

This handout collects a bunch of practice questions for the midterm exam. *Don't worry: The actual midterm will not have this many problems, and some of the problems here would be too long for the midterm.* (The exam is one class period long – 50 minutes.)

A separate handout, which will eventually be available online, includes answers for the problems here. Do the problems before looking at the answers.

Some questions ask you to write code that produces a certain value. In that case, your code may be multiple lines (multiple statements), but the last statement in your code should be an expression whose value is the value the question asks for.

### Problem 1 Whimsical Warmups

- Write Python code that defines a function named `is_long`. The function should take a single argument, a list named `the_list`, and return a boolean value. `is_long` should return `True` if `the_list` has length greater than or equal to 42, and `False` otherwise.
- Now suppose that we have defined `is_long` as described above (even if you didn't do that part). Suppose we have also defined two lists named `some_things` and `other_things`, respectively. Write Python code that produces the value `True` if the concatenation of `some_things` and `other_things` (that is, the list you get when you append `other_things` to `some_things`) is long (that is, has length greater than or equal to 42), and `False` otherwise.
- Suppose we have a table called `students` with a single column labeled "names". We would like to add a column containing abbreviated versions of the names to this table. Let's call it "abbrev nms". To abbreviate a name, we take the first 4 characters in it, including any spaces in the name. (If it's already less than 4 characters, we take the whole thing.) Assuming `students` has already been defined with the table "names", write Python code that adds this new column (again, with name "abbrev nms") to `students`. *There is a (minor) restriction:* You may use one `for` loop in your code *or* call the Table method `apply` once, but not both.
- Write code that does the same thing as in the previous part, but do it in the way you didn't do last time – use `for` if you used `apply`, and vice versa. If you used neither `for` nor `apply` in the previous part, write code that (sensibly) makes use of a `for` loop or of the method `apply` (but not both).

#### Answer:

- ```
def is_long(the_list):
    return len(the_list) >= 42
```

If you're not comfortable with directly returning the boolean value as above (which looks weird to some people, at first), you could write something more verbose:

```
def is_long(the_list):
    if len(the_list) >= 42:
        return True
    else:
        return False
```

- (b) `is_long(some_things + other_things)`
- (c) 

```
def abbreviate(name):
    ABBREVIATION_SIZE = 4
    if len(name) <= 4:
        return name
    else:
        return name[0:4]

students["abbrv nms"] = students.apply(abbreviate, "names")
```
- (d) 

```
ABBREVIATION_SIZE = 4
abbreviated_names = []
names = students["names"]
for name_index in np.arange(students.num_rows):
    name = names[name_index]
    if len(name) <= 4:
        abbreviated_name = name
    else:
        abbreviated_name = name[0:4]
    # abbreviated_names.append(abbreviated_name) would be better style
    # here, but we don't expect you to remember the append method on the
    # exam.
    abbreviated_names = abbreviated_names + [abbreviated_name]

students["abbrv nms"] = abbreviated_names
```

## Problem 2 Python Parody

- (a) Describe, in concise English, what the following Python statement does:

```
def mystery(the_list):
    another_list = []
    for elt in the_list:
        another_list = [elt] + another_list
    return another_list
```

- (b) Describe, in concise English, what the following Python statement does:

```
def science(the_list, the_function):
    another_list = []
    for elt in the_list:
        if the_function(elt):
            another_list = another_list + [elt]
    return another_list
```

- (c) Suppose we execute both of the above statements and then this code:

```
def three_thousand(a_string):
    if a_string == "snakes":
```

```
        return False
    elif a_string == "spiders":
        return False
    else:
        return True

theater = ["watch", "out", "for", "snakes"]
mystery(science(theater, three_thousand))
```

Describe, in concise English, the value of the last expression.

**Answer:**

- (a) The Python statement *defines a function named* `mystery` that reverses a list. Given a list, the defined function produces a new list in the opposite order of the original list.
- (b) The Python statement *defines a function named* `science` that filters a list. Given a list and a function, the defined function produces a new list containing a subset of the original list. Each element is included if a given function returns `True` when called with that element as its argument.
- (c) The value is a list of strings: “for”, “out”, and “watch”, in that order. (Calling `science` on `theater` removes “snakes”, and then `mystery` reverses the remaining elements.)

