have on average 11.192848020434228 rings.
7 large_rings=[d['Rings'] for d in dataset if d['Length']>0.524]
8 small_rings=[d['Rings'] for d in dataset if d['Length']<=0.524]
9 sum1=0.0
10 for k in range(len(small_rings)):
11     sum1=sum1+small_rings[k]
12 ageSmall=sum1/len(small_rings)
13 sum1=0.0
14 for k in range(len(large_rings)):
15     sum1=sum1+large_rings[k]
16 ageLarge=sum1/len(large_rings)
17 # Change variable name if necessary
18 print('Small Abalones have on average', ageSmall, 'rings.')
19 print('Large Abalones have on average', ageLarge, 'rings.')
```

↩ Small Abalones have on average 8.315645514223196 rings.
Large Abalones have on average 11.192848020434228 rings.

▼ Part 3: Data Visualizations

In this course, we learned about [Matplotlib](#), a "Python 2D plotting library which produces publication quality figures in a variety of hardcopy formats and interactive environments across platforms". There are a [variety of figures](#) we can make with Matplotlib, and in conjunction with NumPy, becomes a powerful and versatile skillset.

In lectures, we covered the basics of line plots, histograms, scatter plots, bar plots, and box plots. Let's see a few below:

```
1 import matplotlib.pyplot as plt
2 from matplotlib import colors
3 import numpy
4 from collections import defaultdict
```

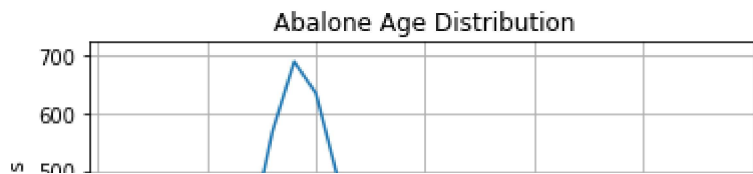
▼ Line Plots

Line plots show the change in data over time. The example Line Plot below plots the change in density of rings (i.e. the distribution of rings). **Note that a line plot is not necessarily the best way to show this data; it doesn't deal with a trend!** Use a histogram (next step) to better showcase this data.

```
1 # Parse out Rings column from dataset
2 rings = [d['Rings'] for d in dataset]
3 rings.sort()
4
5 # Count number of abalones with each number of rings with defaultdict
6 abalone_rings = defaultdict(int)
7 for r in rings:
8     abalone_rings[r] += 1
9 X = list(abalone_rings.keys())
10 Y = list(abalone_rings.values())
11
12 # Customize plot
13 plt.gca().set(xlabel='Rings', ylabel='Number of Abalones',
14               title='Abalone Age Distribution')
15 plt.grid()
16
17 # Show the plot of Rings vs Number of Abalones
18 plt.plot(X, Y)
```



[<matplotlib.lines.Line2D at 0x7f1ddae12978>]



▼ Histograms

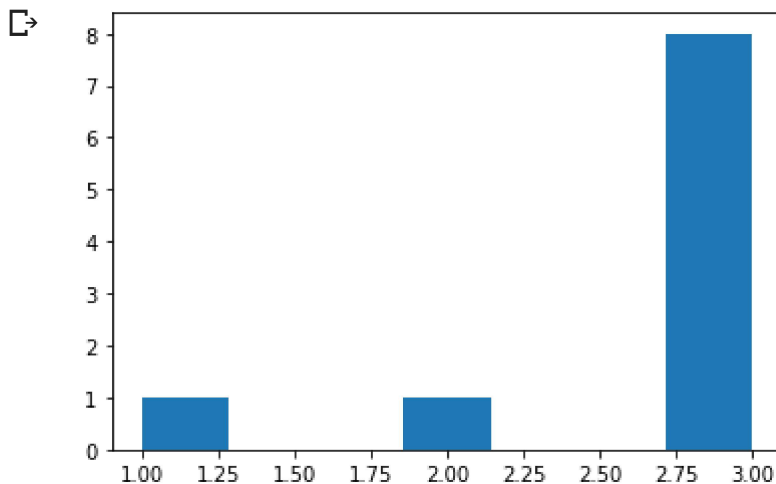
Histograms show the distribution of numeric continuous variables with central tendency and skewne **line plot data from above, plot a histogram showing the distribution of abalone age.** Feel free to exp on your own to customize your histogram and the following visualizations.



