



Data Mining Homework 2

Alessandro Monteleone 1883922

Academic Year: 2023/2024

Contents

1	Introduction	2
2	Problem 1	2
3	Problem 2	2
4	Problem 3	3
4.1	Shingling	3
4.2	MinwiseHashing	3
4.3	LSH	3
4.4	CompareWithLSH	3
4.5	CompareWithJaccard	3
4.6	Report	4
5	Problem 4	4
5.1	CompareWithLSHSpark	5
5.1.1	Shingling	5
5.1.2	Minwise-hashing	5
5.1.3	LSH	5
5.2	CompareWithJaccardSpark	5
5.3	Report	6
6	Guide to use the application	6
7	Images	6

1 Introduction

I have decided to realize this homework as a flask application in order to have a better view on the results. Now I will explain one by one all the solutions to the problems, at the end of the document there will be images taken from the application showing the results and there will also be a little guide to use the application. All the python files are placed in the folder named "server"

2 Problem 1

The functions involved in this peroblem are in `amazon_product_retrieval.py`, `amazon_product_analysis.py` and in `utils_and_classes.py`.

I have decided to use the library `lxml` to do scraping because it works well also with bigger data, this is not the case but it might be useful in future.

My code does the following for the part of retrieving data:

- it iterates on the number of pages to retrieve (I have chosen 10) and each time sends a request for a page
 - if the request is refused it tries again to retrieve the page
- if the response is ok then it finds all the products in the page looking for a certain class in the elements of the html
- for each item it builds an object `Product` which upon creation search inside the html of the element looking for the classes of the interesting informations (description,url,number of reviews, number of stars,primeness) and writes the informations in a tsv file (`amazon_products_gpu.tsv`)

All the classes used for doing scraping are saved as constants inside `utils_and_classes.py` since amazon could decide to change them and this would make this code not working anymore, instead in this way it is sufficient to change the constants to make it working.

For the analysis part let's say something about every requested feature:

- for the price ranges i have plotted the maximum and the minimum price for all the products which include in the description the word describing the category. For the categories I have chosen the series of the gpus like 3080,3090 etc.
- To evaluate the products throught their ratings I have chosen to compute a score based on the ratings times the logarithm of the number of reviews, in this way products with lots of reviews are not advantaged too much
- I have analyzed the primeness looking at the mean price and the mean ratings of the product with it and the ones without it and I have noticed that there is not a big difference between them
- for the top 10 products in term of price or ratings I have simply plotted the table ordered by the required term

All the result of the analysis are visible in the menu "DEFAULT QUERIES"

3 Problem 2

For this problem the code is inside `index_building.py`, `amazon_product_analysis.py` and in `utils_and_classes.py`. For building the index in my code I iterate on the descriptions, for each description I use a class for preprocess the description, this class removes the punctuation and the stopwords and tokenize the description. After preprocessing the list of token is iterated to count the number of times that each word appear in the list, when the iteration finishes the dictionary word:num of appearences is stored in a list since it is useful to compute the TFIDF later. Then for each word in the dictionary built in the iteration it appends the tuple (product id, num of appearences of the word) to the list of the same word in the inverted index. It also computes the normas of each description, to do so it iterates on the list of word:num of appearences computed in previous step and for each description compute the TFIDF of the words (using also tthe inverted index), sums them squared and computes the squared root of the result and put the final result in a list. In this way I can get the norma of a description simply indicating the id of the product.

Now it stores the inverted index and the normas in a json in order to not to have to recompute it every time.

When it is required to compute cosine similarity the index and the normas are retrieved and the similarity is computed as follows:

- the query is preprocessed as the descriptions were
- for each word in the query using the inverted index it is retrieved each description that contains it and it is placed inside a dictionary as id description: cos similarity, each time that the same description appears in the list of any of the words of the query its cosine similarity is updated
- when the iteration on the query has finished all the descriptions with their cosine similarities are put inside an heap
- the number of required descriptions is then retrieved from the heap

The cosine similarity function is available in the menu "QUERY MENU"

4 Problem 3

The code for this problem is inside `hashing_techniques.py` and `utils_for_hashing.py`.

I have realized this problem using 5 classes, I will explain each one.

4.1 Shingling

this class is used to compute the shingles of a string, it takes in input an hash function and the length of the shingles (as k) to be initialized, then it exposes a function to compute the shingles. In this function it iterates on the string that takes in input and for each slice of length k it applies the hash function to it and put the result in a list. At the end of the iterations it returns the list.

4.2 MinwiseHashing

This class is built giving in input to it an integer t and a list of hash functions. The class exposes a function that given a list of shingles returns a list of hashed values of length t, in fact for t times it apply one of the hash functions (possibly always different) to every shingle of the list taking the minimum at each iteration and putting it in another list. After t iterations it returns the list of hashed values.