### Problem 3 Sloppy Syntax

Below are several snippets of Python code, and some contain common bugs. Using your best judgement (and careful reading), determine which ones have bugs – that is, which ones don’t do what the author probably intended. For the ones with bugs, write a fixed version of the code. If there are no bugs, write “OK”.

Assume that the usual imports (for example, `import numpy as np` and `from datascience import *`) have already been executed. Also assume that a table named `marathon_data` has been created and contains columns named “Time (seconds)” and “Name”.

1. 

```
finish_times = marathon_data["Time (seconds)"]
average_finish_time = np.mean(finish_times)
```
2. 

```
marathon_data.sort("Time (seconds)", descending=False)
fastest_runner_name = marathon_data["Name"][0]
```
3. 

```
QUALIFYING_TIME = 8280
qualifying_runners = marathon_data.where("Time (seconds)" <= QUALIFYING_TIME)
qualifying_runner_names = qualifying_runners["Name"]
```
4. 

```
QUALIFYING_TIME = 8280
def create_message(finish_time):
    if finish_time <= QUALIFYING_TIME:
        return "Congratulations, you qualified for the 2016 US Olympic team!"
    else:
        return "Better luck in 2020! :-("
marathon_data["Messages"] = marathon_data.apply(create_message, "Time (seconds)")
```

**Answer:**

1. OK. (The first line assigns the name `finish_times` to an array of finish times, a column from `marathon_data`. The second line computes the average.)

```
2.    marathon_data = Table.read_table("marathons.csv")
       sorted_marathon_data = marathon_data.sort("Finish time", descending=False)
       fastest_person_name = sorted_marathon_data["Name"][0]
```

(The Table method `sort` produces a new table and doesn't modify the table we call it on. So in the original code, the second line does nothing, and the table doesn't end up sorted.)

```
3.    QUALIFYING_TIME = 8280
       qualifying_runners = marathon_data.where(marathon_data["Time (seconds)"] <= QUALIFYING_TIME)
       qualifying_runner_names = qualifying_runners["Name"]
```

(The expression in the original code `"Time (seconds)" <= QUALIFYING_TIME` compares a string to a number, which makes no sense. Instead we want to compare each element of the "Time (seconds)" column of `marathon_data` to `QUALIFYING_TIME`, get an array of booleans, and pass that to `where`.)

4. OK. (The code creates a new column in `marathon_data` with a congratulatory or consolatory message for each runner, depending on whether the runner's marathon time was short enough to qualify.)

#### Problem 4 Fortuitous Function

An analyst was working with a dataset in a table called `the_table`. The data required some manipulation before she could work with them. Specifically, `the_table` included a column called `"percent complete"` listing percentages as *strings*, but the analyst needed to work with *numerical proportions* (that is, with floating-point numbers that are proportions rather than percentages). The analyst found two kinds of strings in the data: a number followed by a percent sign (like "95%" or "0.1%"), or a number followed by " percent", like "95 percent" or "0.1 percent". She wrote a function to convert a string in either format to a *proportion* (a number) and used it to make a new column of proportions using `apply`. (For example, if she had named her function `foo`, she would have written `the_table["proportion complete"] = the_table.apply(foo, "percent complete")`.)

- (a) Describe the function you would write in this situation, using 3-4 English sentences. What is its name, what is its signature, what kind of thing does it output, and what does it do?
- (b) Write Python code that defines your function. *Hint:* The function `float` may be useful.

**Answer:** Note that our answer is just an example of a fully-correct answer. For example, you probably didn't come up with exactly the same function name as we did.

1. The function is named `convert_percentage`. It takes a single argument, a string named `percentage_string`. The argument is expected to contain a percentage followed by either "%" or " percent". The function returns a number, the proportion represented by that string.

```
2.    def convert_percentage(percentage_string):
       if percentage_string[-1] == "%":
           end_offset = len("%")
       else:
           end_offset = len(" percent")
       string_without_suffix = percentage_string[0:len(percentage_string)-end_offset]
       percentage = float(string_without_suffix)
       return percentage / 100
```

## Problem 5 Senior Sample

A simple random sample of voters taken from the U.S. voting population is classified by senior citizen status (yes, no) and political party affiliation (Republican, Democratic, other). The sampled voters' names are removed and replaced by ID numbers 1, 2, 3, etc. The data are entered into a table called `voters`, each row corresponding to one voter. The table has three columns, the first of which contains the ID numbers in increasing order and is called `ID`. The second is called `sen_cit` and contains "yes" or "no" in each row depending on whether the voter is a senior citizen or not. The third column is called `party` and contains the party affiliations "R," "D," and "O."

