

COMS 435 PA2 Report

Brad Warren

184030907

MinHash:

The procedure I used to collect all the terms of within every single one of the documents started with constructing a HashMap with strings as the keys and integers as the values. I made an array of files and with this array of files I made a for loop that goes through every file and runs a method called `getTerms()`. The method directly adds a term to my HashMap if `containsKey(term)` returns false and the length of the term is greater than 2 and not equal to "the" as well as does not have random punctuations like (">", "<", "*", etc. to decrease number of terms). The integer assigned to each term was done in order of which the terms were found. Thus, if the first word in the first document was "hello" then hello would be at index 1 on our HashMap. Then I added all the terms to an Array list in order and the order they came in was their number. In order to construct random permutations, I created a 2d integer array called `permutations` the size of the [number of permutations] [number of terms]. I then constructed a for loop over the number of permutations and used a method called `fisherYatesShuffle` where I implemented fisher Yates Shuffle Algorithm.

MinHashAccuracy:

Running the program and results:

	Permutations: 400	Permutations: 600	Permutations: 800
Threshold: 0.04	12765	5302	4490
Threshold: 0.07	528	515	500
Threshold: 0.09	422	400	270

From this data above we can conclude that out of 1004004 pairs from the space.zip getting this small amount of difference between approximate and exact jaccard is great. We can also conclude that the more permutations we use, the more closely related approximate and exact jaccard are. Also, we can determine that not many pairs have substantial differences in approximate and exact jaccard (seen as threshold decreases). The higher the threshold the smaller the count of pairs which is good.

MinHashTime:

With 600 permutations on the files from space.zip, the runtimes are as follows:

```
Time taken to compute exact Jaccard between each document = 31631.0 ms
Time taken to compute approximate Jaccard between each document = 1298.0 ms
```

We can see that the total time taken to compute the exact Jaccard is 30,333 ms slower than the time it takes to compute approximate Jaccard between each document.

NearDuplicates:

For this I used subset of PA2.zip of 1,610 items, permutations of 600 for every test a test for .9 threshold and .7 threshold for each file tested. The outputs are below:

Space-1000.txt - .7 threshold:

```
space-1000.txt.copy4
space-1000.txt.copy1
space-1000.txt.copy6
space-1000.txt.copy5
space-1000.txt.copy2
space-1000.txt.copy3
space-918.txt.copy3
space-1000.txt.copy7
space-894.txt.copy7
space-805.txt.copy2
space-838.txt.copy2
space-840.txt.copy4
space-983.txt.copy5
space-937.txt.copy6
space-877.txt.copy7
space-897.txt.copy4
space-975.txt.copy4
space-926.txt.copy2
space-855.txt.copy1
space-951.txt.copy3
space-988.txt.copy6
space-874.txt.copy5
space-949.txt.copy6
space-943.txt.copy2
space-983.txt.copy2
space-937.txt.copy7
space-845.txt.copy3
space-872.txt.copy7
space-918.txt.copy1
space-821.txt.copy3
```

- .9 threshold:

```
space-1000.txt.copy5
space-1000.txt.copy2
space-1000.txt.copy3
space-1000.txt.copy1
space-1000.txt.copy6
space-1000.txt.copy7
space-1000.txt.copy4
space-880.txt.copy1
space-964.txt.copy5
space-806.txt.copy1
space-839.txt
space-839.txt.copy5
space-952.txt.copy7
space-802.txt.copy2
space-823.txt.copy6
space-823.txt.copy2
space-997.txt.copy7
space-967.txt.copy4
space-967.txt.copy1
space-901.txt.copy7
space-824.txt.copy1
space-824.txt.copy7
space-824.txt
space-808.txt.copy7
```

space-800.txt.copy2
space-800.txt.copy5
space-800.txt.copy4
space-800.txt.copy6
space-800.txt.copy3
space-800.txt.copy1
space-800.txt.copy7
space-965.txt.copy1
space-960.txt.copy4
space-815.txt.copy3
space-855.txt.copy6
space-993.txt.copy7
space-847.txt.copy2
space-907.txt.copy4
space-877.txt.copy2
space-961.txt.copy2
space-961.txt.copy3
space-961.txt
space-961.txt.copy7
space-835.txt.copy2
space-835.txt.copy1
space-835.txt
space-875.txt.copy4
space-839.txt.copy1
space-817.txt.copy7
space-831.txt.copy6
space-809.txt.copy2
space-996.txt.copy5
space-852.txt.copy4
space-976.txt
space-976.txt.copy1
space-976.txt.copy2
space-976.txt.copy5
space-976.txt.copy4
space-976.txt.copy3
space-900.txt.copy3
space-902.txt.copy5
space-939.txt.copy6
space-939.txt.copy1
space-939.txt.copy7
space-820.txt.copy5