```

1 # Complete this cell with a histogram of abalone age distribution
2
3 # Flatten distribution list into frequency distribution
4 age_freq = []
5 for key in abalone_rings.keys():
6     for i in range(0, abalone_rings.get(key)):
7         age_freq.append(key)
8 k=age_freq[:10]
9 # Plot your histogram here
10 plt.hist(k,bins=7)
11 plt.show()

```



▼ Scatter Plots

Scatter plots show the strength of a relationship between two variables (also known as correlations). *Simple Statistics*, we see that larger abalones tend to be larger, at least from a numbers perspective. **this is actually true by creating a scatter plot showing the relationship between Rings and Length.**

On Your Own: Read up on `scipy` and how you can calculate and graph the correlation as well.

```

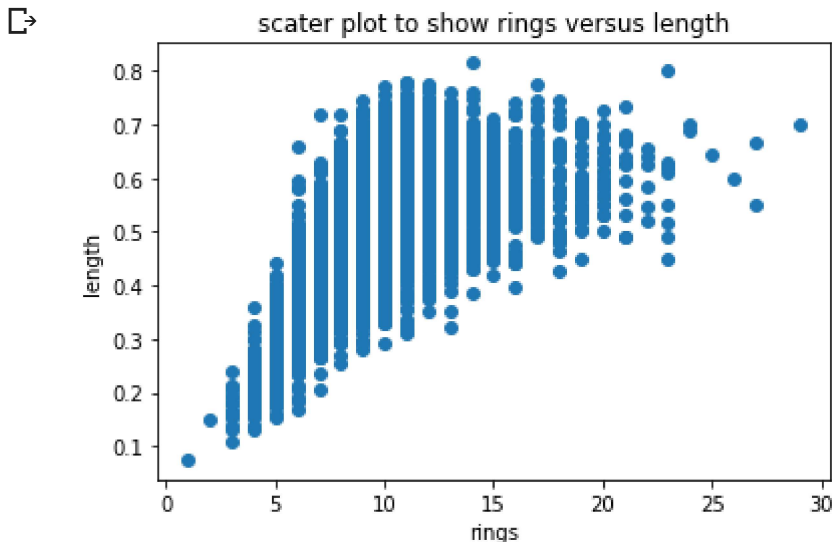
1 # Complete this cell with a scatter plot of age vs length

```

```

2 rings = [d['Rings'] for d in dataset]
3 length = [d['Length'] for d in dataset]
4 plt.scatter(rings,length)
5 plt.xlabel("rings")
6 plt.ylabel("length")
7 plt.title("scater plot to show rings versus length")
8 plt.show()
9

```



▼ Bar Plots

Bar plots are great for comparing categorical variables. There are a few subtypes of bar plots, such a bar chart or stacked bar chart. Since we have the `Sex` field to play with, we can compare data across abalones. Below is a simple stacked bar chart comparing the `Sex` category with the `Shucked Weight` **a bar chart of your choice of data.**

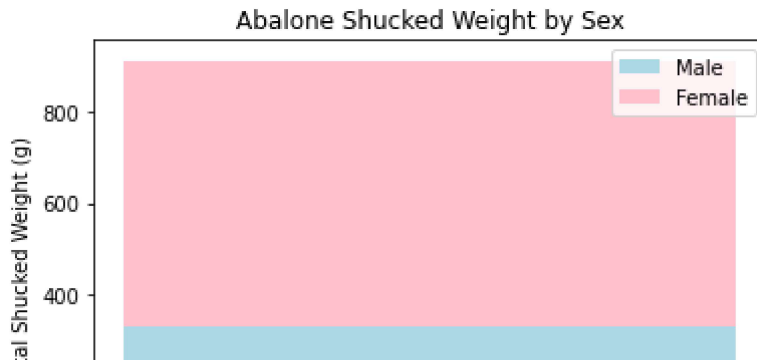
You may refer to the cell below to parse out fields by sex.

