

Complete the following exercises. Each question is a single codeblock with the question at the top as a comment.

- (1) Import numpy as np, matplotlib.pyplot as plt, and pandas as pd. Set matplotlib to inline display.
- (2) Read the Fisher Iris Data in as a variable called df.
- (3) Display the first five lines of df using head.
- (4-7) - Plot a histogram of each of the four variables (leaf/sepal length/width). Adjust the number of bins to produce a resolution you find instructive/visually satisfying. Use xlabel and ylabel to label the axes.
- (8) - Save this figure use savefig.
- (9) Using groupby, group the data by species
- (10) Using the groups, create a bar graph of the mean sepal length by species. Label the axes.
- (11) Use subplots to make a 2x2 pane figure, with each pane a scatterplot of your choosing. Make sure axes are labeled!
- (12) Use tail to print out the last 10 rows of df.
- (13) Use a combination of np.sum() and a selector (e.g., (df['x'] > 3), count how many rows have a petal length greater than 40.
- (14) How many have a sepal width below the median sepal width?
- (15) Use groups to report the mean values for each variable.
- (16) Calculate a new variable based on existing variables, and create a histogram of its distribution.
- (17) Export the data with the new calculated variable as an Excel file. Open it in Excel, and comment on its appearance. Look at the documentation of to_excel and figure out how to make it so the index values (line numbers) don't get exported.
- (18-25) Download the prepared election data

(https://github.com/thomaspingel/geodata/blob/master/election/state_election_data_1976-2016.csv), load it as a new dataframe, and use some of the commands above to explore the dataset. *Each question # should correspond to a single codeblock that examines and sheds light on the dataset.* Focus not only on executing working code, but usefully exploring the data.

```
In [1]: # Question 1
import numpy as np
import matplotlib.pyplot as plt
import pandas as pd
```

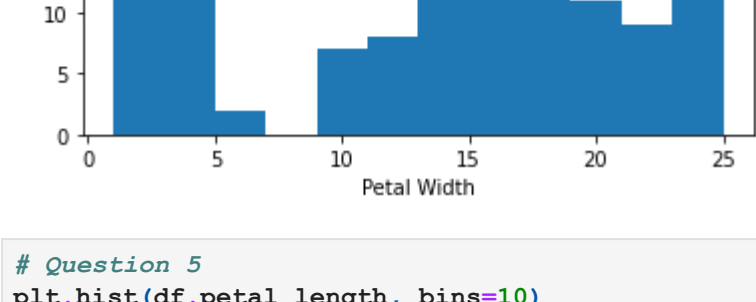
```
In [9]: # Question 2
df = pd.read_csv('fisher_iris_data.csv.csv')
```

```
In [5]: # Question 3
df.head()
```

```
Out[5]:
```

	species	petal_width	petal_length	sepal_width	sepal_length
0	setosa	2	14	33	50
1	setosa	2	10	36	46
2	setosa	2	16	31	48
3	setosa	1	14	36	49
4	setosa	2	13	32	44

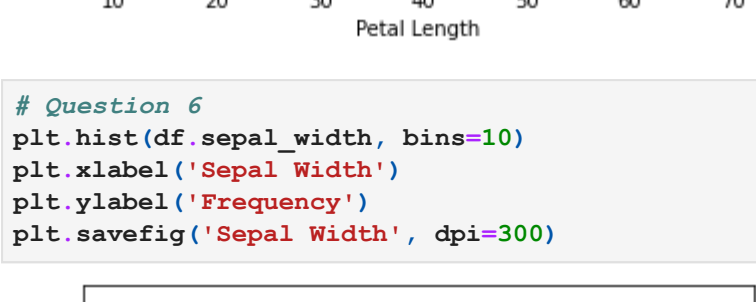
```
In [59]: # Question 4
plt.hist(df.petal_width, bins=12)
plt.xlabel('Petal Width')
plt.ylabel('Frequency')
plt.savefig('Petal Width', dpi=300)
```



