

Report - COVID19 papers (CORD-19)

```
In [1]: # Itamar Fradkin - 312531864
# Ron Boxer - 321219088

first you need to download the files from Here
Note- After downloading the zip file -
Unzip - metadata.csv & pdf_json folder (Keep the sturcted folder((document_parses))).
Insert this files into 'data' folder & create new folder named- 'minidataset' in 'document_parses' folder

We created several utils which can be found in 'utils' folder.

let's load and take a look on the metadata.

In [2]: import json
import os
import cotools
import shutil
import pandas as pd
import numpy as np

In [3]: DATASET_PATH = 'data/'
METADATA_PATH = 'data/metadata.csv'
MINI_DATASET_JSON_PATH = 'data/document_parses/minidataset/'

In [4]: import utils.dataset as ut
df = ut.load_metadata()

C:\Users\ronbo\anaconda3\lib\site-packages\Python\core\Interactiveshell.py:3338: DtypeWarning: Columns (1,4,5,6,13,14,15,16) have mixed types.Specify dtype option on import or set low_memory=False.
if (await self.run_code(code, result, async_asy))):

In [5]: df

Out[5]:
```

	cord_uid	source_x	title	doi	pmcid	pubmed_id	license	abstract	publish_time	authors	journal	mag_id	who_covidence_id	arxiv_id	pdf_json	
	sha															
d7ea1370848eac24ae421c20fc30a69d9fa6566e	9w762whh	Elsevier; PMC; WHO	Patients With Myocarditis Associated With COVID...	10.1016/j.jacc.2022.02.004	PMC8958986		NaN	els-covid	NaN	2022-04-05	Rosner, Carolyn M.; Melary; Saeed, Ibr...	Journal of the American College of Cardiology	NaN	NaN	NaN	document_parses/pdf_json/d7ea1370848eac24ae4
e4494dc889caade275c0fc401e7e8a3d9ae63	6lqjmq0	Elsevier; Medline; PMC; WHO	Electrochemical immunosensor for rapid and hig...	10.1016/j.talanta.2022.123211	PMC8730781	3.49993e+07	no-cc	A label-free electrochemical biosensing approx...	2022-04-01	Mehmandoust, Mohammad; Gumsu, Z. Pinar; Soyлак...	Talanta	NaN	NaN	NaN	document_parses/pdf_json/e4494dc889caade275c0fc401e7e8a3d9ae63	
12bc48a0b3196eef7e887807130b0853e8d47190	vmsyu6w7	Elsevier; Medline; PMC; WHO	Entrepreneurial Interventions for crisis manag...	10.1016/j.jdr.2022.102830	PMC8865137	3.5229e+07	no-cc	This article investigates both the negative an...	2022-04-01	Krishnan, Commander S. Navaneetha; Ganesh, L.S...	Int J Disaster Risk Reduct	NaN	NaN	NaN	document_parses/pdf_json/12bc48a0b3196eef7e	
fea63a0a7eb45059214491c2cb8947fa98b5f65c	bx49s7f8	Medline; PMC; WHO	A Qualitative Examination of COVID-19's Impact...	10.1177/08901171211053845	PMC8851044	3.4787e+07	no-cc	PURPOSE: The COVID-19 pandemic is correlated w...	2022-03-31	Bramon, Grace Ellen; Mitchell, Sophia; Ray, M...	Am J Health Promot	NaN	NaN	NaN	document_parses/pdf_json/fea63a0a7eb45059214491c2cb8947fa98b5f65c	
631cf2b3762b36ecfe68675b08c34351fb09	aagrsfth	Elsevier; PMC; WHO	Population Mobility and Socioeconomic Indicato...	10.1016/j.jid.2021.12.061	PMC8884814		NaN	Purpose To explore the extent that socioeconomic...	2022-03-31	Marwah, A.; Moineddin, R.; Thomas, R...	International Journal of Infectious Diseases	NaN	NaN	NaN	document_parses/pdf_json/631cf2b3762b36ecfe68675b08c34351fb09	
...
cbd8c8f742b75a032e1b041da900a6a51623cc79	htwp70hc	Medline; PMC; WHO	Comparison of chest CT severity scoring system...	10.1007/s00330-021-08432-5	PMC8760133	3.50318e+07	no-cc	PURPOSE: To compare the diagnostic performance...	2022-01-15	Elmokadem, Ali H.; Mourir; Ahmad M.; Ramadan, ...	Eur Radiol	NaN	NaN	NaN	document_parses/pdf_json/cbd8c8f742b75a032e1b041da900a6a51623cc79	
5bacfca63b11318e9788197ca3005a1135138956	4jqc0ny	PMC	Regarding the articles on home sponometry	10.1016/j.jcd.2022.01.002	PMC8760931	35042654	no-cc	NaN	2022-01-15	Curley, Rachael; Campbell, Michael J; Walters...	J Cyst Fibros	NaN	NaN	NaN	document_parses/pdf_json/5bacfca63b11318e9788197ca3005a1135138956	
637abda16edcb7c7245861431108ff825aeaa80f	7221q12	Medline; PMC	Missed opportunities to identify cryptococcosi...	10.1177/20499361211066363	PMC8771738	3.50703e+07	cc-by-nc	SARS-CoV-2 may activate both innate and adapt...	2022-01-15	Chastain, Daniel B.; Henaao-Martinne; Andrés F...	Ther Adv Infect Dis	NaN	NaN	NaN	document_parses/pdf_json/637abda16edcb7c7245861431108ff825aeaa80f	
49ae4f05cc3d5dadcd9see8d7df2a6cbeafcb8fb	33c5xm2	Medline; PMC	Using Simulation Optimization to Solve Patient...	10.3390/healthcare10010164	PMC8775607	3.50523e+07	cc-by	This study investigates patient appointment sc...	2022-01-15	Chen, Ping-Shun; Chen, Gary Yu-Hsin; Liu, Li-W...	Healthcare (Basel)	NaN	NaN	NaN	document_parses/pdf_json/49ae4f05cc3d5dadcd9see8d7df2a6cbeafcb8fb	
a4bb5b7bf9cf309d443cf439d29502cc30e1a5	tbht2hc	Elsevier; Medline; PMC; WHO	Livelihood challenges and healthcare-seeking b...	10.1016/j.aquaculture.2021.737348	PMC8414286	3.44939e+07	no-cc	The outbreak of coronavirus disease (COVID-19)...	2022-01-15	Hossain, Md. Tanvir; Lima, Taposhi Rabya; Ela...	Aquaculture	NaN	NaN	NaN	document_parses/pdf_json/a4bb5b7bf9cf309d443cf439d29502cc30e1a5	

