

Data Mining Assignment 1

Zihao Xu
zihxu@kth.se

Minchong Li
mincli@kth.se

Date: November 8, 2023

1 Short Summary of our Solutions

1.1 Introduction

This homework is aimed at implementing the stages of finding textually similar documents based on Jaccard similarity using the shingling, minHashing, LSH (locality-sensitive hashing) techniques and corresponding algorithms.

The homework is implemented with Python.

1.2 Dataset Construction

We use 2 datasets for the assignment.

The first dataset is part of the text from the davisWiki corpus for this assignment. We selected 10 files for this assignment, which were retrieved based on the following keywords respectively: "zombie", "restaurant", "computer".

Another dataset is a series of texts related generated by ChatGPT. These texts are generated by telling ChatGPT to express the same concept in different ways. These texts are more semantically related.

Dataset can be download from [here](#).

1.3 Shingling

The Shingling part takes text as input and returns k-shingles (where we set $k = 5$) as a Set. In this part, in order to minimize conflicts, we use built-in Python hash functions for mapping, and the

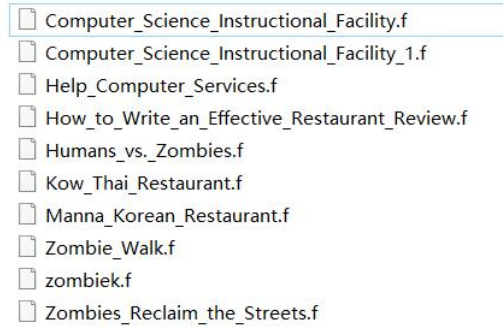


Figure 1: Selected files from davisWiki corpus

mapped results are further processed to reduce the sparsity of the results. Detailed implementation can be seen in `shingling.py`

1.4 Compare Sets

We calculate the Jaccard similarity of different sets of k-shingles. Detailed implementation can be seen in `compareSets.py`

1.5 minHashing

The minHashing part generates a vector (length defined by the user) by applying the min-hashing algorithm with a group of random hash functions. Detailed implementation can be seen in `minHashing.py`

1.6 Compare Signatures

This part compares signatures similarity to estimate the Jaccard similarity between two signatures representing their original documents. Detailed implementation can be seen in `compareSignatures.py`

1.7 LSH (locality-sensitive hashing)

The LSH part takes as input the signature matrix M representing documents, a similarity threshold t , and the number of bands. This function effectively generates pairs of documents with a similarity greater than t and returns them in a set format. The function identifies candidate pairs

by partitioning the hashed signatures into bands and treats two documents as candidates in case of collisions. Detailed implementation can be seen in `lsh.py`

2 How to Build and Run

Firstly, **download the dataset** and unzip it into the assignment1 root.

Secondly, install the requirements:

```
pip install -r requirements.txt
```

Then, use the following command to run the program:

```
cd ./src/  
python main.py ../corpus --k 5 --n 500 -m 1000000007 -t 0.01 -b 5  
python main.py ../corpus_test --k 5 --n 500 -m 1000000007 -t 0.1 -b 20
```

Arguments explanation:

- required positional argument: path of corpus.
- `--k`: Length of a single shingle (k-shingle)
- `--n`: Number of hash functions in signature
- `-m` / `--mod`: Maximum number of hash
- `-t` / `--threshold`: Threshold of being similar
- `-b` / `--band`: Band width in LSH

