

- **Description of the system, design tradeoffs, questions you had and how you resolved them, etc. List the software libraries you used, and for what purpose.**

This system builds on top of previous projects and uses the same design as far as the indexer is concerned. The indexer was modified to create a new document vectors file. The whole project is split into the indexer and retrieval packages, and the retrieval package houses the cluster package. One of the design choices I made was to read the entire document vectors file into memory to be able to access it easily instead of compressing or reading it part-by-part. I also decided to manually dispatch (use switch case) for scoring method rather than using abstract classes and subclasses. Another design choice I made was to store the modified document vectors (with tf-idf weights) in the cluster as opposed to the ones containing pure term counts. I used google's gson jar to easily read/write the document vectors to disk.

- Explain what happens as the threshold value increases. Table the clustering results (number of clusters, size of each) for all of the threshold values.

As the threshold value increases, so does the possibility of no existing cluster matching the threshold requirement, leading to an increase in creation of new clusters. This is what we observe as the threshold value increases, the number of clusters also increases

[illegible]

- **Explain the difference between the four linking strategies. What behavior would you expect from each of them. Compare to the results you obtained using mean.**

Single linking strategy: Candidate document's similarity is calculated against every document in the cluster, and the minimum value encountered is considered to be the representative value for that cluster in comparison to other clusters. Clusters will be more stringy and the distance between the farthest points will be large as the cluster grows in size as compared to mean linkage

Complete linking strategy: Candidate document's similarity is calculated against every document in the cluster, and the maximum value encountered is considered to be the representative value for that cluster in comparison to other clusters. Clusters will be denser as compared to mean linkage

Average linking strategy: Candidate document's similarity is calculated against every document in the cluster, and the average of all values encountered is considered to be the representative value for that cluster in comparison to other clusters. The results will be quite similar to those of mean linkage

Mean linking strategy: Candidate document's similarity is calculated against the mean document vector