Space-800.txt – threshold .7:

space-800.txt.copy2
space-800.txt.copy3
space-800.txt.copy1
space-800.txt.copy5
space-800.txt.copy7
space-800.txt.copy4
space-800.txt.copy6
space-977.txt.copy1
space-977.txt.copy2
space-977.txt.copy3
space-842.txt.copy5
space-817.txt.copy3
space-840.txt.copy7
space-862.txt.copy2
space-842.txt.copy2
space-994.txt.copy4
space-934.txt.copy6
space-934.txt

threshold .9:

space-831.txt.copy6
space-831.txt.copy1
space-942.txt.copy6
space-838.txt.copy5
space-900.txt.copy3
space-831.txt.copy2
space-831.txt.copy4
space-831.txt.copy3
space-907.txt.copy2
space-831.txt.copy7
space-932.txt.copy1
space-932.txt.copy6
space-831.txt.copy5
space-932.txt
space-974.txt.copy6
space-932.txt.copy2
space-969.txt.copy4
space-947.txt.copy2
space-847.txt.copy3
space-853.txt.copy5
space-823.txt.copy5
space-882.txt.copy1
space-923.txt.copy3
space-923.txt.copy4
space-923.txt.copy5
space-923.txt.copy2
space-923.txt.copy7
space-923.txt.copy1
space-923.txt.copy6
space-923.txt
space-968.txt.copy1
space-828.txt.copy4
space-928.txt.copy1
space-866.txt.copy1
space-866.txt

Space-831.txt – threshold .7:

space-831.txt.copy1
space-831.txt.copy7
space-831.txt.copy2
space-831.txt.copy4
space-831.txt.copy6
space-831.txt.copy5
space-831.txt.copy3
space-853.txt.copy3
space-913.txt.copy1
space-855.txt.copy7
space-855.txt
space-855.txt.copy4
space-855.txt.copy2
space-847.txt
space-847.txt.copy4
space-840.txt.copy7
space-883.txt.copy4
space-938.txt.copy3
space-967.txt.copy2
space-864.txt.copy6
space-958.txt.copy3
space-828.txt
space-996.txt.copy1

- threshold .9:

space-866.txt.copy6
space-846.txt.copy3
space-866.txt.copy5
space-866.txt.copy7
space-866.txt.copy1
space-866.txt.copy3
space-866.txt.copy4
space-929.txt.copy1
space-866.txt.copy2
space-990.txt.copy6
space-844.txt.copy4
space-993.txt.copy7
space-993.txt.copy6
space-993.txt.copy3
space-993.txt
space-851.txt.copy6
space-970.txt.copy7
space-947.txt.copy4
space-867.txt.copy7
space-867.txt.copy1
space-867.txt
space-867.txt.copy4
space-813.txt.copy5
space-954.txt.copy2
space-881.txt.copy6
space-801.txt.copy6
space-863.txt.copy3
space-863.txt.copy2
space-863.txt.copy5
space-863.txt
space-863.txt.copy7
space-863.txt.copy1

Space-866.txt – threshold .7:

space-866.txt.copy7
space-866.txt.copy6
space-866.txt.copy3
space-866.txt.copy1
space-916.txt.copy1
space-866.txt.copy4
space-866.txt.copy2
space-866.txt.copy5
space-923.txt.copy4
space-831.txt.copy6
space-816.txt.copy2
space-981.txt.copy1
space-800.txt
space-800.txt.copy5
space-800.txt.copy6
space-896.txt.copy7
space-916.txt
space-916.txt.copy6
space-945.txt.copy6
space-985.txt.copy7
space-940.txt.copy7

- threshold .9:

space-901.txt.copy6
space-901.txt.copy4
space-901.txt.copy1
space-901.txt.copy2
space-901.txt.copy3
space-901.txt.copy5
space-901.txt.copy7
space-885.txt.copy7
space-844.txt.copy4
space-915.txt.copy7
space-941.txt.copy2
space-941.txt.copy3
space-941.txt.copy4
space-941.txt
space-848.txt.copy7
space-864.txt.copy3
space-864.txt
space-864.txt.copy7
space-957.txt.copy2
space-945.txt
space-945.txt.copy5
space-945.txt.copy4
space-988.txt.copy3
space-913.txt.copy6
space-958.txt.copy7
space-900.txt.copy4
space-807.txt
space-807.txt.copy5
space-807.txt.copy2
space-807.txt.copy3
space-807.txt.copy4
space-947.txt.copy2
space-807.txt.copy1
space-807.txt.copy6
space-800.txt.copy4
space-889.txt.copy7
space-831.txt.copy5
space-837.txt.copy3
space-892.txt.copy4