4.3 LSH

This class is built giving in input to it 2 integers r and b and a list of hash functions. This class exposes a function that given a list of elements represented as list of objects of length b*r (if the length is different then it raises an error) returns a set of pairs which are the nearest according to LSH. The function iterates on fragments of the lists, these fragments has length r and therefore there are b fragments for each element, when it is considering the i-th fragment it iterates on the elements, it takes the i-th fragment, hashes it and puts it in a dictionary of kind hash:list of elements with same hash, if there are elements inside the list for the hash of the actual element for each pair actual element, element in list it builds a tuple in which the lower element is always the first in the pair and then add the tuple to the resulting set. In this way there are no duplicates in the result and it is returned when the iterations finishes.

4.4 CompareWithLSH

This class takes in input an instance of each of the previous classes to be initialized. This class simply put together all the 3 above and keep trace of the execution time.

4.5 CompareWithJaccard

This class takes in input an instance of the class Shingles and a value for the threshold to be initialized. The class exposes a method for computing the nearest pairs of elements according to the Jaccard similarity and the given threshold. this method is preatty easy, as first it computes the shingles then it start iterating on all the possible pairs of elements, for each pair it compute the jaccard similarity and if it is higher than the threshold then the pair is put in the result set. At the end the function returns the result set.

The function also keep trace of the execution time.

4.6 Report

This is another class but it is made only to execute the comparing classes and to build a report with the required informations. The fundamental things to say for this class are the coiches for the parameters and the results of the experiment.

- for the hash functions it was used the family showed in the questions file, each function was generated with a different integer
- the values for b and r that I have chosen are respectively 7 and 10, that's because I have tried lots of values in Geogebra and it seems to me that these are the best for the given threshold, the resulting approximation function is showed in figure [9]

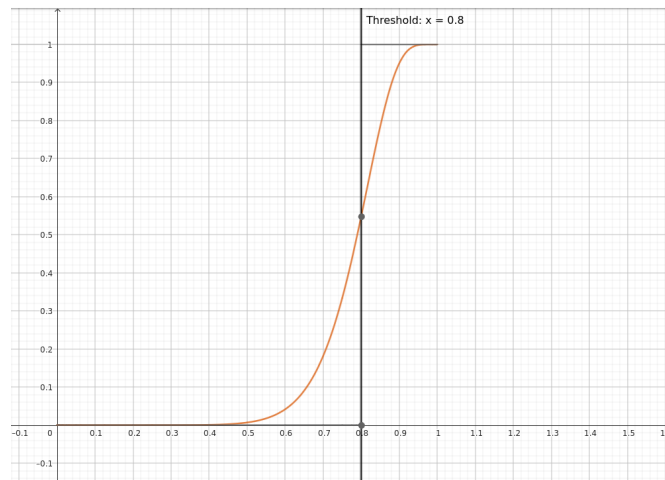


Figure 1: LSH approximation

For which concerns the results I have found that LSH finds 2 false positives and 0 false negative which is pretty good but it is surely due by the kind of data, in first instance the descriptions of the products are short, they can't be compared with real documents, then a lot of products are duplicates, this is due to the fact that often amazon puts some products at the beginning of the page since they are recommended but they are not removed from their original place. Also the execution time is wasted by the bad data, in fact the LSH should be faster since it avoids to compare a lot of pairs of documents but in this case it is the opposite, this is probably caused also by the usage of cryptographic hash functions that are slower than normal hash functions. By the way the bare results are the following:

- Number of duplicates for LSH: 442
- Number of duplicates for Jaccard similarity: 440
- Size of intersection between Jaccard and LSH results: 440
- Time required for the execution of LSH algorithm: 5574
- Time required for the execution of Jaccard similarity algorithm: 1047

All the results of the report are also available in the menu "NEAREST PRODUCTS REPORTS" and the results of the nearest descriptions both for LSH and Jaccard similarity are available in the same menu

5 Problem 4

All the code for this problem is inside `hashing_techniques_with_spark.py`, `hashing_techniques.py` and `utils_for_hashing.py`. I have realized 2 classes to solve the problem, I will explain them separately

5.1 CompareWithLSHSpark

This class is built for comparing the descriptions using LSH and it is implemented with Spark. To be instantiated this class needs 3 integers: the length of the shingles (k), the number of bands (b) and the number of rows (r), then it needs also the hash functions that have to be applied in the steps of the algorithm. The function that this class exposes takes in input a RDD object containing all the descriptions and apply the algorithm (As in previous exercise the function also keep trace of the execution time). Let's explain that through its steps.