```

1 # Example Stacked Bar Chart - Comparisons Between Sexes
2 Mweight_sh = sum([d['Shucked Weight'] for d in dataset if d['Sex'] is 'M'])
3 Fweight_sh = sum([d['Shucked Weight'] for d in dataset if d['Sex'] is 'F'])
4 index = [1]
5 p1 = plt.bar(index, Mweight_sh, color='lightblue')
6 p2 = plt.bar(index, Fweight_sh, bottom=Mweight, color='pink')
7 plt.gca().set(title='Abalone Shucked Weight by Sex', ylabel='Total Shucked Weight (g)');
8 plt.xticks([])
9 plt.legend((p1[0], p2[0]), ('Male', 'Female'))
10 plt.show()

```

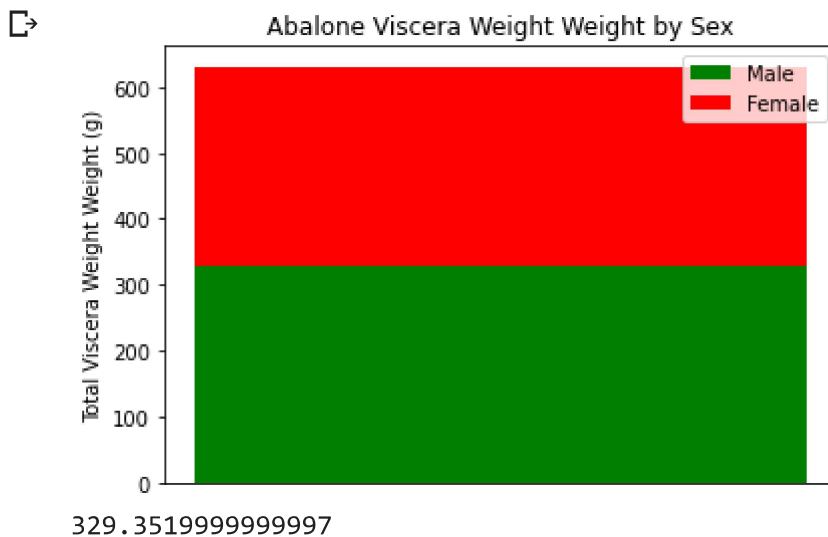




```

1 # Complete this cell with your choice of data
2 # Example Stacked Bar Chart - Comparisons Between Sexes
3 Mweight_vis = sum([d['Viscera Weight'] for d in dataset if d['Sex'] is 'M'])
4 Fweight_vis = sum([d['Viscera Weight'] for d in dataset if d['Sex'] is 'F'])
5 index = [1]
6 p1 = plt.bar(index, Mweight_vis, color='green')
7 p2 = plt.bar(index, Fweight_vis, bottom=Mweight, color='red')
8 plt.gca().set(title='Abalone Viscera Weight Weight by Sex', ylabel='Total Viscera Weight (g)')
9 plt.xticks([1])
10 plt.legend((p1[0], p2[0]), ('Male', 'Female'))
11 plt.show()

```



▼ Box Plots

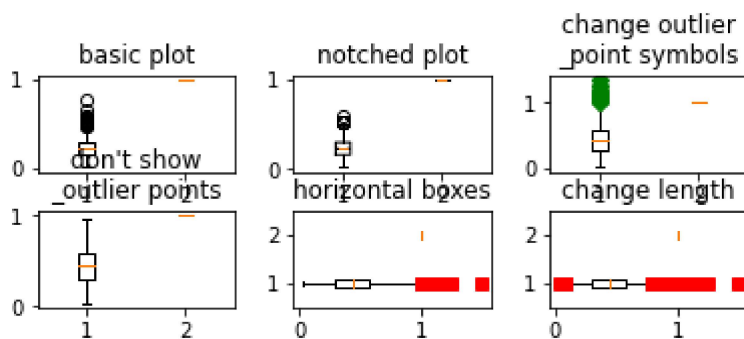
Box plots are useful for comparing distributions of data and are commonly found in research papers. A central portion of a box plot represents 50% of the data, and there are versions where you can mark outliers and extremes. We have the distribution of rings already from the line plot example under the variable name `rings`, assuming you haven't modified it. **Find the distribution of another field of your choice and create one box plot with both of these fields.**

Hint: You can plot multiple box plots with the command `plt.boxplot([plot1, plot2, ..., plotn])` or `subplots()` to draw multiple separate plots at the same time. See [this matplotlib example](#) for more.

