

CS/IT 200

Lab 8: Autocomplete

Recommended Due Date: Thursday, November 5, 11:59pm

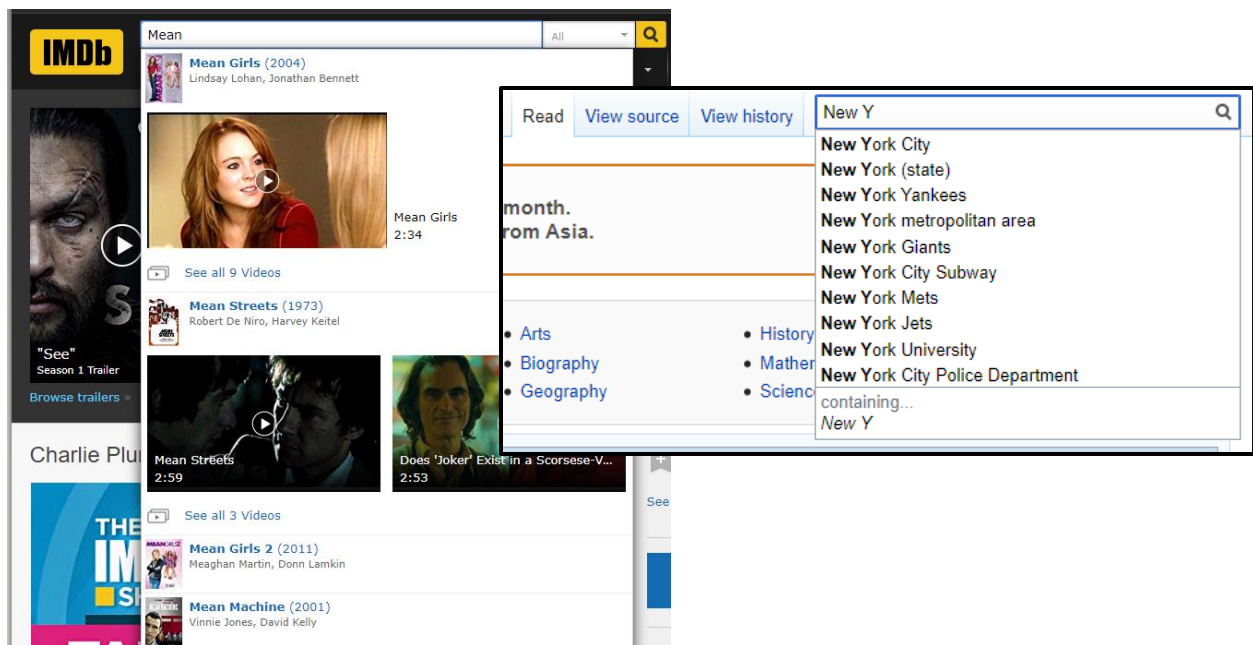
Submission

- Submit all code in a single zip file to the Lab 8 assignment folder on Kodiak.

Goal: Use binary search trees to write a program to implement autocomplete for a given set of N strings and positive weights. That is, given a prefix, find all strings in the set that start with the prefix, and present them in descending order of weight.

Problem

Autocomplete is an important feature of many modern applications. As the user types, the program predicts the complete *query* (typically a word or phrase) that the user intends to type. Autocomplete is most effective when there are a limited number of likely queries. For example, the Internet Movie Database uses it to display the names of movies as the user types; search engines use it to display suggestions as the user enters web search queries; cell phones use it to speed up text input.



In these examples, the application predicts how likely it is that the user is typing each query and presents to the user a list of the top-matching queries, in descending order of weight. These weights are determined by historical data, such as box office revenue for movies, frequencies of search queries from other Google users, or the typing history of a cell phone user. For the purposes of this assignment, you will have access to a list of all possible queries and associated weights (and these queries and weights will not change).

Acknowledgement: This assignment is adapted from *Autocomplete Me*, created by Matthew Drabick and Kevin Wayne, Princeton University, and posted on EngageCSEdu (submitted by Ananda Gunawardena).