3 Results

Results of davisWiki corpus:

```
Jaccard sim of 5 shinglings:  
[('Kow_Thai_Restaurant.f' & Manna_Korean_Restaurant.f', 0.14955127765979312), ('Help_Computer_Services.f' & Manna_Korean_Restaurant.f', 0.12891655951817013), ('Help_Computer_Services.f' &  
Kow_Thai_Restaurant.f', 0.12051337829018925), ('Help_Computer_Services.f' & How_to_Write_an_Effective_Restaurant_Review.f', 0.1118691918597812), ('How_to_Write_an_Effective_Restaurant_  
Review.f' & Manna_Korean_Restaurant.f', 0.1067111296373489), ('How_to_Write_an_Effective_Restaurant_Review.f' & Kow_Thai_Restaurant.f', 0.10430157261794634), ('Computer_Science_Instructi  
onal_Facility.f' & How_to_Write_an_Effective_Restaurant_Review.f', 0.08535926526202053), ('Computer_Science_Instructional_Facility.f' & Help_Computer_Services.f', 0.08534024947511425), (  
'Computer_Science_Instructional_Facility.f' & Kow_Thai_Restaurant.f', 0.08036474164133739), ('Computer_Science_Instructional_Facility.f' & Manna_Korean_Restaurant.f', 0.07803121248499399  
) , ('Help_Computer_Services.f' & Humans_vs._Zombies.f', 0.0754584908301834), ('Humans_vs._Zombies.f' & Kow_Thai_Restaurant.f', 0.0750310902307586), ('Computer_Science_Instructional_Facil  
ity.f' & Humans_vs._Zombies.f', 0.07184399589462881), ('How_to_Write_an_Effective_Restaurant_Review.f' & Humans_vs._Zombies.f', 0.07126294635956099), ('Help_Computer_Services.f' & Zombies_  
_Reclaim_the_Streets.f', 0.0688385269121813), ('Computer_Science_Instructional_Facility.f' & Zombies_Reclaim_the_Streets.f', 0.0670147954743255), ('Humans_vs._Zombies.f' & Zombie_Walk.f'  
, 0.0667010943780582), ('Zombies_Reclaim_the_Streets.f' & Zombie_Walk.f', 0.06656017039403621), ('Humans_vs._Zombies.f' & Zombies_Reclaim_the_Streets.f', 0.06463233735961009), ('Kow_Thai_  
_Restaurant.f' & Zombies_Reclaim_the_Streets.f', 0.06358703214136635), ('How_to_Write_an_Effective_Restaurant_Review.f' & Zombies_Reclaim_the_Streets.f', 0.060899437851342914), ('Humans_  
vs._Zombies.f' & Manna_Korean_Restaurant.f', 0.05742662982600254), ('Manna_Korean_Restaurant.f' & Zombies_Reclaim_the_Streets.f', 0.056979198070545675), ('Computer_Science_Instructional_  
Facility.f' & Zombie_Walk.f', 0.05487437185929648), ('Kow_Thai_Restaurant.f' & Zombie_Walk.f', 0.05425631431244153), ('Help_Computer_Services.f' & Zombie_Walk.f', 0.05343028770154916), ('  
How_to_Write_an_Effective_Restaurant_Review.f' & Zombie_Walk.f', 0.05008880994671403), ('Manna_Korean_Restaurant.f' & Zombie_Walk.f', 0.040126632370399686), ('zombie.f' & Zombies_Reclaim_  
_the_Streets.f', 0.02827586206896516), ('Computer_Science_Instructional_Facility.f' & zombie.f', 0.02648845686512758), ('Kow_Thai_Restaurant.f' & zombie.f', 0.025375268048606146), ('H  
ow_to_Write_an_Effective_Restaurant_Review.f' & zombie.f', 0.024737945492662474), ('Humans_vs._Zombies.f' & zombie.f', 0.023301608139153267), ('Help_Computer_Services.f' & zombie.f', 0  
.022294725394236), ('zombie.f' & Zombie_Walk.f', 0.01793504604944256), ('Manna_Korean_Restaurant.f' & zombie.f', 0.01405857740585774)]
```

Figure 2: Result of davisWiki corpus

Results of texts generated by ChatGPT:

```

Signatures sim:
[('Kow_Thai_Restaurant.f & Manna_Korean_Restaurant.f', 0.168), ('Help_Computer_Services.f & Kow_Thai_Restaurant.f', 0.13), ('Help_Computer_Services.f & How_to_Write_an_Effective_Restaurant_Review.f', 0.126), ('Help_Computer_Services.f & Manna_Korean_Restaurant.f', 0.124), ('How_to_Write_an_Effective_Restaurant_Review.f & Kow_Thai_Restaurant.f', 0.116), ('How_to_Write_an_Effective_Restaurant_Review.f & Manna_Korean_Restaurant.f', 0.112), ('Computer_Science_Instructional_Facility.f & Help_Computer_Services.f', 0.104), ('Computer_Science_Instructional_Facility.f & Kow_Thai_Restaurant.f', 0.098), ('Computer_Science_Instructional_Facility.f & Manna_Korean_Restaurant.f', 0.098), ('Computer_Science_Instructional_Facility.f & How_to_Write_an_Effective_Restaurant_Review.f', 0.086), ('Manna_Korean_Restaurant.f & Zombies_Reclaim_the_Streets.f', 0.082), ('Computer_Science_Instructional_Facility.f & Humans_vs_Zombies.f', 0.08), ('Help_Computer_Services.f & Humans_vs_Zombies.f', 0.078), ('Humans_vs_Zombies.f & Kow_Thai_Restaurant.f', 0.078), ('How_to_Write_an_Effective_Restaurant_Review.f & Humans_vs_Zombies.f', 0.074), ('Computer_Science_Instructional_Facility.f & Zombies_Reclaim_the_Streets.f', 0.074), ('Help_Computer_Services.f & Zombies_Reclaim_the_Streets.f', 0.068), ('Kow_Thai_Restaurant.f & Zombies_Reclaim_the_Streets.f', 0.068), ('Humans_vs_Zombies.f & Zombies_Reclaim_the_Streets.f', 0.068), ('Humans_vs_Zombies.f & Manna_Korean_Restaurant.f', 0.066), ('How_to_Write_an_Effective_Restaurant_Review.f & Zombies_Reclaim_the_Streets.f', 0.054), ('Humans_vs_Zombies.f & Zombie_Walk.f', 0.054), ('Zombies_Reclaim_the_Streets.f & Zombie_Walk.f', 0.052), ('Kow_Thai_Restaurant.f & Zombie_Walk.f', 0.05), ('Computer_Science_Instructional_Facility.f & Zombie_Walk.f', 0.048), ('Help_Computer_Services.f & Zombie_Walk.f', 0.044), ('How_to_Write_an_Effective_Restaurant_Review.f & zombie.f', 0.04), ('How_to_Write_an_Effective_Restaurant_Review.f & Zombie_Walk.f', 0.038), ('Manna_Korean_Restaurant.f & Zombie_Walk.f', 0.034), ('Computer_Science_Instructional_Facility.f & zombie.f', 0.034), ('Humans_vs_Zombies.f & zombie.f', 0.03), ('Help_Computer_Services.f & zombie.f', 0.03), ('Kow_Thai_Restaurant.f & zombie.f', 0.028), ('zombie.f & Zombies_Reclaim_the_Streets.f', 0.026), ('Manna_Korean_Restaurant.f & zombie.f', 0.022), ('zombie.f & Zombie_Walk.f', 0.016)]