Space-901.txt – threshold .7:

space-901.txt.copy4
space-901.txt.copy7
space-901.txt.copy1
space-901.txt.copy6
space-901.txt.copy5
space-901.txt.copy2
space-868.txt.copy5
space-868.txt.copy6
space-901.txt.copy3
space-844.txt.copy4
space-978.txt.copy6
space-812.txt.copy2
space-835.txt.copy2
space-835.txt.copy7
space-892.txt.copy3
space-892.txt.copy6
space-832.txt.copy6
space-800.txt
space-800.txt.copy5
space-808.txt.copy1
space-846.txt.copy4
space-870.txt.copy2
space-870.txt
space-870.txt.copy1
space-876.txt.copy5

- threshold .9:

space-944.txt.copy6
space-944.txt.copy2
space-944.txt.copy5
space-944.txt.copy4
space-944.txt.copy3
space-944.txt.copy1
space-944.txt.copy7
space-945.txt.copy5
space-982.txt.copy3
space-982.txt.copy2
space-982.txt.copy5
space-982.txt.copy7
space-982.txt.copy6
space-982.txt
space-894.txt.copy1
space-831.txt.copy6
space-827.txt.copy7
space-812.txt.copy2
space-812.txt.copy5
space-812.txt.copy3
space-812.txt
space-812.txt.copy1
space-812.txt.copy7
space-933.txt.copy7
space-933.txt.copy1
space-933.txt.copy6
space-933.txt
space-933.txt.copy3
space-933.txt.copy4
space-933.txt.copy2
space-813.txt.copy4
space-959.txt.copy3
space-959.txt
space-959.txt.copy7
space-959.txt.copy6
space-959.txt.copy1
space-834.txt.copy5
space-880.txt.copy5
space-889.txt.copy5

Space-944.txt – threshold .7:

space-944.txt.copy6
space-944.txt.copy7
space-891.txt.copy7
space-944.txt.copy2
space-944.txt.copy1
space-944.txt.copy5
space-944.txt.copy3
space-944.txt.copy4
space-879.txt.copy5
space-914.txt.copy4
space-988.txt.copy1
space-831.txt.copy1
space-835.txt.copy4
space-982.txt.copy4
space-942.txt.copy4
space-942.txt.copy6
space-958.txt
space-942.txt.copy3
space-958.txt.copy7
space-958.txt.copy6
space-799.txt.copy7
space-926.txt.copy3
space-990.txt.copy4
space-960.txt
space-960.txt.copy2
space-960.txt.copy6
space-916.txt.copy4
space-925.txt.copy3
space-925.txt
space-867.txt.copy7

- threshold .9:

space-961.txt.copy3
space-961.txt.copy4
space-961.txt.copy1
space-961.txt.copy7
space-968.txt.copy4
space-961.txt.copy5
space-961.txt.copy2
space-961.txt.copy6
space-847.txt.copy3
space-960.txt.copy7
space-893.txt.copy3
space-893.txt
space-893.txt.copy1
space-893.txt.copy6
space-996.txt.copy6
space-915.txt.copy6
space-869.txt.copy7
space-873.txt.copy1
space-998.txt
space-998.txt.copy1
space-998.txt.copy6
space-998.txt.copy7
space-998.txt.copy5
space-933.txt.copy7
space-933.txt.copy1
space-933.txt.copy6
space-933.txt
space-933.txt.copy3
space-933.txt.copy4
space-933.txt.copy5
space-933.txt.copy2
space-810.txt.copy5
space-912.txt.copy7
space-846.txt.copy3
space-827.txt.copy7
space-939.txt.copy3
space-831.txt.copy6
space-813.txt

Space-961.txt – threshold .7:

space-961.txt.copy5
space-961.txt.copy4
space-961.txt.copy1
space-961.txt.copy6
space-961.txt.copy7
space-961.txt.copy2
space-961.txt.copy3
space-968.txt.copy3
space-968.txt.copy5
space-989.txt.copy1
space-872.txt.copy5
space-851.txt.copy5
space-846.txt.copy1
space-846.txt.copy3
space-846.txt.copy2
space-979.txt
space-969.txt.copy1
space-806.txt.copy1
space-930.txt
space-930.txt.copy6
space-828.txt.copy4
space-956.txt.copy5
space-988.txt.copy7
space-943.txt.copy1
space-946.txt.copy6
space-943.txt.copy3
space-841.txt.copy5
space-843.txt.copy3
space-843.txt
space-843.txt.copy1

- threshold .9:

space-977.txt.copy6
space-977.txt.copy7
space-955.txt.copy5
space-955.txt.copy4
space-977.txt.copy4
space-955.txt.copy6
space-955.txt.copy1
space-955.txt.copy7
space-955.txt
space-977.txt.copy1
space-940.txt.copy5
space-846.txt.copy1
space-940.txt.copy3
space-940.txt.copy4
space-977.txt.copy2
space-940.txt.copy7
space-977.txt.copy3
space-940.txt
space-902.txt.copy3
space-977.txt.copy5
space-839.txt.copy1
space-886.txt.copy5
space-922.txt.copy7
space-887.txt.copy6
space-901.txt.copy4
space-830.txt.copy1
space-890.txt.copy2
space-943.txt.copy5
space-886.txt.copy2
space-886.txt.copy7
space-886.txt
space-813.txt.copy6
space-825.txt.copy2
space-819.txt.copy6
space-881.txt.copy3
space-822.txt.copy6
space-999.txt.copy1
space-999.txt