- Write code that produces (in any clear form) the senior citizen status and party affiliation of Sampled Voter Number 17.
- Write code that produces the party affiliations of all the senior citizens. Say whether your code produces a table, an array, or a list.
- Write code that produces the party affiliation most common among the senior citizens, and the number of senior citizens with that affiliation. You may pick any reasonable output format, but please document it. (For example, if your function produces a string that is the party affiliation and number of people concatenated together (which is not a great idea), you should write a comment that says that's what your function returns.) Be careful: if there is more than one "most common" affiliation, your code should produce the information for all of the affiliations that are tied for most common.

### Answer:

```
(a) voter_seventeen_table = voters.where(voters["ID"] == 17)
    voter_info = [voter_seventeen_table["sen_cit"][0], voter_seventeen_table["party"][0]]
    voter_info
```

```
(b) seniors = voters.where(voters["sen_cit"] == "yes")
    seniors["party"]
```

Our code produces an array. (Columns of tables are arrays.)

```
(c) seniors = voters.where(voters["sen_cit"] == "yes")
    party_counts = seniors.select(["party", "ID len"]).group("party", collect=len)
    party_counts_sorted = party_counts.sort("ID len", descending=True)
    biggest_party_count = party_counts_sorted["ID len"][0]
    biggest_parties = party_counts_sorted.where(party_counts_sorted["ID len"] == biggest_party_count)
    biggest_parties
```

Our code produces a table with a row for each most-common party affiliation; the column `party` gives the affiliation and the column `"ID len"` gives the number of seniors with that affiliation. (It's okay if you forgot that `group` appends `" len"` to the name of the `"ID"` column.)

## Problem 6 Hip Hypotheses

(This problem continues the previous problem.) The surveyors wonder whether there is any relation between party affiliation and being a senior citizen. Help them develop an answer, in the following steps.

- State null and alternative hypotheses as precisely as possible. You might want to review this part after you've done the next part, to make sure that your answers are consistent.
- In order to test your null hypothesis, what kind of statistical test will you perform, and what test statistic will you use? Justify your choices.

- (c) Write Python code that defines a function named `proportion_greater`. It should take two arguments:
- (1) a *list* of statistics (numbers); and
  - (2) a single statistic (a number).

For example, the first argument might be a list of means computed under a null hypothesis, in which case the second argument would be the mean of an observed dataset. It should return the proportion of elements in the list that are greater than the single statistic.

- (d) Write code that tests your null hypothesis, calculates an empirical  $P$ -value, and produces a conclusion (`False` if you reject the null hypothesis and `True` otherwise). You can use a 2.5% cutoff this time. As before, you are free to use any function that you have defined in this homework, but please don't just call functions that have been defined in class.

You don't have to write everything from scratch. We have provided a function to compute the test statistic for a table like `voters` for you. Our code assumes that the test statistic is the total variation distance between the distribution of party affiliations of seniors and the distribution of party affiliations of non-seniors. You may also find the function `proportion_greater`, which you just defined above, useful. (You can assume that it is implemented correctly even if you didn't do the previous part.)

```
# You can ignore this.
def normalize(table, column_name):
    table[column_name] = table[column_name] / sum(table[column_name])

# Takes a table in the same format as the voters table. Returns
# the total variation distance between the distribution of party
# affiliations of seniors and the distribution of party affiliations
# of non-seniors.
def test_statistic(sample):
    proportions = sample.pivot("sen_cit", "party", "ID", collect=len)
    proportions.relabel("yes ID", "senior")
    proportions.relabel("no ID", "non-senior")
    normalize(proportions, "senior")
    normalize(proportions, "non-senior")
    return 0.5 * sum(abs(proportions["senior"] - proportions["non-senior"]))

# Takes a table in the same format as the voters table. Returns
# False if you reject the null hypothesis (as you defined it above), and
# True otherwise.
def test_independence_hypothesis(data):
    # Fill in your hypothesis test code here. This function should
    # return a boolean value as described in the documentation comment
    # above.
```

```
test_independence_hypothesis(voters)
```

**Answer:**

- (a) Null hypothesis: Party affiliation and being a senior citizen are independent and any observed differences between the affiliations of seniors and those of non-seniors are due to random chance. Alternative hypothesis: Party affiliation and being a senior citizen are somehow related.
- (b) Under the null hypothesis, the party affiliation proportions should be about the same for seniors and non-seniors, up to random variation. So our test statistic will be the total variation distance between the distribution of affiliations for seniors and the distribution of affiliations for non-seniors.