```

Figure 3: Result of davisWiki corpus

```

time cost for calculating shingling and minHashing: 63.227
time cost for LSH: 7.199

```

Figure 4: Time cost of davisWiki corpus

```

Jaccard sim of 5 shinglings:
[('1.txt & 8.txt', 0.97), ('1.txt & 7.txt', 0.8253968253968254), ('7.txt & 8.txt', 0.8170347003154574), ('2.txt & 7.txt', 0.3844282238442822), ('1.txt & 2.txt', 0.37411764705882355), ('2.txt & 8.txt', 0.36046511627906974), ('4.txt & 5.txt', 0.2680851063829787), ('5.txt & 6.txt', 0.2463186077643909), ('4.txt & 6.txt', 0.24178712220762155), ('2.txt & 3.txt', 0.23563218390804597), ('3.txt & 7.txt', 0.22072936660268713), ('1.txt & 3.txt', 0.2033271719038817), ('3.txt & 8.txt', 0.2029520295202952), ('3.txt & 4.txt', 0.030534351145038167), ('3.txt & 5.txt', 0.029754204398447608), ('3.txt & 6.txt', 0.02418379685610641), ('2.txt & 4.txt', 0.023415977961432508), ('2.txt & 6.txt', 0.0196078431372549), ('1.txt & 6.txt', 0.016817593790426907), ('6.txt & 8.txt', 0.016795865633074936), ('2.txt & 5.txt', 0.016736401673640166), ('4.txt & 7.txt', 0.01662049861495845), ('1.txt & 4.txt', 0.016282225237449117), ('4.txt & 8.txt', 0.016260162601626018), ('6.txt & 7.txt', 0.015810276679841896), ('1.txt & 5.txt', 0.015193370165745856), ('5.txt & 8.txt', 0.015172413793103448), ('5.txt & 7.txt', 0.012658227848101266)]

Signatures sim:
[('1.txt & 8.txt', 0.978), ('1.txt & 7.txt', 0.826), ('7.txt & 8.txt', 0.826), ('2.txt & 7.txt', 0.376), ('1.txt & 2.txt', 0.354), ('2.txt & 8.txt', 0.344), ('4.txt & 5.txt', 0.282), ('5.txt & 6.txt', 0.258), ('4.txt & 6.txt', 0.252), ('2.txt & 3.txt', 0.246), ('3.txt & 7.txt', 0.228), ('3.txt & 8.txt', 0.224), ('1.txt & 3.txt', 0.222), ('3.txt & 5.txt', 0.024), ('3.txt & 4.txt', 0.02), ('2.txt & 4.txt', 0.018), ('3.txt & 6.txt', 0.016), ('4.txt & 7.txt', 0.016), ('4.txt & 8.txt', 0.014), ('1.txt & 4.txt', 0.014), ('2.txt & 6.txt', 0.014), ('5.txt & 8.txt', 0.012), ('1.txt & 5.txt', 0.012), ('5.txt & 7.txt', 0.01), ('2.txt & 5.txt', 0.01), ('6.txt & 8.txt', 0.008), ('1.txt & 6.txt', 0.008), ('6.txt & 7.txt', 0.008)]

Candidate pairs:{('1.txt', '8.txt'), ('7.txt', '8.txt'), ('1.txt', '7.txt'), ('2.txt', '7.txt')}
Sim pairs:{('1.txt', '8.txt'), ('7.txt', '8.txt'), ('1.txt', '7.txt'), ('2.txt', '7.txt')}

```

Figure 5: Result of ChatGPT generated texts

```

time cost for calculating shingling and minHashing: 4.736
time cost for LSH: 0.549

```

Figure 6: Time cost of ChatGPT generated texts