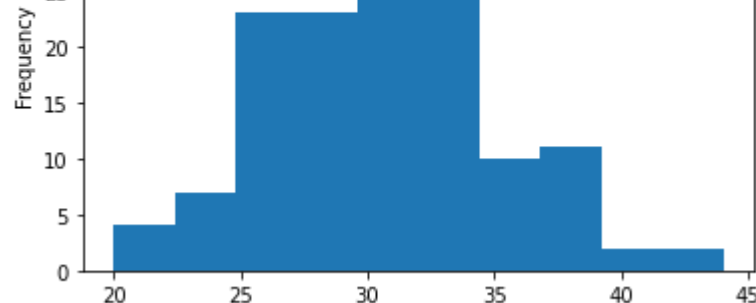
```
In [58]: # Question 5
plt.hist(df.petal_length, bins=10)
plt.xlabel('Petal Length')
plt.ylabel('Frequency')
plt.savefig('Petal Length', dpi=300)
```



```
In [57]: # Question 6
plt.hist(df.sepal_width, bins=10)
plt.xlabel('Sepal Width')
plt.ylabel('Frequency')
plt.savefig('Sepal Width', dpi=300)
```



```
In [56]: # Question 7
plt.hist(df.sepal_length, bins=18)
plt.xlabel('Sepal Length')
plt.ylabel('Frequency')
plt.savefig('Sepal Length', dpi=300)
```



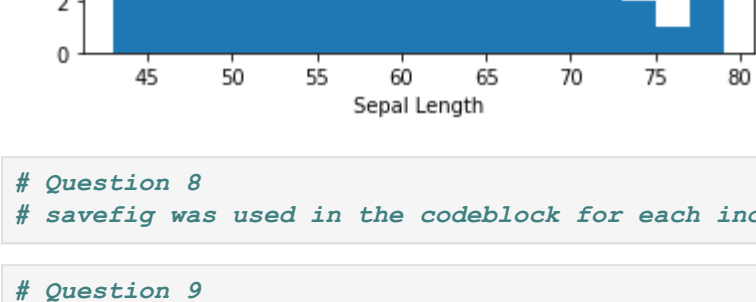
```
In [60]: # Question 8
# savefig was used in the codeblock for each individual histogram
```

```
In [3]: # Question 9
groups = df.groupby(by='species')
groups
```

```
Out[3]: <pandas.core.groupby.generic.DataFrameGroupBy object at 0x00002CF14A81A90>
```

```
In [27]: # Question 10
means = groups['sepal_length'].mean()
means
means.index
plt.bar(means.index, means)
plt.ylabel('Mean Sepal Length')
plt.xlabel('Species')
```

```
Out[27]: Text(0.5, 0, 'Species')
```



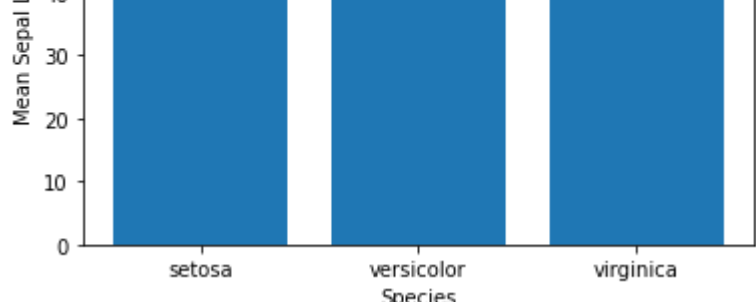
```
In [33]: # Question 11
plt.subplot(2,2,1)
plt.scatter(df.petal_width, df.petal_length)
plt.xlabel('Petal Width')
plt.ylabel('Petal Length')
```

```
plt.subplot(2,2,2)
plt.scatter(df.petal_width, df.sepal_length)
plt.xlabel('Petal Width')
plt.ylabel('Sepal Length')
```

```
plt.subplot(2,2,3)
plt.scatter(df.petal_length, df.sepal_width)
plt.xlabel('Petal Length')
plt.ylabel('Sepal Width')
```

```
plt.subplot(2,2,4)
plt.scatter(df.petal_length, df.sepal_length)
plt.xlabel('Petal Length')
plt.ylabel('Sepal Length')
```

```
plt.subplots_adjust(wspace=.5, hspace=.5)
```



```
In [34]: # Question 12
df.tail(10)
```

```
Out[34]:
```