Space-977.txt – threshold .7:

space-977.txt.copy1
space-977.txt.copy2
space-977.txt.copy5
space-977.txt.copy3
space-977.txt.copy7
space-977.txt.copy4
space-977.txt.copy6
space-954.txt.copy4
space-974.txt.copy5
space-887.txt.copy4
space-887.txt
space-979.txt.copy2
space-891.txt.copy6
space-962.txt.copy3
space-891.txt.copy4
space-909.txt.copy6
space-806.txt.copy5
space-882.txt.copy1
space-830.txt.copy7
space-830.txt.copy3
space-816.txt.copy4
space-905.txt.copy7
space-810.txt.copy4
space-935.txt.copy6
space-935.txt.copy3
space-837.txt.copy6
space-837.txt
space-884.txt.copy7

- threshold .9:

space-985.txt.copy3
space-985.txt.copy5
space-985.txt.copy2
space-985.txt.copy1
space-985.txt.copy6
space-985.txt.copy4
space-985.txt.copy7
space-942.txt.copy3
space-957.txt.copy7
space-905.txt.copy5
space-910.txt.copy4
space-866.txt.copy4
space-953.txt.copy1
space-979.txt.copy4
space-919.txt.copy1
space-838.txt.copy2
space-822.txt.copy3
space-932.txt.copy7
space-966.txt.copy4
space-814.txt.copy3
space-800.txt.copy7
space-860.txt.copy3
space-1000.txt.copy1
space-880.txt.copy5
space-880.txt.copy4
space-880.txt.copy6
space-880.txt.copy7
space-880.txt
space-845.txt.copy3
space-902.txt.copy7
space-877.txt.copy7
space-959.txt.copy1
space-935.txt.copy1
space-833.txt.copy4
space-833.txt.copy1
space-833.txt

Space-985.txt – threshold .7:

- threshold .9:

space-985.txt.copy4
space-985.txt.copy5
space-985.txt.copy2
space-985.txt.copy7
space-985.txt.copy6
space-985.txt.copy1
space-985.txt.copy3
space-854.txt.copy7
space-840.txt.copy6
space-942.txt.copy1
space-942.txt.copy7
space-942.txt.copy4
space-912.txt.copy1
space-898.txt.copy5
space-893.txt.copy5
space-976.txt.copy3
space-941.txt.copy6
space-975.txt.copy7
space-983.txt.copy2

```

space-838.txt.copy4
space-993.txt.copy7
space-993.txt.copy6
space-993.txt.copy2
space-993.txt.copy5
space-961.txt.copy5
space-993.txt.copy1
space-993.txt.copy3
space-993.txt.copy4
space-943.txt.copy1
space-943.txt.copy7
space-943.txt.copy5
space-943.txt
space-968.txt.copy7
space-968.txt.copy1
space-968.txt.copy4
space-968.txt.copy3
space-968.txt.copy5
space-968.txt
space-825.txt.copy1
space-825.txt.copy7
space-825.txt.copy5
space-825.txt.copy2
space-825.txt.copy3
space-825.txt
space-818.txt.copy7
space-833.txt.copy1
space-872.txt.copy4
space-854.txt.copy1
space-854.txt.copy6
space-854.txt

```

Space-993.txt – threshold .7:

```

space-993.txt.copy6
space-993.txt.copy3
space-993.txt.copy2
space-993.txt.copy1
space-993.txt.copy5
space-919.txt.copy7
space-801.txt.copy3
space-993.txt.copy7
space-993.txt.copy4
space-974.txt.copy2
space-974.txt.copy4
space-974.txt.copy6
space-806.txt.copy5
space-988.txt.copy1
space-988.txt.copy7
space-883.txt.copy2
space-973.txt.copy2
space-973.txt
space-870.txt.copy7
space-848.txt.copy2
space-985.txt.copy4
space-807.txt.copy7
space-990.txt.copy5
space-990.txt.copy6
space-990.txt

```

- threshold .9:

From these outputs I can conclude that somewhere along in the process of getting the nearDuplicate documents my data is inconsistent. This is because when I run a threshold of .9 we should only get the outputs of the copies, however this is not the case. I believe it may have something to do with how I get the permutations. But I am unsure. I have looked into this for the last week of the project, but I have yet to come to any conclusions. This data is without removing false positives. Also, each output does have all of the copies of the given file.