

Mappings

Introduction

In this lab, we are going to write Python code to read the complete works of William Shakespeare and then report the word frequencies of the words used in his works.

Objectives

The purpose of this lab is to help you:

1. gain familiarity with class inheritance
2. learn the concept of doubling array size
3. implement and use the basic operations of a mapping (dictionary)
4. learn how to handle a file and practice some string operations
5. learn what to do when we need to sort a dictionary according to the values
6. learn how to implement hash functions and in particular good ones

The Complete Works of William Shakespeare

The text file that has the complete works of William Shakespeare is provided with the skeleton code. This text file has a special format. That is, every word or symbol is followed by a space. This makes the job of splitting words easier. This text file has all the works of Shakespeare, and it is very long. The size of the text file is about 4.5 Mbytes. The total number of lines is 129,107, the total number of words and symbols is 980,637, and the total number of characters is 4,538,523.

Here is the beginning of the text file:

A MIDSUMMER-NIGHT'S DREAM

Now , fair Hippolyta , our nuptial hour
Draws on apace : four happy days bring in
Another moon ; but O ! methinks how slow
This old moon wanes ; she lingers my desires ,
Like to a step dame , or a dowager
Long withering out a young man's revenue .

Four days will quickly steep themselves in night ;
Four nights will quickly dream away the time ;

And then the moon , like to a silver bow
New-bent in heaven , shall behold the night
Of our solemnities .

Go , Philostrate ,
Stir up the Athenian youth to merriments ;
Awake the pert and nimble spirit of mirth ;
Turn melancholy forth to funerals ;
The pale companion is not for our pomp .

Hippolyta , I woo'd thee with my sword ,
And won thy love doing thee injuries ;
But I will wed thee in another key ,
With pomp , with triumph , and with revelling .

We will read this text file and count the number of times each of the words and symbols appear in the text file.

Part 1: Solving the problem using Python dictionary (in-lab)

We first want to solve the problem using Python dictionary. To this end you are required to write a function called `getTokensFreq` that receives as a parameter a file name and returns a dictionary that includes all words as keys and their frequencies as the values.

Further, you are required to write a function called `getMostFrequent` that receives as parameters a dictionary `d` and the number of required frequent tokens `k` . The function returns a list of tuples of the `k` frequent key-value pairs in the dictionary.

Initially, read the text file into a string and split the string into symbols and words. And then loop through the list of symbols and words and use a dictionary to count the number of times these symbols and words appeared in the list. After that, use sorted function to sort the dictionary according to the values and save the result into a list. Lastly, return the top `k` frequently used tokens and their corresponding frequencies.

Example:

Assume that the following text is stored in `f.txt` :

"I felt happy because I saw the others were happy and because I knew I should feel happy, but I wasn't really happy."

```
d = getTokensFreq('f.txt')
```

```

print(d)
# prints {'i': 5, 'felt': 1, 'happy': 4, 'because': 2, 'saw': 1,
# 'the': 1, 'others': 1, 'were': 1, 'and': 1, 'knew': 1, 'should': 1,
# 'feel': 1, ',': 1, 'but': 1, 'was': 1, 'not': 1, 'really': 1, '.': 1}

freq = getMostFrequent(d, 5)
print(freq)
# prints [('i', 5), ('happy', 4), ('because', 2), ('felt', 1), ('saw', 1)]

```

Part 2: Solving the problem using HashMapping

Your job now is to implement a class `HashMapping` and use this class to replace the dictionary above and repeat the above operations and obtain the same results. Note that you are provided with classes `Entry`, and `ListMapping`. You will need these two to implement the class `HashMapping`.

The `Entry` class is to store an entry of key-value pair. It has two attributes, `key` and `value`. `ListMapping` stores all the entries in a list. Its important methods are to get a value associated to a given key and to set a new value to a given key.

The `ListMapping` takes linear time because it has to iterate through the list. We could make this faster if we had many short lists instead of one large list. Then, we just need to have a quick way of knowing which short list to search or update.

We're going to store a list of `ListMappings`. The size of this list will be `1000`. For any key `k`, we want to compute the index of the right `ListMapping` for `k`. We often call these `ListMapping`s buckets. This term goes back to the idea that you can quickly group items into buckets. Then, when looking for something in a bucket, you can check all the items in there assuming there aren't too many.

