

# hw02

March 28, 2020

## 1 Homework 2: Arrays and Tables

**Recommended Reading:** \* [Data Types](#) \* [Sequences](#) \* [Tables](#)

Please complete this notebook by filling in the cells provided. Before you begin, execute the following cell to load some necessary information for the assignment.

Homework 2 is due Thursday, 1/30 at 11:59pm. Start early so that you can come to office hours if you're stuck. Check the website for the office hours schedule. Late work will not be accepted as per the policies of this course.

```
[3]: # Don't change this cell; just run it.
```

```
import numpy as np
from datascience import *
```

### 1.1 1. Creating Arrays

**Question 1.1.** Make an array called `weird_numbers` containing the following numbers (in the given order):

1. -2
2. the sine of 1.2
3. 3
4. 5 to the power of the cosine of 1.2

*Hint: sin and cos are functions in the math module.*

```
[4]: # Our solution involved one extra line of code before creating
# weird_numbers.
import math
weird_numbers = make_array (-2, math.sin(1.2), 3, 5**math.cos(1.2))
weird_numbers
```

```
[4]: array([-2.          ,  0.93203909,  3.          ,  1.79174913])
```

**Question 1.2.** Make an array called `book_title_words` containing the following three strings: "Eats", "Shoots", and "and Leaves".

```
[5]: book_title_words = make_array ("Eats", "Shoots", "Leaves")
book_title_words
```

```
[5]: array(['Eats', 'Shoots', 'Leaves'], dtype='<U6')
```

Strings have a method called `join`. `join` takes one argument, an array of strings. It returns a single string. Specifically, the value of `a_string.join(an_array)` is a single string that's the **concatenation** ("putting together") of all the strings in `an_array`, **except** `a_string` is inserted in between each string.

**Question 1.3.** Use the array `book_title_words` and the method `join` to make two strings:

1. "Eats, Shoots, and Leaves" (call this one `with_commas`)
2. "Eats Shoots and Leaves" (call this one `without_commas`)

*Hint:* If you're not sure what `join` does, first try just calling, for example, `"foo".join(book_title_words)`.

```
[6]: with_commas = ("Eats, Shoots, and Leaves")
    without_commas = ("Eats Shoots and Leaves")

    # These lines are provided just to print out your answers.

    print (' , '.join(book_title_words))

    print (' '.join(book_title_words) )
```

```
Eats, Shoots, Leaves
Eats Shoots Leaves
```

## 1.2 2. Indexing Arrays

These exercises give you practice accessing individual elements of arrays. In Python (and in many programming languages), elements are accessed by *index*, so the first element is the element at index 0.

**Question 2.1.** The cell below creates an array of some numbers. Set `third_element` to the third element of `some_numbers`.

```
[7]: some_numbers = make_array(-1, -3, -6, -10, -15)

    third_element = some_numbers.item(2)
    third_element
```

```
[7]: -6
```

**Question 2.2.** The next cell creates a table that displays some information about the elements of `some_numbers` and their order. Run the cell to see the partially-completed table, then fill in the missing information in the cell (the strings that are currently "???" ) to complete the table.

```
[8]: elements_of_some_numbers = Table().with_columns(
    "English name for position", make_array("first", "second", "third",
    → "fourth", "fifth"),
    "Index", make_array("0", "1", "2", "3", "4"),
    "Element", some_numbers)
    elements_of_some_numbers
```

```
[8]: English name for position | Index | Element
    first                    | 0    | -1
```

second	1	-3
third	2	-6
fourth	3	-10
fifth	4	-15

**Question 2.3.** You'll sometimes want to find the *last* element of an array. Suppose an array has 142 elements. What is the index of its last element?

```
[9]: index_of_last_element = 141
```

More often, you don't know the number of elements in an array, its *length*. (For example, it might be a large dataset you found on the Internet.) The function `len` takes a single argument, an array, and returns the length of that array (an integer).

