

# cse13s asgn7 WRITEUP.pdf

Lucas Lee; CruzID: luclee

3/11/2022

## 1 Program Details

This program will use Hash tables, Bloom Filters, and Bitvectors to find k most likely authors to a given text.

## 2 Program issues, Hash table/Bloom filter analysis

My program produces the same order of authors as the example binary when I run the program with the normal database, but it runs very slowly compared to the example. This is because my program uses an insertion method rather than the heap method which makes my program run at  $O(n)$ . Other than this, my program's excessively long run time could be because I may be performing extra hash table lookups throughout the process of computing the distances between the texts and/or computing the Texts.

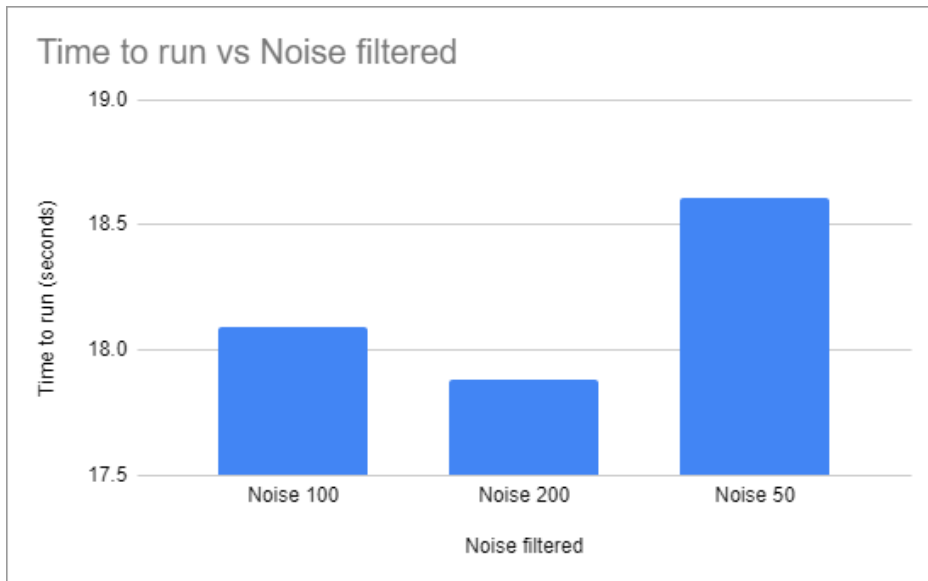
In terms of the smaller databases, the program still runs a bit slow, but the difference is much less than the example binary. For the large databases, it takes about 3 minutes and a half, while running it on the small database is about 20 seconds to run. The accuracy on smaller text files is a bit off, but the order is generally correct with the larger database. There are some files in the small database that sometimes display the wrong ordering, but the only thing that it changes is that it swaps two different texts. The distances for these texts may have something to do with my lookup function, as when I would change the function for testing purposes, it would change the distance to be closer between those two texts. I believe this is when you run the small database to display 10 matches, and the two swapped authors are Henry James and H. G. Wells. Otherwise I have not seen any differences with the result ordering for my program and the example binary.

In the specifics of this program, we had to use bloom filters alongside the hash tables, and what made this program run much slower was the linear probing that we had to implement inside of our hash tables to look for specific words. Since we are hashing the words somewhat randomly throughout the hash table that has  $2^{19}$  indices, if a text is larger then the linear probing will take a long time. The bloom filter only speeds up the process of lookups in the event that a word is not inside of the hash table, otherwise it has to perform another hash table lookup, which makes the program run longer. In larger texts, this is going to make the process much slower, since many more bits are being set in the bloom filter for each word, there is a higher probability for false positives since the bit vector will be more full of set bits. This could also be why the program's run time is extremely long for the larger databases.

### 3 Results when changing noise words

When changing the noise limit, the speed at which the identify binary changes, as well as many of the results.

If you increase the noise limit to around 200, meaning you filter out double the amount of words, the program runs a bit faster than it did before but not by a large difference. In the small database it makes it run a fraction of a second faster because it does not have to check the texts for as many words. For the default database, it decreases the run time by a couple seconds due to the text having much less words since there are much more texts to be created.