27000 rows × 18 columns

As we only need some of the data, we move a small cohort into an entirely separate folder. (It will take a while :))

```
In [6]: ut.move_files(df)

Now we will use outsource package - 'cord-19-tools' to help us load the papers effiantly.

In [7]: pip install cord-19-tools

Requirement already satisfied: cord-19-tools in c:\users\ronbo\anaconda3\lib\site-packages (0.3.3)
Requirement already satisfied: xmltodict in c:\users\ronbo\anaconda3\lib\site-packages (from cord-19-tools) (0.12.0)
Note: you may need to restart the kernel to use updated packages.

In [8]: # Load papers
import cotools
data = cotools.PaperSet(MINI_DATASET_JSON_PATH)

In [9]: covid_df = ut.load_papers_into_df(data, size=20000)
covid_df.head()

Out[9]:
```

	title	paper_id	abstract	body_text
0	Recruitment methods and yield rates in a clini...	000e6f401f046b240cc8312c9ef41914bfb9e06	Background: Although the prevalence of hyperten...	The prevalence of hypertension is increasing w...
1	Impact of the COVID-19 Pandemic on Daily Life...	0012308e8c02792b907c31935868f1abff2c59b	Using a mixed methods design, this study aimed...	travel restrictions. These measures aimed to p...
2	Colchicine Against SARS-CoV-2 Infection: What ...	0015cecc2299c3bdb9bd0e0b4b38ebdcca716f	Coronavirus disease 2019 caused by the severe ...	Systemic inflammation is the hallmark of coron...
3	Use and misuse of prescription stimulants by u...	001716d1b1c1e9d9c1c8b3751e6203712ddfd24	Background: Misuse of prescription stimulants u...	Global Drug Survey, an increase of neuroenhan...
4	A Neural Phillips Curve and a Deep Output Gap	001746a8ab396aa3a7bdc0b45ba75f26ec860fa9	Many problems plague the estimation of Philip...	Few equations are as central to modern macroec...

Part A - Paper Similarity using gzip as our compressing method on text data only
Normalized compression distance (NCD) is a way of measuring the similarity between two objects, be it two documents, two letters, two emails, two music scores, two languages, two programs, two pictures, two systems, two genomes, to name a few. Such a measurement should not be application dependent or arbitrary. A reasonable definition for the similarity between two objects is how difficult it is to transform them into each other. [Further information can be found on this blog post](#)

If two objects compress better together than separately, it means they share common patterns and are similar!!
here is the formula-



