

## Part 1: Parse the Corpus

1.

There are four main steps in this part: remove stopwords, stemming words, calculate TFIDF for each word in each document, calculate cosine similarities.

### Removing stopwords:

Given a list of documents, find all the unique words and store them in a *bag\_of\_words*

For each word in the *bag\_of\_words*, count the number of documents that contain such words, denoted as *count*

Calculate the idf using  $\log(N / \text{count})$ , where N is the total number of documents in the corpus.

The output shows a large number of words that only occur in one or two documents, which results in a very large idf score ( $\log(N/1) = 7.9$ ). Hence, I chose all the words that have idf larger than 95% of the highest idf score as my stopwords.

```
['circuitry', 7.914983005848394), ('localized', 7.914983005848394), ('appropriation', 7.914983005848394), ('812000', 7.914983005848394), ('tbody', 7.914983005848394), ('broiler', 7.914983005848394), ('woods', 7.914983005848394), ('6722527', 7.914983005848394), ('hanging', 7.914983005848394), ('sustaining', 7.914983005848394), ('diseasecontrol', 7.914983005848394), ('saddams', 7.914983005848394), ('centalized', 7.914983005848394), ('designing', 7.914983005848394), ('metalurgicos', 7.914983005848394), ('nawg', 7.914983005848394), ('perforations', 7.914983005848394), ('chinh', 7.914983005848394), ('sickline', 7.914983005848394), ('kids', 7.914983005848394), ('spotty', 7.914983005848394), ('hurricaneforce', 7.914983005848394), ('wineapple', 7.914983005848394), ('46864', 7.914983005848394), ('1001900', 7.914983005848394), ('hodler', 7.914983005848394), ('dnj', 7.914983005848394), ('mitsuya', 7.914983005848394), ('populations', 7.914983005848394), ('163501', 7.914983005848394), ('ninety-six', 7.914983005848394), ('locked', 7.914983005848394), ('locker', 7.914983005848394), ('118800', 7.914983005848394), ('13050', 7.914983005848394), ('absolute', 7.914983005848394), ('sdwkl', 7.914983005848394), ('cutback', 7.914983005848394), ('nazmi', 7.914983005848394)
```

Words with low idf ( $\log(N/300) = 2.197$ ) will also be removed because it indicates that this word occurs in most of the documents. Therefore, I chose all words that have idf smaller than 1.05% of the lowest idf score as my stopwords.

```
(2.2837712240270287), ('five', 2.2801934026791444), ('international', 2.255500790088773), ('due', 2.2520225257124484), ('dlr', 2.2451020828678745), ('march', 2.2416597386769017), ('group', 2.2416597386769017), ('sales', 2.2280076495085743), ('over', 2.2012502003390253), ('per', 2.1913979038960134), ('shares', 2.1848832228748196), ('told', 2.17519009366916), ('years', 2.171979818038912), ('they', 2.1592407922614822), ('more', 2.1497919030635497), ('three', 2.146662010054622), ('1987', 2.1435418827183783), ('after', 2.1250228349511406), ('shr', 2.118925255083022), ('market', 2.10684051586795), ('than', 2.100852474023328), ('other', 2.0831005285648776), ('bank', 2.0831005285648776), ('stock', 2.071438588170344), ('first', 2.04568609271462), ('share', 2.04568609271462), ('up', 1.9800888102288066), ('may', 1.9591456363835633), ('been', 1.9037158314442326), ('april', 1.8818967840495928), ('two', 1.8327640954719482), ('or', 1.8327640954719482), ('1986', 1.8214132358032586), ('had', 1.8191584434161694), ('billion', 1.8124244112348253), ('co', 1.7968858078070462), ('about', 1.7794181147666552), ('net', 1.7707973717227483), ('vs', 1.737038891797794), ('one', 1.724667599995247), ('he', 1.724667599995247), ('also', 1.7104252432797038), ('last', 1.6904245765730344), ('but', 1.6650077635889111), ('this', 1.6496817931106842), ('new', 1.6477824573070319), ('have', 1.6345871668881903), ('cts', 1.6234138662900741), ('us', 1.5996250043260594), ('were', 1.5870462221191992), ('are', 1.571102571220632), ('would', 1.5623536095288277), ('corp', 1.4933607380418763), ('not', 1.448838281610775), ('which', 1.4167008563719603), ('as', 1.397311732936119), ('inc', 1.3696333455139746), ('year', 1.2595426554807467), ('company', 1.2544078560087082), ('pct', 1.2179487581819097), ('has', 1.2117948926075315), ('a', 1.197178310824703), ('that', 1.1959698514631345), ('was', 1.182772299381188), ('at', 1.0585210212538076), ('with', 1.0213266512457588), ('be', 1.0213266512457588), ('will', 1.0122402686898013), ('from', 0.9893878087379261), ('on', 0.9423767545466408), ('by', 0.9377016642176469), ('dlrs', 0.9367692632176958), ('is', 0.9028677115420144), ('its', 0.891224051109951), ('min', 0.8570850684365375), ('it', 0.6192479330991123), ('for', 0.5453822853219851), ('in', 0.42855279081682867), ('a', 0.37223946047984396), ('to', 0.33428325362383127), ('and', 0.3150810466389545), ('said', 0.23619400764924056), ('of', 0.2315793247945681), ('the', 0.2210456803912245)]
```