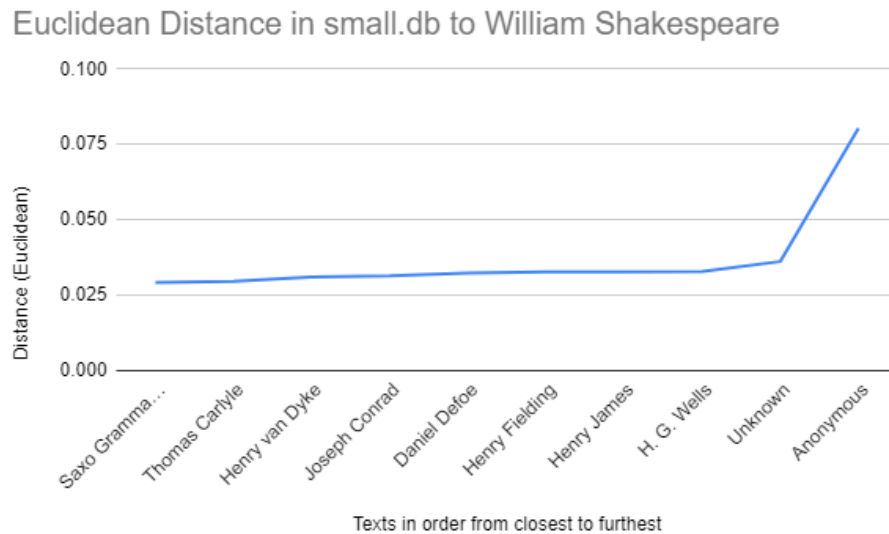
If you decrease the noise limit to about 50, you are leaving in half as many words that would normally be filtered out, making the program run slower. For the smaller database, this does not matter as much, most likely because the texts in the smaller database use the 50 words, that are not filtering out, very rarely. The run time changes for the small database by about half a second, while for the larger database, it runs for about an extra 15 seconds.

### 4 Comparing large vs small texts

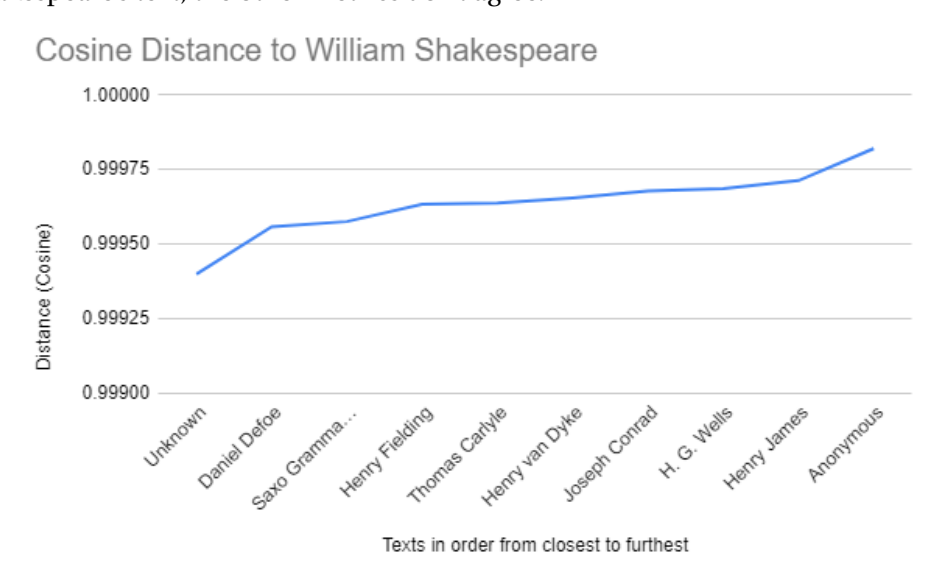
My program correctly is able to identify the closest authors using small texts. For extremely large texts, the amount of words is so large that the numbers calculated for the distances are very close. The distance values are sometimes extremely close when comparing two different large texts because of the amount of words in those texts. In a large text, the text creation process will take much longer because of the hash tables and bloom filters for those texts. In the Program issues section, the reason for this longer run time and creation is stated for these larger texts. For smaller texts, this will not matter as much because there are much less set bits inside of the bit vector for the bloom filter to check, leading it to give false positives less often.