5.1.1 Shingling

To produce the shingles as first it applies a flat map that returns set of tuples (id,hash(shingle)), each tuple is produced dividing each element of the RDD into its shingles of length k

5.1.2 Minwise-hashing

To produce the signature of the descriptions in the RDD minwise hashing is executed in 2 steps:

- in step 1 it is applied again a flat map to the shingles and it is returned set of tuples ((id, num hash), hashed element), each tuple is produced applying to each element the t hashes (if you apply i-th hash then num hash = i)
- in step 2 it is applied a reduce by key function, since the key is (id, num hash) then all the elements of the same document to which it was applied the same hash are reduced and the reduce function is minimum. In this way the result is elements of kind ((id, num hash), hashed element) but now there is only 1 element or each pair (id, num hash) and the "hashed element" is the minimum

5.1.3 LSH

To apply LSH more steps are needed:

- in step 1 each element is mapped to its band, the resulting is a tuple of kind (id, num band), array with the element placed in its position inside the signature of the document (id). For each element it is computed the band and the resulting position inside an array of length r in which all the elements are null except for the hashed element that is placed in the computed position
- In step 2 a reduce by key function is called and the function applied is a merging of arrays so that every element of the same document that belongs to the same band is merged in one single array
- In step 3 every element is mapped to a tuple ((num band, hash(element)), id) where hash(element) is the hash function applied to the array
- In step 4 it is used a reduce by key function which has to reduce the elements with same (num band, hash(element)) to a single set of ordered pairs, in this way now each element of our RDD is like a bucket containing all the pairs of elements which had same (num band, hash(element))
- In step 5 we simply filter all the elements which haven't a set as value, since there could be pairs (num band, hash(element)) that are not repeated in the RDD
- In step 6 we compute the result applying a reduce function which merge all the sets together

5.2 CompareWithJaccardSpark

This class is built for comparing descriptions using Jaccard similarity and it is implemented with Spark.

To be instantiated this class needs an integer k which represents the length of the shingles, a float which is the threshold that tells us which descriptions are considered near duplicates and which not and a hash function to be applied during shingling.

As the other class also this exposes a function to compute the nearest descriptions using Jaccard similarity and it also keep trace of the execution time. The first part of the algorithm is shared with the previous class (Shingling part 5.1.1) the real algorithm is executed after and here it is explained in few steps:

- In step 1 the shingles of a same document are put together using a reduce by key function, the result is for each document a set of shingles

- In step 2 it is applied a flat map and for each element it is created a set of elements with the same value (the set) but different id, the id is a ordered pair that represent the comparing to be executed in the next step
- In step 3 it is applied a reduce by key function that for each pair of elements with the same key compute the Jaccard similarity and return True if it is higher than the threshold, false otherwise. So now the elements are ((id1,id2) are near?)
- In step 4 it is simply applied a filter on the tuples that are near duplicates (are near? = true)
- In step 5 the result is computed through a reduce function that puts all the pairs in a set

5.3 Report

The class report is also used in this case to execute the 2 classes and build the report. From the result we can notice that the resulting sets are the same as the ones computed by the code which doesn't use Spark so same number of elements in each result and same intersections between them. What change is the execution time which is anyways higher than the one of the algorithm which doesn't use Spark. It is probably due to the fact that the data are not so suitable to be treated in a distributed way or to the fact that my machine has only 4 cores so the parallelization is not so good. Anyway we can notice that beside the Jaccard similarity algorithm execution time has been more than duplicated the execution time of LSH algorithm is increased only of nearly 50%, then we could say that the LSH algorithm has shown to be more suitable for parallel execution. These are the results of the execution times:

- Time required for the execution of LSH algorithm: 7873
- Time required for the execution of Jaccard similarity algorithm: 2888

All the results of the report are also available in the menu "NEAREST PRODUCTS REPORTS" and the results of the nearest descriptions both for LSH and Jaccard similarity are available in the same menu

6 Guide to use the application

The application is made by 2 fundamental parts the script for installing the required python libraries install-requirements.sh and the application: app.py. This one is inside the folder "server" and has to be executed using the command python3 app.py (easy) now it should launch a server available on local host at url http://127.0.0.1:5000. If you click on the link (after launching the app) you should find all the result and the functions visualized in a better way (nothing spectacular but way better than a txt file), you can navigate through menus and try all the functions.