Random variation will make the test statistic non-zero even if the null hypothesis is true. So we would like to be able to generate data under the assumption that the null hypothesis is true, and see how large the test statistic typically will be under that hypothesis. To do that, we can use a *permutation test*: Under the null, any permutation of the party affiliations has the same distribution, so we can permute the party affiliations in our data to generate new samples under the null.

We will repeatedly permute the party affiliations and compute the test statistic for each permuted data set. Then we will look at the distribution of test statistics. Qualitatively, if the test statistic we actually observed is “in the tail” of this histogram (that is, larger than most of the test statistics we would see if the null were true), then that is evidence against the null hypothesis. Quantitatively, we could compute an empirical  $P$ -value – the proportion of test statistics larger than the one we observed – and decide to reject the null if that  $P$ -value is smaller than, say, 0.025.

- (c)
- ```
def proportion_greater(statistics_under_null, actual_statistic):
    num_greater_simulated_statistics = np.count_nonzero(np.array(statistics_under_null) > actual_statistic)
    return num_greater_simulated_statistics / len(statistics_under_null)
```
- (d)
- ```
def test_independence_hypothesis(data):
    # Fill in your hypothesis test code here.
    n = data.num_rows
    actual_statistic = test_statistic(data)
    statistics_under_null = []
    NUM_SIMULATIONS = 10000
    simulation_data = data.select(["party", "sen_cit", "ID"])
    for simulation_index in np.arange(NUM_SIMULATIONS):
        permuted_party = simulation_data.select("party").sample(n, with_replacement=False)["party"]
        simulation_data["party"] = permuted_party
        simulated_test_statistic = test_statistic(simulation_data)
        statistics_under_null = statistics_under_null + [simulated_test_statistic]
    p = proportion_greater(statistics_under_null, actual_statistic)
    return p >= 0.025
```

```
test_independence_hypothesis(voters)
```

## Problem 7 Bin Boon

Suppose we would like to make a histogram of some data. Figuring out good bins for a histogram can be hard, so we decide to use the bins we used last time we analyzed a similar dataset. But our existing bins don't quite cover all of our data; the maximum is too low and the minimum is too high. We would like to modify the bins so that they do.

We are going to define a function called `stretch_bins`. This function shifts and rescales an increasing array of numbers to have a new minimum and maximum, stretching the middle values proportionally. It takes three arguments:

- (1) an array (of length at least 2) containing floating-point numbers in increasing order (the kind of thing we might pass as bins to the `hist` function, like `np.array([-1.0, 1.5, 2.0, 4.0])`);
- (2) a number, the new minimum value of the array; and
- (3) a number, the new maximum value of the array.

`stretch_bins` returns a new array whose first value is the new minimum value and whose last value is the new maximum value. The array is still in increasing order, and the entries in the middle are rescaled so that the ratios of the distances between consecutive numbers are the same as in the original array.

For example, `stretch_bins(np.array([-1.0, 1.5, 2.0, 4.0]), 1.0, 11.0)` should return an array equal to `np.array([1.0, 6.0, 7.0, 11.0])`.

- (a) What is the value of `stretch_bins(np.arange(-1.0, 4.0, 1.0), -5.0, 15.0)`?
- (b) What is the value of `stretch_bins(np.array([3.0, 4.0, 6.0, 7.0]), -4.0, -2.0)`?
- (c) Write Python code that defines `stretch_bins`.

**Answer:**

- (a) `np.array([-5.0, 0.0, 5.0, 10.0, 15.0])`
- (b) `np.array([-4.0, -3.5, -2.5, -2.0])`
- (c)

```
def stretch_bins(bins, new_min, new_max):
    new_range = new_max - new_min
    old_range = bins[-1] - bins[0]
    rescaling_factor = new_range / old_range
    rescaled_bins = rescaling_factor * bins
    shift_amount = new_min - rescaled_bins[0]
    shifted_bins = rescaled_bins + shift_amount
    return shifted_bins
```