

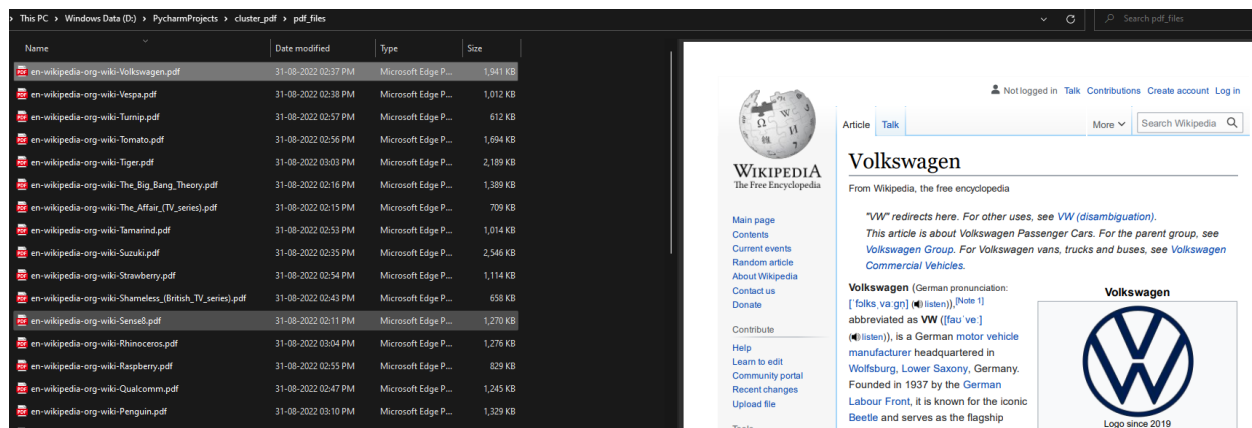
Unsupervised Machine Learning Course Project: Automatic Folder Organization

Objective:

A clustering model was proposed to divide a folder containing multiple PDF documents into different folders, where each folder has PDF files corresponding to a topic. The benefit of this application is that a user can reorganize a folder automatically based on the text present inside them.

Data:

84 Wikipedia pages were scraped and saved as a PDF (using webtopdf.com) in a folder (as shown below in screenshot). It consisted of pages of TV series/movies, automobile companies, fruits and vegetables, birds and animals, software companies, semiconductor companies. Each webpage when converted to PDF resulted in multiple pages, but only the text from the first page from each PDF was chosen for clustering the document.



Data Exploration and Feature Engineering:

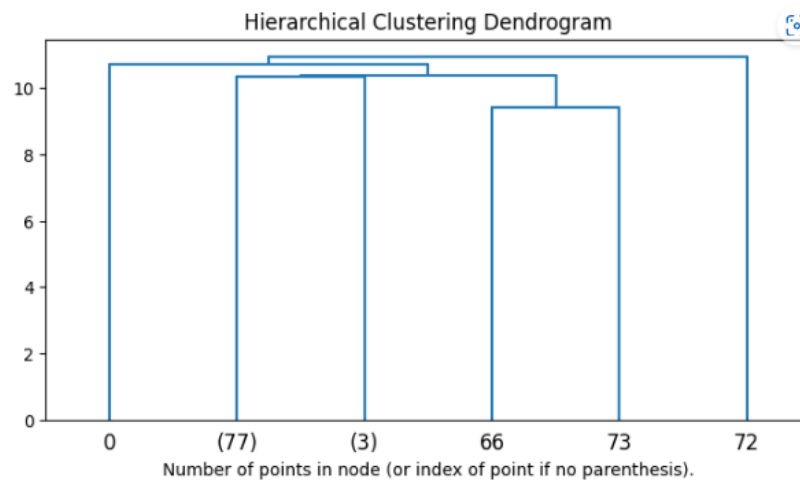
Only the first page of each PDF was used to reduce the time taken for clustering the PDF files and move into a new folder. For each first page in a PDF, the stop words were removed, the remaining words were then lemmatized and TD-IDF (Term Frequency-Inverse Document frequency) vectors were computed. Only unigrams were chosen for creating the TF-IDF vector as presence of same bigram in documents of same topic was less common. The max number of features for the TD-IDF vectorizer was chosen as 1024.