## 5 Comparing the different metrics

In my program, all three of the metrics correctly generate the exact author, if given an author that exists in that database. Otherwise, my program, for the most part (in program issues) mimics the output from the example binary for each metric. Running my program using the parameters in the Assignment document, I am given the same list of top 10 authors for William Shakespeare. In the next couple figures, they will use the small database instead to compare and contrast results.

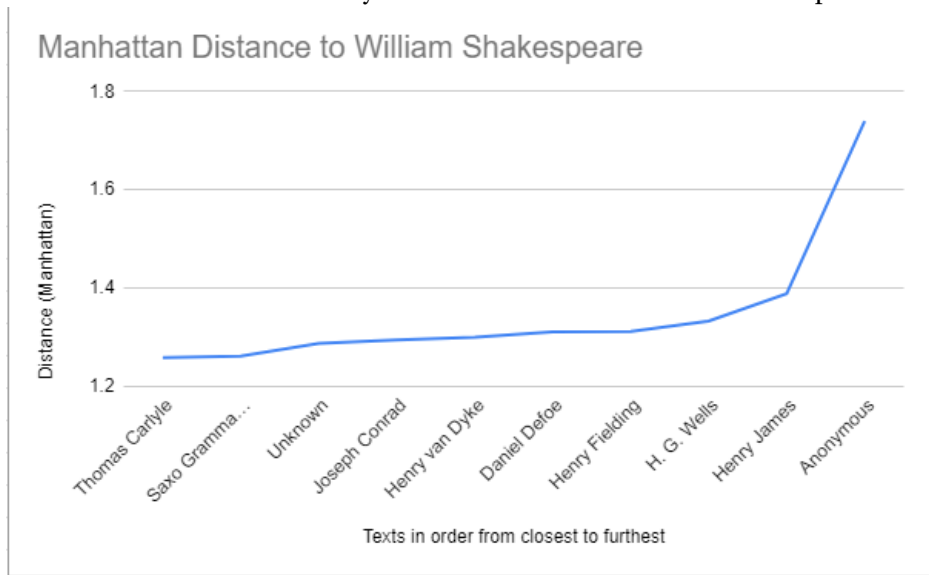


In the above graph is the nearest neighbors in small.db using the Euclidean distance from William Shakespeare's text, which matches the order of the example binary. While the Euclidean distance says that Saxo Grammaticus is the most likely author out of the small.db authors to be the one to write William Shakespeare's text, the other metrics don't agree.



The above graph shows the Cosine distance of the texts in small.db to William Shakespeare's text. The calculations essentially generate a larger sum based on the similarity of the texts and subtract that

sum from 1. In the case for this calculation, the numbers are much higher than those of the Euclidean distance, and they also produce a different result. According to the Cosine distance calculation, the Unknown author is the most likely author to be the one to write Shakespeare's text.



The above graph shows the Manhattan distance of the texts in small.db to Shakespeare's text once again. These calculations are even larger, since the calculation simply gets the sum of the difference between vectors. This metric claims that Thomas Carlyle is the most likely author of Shakespeare's text. All 3 of these metrics agree on which text is the most unlikely text, which is the anonymous text, but only Manhattan and Euclidean agree that the author is most likely either Thomas Carlyle or Saxo Grammaticus, since their values are extremely close in distance. Overall, all three of these distance calculations give different values and approximate the distance very differently.

## 6 Other Comments

This program took a long time to compile in Valgrind, even for the small database, so I did not check this program with the large database through Valgrind, but I would assume that if the program is able to run without memory leaks in the smaller database, that it would not cause a problem checking for leaks in a larger database. I also tried to decrease the floating point round off errors using long doubles instead of doubles, but since the specified return values for the given functions are in doubles, I was not able to decrease the rounding error by a lot. While there is still a difference in the values for my program and the example binary executable, the results in terms of most likely authors is mostly the same.