7 Images

DataMining HW 2 Server						
Top 10 products by price						
	description	price	prime	url	stars	reviews
100	Intel Core i7-10700K Processor 8 Cores 16 Threads 12.5 MB Cache	399.00	Prime	https://www.amazon.it/Intel-Core-i7-10700K-Processor/dp/B087RBYDQC/ref=sr_1_97?keywords=cpu&qid=1699106766&s=pc-rank	4.5	100
247	Intel Core i7-10700K Processor 8 Cores 16 Threads 12.5 MB Cache	399.00	Prime	https://www.amazon.it/Intel-Core-i7-10700K-Processor/dp/B087RBYDQC/ref=sr_1_204?keywords=cpu&qid=1699106766&s=pc-rank	4.5	100
248	Intel Core i7-10700K Processor 8 Cores 16 Threads 12.5 MB Cache	399.00	Prime	https://www.amazon.it/Intel-Core-i7-10700K-Processor/dp/B087RBYDQC/ref=sr_1_263?keywords=cpu&qid=1699106766&s=pc-rank	4.5	100
249	Intel Core i7-10700K Processor 8 Cores 16 Threads 12.5 MB Cache	399.00	Prime	https://www.amazon.it/Intel-Core-i7-10700K-Processor/dp/B087RBYDQC/ref=sr_1_112?keywords=cpu&qid=1699106766&s=pc-rank	4.5	100
250	Intel Core i7-10700K Processor 8 Cores 16 Threads 12.5 MB Cache	399.00	Prime	https://www.amazon.it/Intel-Core-i7-10700K-Processor/dp/B087RBYDQC/ref=sr_1_49?keywords=cpu&qid=1699106766&s=pc-rank	4.5	100
251	Intel Core i7-10700K Processor 8 Cores 16 Threads 12.5 MB Cache	399.00	Prime	https://www.amazon.it/Intel-Core-i7-10700K-Processor/dp/B087RBYDQC/ref=sr_1_89?keywords=cpu&qid=1699106766&s=pc-rank	4.5	100
252	Intel Core i7-10700K Processor 8 Cores 16 Threads 12.5 MB Cache	399.00	Prime	https://www.amazon.it/Intel-Core-i7-10700K-Processor/dp/B087RBYDQC/ref=sr_1_93?keywords=cpu&qid=1699106766&s=pc-rank	4.5	100
253	Intel Core i7-10700K Processor 8 Cores 16 Threads 12.5 MB Cache	399.00	Prime	https://www.amazon.it/Intel-Core-i7-10700K-Processor/dp/B087RBYDQC/ref=sr_1_86?keywords=cpu&qid=1699106766&s=pc-rank	4.5	100

Figure 2: Top 10 products by price

[illegible]

Figure 3: Top 10 products by rating

[illegible]

Figure 4: Top 10 products by rating weighted by the number of reviews

Return to fixed queries menu

DataMining HW 2 Server

Analyze primeness

average price		average price prime products	average stars	average stars prime products
0.247.9	232.8	4.38	4.2	

Figure 5: Analysis of primeness with respect to price and ratings

DataMining HW 2 Server

Price range for different series of gpu

	categories	min	price	max	price
0	4090	17.99	2119.00		
1	4080	19.99	1739.99		
2	7900	1204.38	1204.38		
3	4070	22.86	1079.23		
4	3090	8.90	3251.73		
5	3080	167.63	1015.68		
6	6800	5750.00	5750.00		
7	3070	502.00	931.90		
8	2080	31.98	31.98		
9	6700	377.90	407.89		
10	4060	346.00	563.00		
11	3060	318.00	377.89		
12	6650	285.54	394.57		
13	7600	316.58	361.47		
14	6600	217.00	253.00		
15	1080	15.95	147.30		
16	2060	307.00	369.88		
17	0500	258.62	317.20		
18	1660	233.85	233.85		
19	590	19.99	205.99		
20	6500	179.65	179.65		
21	2060	307.00	369.88		
22	580	109.90	148.75		
23	1650	179.90	189.96		
24	680	109.76	109.76		
25	780	800.42	800.42		

