

Name: Dan Schumacher

abc123: hdd249

Remember, it is okay to talk to your classmates about the homework, but do not share solutions. You should complete the homework on your own. Below, please list all of your peers ("collaborators") that you discussed the homework with.

Colaborators:

- Name 1
- Name 2
- Name 3

Problem Set 1 (PS 1)

Complete the 4 functions as directed to:

- validate the format of a string
- compute probabilities and store them in a lookup table (dictionary)
- generate a string deterministically based on a given string

Each function is worth 1 to 4 points, for a total of 10 points. Other assignments will be more complicated, and worth more points, for this I just want to make sure you are in the correct course. I care mainly about correctness of the implementation, but reserve the right to deduct points for code that is difficult to read.

This assignment will test your skills at implementing algorithms in Python from scratch. If this is your first experience with traditional programming assignments, then it may take some time to figure everything out. Please schedule time to meet with me if you need help. I recommend you start early. For others with significant programming experience, it will not take long.

Submission Instructions

After completing the exercises below, generate a PDF *or* HTML of the code **with** outputs. Then, create a zip file containing both the completed exercise and the generated PDF/HTML. You are **required** to check the PDF/HTML to make sure all the code **and** outputs are clearly visible and easy to read. If your code goes off the page, you should reduce the line size. I generally recommend not going over 80 characters.

Finally, name the zip file using a combination of the assignment ID and your name, e.g., ps1_rios.zip. Submit your homework on Blackboard.

Preliminary Information

For the homework exercises, you will need to write code in the specified functions. The variable specified in the open and close parenthesis "my_function(**variable_name**)" of the function should **NOT** be reassigned. They will hold important information---such as the string you should process---which is required to complete each task. Your job is to write code that returns the required value for each exercise.

Please refer to the example below to understand how functions work in python:

```
def my_function(my_variable):
    # Your code goes here
    print("The parameter is {}".format(my_variable))
    new_variable = my_variable + 100
    return new_variable

# my_variable in my_function() will be assigned 42.
return_value = my_function(42) # return value will be assigned the value
given by the "return" in my_function()
print("The function returned the value {}".format(return_value))
```

If you run the code above, you will get the following output:

```
The parameter is 42
The function returned the value 142
```

Please note that you do **NOT** need to create new functions. The functions have already been specified. You simply need to fill them out with code.

Exercise 1 (4 points)

Checks whether the string is a valid employee ID. An employee ID is valid if and only if it consists only of 6-10 alphabetic characters (letters), followed by 2 numeric digits. Implement the validation **without** any external libraries (i.e., do not use the "re" package). You may use the .isalpha() and .isdigit() methods, variables, standard data structures such as lists, if statements, loops, etc.

```
In [1]: # Example use of .isalpha() and .isdigit(). This code is just to help, you can delete
print("32".isdigit()) # if the string is a digit, .isdigit() will return True, otherwise
print("lkjsdaf".isdigit()) # Here is an example of .isdigit() returning False
print("lskdjf".isalpha()) # In this case we can check if all characters are "alphabetic"
print("ab1".isalpha()) # Here is an example of .isalpha() resulting in a False
print(len("ALKJDSLKFJDS")) # You can check how many characters are in a string by using len()

True
False
True
False
12
```

to summarize

1. 1-6 (up to 10) alpha
2. end with 2 digits
 - My plan, is to check last 2 in string for numeric. Pop them then check len() and alpha of the rest

```
In [2]: # Write code here. Do NOT delete this cell.
def validate(s):

    # RULE 1: Last 2 chars must be digits
    rule1 = s[::-1][0:2].isdigit()

    # We are done with those digits so Lets chop them (is this bad practice?)
    stop = len(s) - 2
    s = s[0:stop]

    # RULE 2: Length of remainder must be between 6 and 10 (inclusive)
    if len(s) >= 6 and len(s) <= 10:
        rule2 = True
    else:
        rule2 = False

    # RULE 3: remainder of s must be alpha!
    rule3 = s.isalpha()

    #check that all rules are true and return a T/F value
    if rule1 and rule2 and rule3:
        return True
    else:
        return False
```

The lines below give example inputs and correct outputs using asserts, and can be run to test the code. Passing these tests is necessary, but **NOT** sufficient to guarantee your implementation is correct. You may add additional test cases, but do not remove any tests.

```
In [3]: #teacher asserts
assert(validate('AbcdEf00') == True)
assert(validate('$0RQLpCHz49') == False)

#ones I added the check
assert(validate('$asldf5hlaksjdh!!') == False)
assert(validate('AaBb123') == False)
assert(validate('$AaBb!!123') == False)
assert(validate('$AaBbCc123') == False)
assert(validate('$AaBbCcd123') == False)
assert(validate('$0RQLpCHz49') == False)

print("Asserts Completed Successfully")
```

Asserts Completed Successfully

Exercise 2 (3 points)

Given a sequence of the DNA bases {A, C, G, T}, stored as a string, returns a probability table in a data structure such that one base can be looked up to get the probability p(base). More

specifically, write code to return a dictionary such that the value for each key is the probability of that key occurring in the string. For example, given the string "abb", $p(a)$ is equal to $\frac{1}{3}$ and $p(b)$ is $\frac{2}{3}$. Simply, divide the number of times each character appears by the total number of characters. So, the code should return the following dictionary:

```
{"a": 0.33333, "b": 0.66666}
```

