

```

1 # Week 12 - Assessed exercises
2
3 # This week we learnt some advanced data manipulation methods, about APIs and
4 # about webscraping. In this last set of assessed exercises you must complete the
5 # Brightspace quiz 'W12 - Assessed exercises' and submit a .py file with the
6 # code you used to answer the questions in your quiz. Each question is worth
7 # 0.5 marks and it is either correct (full marks) or incorrect (0 marks).
8
9 # This template file contains code that will help you answer the questions in
10 # the quiz.
11
12 # Q1 and Q2 are based on the advanced data manipulation section. You will need to
13 # use the titanic dataset which is part of the seaborn package and can be loaded
14 # using the following commands
15 import seaborn as sb
16 import re
17 from bs4 import BeautifulSoup
18 import requests
19 import matplotlib.pyplot as plt
20 import datetime
21 import wldata as wbd
22 from pandas import Series, DataFrame
23 import pandas as pd
24 import numpy as np
25 titanic = sb.load_dataset('titanic')
26 titanic
27 # The questions involve using the groupby function and applying functions to that
28 # grouped object. For questions involving the interquartile range, use the function
29 # from lecture_12_code.py
30
31
32 # Q1
33 fare_group = titanic.fare.groupby([titanic.sex, titanic.survived])
34 q1 = fare_group.mean()
35 print('The average fare of male passengers who did not survive is: {}'.format(
36     round(q1[2], 0)))
37 # Q2
38 survived_group = titanic.survived.groupby([titanic.sex, titanic.pclass])
39 q2 = survived_group.mean()
40 print('The percentage of female first class passengers (pclass=1) survived:
41 {:.2f}'.format(
42     q2[0]))
43 # Q3 to Q5 relate to the World Bank API. You will be asked to search for indicator
44 # and country codes in Q3 and Q4. In Q5 you will need to extract data from the
45 # the World Bank for a particular indicator, country and year
46
47 # Q3
48 age_group = titanic.groupby(['class', 'sex', 'survived'])
49 quant_3 = age_group.age.quantile(q=0.75)[9]
50 quant_1 = age_group.age.quantile(q=0.25)[9]
51
52 print(' the interquartile range for the age of female third class passengers
53 (class=3) who survived (survived=1) is : {}'.format(
54     round(quant_3-quant_1, 0)))
55
56 # Q4
57 # the indicator code for "Taxes on exports (% of tax revenue)" is GC.TAX.EXPT.ZS.
58
59 # Q5
60 indicator1 = {

```

```

59     'GC.TAX.YPKG.RV.ZS': 'the taxes on income, profits and capital gains, as a % of
revenue, for the Egypt in 2007'}
60 data_date1 = (datetime.datetime(2007, 1, 1), datetime.datetime(2007, 12, 31))
61 data1 = wbd.get_dataframe(indicator1, 'EGY', data_date1)
62 data1
63 # The taxes on income, profits and capital gains, as a % of revenue, for the Egypt
in 2007 is about 28%.
64
65 # Q6 to Q8 relate to webscraping and uses the Spotify weekly charts. You will need
66 # to import BeautifulSoup and the requests package
67 # The below code loads the data from the Spotify weekly charts for the week
68 # 2017-06-30 to 2017-07-07, and uses BeautifulSoup to parse the html.
69 spotify = requests.get(
70     'https://spotifycharts.com/regional/global/weekly/2017-06-30--2017-07-07')
71 soup = BeautifulSoup(spotify.text, "html.parser")
72 # The following commands extract the information related to the tracks and removes
73 # the html tags
74 track = soup.find_all('td', class_="chart-table-track")
75 tracks = [x.text.strip() for x in track]
76 # Q6 asks you to search through tracks to find the number of times a particular
77 # artist appears in this weekly chart
78 count1 = 0
79 for i in tracks:
80     if 'Justin Bieber' in i:
81         count1 += 1
82 count1
83 # Justin Bieber appeared 5 times
84
85 # The following commands extract the information related to the number of plays,
86 # removes the html tags and commas, and converts the value to an integers
87 play = soup.find_all('td', class_="chart-table-streams")
88 plays = [int(x.text.strip().replace(',', '')) for x in play]
89
90 # Q7 asks you to perform some statistical analysis on these numbers
91 count2 = pd.Series(plays).mean()
92 print('the mean number of plays for the 200 songs in the 2017-06-30 to 2017-07-07 is
{}'.format(round(count2,0)))
93
94
95 # Q8 asks you to load the charts for a different week and determine how many of
96 # the songs from the original week 2017-06-30 to 2017-07-07 are still in the
97 # charts at this later week. To load in the data for the new week, change the
98 # date range in the url to the date range specified in your question.
99 spotify2 = requests.get(
100     'https://spotifycharts.com/regional/global/weekly/2017-08-25--2017-09-01')
101 soup2 = BeautifulSoup(spotify2.text, "html.parser")
102 # The following commands extract the information related to the tracks and removes
103 # the html tags
104 track2 = soup2.find_all('td', class_="chart-table-track")
105 tracks2 = [x.text.strip() for x in track2]
106
107 count3 = 0
108 for i in tracks:
109     if i in tracks2:
110         count3 += 1
111 count3
112 # 125 songs were still in the chart.
113

```