	species	petal_width	petal_length	sepal_width	sepal_length
140	virginica	16	58	30	72
141	virginica	21	59	30	71
142	virginica	18	56	29	63
143	virginica	23	69	26	77
144	virginica	19	61	28	74
145	virginica	18	63	29	73
146	virginica	22	58	30	65
147	virginica	19	53	27	64
148	virginica	20	50	25	57
149	virginica	24	51	28	58

```
In [36]: # Question 13
np.sum(df.petal_length > 40)
```

```
Out[36]: 85
```

```
In [40]: # Question 14
median = np.median(df.sepal_width)
median
```

```
np.sum(df.sepal_width < median)
```

```
Out[40]: 57
```

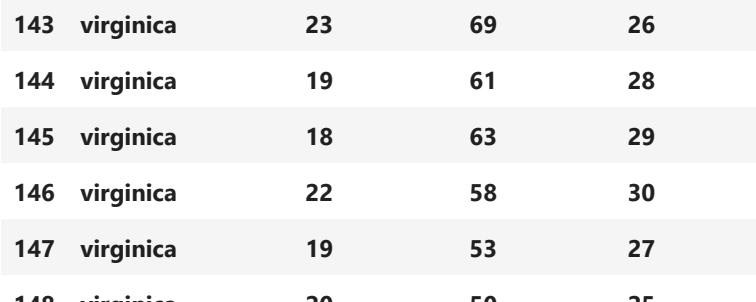
```
In [8]: # Question 15
means = groups.mean()
means
```

```
Out[8]:
```

	petal_width	petal_length	sepal_width	sepal_length
species				
setosa	2.46	14.62	34.28	50.10
versicolor	13.26	43.22	27.64	59.36
virginica	20.06	55.52	29.74	65.88

```
In [63]: # Question 16
df['sepal_thickness'] = df.sepal_length / df.sepal_width
plt.hist(df.sepal_thickness, bins=10)
plt.xlabel('Sepal Thickness')
plt.ylabel('Frequency')
```

```
Out[63]: Text(0, 0.5, 'Frequency')
```



```
In [95]: # Question 17
df.to_excel('fisherdata.xlsx', index=False)
```

```
# The excel file looks just like the dataframe after sepal_thickness was added. The
# index values are unnecessary because excel sheets include them by default,
# so the index needed to be set to False for them to not show up.
```

```
In [15]: # Data for 18-25
url = 'https://raw.githubusercontent.com/thomaspingel/geodata/master/election/state_election_data_1976-2020.csv'
election_df = pd.read_csv(url)
election_df.head()
```

```
Out[15]:
```

	state	state_po	FIPS	gop_1976_votes	dem_1976_votes	totalvotes_1976	gop_1976_prc	dem_1976_prc	gop_minus_dem_prc_1976	gop_minus_dem_prc_2020
0	ALABAMA	AL	1	504070	659170	1182850	42.61	55.73	-13.12	-13.12
1	ALASKA	AK	2	71555	44058	123574	57.90	35.65	22.25	22.25
2	ARIZONA	AZ	4	418642	295602	742719	56.37	39.80	16.57	16.57
3	ARKANSAS	AR	5	267903	498604	767535	34.90	64.96	-30.06	-30.06
4	CALIFORNIA	CA	6	3882244	3742284	7803770	49.75	47.95	1.80	1.80

5 rows x 75 columns

```
In [94]: # Question 18
election_df['voter_turnout_change'] = election_df.totalvotes_2020 - election_df.totalvotes_1976
election_df.head(10)
```

```
# Adds a column to the table showing increase in voter turnout for each state between
# 1976 and 2020
```

```
Out[94]:
```

	state	state_po	FIPS	gop_1976_votes	dem_1976_votes	totalvotes_1976	gop_1976_prc	dem_1976_prc	gop_minus_dem_prc_1976	gop_minus_dem_prc_2020
0	ALABAMA	AL	1	504070	659170	1182850	42.61	55.73	-13.12	-13.12
1	ALASKA	AK	2	71555	44058	123574	57.90	35.65	22.25	22.25
2	ARIZONA	AZ	4	418642	295602	742719	56.37	39.80	16.57	16.57
3	ARKANSAS	AR	5	267903	498604	767535	34.90	64.96	-30.06	-30.06
4	CALIFORNIA	CA	6	3882244	3742284	7803770	49.75	47.95	1.80	1.80