In this case, I removed words that only appeared in one document (high idf) and words that appeared in approximately 300 documents (low idf) as my stopwords.

### Stemming words:

I used PorterStemmer from nltk package to stem words in each document.

### TFIDF:

First of all, I calculated IDF for each word in the same way as above. Get the number of count of each word in each document, find the idf score for such a word. Get TFIDF by multiplying tf with idf, and store the TFIDF score for each word.

**Cosine Similarity:**

For each document, calculate the cosine similarity with every other document in the corpus (including itself).

2.

Packages used:

BeautifulSoup4 for parsing the sgm file. Using findAll to find the NEWID and BODY

Nltk for tokenize document into a list of words, PorterStemmer for stemming words.

## **Part 2: Run Complete and Single-Linkage Clustering**

I used the dendrogram library from *scipy.cluster.hierarchy* package to plot the graph. We can utilize this to plot both single and complete clusters. Each function returns a matrix *z*. *Z* contains four columns, which specifies [cluster, cluster, distance, number of nodes] in each row. I used this matrix *z* to find all clusters for each document.

Since the first and second column in the matrix *z* represents the clusters, and each document is also a cluster itself. I put all documents under the same cluster into a list and the key is the corresponding clusterID. In this case, I union two clusters, which already in the dictionary, by grouping all the documents under these two child clusters.

One thing that I need to change is that the highest cluster number in the matrix represents the largest cluster, whereas the smallest cluster number is required to be the largest cluster in this project. I changed the key (clusterID) to meet the requirement.

I did the same thing for complete clustering.

One interesting thing is that documents have much more common top levels clustering in single.txt than in complete.txt. This would result in a higher and thinner single clustering.

However, there are less common top level clustering in complete.txt, which means that the dendrogram will be lower and wider. This happens because single clustering uses the nearest point algorithm, whereas the complete clustering uses the farthest point algorithm.

## Part 3: Evaluation

### 1. Number of clusters

If we assume that each cluster only has two children in the dendrogram, then the total number of clusters should be  $2N-1$ , where  $N$  is the number of documents in the corpus (each document is a cluster itself).

### 2. Evaluation Metric Overview

The way I derive the evaluation metric is to utilize the `<TOPIC>` tag in the corpus. Each document can have one or more topics, and documents with the same topic should be clustered together.

### 4. Why this is a good measure

I think it is a good measure because topic is a way of clustering documents. If two documents have more common topics, these two documents should have more common clusters.

### 3. Calculate # of common topics (denoted as `common_topic_score`)

First, I extracted all the unique topics in the corpus. There should be one or more documents containing such topic. If two documents are having the same topic, I updated the # of common topics by one each time. I calculated all # of common topics for every two documents by looping through all the unique topic, and this yields a 2D array.

### 4. Calculate # of common clusters (normalized) (denoted as `clustering_score`)

Second, for each two documents, union their clusters and get the total number. Additionally, find the number of intersecting clusters of these two documents. Using the number of intersecting divided by the number of union should give us a normalized score.

### 5. Get the relationship between `common_topic_score` : `clustering_score`

In the above steps, we get the number # of common topics and # of common clusters (normalized). However, different # of common topics might have a list of different `clustering_score`. There might exist different cluster\_score for a same `common_topics_score`. (e.g: {1: [0.2223, 0.3334, 0.007]}). In this case, I chose to use the median of the list because the median is not affected by outliers.

6.

