

assignment1

March 21, 2021

1 Assignment 1

For this assignment you are welcomed to use other regex resources such a regex “cheat sheets” you find on the web.

Before start working on the problems, here is a small example to help you understand how to write your own answers. In short, the solution should be written within the function body given, and the final result should be returned. Then the autograder will try to call the function and validate your returned result accordingly.

```
[5]: def example_word_count():  
    # This example question requires counting words in the example_string below.  
    example_string = "Amy is 5 years old"  
  
    # YOUR CODE HERE.  
    # You should write your solution here, and return your result, you can  
    →comment out or delete the  
    # NotImplementedError below.  
    result = example_string.split(" ")  
    return len(result)  
  
    #raise NotImplementedError()
```

1.1 Part A

Find a list of all of the names in the following string using regex.

```
[7]: import re  
def names():  
    simple_string = """Amy is 5 years old, and her sister Mary is 2 years old.  
    Ruth and Peter, their parents, have 3 kids."""  
  
    pattern = re.compile("[A-Z]\w*")  
    names = pattern.findall(simple_string)  
  
    return names  
    # YOUR CODE HERE  
    raise NotImplementedError()
```

```
[8]: assert len(names()) == 4, "There are four names in the simple_string"
```

1.2 Part B

The dataset file in `assets/grades.txt` contains a line separated list of people with their grade in a class. Create a regex to generate a list of just those students who received a B in the course.

```
[22]: import re
def grades():
    with open ("assets/grades.txt", "r") as file:
        grades = file.read()
        pattern = re.compile("[-\w ]*:\sB")
        matches=pattern.findall(grades)

        student_B = []

        for i in range(len(matches)):
            student_B.append(matches[i][:3])

        return student_B
        # YOUR CODE HERE
        raise NotImplementedError()
```

```
[20]: assert len(grades()) == 16
```

1.3 Part C

Consider the standard web log file in `assets/logdata.txt`. This file records the access a user makes when visiting a web page (like this one!). Each line of the log has the following items: * a host (e.g., '146.204.224.152') * a user_name (e.g., 'feest6811' **note: sometimes the user name is missing! In this case, use '-' as the value for the username.**) * the time a request was made (e.g., '21/Jun/2019:15:45:24 -0700') * the post request type (e.g., 'POST /incentivize HTTP/1.1' **note: not everything is a POST!**)

Your task is to convert this into a list of dictionaries, where each dictionary looks like the following:

```
example_dict = {"host": "146.204.224.152",
                "user_name": "feest6811",
                "time": "21/Jun/2019:15:45:24 -0700",
                "request": "POST /incentivize HTTP/1.1"}
```

```
[11]: import re
def logs():
    with open("assets/logdata.txt", "r") as file:
        logdata = file.read()

    logs = []
    #to pick host pattern
    pattern = re.compile("(?:[0-9]{1,3}\.){3}[0-9]*\d")
```

```

host = pattern.findall(logdata)

#to pick user_name pattern
pattern = re.compile("-\s+[\w-]*")
user_name = pattern.findall(logdata)

#to pick time pattern
pattern = re.compile("\[(.*?)\]")
time = pattern.findall(logdata)

#to pick request pattern
pattern = re.compile("\"(.*?)\"")
request = pattern.findall(logdata)

#to clean some unwanted character and whitespace
user_name_list = []
for x in user_name:
    user_name_list.append( x.replace("-", "", 1).replace(" ", ""))

for i in range(len(user_name_list)):
    one_item = {'host': host[i], 'user_name': user_name_list[i], 'time':
→time[i], 'request': request[i]}
    logs.append(one_item)

return logs
raise NotImplementedError()

```

```
[12]: assert len(logs()) == 979
```

```

one_item={'host': '146.204.224.152',
  'user_name': 'feest6811',
  'time': '21/Jun/2019:15:45:24 -0700',
  'request': 'POST /incentivize HTTP/1.1'}
assert one_item in logs(), "Sorry, this item should be in the log results,
→check your formating"

```

```
[ ]:
```