Hierarchical Agglomerative Clustering

Since the number of clusters for a given set of documents is not known and the number of PDFs corresponding to a cluster may not be even, hierarchical clustering was chosen. It is also scalable in terms of large number of documents to be clustered and large number of clusters which might be present.

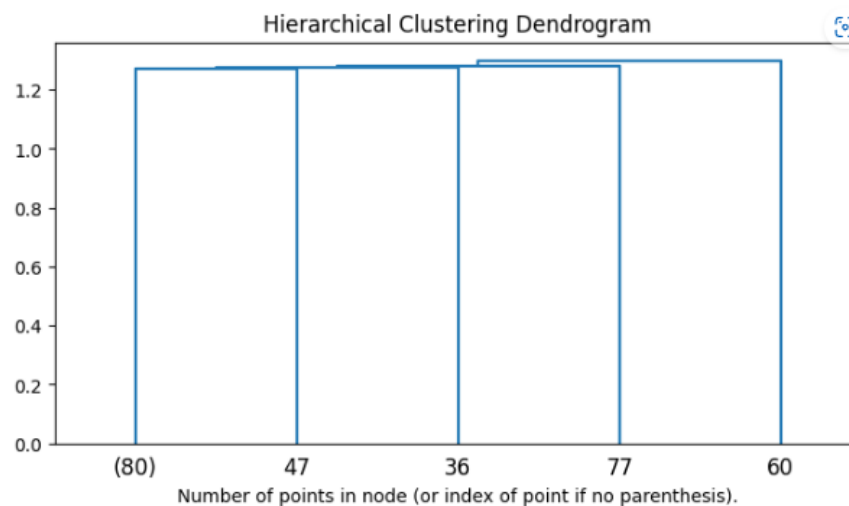
Single Linkage with L1 distance metric

From the dendrogram, we saw that it **failed to cluster the document properly** as most (77) of the documents fell into a single cluster.



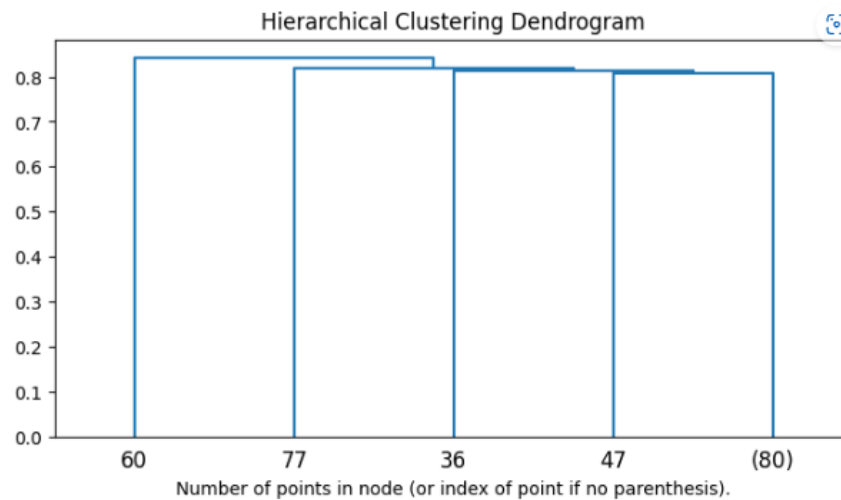
Single Linkage with L2 distance metric

From the dendrogram, we saw that it **failed to cluster the document properly** as most (80) of the documents fell into a single cluster.



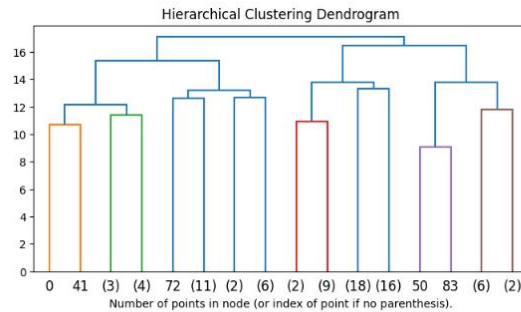
Single Linkage with cosine distance metric

From the dendrogram, we saw that it **failed to cluster the document properly** as most (80) of the documents fell into a single cluster.