The performance of autocomplete functionality is critical in many systems. For example, consider a search engine which runs an autocomplete application on a server farm. According to one study, the application has only about 50ms to return a list of suggestions for it to be useful to the user. Moreover, in principle, it must perform this computation *for every keystroke typed into the search bar, and for every user!* (You won't have to handle each individual keystroke.)

In this assignment, you will implement autocomplete by using binary search trees to find the set of queries that start with a given prefix; and sorting the matching queries in descending order by weight.

Part I – Autocomplete Term

In `autocomplete.py`, write a class `Term` that represents an autocomplete term: a string query and an associated positive real-valued (float) weight. You must implement the following methods for the `Term` class:

- A constructor that takes a query and a weight as arguments. The constructor should raise a `TypeError` if query is `None`. It should raise a `ValueError` if the weight is negative.
- A `get_weight()` method that returns the weight.
- The `__eq__` method, which should return `True` if the two `Term`'s queries are equal.
- The `__lt__` method (which represents the `<` operator), which should return `True` if this `Term`'s query comes before the other `Term`'s query in lexicographic order. (Lexicographic order is based on the ASCII order of characters. That is, `A < Z < a < z`. Lexicographic order is the default ordering for Python strings.)
- The `__str__` method, which should return a string containing this `Term`'s weight and query.

Part II – Autocomplete

In this part, you will implement the functions that provide autocomplete functionality using the `Term` class and a binary search tree. To do so, you will add each `Term` to the binary search tree (BST), use the tree's methods to find the set of `Terms` that start with a given prefix, and sort the results in descending order by weight.

You will organize your program into the following functions. Add these functions to `autocomplete.py`. (They should not be part of the `Term` class.)

- `build_tree(filename)` which takes a filename for the file containing weights and queries. This function should create a `Term` for each query/weight and add them to a `TreeMap` (our BST from class). Return that BST.
 - Our input files are described below. Because they are in Unicode format, you will need to open them with the following command:

```
file = open(filename, 'r', encoding='utf-8')
```

- `all_matches(tree, prefix)` which takes a BST and a string prefix. The function should return a list of `Terms` in the tree that start with the given prefix. That list should be sorted by the `Term`'s weight in descending order (highest weights first).
 - Your solution should not loop over the entire tree. Carefully study the functions provided in `TreeMap` that might help with this task.
 - Python's built-in `sort()` function for lists can be told how to sort a list by using the `key` keyword argument. You can also specify descending sort using the `reverse` keyword argument. If you followed the naming conventions above, you should be able to sort a list of `Terms` (`termlist`) by their weight in descending order as follows:

```
termlist.sort(key=Term.get_weight, reverse=True)
```

- `main()` should do the following:
 - Ask the user for a file name and build the tree from that file.
 - Ask the user for a number of results to show.
 - Ask the user for a search term and display the appropriate number of results. The program should continue this step until the user enters a blank search term (by just pressing enter).
 - Time how long each query took to complete. (This should only be based on how long it took for `all_matches()` to execute.)

```
import timeit
start = timeit.default_timer()
# Code that you want to time
end = timeit.default_timer()
total_time_in_sec = end-start
# Note this is seconds, not ms
```

Input Format: We provide two sample input files for testing. Each file consists of an integer `N` followed by `N` pairs of query strings and positive weights. There is one pair per line, with the weight and string separated by a tab. The query strings are in Unicode and can contain any Unicode characters (including spaces).

- The file `wiktionary.txt` contains the 10,000 most common words in Project Gutenberg, with weights equal to their frequencies.
- The file `cities.txt` contains nearly 100,000 cities, with weights equal to their populations.

Here are a few sample runs from our programs (user input in blue):

Enter file: <code>wiktionary.txt</code> Enter number of results to show: <code>5</code> Enter search: <code>auto</code> 619695.0 automobile 424997.0 automatic Query took 0 ms	Enter file: <code>cities.txt</code> Enter number of results to show: <code>7</code> Enter search: <code>M</code> 12691836.0 Mumbai, India 12294193.0 Mexico City, Distrito Federal, Mexico
---	---