Figure 6: Price range for different series of gpus

DataMining HW 2 Server										
Result of query: 3080 mvidia										
Insert query										
		description	price	price				url	status	
22	ms	MSI RTX 3080 NVIDIA GeForce RTX 3080 10 GB Edition Gaming Graphics, 12GB GDDR6X 384-bit 10	N/A	False				https://www.amazon.de/B085Y7J7-GeForce-DesktopPC-Complataq/gp/product/B085Y7J7-GeForce-DesktopPC-Complataq?pf_rd_p=1599525292	ms	1200
23	ms	GeForce RTX 3080 10 GB NVIDIA Ampere, ROG STRIX RTX3080 10 GB 25-Gigamem	N/A	False				https://www.amazon.de/B085Y7J7-GeForce-DesktopPC-Complataq/gp/product/B085Y7J7-GeForce-DesktopPC-Complataq?pf_rd_p=1599525292	ms	1200
24	ms	GeForce RTX 3080 10 GB NVIDIA Ampere, ROG STRIX RTX3080 10 GB 25-Gigamem	N/A	False				https://www.amazon.de/B085Y7J7-GeForce-DesktopPC-Complataq/gp/product/B085Y7J7-GeForce-DesktopPC-Complataq?pf_rd_p=1599525292	ms	1200
25	ms	GeForce RTX 3080 10 GB NVIDIA Ampere, ROG STRIX RTX3080 10 GB 25-Gigamem	N/A	False				https://www.amazon.de/B085Y7J7-GeForce-DesktopPC-Complataq/gp/product/B085Y7J7-GeForce-DesktopPC-Complataq?pf_rd_p=1599525292	ms	1200
26	ms	GeForce RTX 3080 10 GB NVIDIA Ampere, ROG STRIX RTX3080 10 GB 25-Gigamem	N/A	False				https://www.amazon.de/B085Y7J7-GeForce-DesktopPC-Complataq/gp/product/B085Y7J7-GeForce-DesktopPC-Complataq?pf_rd_p=1599525292	ms	1200
27	ms	GeForce RTX 3080 10 GB NVIDIA Ampere, ROG STRIX RTX3080 10 GB 25-Gigamem	N/A	False				https://www.amazon.de/B085Y7J7-GeForce-DesktopPC-Complataq/gp/product/B085Y7J7-GeForce-DesktopPC-Complataq?pf_rd_p=1599525292	ms	1200
28	ms	GeForce RTX 3080 10 GB NVIDIA Ampere, ROG STRIX RTX3080 10 GB 25-Gigamem	N/A	False				https://www.amazon.de/B085Y7J7-GeForce-DesktopPC-Complataq/gp/product/B085Y7J7-GeForce-DesktopPC-Complataq?pf_rd_p=1599525292	ms	1200
29	ms	GeForce RTX 3080 10 GB NVIDIA Ampere, ROG STRIX RTX3080 10 GB 25-Gigamem	N/A	False				https://www.amazon.de/B085Y7J7-GeForce-DesktopPC-Complataq/gp/product/B085Y7J7-GeForce-DesktopPC-Complataq?pf_rd_p=1599525292	ms	1200
30	ms	GeForce RTX 3080 10 GB NVIDIA Ampere, ROG STRIX RTX3080 10 GB 25-Gigamem	N/A	False				https://www.amazon.de/B085Y7J7-GeForce-DesktopPC-Complataq/gp/product/B085Y7J7-GeForce-DesktopPC-Complataq?pf_rd_p=1599525292	ms	1200
31	ms	GeForce RTX 3080 10 GB NVIDIA Ampere, ROG STRIX RTX3080 10 GB 25-Gigamem	N/A	False				https://www.amazon.de/B085Y7J7-GeForce-DesktopPC-Complataq/gp/product/B085Y7J7-GeForce-DesktopPC-Complataq?pf_rd_p=1599525292	ms	1200
32	ms	GeForce RTX 3080 10 GB NVIDIA Ampere, ROG STRIX RTX3080 10 GB 25-Gigamem	N/A	False				https://www.amazon.de/B085Y7J7-GeForce-DesktopPC-Complataq/gp/product/B085Y7J7-GeForce-DesktopPC-Complataq?pf_rd_p=1599525292	ms	1200
33	ms	GeForce RTX 3080 10 GB NVIDIA Ampere, ROG STRIX RTX3080 10 GB 25-Gigamem	N/A	False				https://www.amazon.de/B085Y7J7-GeForce-DesktopPC-Complataq/gp/product/B085Y7J7-GeForce-DesktopPC-Complataq?pf_rd_p=1599525292	ms	1200
34	ms	GeForce RTX 3080 10 GB NVIDIA Ampere, ROG STRIX RTX3080 10 GB 25-Gigamem	N/A	False				https://www.amazon.de/B085Y7J7-GeForce-DesktopPC-Complataq/gp/product/B085Y7J7-GeForce-DesktopPC-Complataq?pf_rd_p=1599525292	ms	1200
