```python
1  #%%
2  import numpy as np
3  import pandas as pd
4  import sns as sns
5
6  #%% md
7  1. Read the Auto data (5 points)
8      a. use pandas to read the data
9      b. output the first few rows
10     c. output the dimensions of the data
11 #%%
12 data = pd.read_csv("Auto.csv")
13 print(data)
14 #%%
15 print(data.head())
16 print("Dimensions:", data.shape)
17 #%% md
18 2. Data exploration with code (5 points)
19     a. use describe() on the mpg, weight, and year
   columns
20     b. write comments indicating the range and
   average of each column
21 #%%
22 print("MPG column:")
23 print(data["mpg"].describe())
24
25 print("Weight column:")
26 print(data["weight"].describe())
27
28 print("Year column:")
29 print(data["year"].describe())
30 #%% md
31 MPG Column: Range = 37.6 , Mean = 23.45
32 Weight Column: Range = 3527 , Mean = 2977.58
33 Year Column: Range = 12 , Mean = 76.01
34 #%% md
35 3. Explore data types (5 points)
36     a. check the data types of all columns
37     b. change the cylinders column to categorical (
   use cat.codes)
38     c. change the origin column to categorical (don
```

```
38 't use cat.codes)
39     d. verify the changes with the dtypes attribute
40 #%%
41 print(data.dtypes)
42 #%%
43 data["cylinders"] = data["cylinders"].astype("
   category")
44 data["cylinders_cat"] = data["cylinders"].cat.codes
45
46 data["origin"] = data["origin"].astype("category")
47
48 print(data.dtypes)
49 print(data.shape)
50 #%% md
51 4. Deal with NAs (5 points)
52     a. delete rows with NAs
53     b. output the new dimensions
54 #%%
55 print(data.shape)
56 data.dropna(inplace=True)
57 print("Dimensions after deleting rows with NAs:")
58 print(data.shape)
59 #%% md
60 5. Modify columns (10 points)
61     a. make a new column, mpg_high, and make it
   categorical:
62         i. the column == 1 if mpg > average mpg,
   else == 0
63     b. delete the mpg and name columns (delete mpg
   so the algorithm doesn't just learn
64     to predict mpg_high from mpg)
65     c. output the first few rows of the modified
   data frame
66 #%%
67 avg_mpg = data['mpg'].mean()
68 data['mpg_high'] = (data['mpg'] > avg_mpg).astype(
   int)
69
70 data.drop(columns=['mpg', 'name'], inplace=True)
71
72 print(data.head())
```

```
73  #%% md
74  6. Data exploration with graphs (15 points)
75      a. seaborn catplot on the mpg_high column
76      b. seaborn relplot with horsepower on the x
    axis, weight on the y axis, setting hue or
77      style to mpg_high
78      c. seaborn boxplot with mpg_high on the x axis
     and weight on the y axis
79      d. for each graph, write a comment indicating
    one thing you learned about the data
80      from the graph
81  #%%
82  import seaborn as sns
83  import matplotlib.pyplot as plt
84
85  sns.catplot(x="mpg_high", kind="count", data=data)
86  sns.relplot(x="horsepower", y="weight", hue="
    mpg_high", data=data)
87  sns.boxplot(x="mpg_high", y="weight", data=data)
88  plt.show()
89  #%% md
90  Catplot: We can see that there are more vehicles
    that have mpg lower than the average. However the
    numbers are really close.
91  Retplot and Boxplot: we can see that the cars with
     a lower weight has less mpg. This makes sense
    because the engine now requires less amount of
    work to move the vehicle.
92  #%% md
93  7. Train/test split (5 points)
94      a. 80/20
95      b. use seed 1234 so we all get the same
    results
96      c. train /test X data frames consists of all
    remaining columns except mpg_high
97      d. output the dimensions of train and test
98  #%%
99  from sklearn.model_selection import
    train_test_split
100
101 np.random.seed(1234)
```

```python
102 X = data.drop('mpg_high', axis=1)
103 y = data['mpg_high']
104
105 X_train, X_test, y_train, y_test =
    train_test_split(X, y, test_size=0.2)
106
107 print("Train:", X_train.shape, y_train.shape)
108 print("Test:", X_test.shape, y_test.shape)
109 #%% md
110 8. Logistic Regression (10 points)
111     a. train a logistic regression model using
    solver lbfgs
112     b. test and evaluate
113     c. print metrics using the classification
    report
114 #%%
115 from sklearn.linear_model import
    LogisticRegression
116 from sklearn.metrics import classification_report
117
118 clasf = LogisticRegression(solver='lbfgs')
119 clasf.fit(X_train, y_train)
120
121 y_pred = clasf.predict(X_test)
122 print(classification_report(y_test, y_pred))
123 #%% md
124 9. Decision Tree (10 points)
125     a. train a decision tree
126     b. test and evaluate
127     c. print the classification report metrics
128 #%%
129 from sklearn.tree import DecisionTreeClassifier
130 from sklearn.metrics import classification_report
131
132 tree = DecisionTreeClassifier(max_depth=3,
    random_state=1234)
133 tree.fit(X_train, y_train)
134 y_pred = tree.predict(X_test)
135
136 print(classification_report(y_test, y_pred))
137 #%% md
```

```
138 10. Neural Network (15 points)
139     a. train a neural network, choosing a network
    topology of your choice
140     b. test and evaluate
141     c. train a second network with a different
    topology and different settings
142     d. test and evaluate
143     e. compare the two models and why you think
    the performance was same/different
144 #%%
145 from sklearn.neural_network import MLPClassifier
146
147 nn = MLPClassifier(hidden_layer_sizes=(20,),
    max_iter=700, random_state=1234)
148 nn.fit(X_train, y_train)
149 y_pred = nn.predict(X_test)
150 print("Neural Network 1:")
151 print(classification_report(y_test, y_pred))
152
153 nn2 = MLPClassifier(hidden_layer_sizes=(20, 10),
    max_iter=1000, random_state=1234)
154 nn2.fit(X_train, y_train)
155 y_pred2 = nn2.predict(X_test)
156 print("Neural Network 2:")
157 print(classification_report(y_test, y_pred2))
158 #%% md
159 The two neural networks have very similar results
    . The First one is predicting better results for
    class 1 while the second one is producing better
    results for class 0. The result are same because
    the topology being very similar too.
160 #%% md
161 11. Analysis (15 points)
162     a. which algorithm performed better?
163     It looks like the logistic algorithm performed
    the best.
164     b. compare accuracy, recall and precision
    metrics by class
165     By comparing all three we can still see that
    logistic expression is giving us the best results.
166     c. give your analysis of why the better-
```

```
166 performing algorithm might have outperformed
167     the other
168     We can see that Logistic expression is giving
    us the best accuracy. This is because simple
    algorithm like logistic expression is suited for
    this problem
169     d. write a couple of sentences comparing your
    experiences using R versus sklearn. Feel
170     free to express strong preferences.
171     Python is definitely easier for me because I
    did a similar but more complexive project like
    this last semester which is why I had some
    previous experience. However in R, it is my first
    time trying it so it was quite difficult for me at
     the beginning.
```