You may use the collections module, but no other libraries.

```
In [4]: def dna_prob1(seq):
        # make sure that capitalization is consistent
        seq = seq.upper()

        # construct a dictionary blank DNA dictionary
        dnaDict = {'A':0, 'C':0, 'G':0, 'T':0}

        # iterate through the string and count a,c,g, and t values
        for character in seq:
            if character == 'A':
                dnaDict['A'] +=1

            if character == 'C':
                dnaDict['C'] +=1

            if character == 'G':
                dnaDict['G'] +=1

            if character == 'T':
                dnaDict['T'] +=1

        # turn dictionary values into ratios.
        for key in dnaDict:
            dnaDict[key] = dnaDict[key]/len(seq)

        #return our ratio-dictionary
        return(dnaDict)
```

The lines below give example inputs and correct outputs using asserts, and can be run to test the code. Passing these tests is **NOT** sufficient to guarantee your implementation is correct. You may add additional test cases, but do not remove any tests.

```
In [5]: tbl = dna_prob1('ATCGATTGAGCTCTAGCG')
        assert(tbl['T'] == 5/18)
        assert(tbl['G'] == 5/18)
        assert(tbl['C'] == 4/18)
        print("Asserts Completed Successfully")
```

Asserts Completed Successfully

Exercise 3 (3 points)

Given a string representing a sequence of DNA bases, returns the paired sequence, also as a string, where A is always paired with T and C with G, i.e., replace A with T, T with A, ...

Do not use any libraries.

Hint: this can be done in one line. (More than one line is okay too.)

```
In [6]: def dna_bp(seq):
        #dictionary
        conversion = {'A':'T',
                      'T':'A',
                      'C':'G',
                      'G':'C'}

        #for each item in seq, use dictionary above to grab opposite dna base.
        dna2 = ''.join([conversion[character] for character in seq])

        return(dna2)
```

The lines below give example inputs and correct outputs using asserts, and can be run to test the code. Passing these tests is **NOT** sufficient to guarantee your implementation is correct. You may add additional test cases, but do not remove any tests.

```
In [7]: assert(dna_bp('ATCGATTGAGCTCTAGCG') == 'TAGCTAACTCGAGATCGC')
        print("Asserts Completed Successfully")
```

Asserts Completed Successfully

```
In [8]: #to do this in one line... although I think it is ugly
        def dna_bp(seq): return ''.join(['A':'T', 'T':'A', 'C':'G', 'G':'C'][character] for char
```

```
In [9]: assert(dna_bp('ATCGATTGAGCTCTAGCG') == 'TAGCTAACTCGAGATCGC')
        print("Asserts Completed Successfully")
```

Asserts Completed Successfully

Extra Credit (2 points)

Given a sequence of the DNA bases {A, C, G, T}, stored as a string, returns a conditional probability table in a data structure such that one base (b1) can be looked up, and then a second (b2), to get the probability $p(b2 | b1)$ of the second base occurring immediately after the first. (Assumes the length of seq is ≥ 3 , and that the probability of any b1 and b2 which have never been seen together is 0. Ignores the probability that b1 will be followed by the end of the string character.)

Here is an example:

GAAAGG

$P(A | G)$ = probability of A given the letter G occurred directly before it = $1/2 = 0.5$. The denominator is 2 instead of 3 because the last G does not have a letter after it. The denominator should be the number of times a letter appears with another letter directly after it.

$P(G | G)$ = probability of G given the letter G occurred directly before it = $1/2 = 0.5$

$P(G | A)$ = probability of G given the letter A occurred directly before it = $1/3 = 0.33$

$P(A | A)$ = probability of A given the letter A occurred directly before it = $2/3 = .66$

The dictionary returned by dna_prob2 should be:

```
{ 'A': { 'A': 2/3, 'G': 1/3, 'T': 0, 'C': 0 }, 'G': { 'A': 1/2, 'G': 1/2, 'T': 0, 'C': 0 }, 'T': { 'A': 0, 'T': 0, 'G': 0, 'C': 0 }, 'C': { 'A': 0, 'G': 0, 'T': 0, 'C': 0 } }
```

Notice this is a "nested" dictionary, where the value of the outer dictionary is also a dictionary. All of the keys are strings.

Note: 2 extra points adds 20% to your final homework grade. This can improve your grade substantially. So, I recommend working on this, if you are worried about the quizzes and Midterm.

```
In [10]: def dna_prob2(seq):
    dnaDict = { 'A': { 'A': 0, 'G': 0, 'T': 0, 'C': 0 },
                'G': { 'A': 0, 'G': 0, 'T': 0, 'C': 0 },
                'T': { 'A': 0, 'T': 0, 'G': 0, 'C': 0 },
                'C': { 'A': 0, 'G': 0, 'T': 0, 'C': 0 }

    for x,y in zip(seq, seq[1::]):
        dnaDict[x][y] += 1

    # turn those counts into probabilities.
    for bigkey in dnaDict:

        # the total number of A,G,T, or C respectively
        total = sum(dnaDict[bigkey].values())

        for smallkey in dnaDict[bigkey]:
            # the proceeding of A,G,T, or Cs
            count = dnaDict[bigkey][smallkey]

            #can't divide by zero!
            if total > 0:
                #update that slot to be the ratio rather than count
                dnaDict[bigkey][smallkey] = count/total

    return dnaDict
```

The lines below give example inputs and correct outputs using asserts, and can be run to test the code. Passing these tests is **NOT** sufficient to guarantee your implementation is correct. You may add additional test cases, but do not remove any tests.

```
In [11]: tbl = dna_prob2('ATCGATTGAGCTCTAGCG')
assert(tbl['T']['T'] == 0.2)
assert(tbl['G']['A'] == 0.5)
assert(tbl['C']['G'] == 0.5)
print("Asserts Completed Successfully")
```

Asserts Completed Successfully