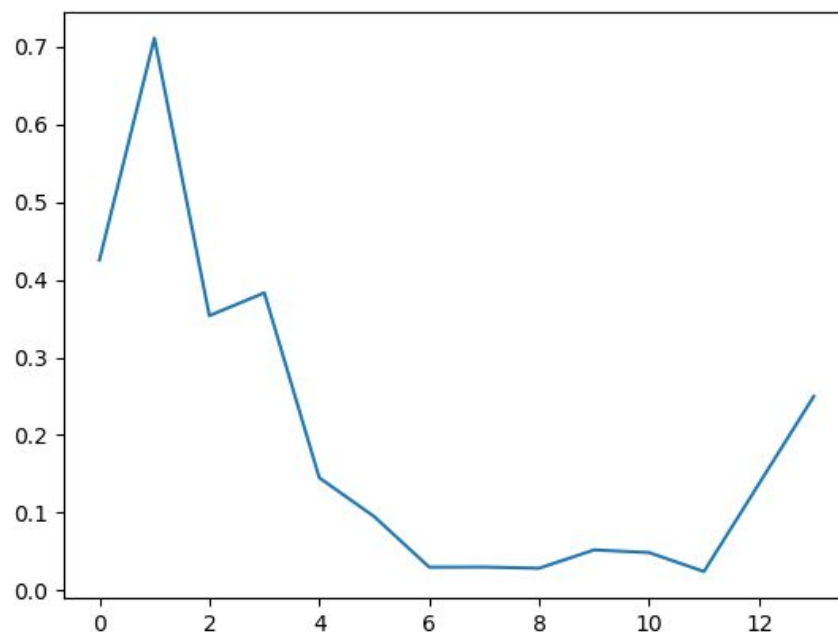
Single Percentage:

{0: 0.4253968253968254, 1: 0.7110608455189744, 2: 0.35374149659863946, 3: 0.383133637313193, 4: 0.1450653983353151, 5: 0.09512292164007204, 6: 0.030017921146953404, 7: 0.030175337971186107, 8: 0.02857142857142857, 9: 0.0521759697256386, 10: 0.04878048780487805, 11: 0.024390243902439025, 13: 0.25}

Best-Score achieved: 0.7110608455189744

The x-axis is the common\_topic\_score

The y-axis is clustering\_score



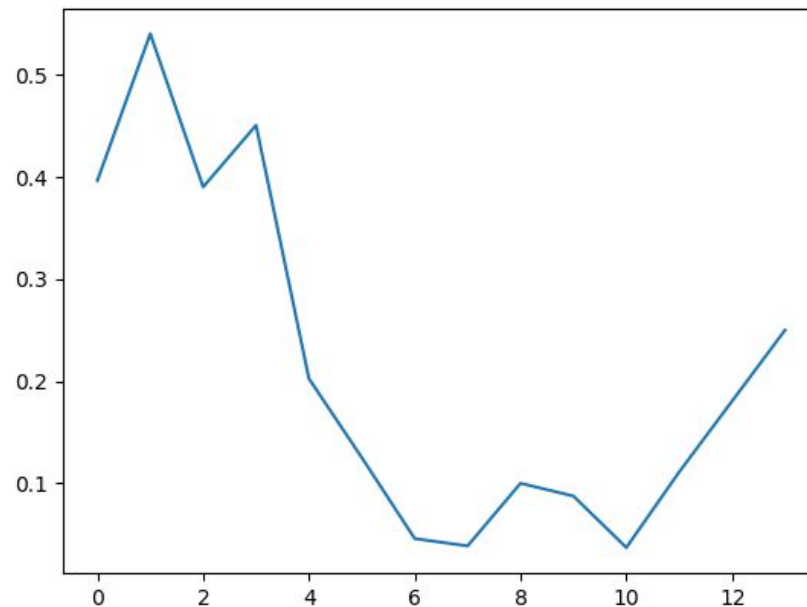
Complete Percentage:

{0: 0.39644970414201186, 1: 0.5403225806451613, 2: 0.3902439024390244, 3: 0.450713453698311, 4: 0.20253164556962025, 5: 0.12508218277449046, 6: 0.04583189207886544, 7: 0.0386302294197031, 8: 0.1, 9: 0.0875, 10: 0.037037037037037035, 11: 0.11111111111111111, 13: 0.25}

Best-Score achieved: 0.5403225806451613

The x-axis is the common\_topic\_score

The y-axis is clustering\_score



As we can see above, I think the clustering is not very good. The expected result should be a positive relation between `common_topic_score` and `clustering_score`. As the number of `common_topic` increases, the number of common clusters should also be increased. However, the graph does not show an explicit positive relationship.

#### 6. Discuss the different in performance between single-linkage and complete-linkage

The single-linkage results in a longer time than complete-linkage does.

I think the reason is that the height of single-linkage is higher than complete-linkage. Just like a binary tree, if the depth is large, it would result in a long time to get to the bottom. On the other hand, complete-linkage has a small depths and it spreads more. Hence, complete-linkage would result a better performance.