35	ms	GeForce RTX 3080 10 GB NVIDIA Ampere, ROG STRIX RTX3080 10 GB 25-Gigamem	N/A	False				https://www.amazon.de/B085Y7J7-GeForce-DesktopPC-Complataq/gp/product/B085Y7J7-GeForce-DesktopPC-Complataq?pf_rd_p=1599525292	ms	1200
36	ms	GeForce RTX 3080 10 GB NVIDIA Ampere, ROG STRIX RTX3080 10 GB 25-Gigamem	N/A	False				https://www.amazon.de/B085Y7J7-GeForce-DesktopPC-Complataq/gp/product/B085Y7J7-GeForce-DesktopPC-Complataq?pf_rd_p=1599525292	ms	1200
37	ms	GeForce RTX 3080 10 GB NVIDIA Ampere, ROG STRIX RTX3080 10 GB 25-Gigamem	N/A	False				https://www.amazon.de/B085Y7J7-GeForce-DesktopPC-Complataq/gp/product/B085Y7J7-GeForce-DesktopPC-Complataq?pf_rd_p=1599525292	ms	1200
38	ms	GeForce RTX 3080 10 GB NVIDIA Ampere, ROG STRIX RTX3080 10 GB 25-Gigamem	N/A	False				https://www.amazon.de/B085Y7J7-GeForce-DesktopPC-Complataq/gp/product/B085Y7J7-GeForce-DesktopPC-Complataq?pf_rd_p=1599525292	ms	1200
39	ms	GeForce RTX 3080 10 GB NVIDIA Ampere, ROG STRIX RTX3080 10 GB 25-Gigamem	N/A	False				https://www.amazon.de/B085Y7J7-GeForce-DesktopPC-Complataq/gp/product/B085Y7J7-GeForce-DesktopPC-Complataq?pf_rd_p=1599525292	ms	1200
40	ms	GeForce RTX 3080 10 GB NVIDIA Ampere, ROG STRIX RTX3080 10 GB 25-Gigamem	N/A	False				https://www.amazon.de/B085Y7J7-GeForce-DesktopPC-Complataq/gp/product/B085Y7J7-GeForce-DesktopPC-Complataq?pf_rd_p=1599525292	ms	1200
41	ms	GeForce RTX 3080 10 GB NVIDIA Ampere, ROG STRIX RTX3080 10 GB 25-Gigamem	N/A	False				https://www.amazon.de/B085Y7J7-GeForce-DesktopPC-Complataq/gp/product/B085Y7J7-GeForce-DesktopPC-Complataq?pf_rd_p=1599525292	ms	1200
42	ms	GeForce RTX 3080 10 GB NVIDIA Ampere, ROG STRIX RTX3080 10 GB 25-Gigamem	N/A	False				https://www.amazon.de/B085Y7J7-GeForce-DesktopPC-Complataq/gp/product/B085Y7J7-GeForce-DesktopPC-Complataq?pf_rd_p=1599525292	ms	1200
43	ms	GeForce RTX 3080 10 GB NVIDIA Ampere, ROG STRIX RTX3080 10 GB 25-Gigamem	N/A	False				https://www.amazon.de/B085Y7J7-GeForce-DesktopPC-Complataq/gp/product/B085Y7J7-GeForce-DesktopPC-Complataq?pf_rd_p=1599525292	ms	1200
44	ms	GeForce RTX 3080 10 GB NVIDIA Ampere, ROG STRIX RTX3080 10 GB 25-Gigamem	N/A	False				https://www.amazon.de/B085Y7J7-GeForce-DesktopPC-Complataq/gp/product/B085Y7J7-GeForce-DesktopPC-Complataq?pf_rd_p=1599525292	ms	1200
45	ms	GeForce RTX 3080 10 GB NVIDIA Ampere, ROG STRIX RTX3080 10 GB 25-Gigamem	N/A	False				https://www.amazon.de/B085Y7J7-GeForce-DesktopPC-Complataq/gp/product/B085Y7J7-GeForce-DesktopPC-Complataq?pf_rd_p=1599525292	ms	1200
46	ms	GeForce RTX 3080 10 GB NVIDIA Ampere, ROG STRIX RTX3080 10 GB 25-Gigamem	N/A	False				https://www.amazon.de/B085Y7J7-GeForce-DesktopPC-Complataq/gp/product/B085Y7J7-GeForce-DesktopPC-Complataq?pf_rd_p=1599525292	ms	1200
47	ms	GeForce RTX 3080 10 GB NVIDIA Ampere, ROG STRIX RTX3080 10 GB 25-Gigamem	N/A	False				https://www.amazon.de/B085Y7J7-GeForce-DesktopPC-Complataq/gp/product/B085Y7J7-GeForce-DesktopPC-Complataq?pf_rd_p=1599525292	ms	1200