Enter search: <code>comp</code> 13315900.0 <code>company</code> 7803980.0 <code>complete</code> 6038490.0 <code>companion</code> 5205030.0 <code>completely</code> 4481770.0 <code>comply</code> Query took 0 ms Enter search: <code>the</code> 5627187200.0 <code>the</code> 334039800.0 <code>they</code> 282026500.0 <code>their</code> 250991700.0 <code>them</code> 196120000.0 <code>there</code> Query took 0 ms Enter search:	10444527.0 <code>Manila, Philippines</code> 10381222.0 <code>Moscow, Russia</code> 3730206.0 <code>Melbourne, Victoria, Australia</code> 3268513.0 <code>Montréal, Quebec, Canada</code> 3255944.0 <code>Madrid, Spain</code> Query took 150 ms Enter search: <code>Al M</code> 431052.0 <code>Al Maḥallah al Kubrá, Egypt</code> 420195.0 <code>Al Maṣṣūrah, Egypt</code> 290802.0 <code>Al Mubarras, Saudi Arabia</code> 258132.0 <code>Al Mukallā, Yemen</code> 227150.0 <code>Al Minyā, Egypt</code> 128297.0 <code>Al Manāqil, Sudan</code> 99357.0 <code>Al Maṭariyah, Egypt</code> Query took 0 ms Enter search:
---	--

Part III – Autocomplete with a Hash Table

Make a copy of `autocomplete.py` and save it as `autocomplete2.py`. In this copy, make the following changes:

- Change `build_tree()` to `build_map()`. Instead of returning a binary search tree, it should return a `ProbeHashMap` (our linear probing hash table from the `hash_tables.py` module) containing the same data as before. Make sure you update any calls to `build_tree()` to say `build_map()`.
- Re-write `all_matches()` so that it uses the hash table of `Terms` instead of the tree.

Part IV – Data Collection

Select eight search queries to run on both versions of our autocomplete code. Use the `cities.txt` input for these searches. Make a table that shows the query, how long it took in our BST version of our algorithm, and how long it took with the hash table version our algorithm.

Answer the following questions:

1. Overall, which version of our algorithm was faster, the one that used `TreeMap` (the BST) or `ProbeHashMap` (the hash table)?
2. Is this what you expected? Why or why not?
3. Is there any alternative to `TreeMap` or `ProbeHashMap` that you would expect to be superior to both? If so, why?

What to Submit:

- Your code (`autocomplete.py`, `autocomplete2.py` and any other necessary modules, including modules given as class materials)
- Your table and answers from Part IV

Rubrics

For this assignment, I've broken rubrics down by part more clearly. The overall rubric is at the top, while rubrics for each part follow.

Overall Rubric

Grade	Requirements
4	A 4 in Part II and Part IV; At least a 3 in other parts
3	At least a 3 in Part II and Part IV; At least a 2 in other parts
2	At least a 2 in Part II; At least a 1 in other parts
1	At least a 1 in all parts
0	

Part I Rubric

Grade	Requirements
4	All functions in Term behave as expected, including raising errors when expected.
3	Incorrectly raises error, or minor bug in <code>__eq__</code> or <code>__lt__</code> methods.
2	Two of the bugs listed in (3)
1	Multiple bugs, but the class is usable.
0	Syntax or other errors prevent creation and use of the Term class.

Part II Rubric

Grade	Requirements
4	Correctly builds BST; efficiently locates matching terms and returns them in sorted order; <code>main()</code> function behaves as expected, including user input and timing
3	One of the following issues: BST key/value is incorrect; Finds matching terms inefficiently; Returns correct results but not sorted; Shows incorrect number of results; Does not allow for repeated queries
2	Two of the issues above, or incorrect results from <code>all_matches()</code>
1	Up to four of the issues above possibly including incorrect results from <code>all_matches()</code>
0	Either function does not complete due to errors

Part III Rubric

Grade	Requirements
-------	--------------

4	build_map() and all_matches() correctly updated to use a hash table.
3	Minor error in one function
2	--
1	Major error in one function
0	Code cannot be used to draw comparisons with Part II results

Part IV Rubric

Grade	Requirements
4	All required data is present and questions are answered correctly.
3	All required data is present with one error in answers
2	All required data is present with multiple errors in answers
1	All required data present without questions answered, or missing data for one version with correct questions.
0	Missing data and incorrect questions