**Question 2.4.** The cell below loads an array called `president_birth_years`. The last element in that array is the most recent birth year of any deceased president as of 2017. Assign that year to `most_recent_birth_year`.

```
[10]: president_birth_years = Table.read_table("president_births.csv").column('Birth_
      ↳Year')
      most_recent_birth_year = max(president_birth_years)
      most_recent_birth_year
```

```
[10]: 1917
```

**Question 2.5.** Finally, assign `sum_of_birth_years` to the sum of the first, tenth, and last birth year in `president_birth_years`

```
[11]: sum_of_birth_years = president_birth_years.item(0) + president_birth_years.
      ↳item(9) + most_recent_birth_year
      sum_of_birth_years
```

```
[11]: 5433
```

### 1.3 3. Basic Array Arithmetic

**Question 3.1.** Multiply the numbers 42, 4224, 42422424, and -250 by 157. For this question, **don't** use arrays.

```
[12]: first_product = 42 * 157
      second_product = 4224 * 157
      third_product = 42422424 * 157
      fourth_product = -250 * 157
      print(first_product, second_product, third_product, fourth_product)
```

```
6594 663168 6660320568 -39250
```

**Question 3.2.** Now, do the same calculation, but using an array called `numbers` and only a single multiplication (`*`) operator. Store the 4 results in an array named `products`.

```
[13]: numbers = make_array(42, 4224, 42422424, -250)
      products = numbers * 157
      products
```

```
[13]: array([ 6594, 663168, 6660320568, -39250])
```

**Question 3.3.** Oops, we made a typo! Instead of 157, we wanted to multiply each number by 1577. Compute the fixed products in the cell below using array arithmetic. Notice that your job is really easy if you previously defined an array containing the 4 numbers.

```
[14]: fixed_products = numbers * 1577
      fixed_products
```

```
[14]: array([ 66234, 6661248, 66900162648, -394250])
```

**Question 3.4.** We've loaded an array of temperatures in the next cell. Each number is the highest temperature observed on a day at a climate observation station, mostly from the US. Since they're from the US government agency [NOAA](#), all the temperatures are in Fahrenheit. Convert them all to Celsius by first subtracting 32 from them, then multiplying the results by  $\frac{5}{9}$ . Make sure to **ROUND** each result to the nearest integer using the `np.round` function.

```
[15]: max_temperatures = Table.read_table("temperatures.csv").column("Daily Max_
      ↪Temperature")

      celsius_max_temperatures = np.round ((max_temperatures-32)*5/9)
      celsius_max_temperatures
```

```
[15]: array([-4., 31., 32., ..., 17., 23., 16.])
```

**Question 3.5.** The cell below loads all the *lowest* temperatures from each day (in Fahrenheit). Compute the size of the daily temperature range for each day. That is, compute the difference between each daily maximum temperature and the corresponding daily minimum temperature. **Give your answer in Celsius!** Make sure **NOT** to round your answer for this question!

```
[16]: min_temperatures = Table.read_table("temperatures.csv").column("Daily Min_
      ↪Temperature")

      celsius_temperature_ranges = ((max_temperatures-32) * 5/9) - ((min_temperatures_
      ↪ - 32)* 5/9)
      celsius_temperature_ranges
```

```
[16]: array([ 6.66666667, 10.          , 12.22222222, ..., 17.22222222,
      11.66666667, 11.11111111])
```

## 1.4 4. World Population

The cell below loads a table of estimates of the world population for different years, starting in 1950. The estimates come from the [US Census Bureau website](#).

```
[17]: world = Table.read_table("world_population.csv").select('Year', 'Population')
      world.show(4)
```

<IPython.core.display.HTML object>

The name `population` is assigned to an array of population estimates.

```
[18]: population = world.column(1)
      population
```

```
[18]: array([2557628654, 2594939877, 2636772306, 2682053389, 2730228104,
          2782098943, 2835299673, 2891349717, 2948137248, 3000716593,
          3043001508, 3083966929, 3140093217, 3209827882, 3281201306,
          3350425793, 3420677923, 3490333715, 3562313822, 3637159050,
          3712697742, 3790326948, 3866568653, 3942096442, 4016608813,
          4089083233, 4160185010, 4232084578, 4304105753, 4379013942,
          4451362735, 4534410125, 4614566561, 4695736743, 4774569391,
          4856462699, 4940571232, 5027200492, 5114557167, 5201440110,
          5288955934, 5371585922, 5456136278, 5538268316, 5618682132,
          5699202985, 5779440593, 5857972543, 5935213248, 6012074922,
          6088571383, 6165219247, 6242016348, 6318590956, 6395699509,
          6473044732, 6551263534, 6629913759, 6709049780, 6788214394,
          6866332358, 6944055583, 7022349283, 7101027895, 7178722893,
          7256490011])
```