48	ms	GeForce RTX 3080 10 GB NVIDIA Ampere, ROG STRIX RTX3080 10 GB 25-Gigamem	N/A	False				https://www.amazon.de/B085Y7J7-GeForce-DesktopPC-Complataq/gp/product/B085Y7J7-GeForce-DesktopPC-Complataq?pf_rd_p=1599525292	ms	1200
49	ms	GeForce RTX 3080 10 GB NVIDIA Ampere, ROG STRIX RTX3080 10 GB 25-Gigamem	N/A	False				https://www.amazon.de/B085Y7J7-GeForce-DesktopPC-Complataq/gp/product/B085Y7J7-GeForce-DesktopPC-Complataq?pf_rd_p=1599525292	ms	1200
50	ms	GeForce RTX 3080 10 GB NVIDIA Ampere, ROG STRIX RTX3080 10 GB 25-Gigamem	N/A	False				https://www.amazon.de/B085Y7J7-GeForce-DesktopPC-Complataq/gp/product/B085Y7J7-GeForce-DesktopPC-Complataq?pf_rd_p=1599525292	ms	1200
51	ms	GeForce RTX 3080 10 GB NVIDIA Ampere, ROG STRIX RTX3080 10 GB 25-Gigamem	N/A	False				https://www.amazon.de/B085Y7J7-GeForce-DesktopPC-Complataq/gp/product/B085Y7J7-GeForce-DesktopPC-Complataq?pf_rd_p=1599525292	ms	1200
52	ms	GeForce RTX 3080 10 GB NVIDIA Ampere, ROG STRIX RTX3080 10 GB 25-Gigamem	N/A	False				https://www.amazon.de/B085Y7J7-GeForce-DesktopPC-Complataq/gp/product/B085Y7J7-GeForce-DesktopPC-Complataq?pf_rd_p=1599525292	ms	1200
53	ms	GeForce RTX 3080 10 GB NVIDIA Ampere, ROG STRIX RTX3080 10 GB 25-Gigamem	N/A	False				https://www.amazon.de/B085Y7J7-GeForce-DesktopPC-Complataq/gp/product/B085Y7J7-GeForce-DesktopPC-Complataq?pf_rd_p=1599525292	ms	1200
54	ms	GeForce RTX 3080 10 GB NVIDIA Ampere, ROG STRIX RTX3080 10 GB 25-Gigamem	N/A	False				https://www.amazon.de/B085Y7J7-GeForce-DesktopPC-Complataq/gp/product/B085Y7J7-GeForce-DesktopPC-Complataq?pf_rd_p=1599525292	ms	1200
55	ms	GeForce RTX 3080 10 GB NVIDIA Ampere, ROG STRIX RTX3080 10 GB 25-Gigamem	N/A	False				https://www.amazon.de/B085Y7J7-GeForce-DesktopPC-Complataq/gp/product/B085Y7J7-GeForce-DesktopPC-Complataq?pf_rd_p=1599525292	ms	1200
56	ms	GeForce RTX 3080 10 GB NVIDIA Ampere, ROG STRIX RTX3080 10 GB 25-Gigamem	N/A	False				https://www.amazon.de/B085Y7J7-GeForce-DesktopPC-Complataq/gp/product/B085Y7J7-GeForce-DesktopPC-Complataq?pf_rd_p=1599525292	ms	1200
57	ms	GeForce RTX 3080 10 GB NVIDIA Ampere, ROG STRIX RTX3080 10 GB 25-Gigamem	N/A	False				https://www.amazon.de/B085Y7J7-GeForce-DesktopPC-Complataq/gp/product/B085Y7J7-GeForce-DesktopPC-Complataq?pf_rd_p=1599525292	ms	1200
58	ms	GeForce RTX 3080 10 GB NVIDIA Ampere, ROG STRIX RTX3080 10 GB 25-Gigamem	N/A	False				https://www.amazon.de/B085Y7J7-GeForce-DesktopPC-Complataq/gp/product/B085Y7J7-GeForce-DesktopPC-Complataq?pf_rd_p=1599525292	ms	1200
59	ms	GeForce RTX 3080 10 GB NVIDIA Ampere, ROG STRIX RTX3080 10 GB 25-Gigamem	N/A	False				https://www.amazon.de/B085Y7J7-GeForce-DesktopPC-Complataq/gp/product/B085Y7J7-GeForce-DesktopPC-Complataq?pf_rd_p=1599525292	ms	1200
60	ms	GeForce RTX 3080 10 GB NVIDIA Ampere, ROG STRIX RTX3080 10 GB 25-Gigamem	N/A	False				https://www.amazon.de/B085Y7J7-GeForce-DesktopPC-Complataq/gp/product/B085Y7J7-GeForce-DesktopPC-Complataq?pf_rd_p=1599525292	ms	1200
61	ms	GeForce RTX 3080 10 GB NVIDIA Ampere, ROG STRIX RTX3080 10 GB 25-Gigamem	N/A	False				https://www.amazon.de/B085Y7J7-GeForce-DesktopPC-Complataq/gp/product/B085Y7J7-GeForce-DesktopPC-Complataq?pf_rd_p=1599525292	ms	1200