This means, we want an integer, i.e. the index into our list of buckets. A hash function takes a key and returns an integer. Most classes in python implement a method called `__hash__` that does just this. We can use it to implement a mapping scheme that improves on the `ListMapping`. This is called `HashMapping` which you will be implementing.

Below is a screenshot of results for the top 20 words. Using dictionary, we get the job done in 0.6 seconds. Using `HashMapping`, it takes around 7.5 seconds.

```
1 ('the', 26804)
2 ('and', 24037)
3 ('i', 20041)
4 ('to', 18532)
5 ('of', 16006)
6 ('you', 13833)
7 ('a', 13678)
8 ('my', 12256)
9 ('that', 10718)
10 ('in', 10524)
11 ('is', 9138)
12 ('not', 8450)
13 ('me', 7757)
14 ('it', 7736)
15 ('for', 7538)
16 ('with', 7141)
17 ('be', 6840)
18 ('your', 6744)
19 ('this', 6584)
20 ('his', 6528)
0.6034021377563477
```

```
-----
1 ('the', 26804)
2 ('and', 24037)
3 ('i', 20041)
4 ('to', 18532)
5 ('of', 16006)
6 ('you', 13833)
7 ('a', 13678)
8 ('my', 12256)
9 ('that', 10718)
10 ('in', 10524)
11 ('is', 9138)
12 ('not', 8450)
13 ('me', 7757)
14 ('it', 7736)
15 ('for', 7538)
16 ('with', 7141)
17 ('be', 6840)
18 ('your', 6744)
19 ('this', 6584)
20 ('his', 6528)
7.486926794052124
```

Let's see if you can do a better job implementing `HashMapping` than this.

Part 3: Shakespeare Tokens

In this part you are required to develop a class, called `ShakespeareToken`, that represents a Shakespeare token, which is really just a word from his text. Every token is a string so we can define `ShakespeareToken` class to inherit from Python string class which is called `str`. So the definition of your class will look like:

```
class ShakespeareToken(str):  
    pass
```

This minimal definition of `ShakespeareToken` can be used in the same way as a string. For example, the following code works:

```
[>>> class ShakespeareToken(str):  
[...     pass  
[...  
[>>> s = ShakespeareToken("Hello")  
[>>> len(s)  
5  
[>>> s[2]  
'l'  
[>>> s[2:5]  
'llo'  
[>>> hash(s)  
5125980898720904098  
[>>> t = ShakespeareToken("Hello")  
[>>> s == t  
True  
[>>> u = "Hello"  
[>>> s == u  
True  
[>>> id(s)  
4400775960  
[>>> id(t)  
4400776080  
>>> █
```

You can see that the hash method and the related comparison operator work as well.

Note: If you do not know why the above code works without having this functionality defined explicitly in `ShakespeareToken` class please learn more about inheritance before proceeding further.

Your job in this part is to develop your own hash method for the `ShakespeareToken` class. To do that you need to override the `__hash__` method. In this part the hash method that you are required to define is a “bad” one where it returns the length of the underlying token.

```
s = ShakespeareToken("Hello")
t = ShakespeareToken("Hello")
u = "Hello"

hash(s)      ## returns 5
hash(t)      ## returns 5
```

Part 4:

Write a second implementation of `ShakespeareToken` called `ShakespeareToken2` that improves the `__hash__` method. This method will just add up the values of the letters (o for 'a', 1 for 'b' etc.) in the given string and mod it by m , where m is the bucket size, and return this value.

Let the value of m be $10^9 + 9$.

```
s = ShakespeareToken2("Hello")
print(hash(s))    ## prints 47
```

Next, write a third implementation of `ShakespeareToken` called `ShakespeareToken3` that further improves the `__hash__` method. This method will use the following formula to calculate the hash of a string.

$$\begin{aligned}\text{hash}(s) &= s[0] + s[1] \cdot p + s[2] \cdot p^2 + \dots + s[n-1] \cdot p^{n-1} \mod m \\ &= \sum_{i=0}^{n-1} s[i] \cdot p^i \mod m,\end{aligned}$$

where m is the number of buckets and p is a prime number. Refer to the **Calculation of the hash of a string** section of the following link for more information: <https://cp-algorithms.com/string/string-hashing.html>

Part 5

Will be updated at the earliest.