In this question, you will apply some built-in Numpy functions to this array.

The difference function `np.diff` subtracts each element in an array by the element that precedes it. As a result, the length of the array `np.diff` returns will always be one less than the length of the input array.

The cumulative sum function `np.cumsum` outputs an array of partial sums. For example, the third element in the output array corresponds to the sum of the first, second, and third elements.

**Question 4.1.** Very often in data science, we are interested understanding how values change with time. Use `np.diff` and `np.max` (or just `max`) to calculate the largest annual change in population between any two consecutive years.

```
[19]: largest_population_change = max(np.diff(population))
      largest_population_change
```

```
[19]: 87515824
```

**Question 4.2.** Describe in words the result of the following expression. What do the values in the resulting array represent (choose one)?

```
[20]: np.cumsum(np.diff(population))
```

```
[20]: array([ 37311223,  79143652, 124424735, 172599450, 224470289,
          277671019, 333721063, 390508594, 443087939, 485372854,
          526338275, 582464563, 652199228, 723572652, 792797139,
          863049269, 932705061, 1004685168, 1079530396, 1155069088,
          1232698294, 1308939999, 1384467788, 1458980159, 1531454579,
          1602556356, 1674455924, 1746477099, 1821385288, 1893734081,
          1976781471, 2056937907, 2138108089, 2216940737, 2298834045,
          2382942578, 2469571838, 2556928513, 2643811456, 2731327280,
          2813957268, 2898507624, 2980639662, 3061053478, 3141574331,
          3221811939, 3300343889, 3377584594, 3454446268, 3530942729,
          3607590593, 3684387694, 3760962302, 3838070855, 3915416078,
          3993634880, 4072285105, 4151421126, 4230585740, 4308703704,
          4386426929, 4464720629, 4543399241, 4621094239, 4698861357])
```

- 1) The total population change between consecutive years, starting at 1951.
- 2) The total population change between 1950 and each later year, starting at 1951.

- 3) The total population change between 1950 and each later year, starting inclusively at 1950 (with a total change of 0).

```
[21]: # Assign cumulative_sum_answer to 1, 2, or 3
      cumulative_sum_answer = 2
```

## 1.5 5. Old Faithful

Old Faithful is a geyser in Yellowstone that erupts every 44 to 125 minutes (according to [Wikipedia](#)). People are often told that the geyser erupts every hour, but in fact the waiting time between eruptions is more variable. Let's take a look.

**Question 5.1.** The first line below assigns `waiting_times` to an array of 272 consecutive waiting times between eruptions, taken from a classic 1938 dataset. Assign the names `shortest`, `longest`, and `average` so that the print statement is correct.

```
[22]: waiting_times = Table.read_table('old_faithful.csv').column('waiting')

      shortest = min(waiting_times)
      longest = max(waiting_times)
      average = sum(waiting_times) / len(waiting_times)

      print("Old Faithful erupts every", shortest, "to", longest, "minutes and
      →every", average, "minutes on average.")
```

Old Faithful erupts every 43 to 96 minutes and every 70.8970588235294 minutes on average.

**Question 5.2.** Assign `biggest_decrease` to the biggest decrease in waiting time between two consecutive eruptions. For example, the third eruption occurred after 74 minutes and the fourth after 62 minutes, so the decrease in waiting time was  $74 - 62 = 12$  minutes.

*Hint:* You'll need an array arithmetic function mentioned in the textbook.