```

1 # Complete this cell with multiple box plots
2 Mweight_vis = [d['Viscera Weight'] for d in dataset if d['Sex'] is 'M']
3 Fweight_vis = [d['Viscera Weight'] for d in dataset if d['Sex'] is 'F']
4 Mweight_sh = [d['Shucked Weight'] for d in dataset if d['Sex'] is 'M']
5 Fweight_sh = [d['Shucked Weight'] for d in dataset if d['Sex'] is 'F']
6 fig, axs = plt.subplots(2, 3)
7 # basic plot
8 data=(Mweight_vis,index)
9 axs[0, 0].boxplot(data)
10 axs[0, 0].set_title('basic plot')
11 # notched plot
12 data=(Fweight_vis,index)
13 axs[0, 1].boxplot(data, 1)
14 axs[0, 1].set_title('notched plot')
15
16 # change outlier point symbols
17 data=(Mweight_sh,index)
18 axs[0, 2].boxplot(data, 0, 'gD')
19 axs[0, 2].set_title('change outlier\n_point symbols')
20
21 # don't show outlier points
22 data=(Fweight_sh,index)
23 axs[1, 0].boxplot(data, 0, '')
24 axs[1, 0].set_title("don't show\n_outlier points")
25
26 # horizontal boxes
27 axs[1, 1].boxplot(data, 0, 'rs', 0)
28 axs[1, 1].set_title('horizontal boxes')
29
30 # change whisker length
31 axs[1, 2].boxplot(data, 0, 'rs', 0, 0.75)
32 axs[1, 2].set_title('change length')
33
34 fig.subplots_adjust(left=0.08, right=0.9, bottom=0.5, top=0.9,
35                     hspace=0.4, wspace=0.3)

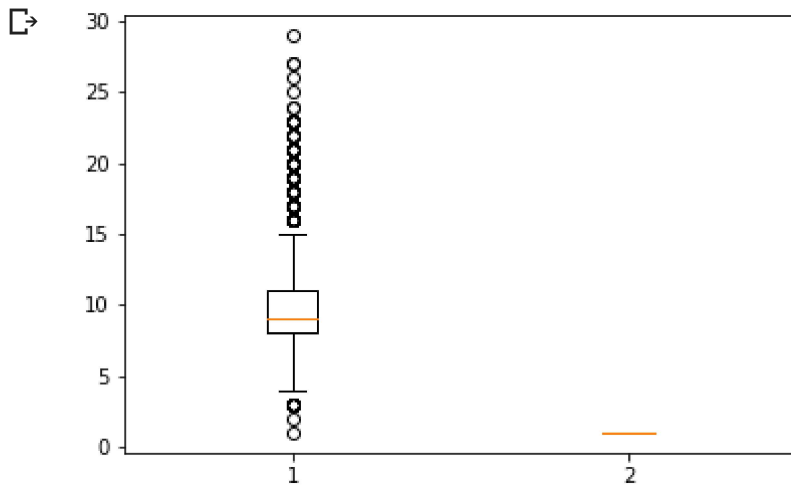
```



```

1 d2=[d['Rings'] for d in dataset]
2 data = [d2,index]
3 plt.boxplot(data)
4 plt.show()

```

All Done!

In this notebook, we covered loading a dataset, simple statistics, basic data visualizations, and web-scraping. These will be immensely helpful as you move forwards in building your skills in data science.

By now, you hopefully feel a little more confident with tackling your final project. It is up to you to find data, build your own notebook, and show others what you have achieved. Best of luck!