10 rows x 79 columns

```
In [47]: # Question 19
np.sum(election_df.dem_2020_votes > election_df.gop_2020_votes)
```

```
# Counts how many states had more votes for Democrat than Republican in 2020
```

```
Out[47]: 26
```

```
In [93]: # Question 20
np.sum(election_df.dem_2020_prc > election_df.dem_2016_prc)
```

```
# Counts how many states had an increase in Democratic vote percentage between 2016 and
# 2020
```

```
Out[93]: 51
```

```
In [92]: # Question 21
election_df['dem_voter_prc_increase_2016_to_2020'] = election_df.dem_2020_prc - election_df.dem_2016_prc
election_df.head(10)
```

```
# Adds a column to the table that shows, from 2016 to 2020, the increase in percentage
# of Democratic votes for each state
```

```
Out[92]:
```

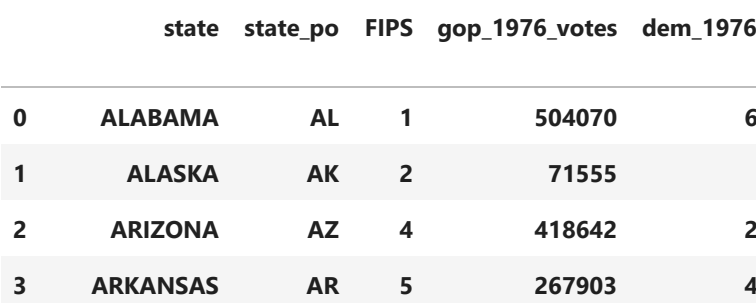
	state	state_po	FIPS	gop_1976_votes	dem_1976_votes	totalvotes_1976	gop_1976_prc	dem_1976_prc	gop_minus_dem_prc_1976	gop_minus_dem_prc_2020
0	ALABAMA	AL	1	504070	659170	1182850	42.61	55.73	-13.12	-13.12
1	ALASKA	AK	2	71555	44058	123574	57.90	35.65	22.25	22.25
2	ARIZONA	AZ	4	418642	295602	742719	56.37	39.80	16.57	16.57
3	ARKANSAS	AR	5	267903	498604	767535	34.90	64.96	-30.06	-30.06
4	CALIFORNIA	CA	6	3882244	3742284	7803770	49.75	47.95	1.80	1.80

10 rows x 79 columns

```
In [40]: # Question 22
plt.scatter(election_df.gop_1976_votes / election_df.totalvotes_1976, election_df.gop_2020_votes / election_df.totalvotes_2020)
plt.xlabel('1976 GOP % of Vote')
plt.ylabel('2020 GOP % of Vote')
plt.title('State GOP Vote 1976 and 2020')
```

```
# Makes a scatterplot showing GOP vote percentage in 1976 and 2020
```

```
Out[40]: Text(0.5, 1.0, 'State GOP Vote 1976 and 2020')
```



```
In [43]: # Question 23
election_means = np.mean(election_df.totalvotes_2020 - election_df.totalvotes_2016)
election_means
```

```
# Calculates the mean increase in voter turnout per state between from 2016 to 2020
```

```
Out[43]: 65798.07843137255
```

```
In [91]: # Question 24
np.sum(election_df.dem_2020_prc > 60)
```

```
# Counts the number of states in which the 2020 Democratic margin of victory was greater
# than 20%
```

```
Out[91]: 7
```

```
In [58]: # Question 25
election_means2 = np.mean(election_df.gop_2020_votes)
election_means2
```

```
# Calculates the mean number of Republican voters per state in 2020
```

```
Out[58]: 1455218.549019608
```