*Hint 2:* The function you use may report positive or negative values. You will have to determine if the biggest decrease corresponds to the highest or lowest value. Ultimately, we want to return the absolute value of the biggest decrease so if it is a negative number, make it positive.

```
[36]: biggest_decrease = abs(min(np.diff(waiting_times)))

      biggest_decrease
```

[36]: 45

**Question 5.3.** If you expected Old Faithful to erupt every hour, you would expect to wait a total of  $60 * k$  minutes to see  $k$  eruptions. Set `difference_from_expected` to an array with 272 elements, where the element at index  $i$  is the absolute difference between the expected and actual total amount of waiting time to see the first  $i+1$  eruptions. *Hint:* You'll need to compare a cumulative sum to a range.

For example, since the first three waiting times are 79, 54, and 74, the total waiting time for 3 eruptions is  $79 + 54 + 74 = 207$ . The expected waiting time for 3 eruptions is  $60 * 3 = 180$ . Therefore, `difference_from_expected.item(2)` should be  $|207 - 180| = 27$ .

```
[37]: difference_from_expected = abs ((np.cumsum(waiting_times))- (60* np.
      ↳arange(1,273)))
      difference_from_expected
```

```
[37]: array([ 19,  13,  27,  29,  54,  49,  77, 102,  93, 118, 112,
          136, 154, 141, 164, 156, 158, 182, 174, 193, 184, 171,
          189, 198, 212, 235, 230, 246, 264, 283, 296, 313, 319,
          339, 353, 345, 333, 353, 352, 382, 402, 400, 424, 422,
          435, 458, 462, 455, 477, 476, 491, 521, 515, 535, 529,
          552, 563, 567, 584, 605, 604, 628, 616, 638, 638, 670,
          688, 706, 711, 724, 746, 742, 761, 772, 774, 790, 790,
          808, 824, 847, 862, 884, 894, 899, 912, 940, 956, 976,
          964, 990, 990, 1020, 1010, 1028, 1031, 1043, 1067, 1082, 1073,
          1095, 1097, 1125, 1114, 1137, 1158, 1145, 1169, 1161, 1187, 1208,
          1223, 1222, 1251, 1270, 1269, 1290, 1280, 1305, 1304, 1331, 1324,
          1333, 1350, 1346, 1374, 1395, 1380, 1402, 1397, 1427, 1412, 1435,
          1431, 1460, 1446, 1468, 1459, 1485, 1478, 1497, 1518, 1518, 1540,
          1557, 1573, 1572, 1592, 1581, 1617, 1610, 1627, 1644, 1649, 1670,
          1681, 1691, 1712, 1745, 1738, 1767, 1752, 1778, 1776, 1794, 1800,
          1816, 1819, 1847, 1839, 1872, 1861, 1858, 1875, 1883, 1904, 1925,
          1938, 1928, 1953, 1967, 1962, 1979, 2002, 2025, 2016, 2034, 2058,
          2044, 2067, 2062, 2083, 2080, 2096, 2120, 2137, 2158, 2185, 2202,
          2193, 2211, 2211, 2233, 2264, 2257, 2275, 2261, 2278, 2302, 2291,
          2314, 2325, 2345, 2334, 2349, 2353, 2369, 2362, 2396, 2391, 2407,
          2397, 2419, 2413, 2428, 2446, 2465, 2483, 2501, 2511, 2530, 2540,
          2534, 2560, 2550, 2580, 2574, 2568, 2585, 2604, 2608, 2623, 2610,
          2636, 2639, 2664, 2686, 2683, 2705, 2712, 2726, 2720, 2743, 2756,
          2769, 2797, 2817, 2828, 2851, 2847, 2866, 2884, 2908, 2906, 2929,
          2912, 2912, 2927, 2948, 2934, 2964, 2950, 2964])
```

**Question 5.4.** If instead you guess that each waiting time will be the same as the previous waiting time, how many minutes would your guess differ from the actual time, averaging over every wait time except the first one.