Complete Linkage with L1 distance metric

The dendrogram for complete linkage is shown below. The distance threshold was then chosen as 14 and the file names for each of the clusters thus obtained was plotted in the form of word cloud. Most of the movies/ TV shows, fruits and vegetables, and companies (in 2 clusters) were clustered separately. **The type of company like semiconductor, IT or automobile company got mixed up in this clustering.**

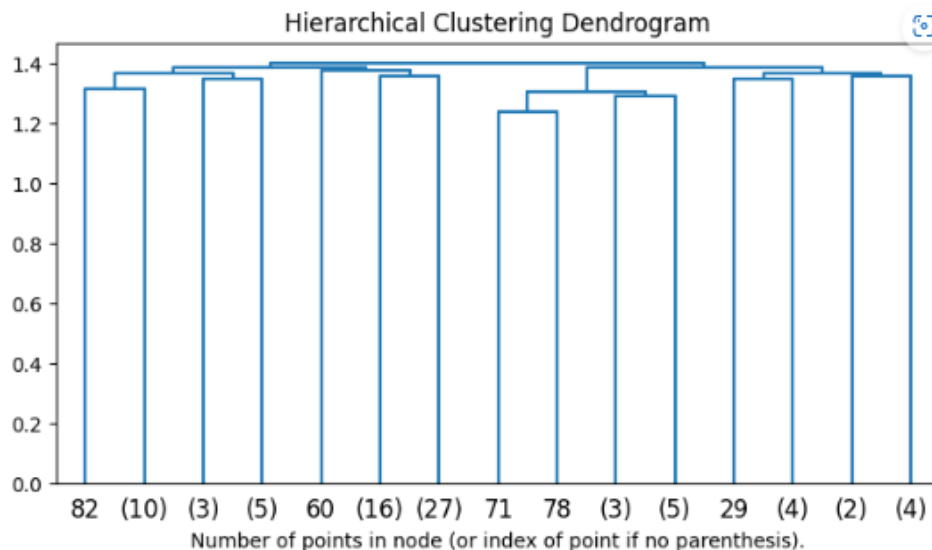


The_Affair_Game_of_Thrones_
The_Big_Bang_Theory_pdf
Friends_pdf
Shameless_
_Young_Sheldon_pdf
Narcos_pdf
House_of_Cards_
American_TV_series
British_TV_series
Sense8_pdf
How_I_Met_Your_Mother_pdf
season_8
TV_series
_Netflix_pdf
_Hero_MotoCorp_pdf
_Meta_Platforms_pdf
_Intel_pdf
Audi_pdf
Bajaj_Auto_p
Honda_pdf
Royal_Enfield_pdf
BMW_pdf
MercedesBenz_pdf
Nvidia_pdf
NSU_Motorenwerke_pdf
_Microsoft_pdf
Cypress_Semiconductor_pdf
_KTM_pdf
HewlettPackard_pdf
_Microelectronic_pdf

Penguin_pdf
Cauliflower_pdf
Banana_pdf
Cattle_pdf
Garlic_pdf
Lychee_pdf
Gooseberry_pdf
Old_World_sparrow_pdf
Okra_pdf
Ginger_pdf
Elephant_pdf
Hippopotamus_pdf
Chili_pepper_pdf
Beetroot_pdf
Carrot_pdf
Advanced_Micro_Devices_pdf
Volkswagen_pdf
Apple_Inc_pdf
Qualcomm_pdf
company_pdf
Google_pdf
Amazon_
Suzuki_pdf
Ducati_Motor_Holding_pdf