62	ms	GeForce RTX 3080 10 GB NVIDIA Ampere, ROG STRIX RTX3080 10 GB 25-Gigamem	N/A	False				https://www.amazon.de/B085Y7J7-GeForce-DesktopPC-Complataq/gp/product/B085Y7J7-GeForce-DesktopPC-Complataq?pf_rd_p=1599525292	ms	1200
63	ms	GeForce RTX 3080 10 GB NVIDIA Ampere, ROG STRIX RTX3080 10 GB 25-Gigamem	N/A	False				https://www.amazon.de/B085Y7J7-GeForce-DesktopPC-Complataq/gp/product/B085Y7J7-GeForce-DesktopPC-Complataq?pf_rd_p=1599525292	ms	1200
64	ms	GeForce RTX 3080 10 GB NVIDIA Ampere, ROG STRIX RTX3080 10 GB 25-Gigamem	N/A	False				https://www.amazon.de/B085Y7J7-GeForce-DesktopPC-Complataq/gp/product/B085Y7J7-GeForce-DesktopPC-Complataq?pf_rd_p=1599525292	ms	1200
65	ms	GeForce RTX 3080 10 GB NVIDIA Ampere, ROG STRIX RTX3080 10 GB 25-Gigamem	N/A	False				https://www.amazon.de/B085Y7J7-GeForce-DesktopPC-Complataq/gp/product/B085Y7J7-GeForce-DesktopPC-Complataq?pf_rd_p=1599525292	ms	1200
66	ms	GeForce RTX 3080 10 GB NVIDIA Ampere, ROG STRIX RTX3080 10 GB 25-Gigamem	N/A	False				https://www.amazon.de/B085Y7J7-GeForce-DesktopPC-Complataq/gp/product/B085Y7J7-GeForce-DesktopPC-Complataq?pf_rd_p=1599525292	ms	1200
67	ms	GeForce RTX 3080 10 GB NVIDIA Ampere, ROG STRIX RTX3080 10 GB 25-Gigamem	N/A	False				https://www.amazon.de/B085Y7J7-GeForce-DesktopPC-Complataq/gp/product/B085Y7J7-GeForce-DesktopPC-Complataq?pf_rd_p=1599525292	ms	1200
68	ms	GeForce RTX 3080 10 GB NVIDIA Ampere, ROG STRIX RTX3080 10 GB 25-Gigamem	N/A	False				https://www.amazon.de/B085Y7J7-GeForce-DesktopPC-Complataq/gp/product/B085Y7J7-GeForce-DesktopPC-Complataq?pf_rd_p=1599525292	ms	1200
69	ms	GeForce RTX 3080 10 GB NVIDIA Ampere, ROG STRIX RTX3080 10 GB 25-Gigamem	N/A	False				https://www.amazon.de/B085Y7J7-GeForce-DesktopPC-Complataq/gp/product/B085Y7J7-GeForce-DesktopPC-Complataq?pf_rd_p=1599525292	ms	1200
70	ms	GeForce RTX 3080 10 GB NVIDIA Ampere, ROG STRIX RTX3080 10 GB 25-Gigamem	N/A	False				https://www.amazon.de/B085Y7J7-GeForce-DesktopPC-Complataq/gp/product/B085Y7J7-GeForce-DesktopPC-Complataq?pf_rd_p=1599525292	ms	1200
71	ms	GeForce RTX 3080 10 GB NVIDIA Ampere, ROG STRIX RTX3080 10 GB 25-Gigamem	N/A	False				https://www.amazon.de/B085Y7J7-GeForce-DesktopPC-Complataq/gp/product/B085Y7J7-GeForce-DesktopPC-Complataq?pf_rd_p=1599525292	ms	1200
72	ms	GeForce RTX 3080 10 GB NVIDIA Ampere, ROG STRIX RTX3080 10 GB 25-Gigamem	N/A	False				https://www.amazon.de/B085Y7J7-GeForce-DesktopPC-Complataq/gp/product/B085Y7J7-GeForce-DesktopPC-Complataq?pf_rd_p=1599525292	ms	1200
73	ms	GeForce RTX 3080 10 GB NVIDIA Ampere, ROG STRIX RTX3080 10 GB 25-Gigamem	N/A	False				https://www.amazon.de/B085Y7J7-GeForce-DesktopPC-Complataq/gp/product/B085Y7J7-GeForce-DesktopPC-Complataq?pf_rd_p=1599525292	ms	1200
74	ms	GeForce RTX 3080 10 GB NVIDIA Ampere, ROG STRIX RTX3080 10 GB 25-Gigamem	N/A	False				https://www.amazon.de/B085Y7J7-GeForce-DesktopPC-Complataq/gp/product/B085Y7J7-GeForce-DesktopPC-Complataq		

Figure 7: Result of a short query using cosine similarity

[illegible]

Figure 9: Some resulting near duplicates computed using LSH