For example, since the first three waiting times are 79, 54, and 74, the average difference between your guess and the actual time for just the second and third eruption would be  $\frac{|79-54|+|54-74|}{2} = 22.5$ .

```
[26]: average_error = sum (abs(np.diff(waiting_times))) / (len (waiting_times)-1)
      average_error
```

```
[26]: 20.52029520295203
```

## 1.6 6. Tables

**Question 6.1.** Suppose you have 4 apples, 3 oranges, and 3 pineapples. (Perhaps you're using Python to solve a high school Algebra problem.) Create a table that contains this information. It should have two columns: "fruit name" and "count". Give it the name fruits.

**Note:** Use lower-case and singular words for the name of each fruit, like "apple".

```
[39]: # Our solution uses 1 statement split over 3 lines.
fruits = Table().with_columns (
    'fruit name', make_array ('apple', 'orange', 'pineapple'),
    'count', make_array (4,3,3))

fruits
```

```
[39]: fruit name | count
apple      | 4
orange     | 3
pineapple  | 3
```

**Question 6.2.** The file `inventory.csv` contains information about the inventory at a fruit stand. Each row represents the contents of one box of fruit. Load it as a table named `inventory`.

```
[43]: inventory = Table.read_table ('inventory.csv')
inventory
```

```
[43]: box ID | fruit name | count
53686  | kiwi       | 45
57181  | strawberry | 123
25274  | apple      | 20
48800  | orange     | 35
26187  | strawberry | 255
57930  | grape      | 517
52357  | strawberry | 102
43566  | peach      | 40
```

**Question 6.3.** Does each box at the fruit stand contain a different fruit?

```
[47]: # Set all_different to "Yes" if each box contains a different fruit or
# to "No" if multiple boxes contain the same fruit
all_different = "No"
all_different
```

```
[47]: 'No'
```

**Question 6.4.** The file `sales.csv` contains the number of fruit sold from each box last Saturday. It has an extra column called “price per fruit (\$)” that’s the price *per item of fruit* for fruit in that box. The rows are in the same order as the `inventory` table. Load these data into a table called `sales`.

```
[50]: sales = Table.read_table("sales.csv")

sales
```

```
[50]: box ID | fruit name | count sold | price per fruit ($)
53686  | kiwi       | 3           | 0.5
57181  | strawberry | 101         | 0.2
25274  | apple      | 0           | 0.8
48800  | orange     | 35          | 0.6
26187  | strawberry | 25          | 0.15
57930  | grape      | 355         | 0.06
```



52357	strawberry	102	0.25
43566	peach	17	0.8

**Question 6.5.** How many fruits did the store sell in total on that day?

```
[53]: total_fruits_sold = sum(sales.column(2))
      total_fruits_sold
```

[53]: 638

**Question 6.6.** What was the store's total revenue (the total price of all fruits sold) on that day?

*Hint:* If you're stuck, think first about how you would compute the total revenue from just the grape sales.

```
[54]: total_revenue = sum (sales.column (2) * sales.column(3))
      total_revenue
```

[54]: 106.85

**Question 6.7.** Make a new table called `remaining_inventory`. It should have the same rows and columns as `inventory`, except that the amount of fruit sold from each box should be subtracted from that box's count, so that the "count" is the amount of fruit remaining after Saturday.

```
[57]: remaining_inventory = Table ().with_columns (
      'box ID', make_array ('53686', '57181', '25274', '48800', '26187', '57930',
      →'52357', '43566'),
      'fruit name', make_array
      →('kiwi', 'strawberry', 'apple', 'orange', 'strawberry', 'grape', 'strawberry', 'peach'),
      'count', make_array ( 45-3, 123-101, 20-0, 35-35, 255-25, 517-355, 102-102,
      →40-17 ))

remaining_inventory
```

```
[57]: box ID | fruit name | count
      53686 | kiwi       | 42
      57181 | strawberry | 22
      25274 | apple     | 20
      48800 | orange    | 0
      26187 | strawberry | 230
      57930 | grape     | 162
      52357 | strawberry | 0
      43566 | peach     | 23
```

## 1.7 7. Submission

Once you're finished, submit your assignment as a .ipynb (Jupyter Notebook) and .pdf (download as .html, then print to save as a .pdf) on the class Canvas site.