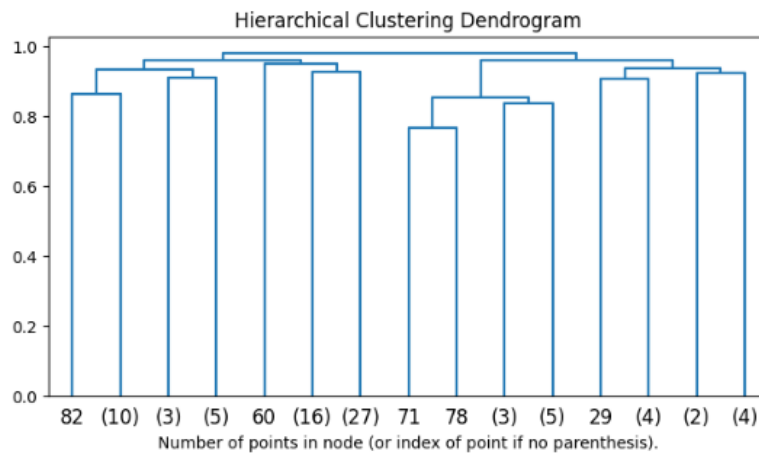
Complete Linkage with L2 distance metric

From the dendrogram, we saw that the range of threshold for which meaningful clusters were getting formed was very small. So choosing a **threshold would be very sensitive**.



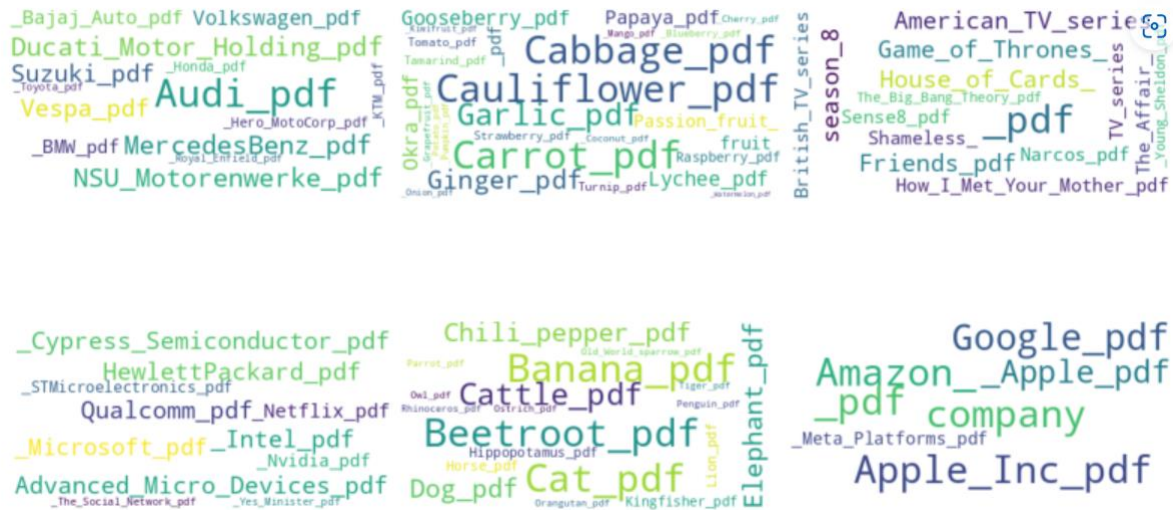
Complete Linkage with cosine distance metric

From the dendrogram, we saw that the range of threshold for which meaningful clusters were getting formed was very small. So choosing a **threshold would be very sensitive**.



WARD Linkage with L2 distance metric

The dendrogram for ward linkage is shown below. The distance threshold was then chosen as 1.5 and the file names for each of the clusters thus obtained was plotted in the form of word cloud. It can be seen that most of the automobile companies, fruits and vegetables, movies/ TV shows, semiconductor companies, fruits and vegetables, and internet companies were clustered separately. **This has the best formed clusters.**



Summary

A clustering approach was used to automatically restructure a folder of PDF documents into folder of PDF files belonging to a topic based on the text present in them. Hierarchical Agglomerative Clustering was successfully used for this task with WARD Linkage with Euclidean distance metric. It was able to cluster various topics present in the dataset including clusters which were fine grained like semiconductor, internet and automobile companies.

Code

Github Link : [CV-and-Neural-Nets-Basic/Automatic Folder Organization at master · kritiksoman/CV-and-Neural-Nets-Basic \(github.com\)](https://github.com/kritiksoman/CV-and-Neural-Nets-Basic)

Future Work

Adding more features like file metadata, support for different file formats, using images embeddings as features are few improvements which can be made to this work. The clustering approach should also have support for handling outliers. All these are missing in the current approach.