

FIT1008 – Intro to Computer Science

Assessed Prac 3 – Weeks 11 and 12

Semester 1, 2018

Objectives of this practical session

To be able to implement and use hash tables in Python.

Note:

- You should provide documentation and testing for each piece of functionality in your code. Your documentation needs to include pre and post conditions, and information on any parameters used.
- Create a new file/module for each task or subtask.
- You are not allowed to use the Python built-in dict methods.
- Name your files task[num] to keep them organised.

Testing

For this prac, you are required to write:

- (1) a function to test each function you implement, and
- (2) at least two test cases per function.

The cases need to show that your functions can handle both valid and invalid inputs. There is no need to test menu functions.

Marks

For this assessed prac, there are a total of 30 marks. There are 10 marks allocated to your understanding of the solutions and your implementations in the prac overall. In addition to these, the marks for each task are listed with the tasks. A marking rubric is available online for you to know and understand how you will be marked.

Task 1 [7 marks]

Implement a complete version of a hash table using Linear Probing to resolve collisions. Include implementations for the following 4 functions:

- `__getitem__(self, key)`: Returns the value corresponding to key in the hash table. Raises a `KeyError` if key does not exist in the hash table. Called by `table[key]`
- `__setitem__(self, key, value)`: Sets the value corresponding to key in the hash table to be value. Raise an exception if the hash

table is full and the key does not exist in the table yet. Called by
`table[key] = value`

- `__contains__(self, key)`: Returns True if key is in the table and False otherwise.
- `hash(self, key)`: Calculates the hash value for the given key. Use the hash function given below.

```

1 def hash_value(self, key):
2     a = 101
3     h = 0
4     for c in key:
5         h = (h * a + ord(c)) % self.table_size
6     return h

```

Task 2 [6 marks]

Download from Moodle the dictionary files `english_small.txt`, `english_large.txt` and `french.txt`.

For each of these dictionaries, time how long it takes to read all words in it into the hash table (i.e. *wall time*). Do this for each of the following `table_size` values: 210000, 209987, 400000 and 399989 and 202361. Present the wall times recorded in a table.

Write a short analysis reporting what values work best and which work poorly. Explain why these might be the case ¹

Task 3 [3 marks]

Modify your hash table implementation to now track the number of collisions, the load, as well as the average probe length. For the latter, it is advisable to track the total probe length in an instance variable. The average probe length is the total probe length divided by the number of items on the table.

Using collisions, probe length and wall time, choose appropriate values of *a* (in your hash function) and `table_size`. You want to find values that perform well across all three files. For this task use a maximum table size of 400000.

You should try at least 10 values, and explain your choice by presenting all data behind your reasoning recorded in a table ².

¹ You can use Ctrl+c to stop execution of your program in case some combination of values takes too long. In total you should consider 15 possibilities (3 files for each table size). In this task you use the word as both, key and data.

² Just for fun we'll collect the best performer in each lab class

CHECKPOINT

(You should reach this point during week 10)

Task 4 [4 marks]

Modify your hash table from the previous tasks to:

- use Quadratic Probing to resolve collisions.
- implement dynamic hashing, by doubling the size of the underlying array (and rehashing) every time the load exceeds $\frac{2}{3}$

Compare the number of collisions, probe length and running time found when loading each dictionary against the best combination found using the Linear Probing hash table.

Task 5 [4 marks]

Implement a hash table which uses Separate Chaining to resolve collisions. It is a good idea to first implement and test the Linked List separately. Compare the performance of Separate Chaining against the linear probe above.

Background

In most large collections of written language, the frequency of a given word in that collection is inversely proportional to its rank in the words. That is to say that the second most common word appears about half as often as the most common word, the third most common word appears about a third as often as the most common word and so on³.

If we count the number of occurrences of each word in such a collection, we can use just the number of occurrences to determine the approximate rank of each word. Taking the number of occurrences of the most common word to be max and the relationship described earlier, we can assume that any word that appears at least $\text{max}/100$ times appears in the top 100 words and is therefore common.

The same can be applied as $\text{max}/1000$ times for the next 900 words, rounding out the top 1000 words, considered to be uncommon. All words that appear less than $\text{max}/1000$ times are considered to be rare. In this prac we have been implementing hash tables and so we will use one here to facilitate a frequency analysis on a given text file.

³ This is known as Zipf's law. You can read more at https://en.wikipedia.org/wiki/Zipf%27s_law

Task 6 [6 marks]

Download some ebooks as text files from <https://www.gutenberg.org/> and use these files as a collection of English. Read these into your Quadratic Probing hash table to count the number of each word in the texts. Considering the data you collected in Task 2, select an appropriate table size for the text files. Use these occurrence counts to construct a table which can be used to look up whether a given word is common, uncommon or rare within written English.

Task 7 – Just for fun [0 marks]

Implement the program in Task 3 using Python's dict class. How does the performance of your own implementations compare to Python's in terms of the wall time? What factors may drive the differences? How close can you get to the reference implementation?