Let's take a paper as an example. Calculate the NCD distance to all other papers and take the 10 closest papers.

```
In [10]: import utils.model as model
paper_id = '88b80f02e54d2ed8c9f05f422b45f4818a6a405b' # Example
k = 10
similarity_df = model.k_similar_papers(paper_id, k, covid_df, txt_to_encode='abstract')
similarity_df

Out[10]:
```

	title	paper_id	abstract	body_text	ncd_distance
14448	Explainable machine learning to predict long-t...	88b80f02e54d2ed8c9f05f422b45f4818a6a405b	Background: Machine learning (ML) model is inc...	The long-term outcome is currently an emergin...	0.037691
14255	Comparative analysis of explainable machine le...	8710f527c4709f1ae2013a4d7688b1ec03f73aa	Background: Machine learning (ML) holds the pr...	underlying structure of the data, while ML mod...	0.777518
2529	Prediction of 3-year risk of diabetic kidney d...	18a795e3045a0f4ab3eb1be1198a805de0e1a	Background: Established prediction models of D...	clinical decision-making. Understanding the ri...	0.795195
9839	Predicting the necessity of oxygen therapy in ...	5d6ad19f8d11971024d1b7fa4a23e6cd7378cd15b	Medical oxygen is a critical element in the tr...	Since the spread of COVID-19 in December 2019...	0.805654
9625	Article 852736 1 (2022) A Promising Preopera...	5b9d41960cd050bc52122130019030e29f1b4ae2	Background: The non-invasive preoperative diag...	Hepatocellular carcinoma (HCC) is one of the m...	0.807205
10705	Mortality Predictive Value of APACHE II and SO...	6621c2bde72488b15a57f9d2fb8613a77461266b	Background: COVID-19 pandemic has become a glo...	In December 2019, severe acute respiratory syn...	0.813725
5295	Journal Pre-proofs Can we reliably automate cl...	32d17d00540281827e8a013e5198e953e7f114f	Background: Building Machine Learning (ML) mod...	What was already known on the topic: Classic...	0.818610
14952	Predicting ionized hyppocalcemia: External vali...	8dd1465054a2ef6c952fedf0ed9af9a61ea1a38ef	Background: Ionized hypocalcemia is common in ...	New York City was the epicenter of the COVID-1...	0.824499
15770	Early changes in laboratory tests predict live...	959423993ae9a583d8f26edf1842128dc086ae06	Background: Most patients with coronavirus dis...	Studies have shown that COVID-19 can affect ...	0.828033
8813	Comparing different machine learning technique...	53d8d480246f9722396c5c6e87c0384a40805ba514	Background: Coronavirus disease 2019 (COVID-1...	Coronavirus disease 2019 (COVID-19) has been a...	0.828685
24656	Comparison of Regression and Machine Learning ...	e997ad775d1eed1188a850c60f050a87c746851	Background: Depression is highly prevalent and...	Depression is the most common psychiatric diso...	0.828889

In this example, we found some similar papers using the comparison method. our paper was in the field of machine learning From the results, we can see that we got papers in the field of machine learning are also represented.

```
Part B - Clustering using Compression method from Part A
Our idea is to make several features based on the compression method. we saw in part A that actually, we can get some similar papers based on this method. so we decided to enrich the idea and create 3 features based on compression on the title, paper body text, and abstract text. we think this logic will find clusters with similar papers.

In [11]: import gzip

def compress(x):
    x = str(x).encode()
    l_x = len(gzip.compress(x))
    return l_x

# Apply Compression method from Part A
train = covid_df.copy()
train['comp_title'] = train['title'].apply(compress)
train['comp_body'] = train['body_text'].apply(compress)
train['comp_abstract'] = train['abstract'].apply(compress)

In [12]: # Pre- proscece before clustering - Normalization
from sklearn.preprocessing import StandardScaler, MinMaxScaler

normalizer = MinMaxScaler()
paper_ids = covid_df['paper_id']
train = train[['paper_id', 'comp_title', 'comp_body', 'comp_abstract']]
train['comp_title'] = normalizer.fit_transform(train[['comp_title']])
train['comp_body'] = normalizer.fit_transform(train[['comp_body']])
train['comp_abstract'] = normalizer.fit_transform(train[['comp_abstract']])
train

Out[12]:
```

	paper_id	comp_title	comp_body	comp_abstract
0	000e6f401f046b240cc8312c9ef41914bfb9e06	0.135011	0.026254	0.079546
1	0012308e8c02792b907c31935868f1abff2c59b	0.092677	0.050131	0.030028
2	0015cecc2299c3bdb9bd0e0b4b38ebdcca716f	0.069794	0.027313	0.050121
3	001716d1b1c1e9d9c1c8b3751e6203712ddfd24	0.127002	0.036918	0.055819
4	001746a8ab396aa3a7bdc0b45ba75f26ec860fa9	0.050343	0.141596	0.166362
...
26979	f1eb944868a31900f661d78835ae1b22655fb22	0.137300	0.021210	0.072404
26980	f1eba2e1ac4562c93da63f732f73cea26327793	0.081236	0.031219	0.047636
26986	m951c1c32133ecdfef80506d630494052b126b	0.077803	0.049735	0.010156
26988	ffac3b29c594a89bc2e2db9a844a496c540ad913c	0.148741	0.022383	0.000658
26990	ffed5cd958b4eb441f517c58f6d04990fe2c7c2	0.141876	0.019907	0.058815

17664 rows × 4 columns

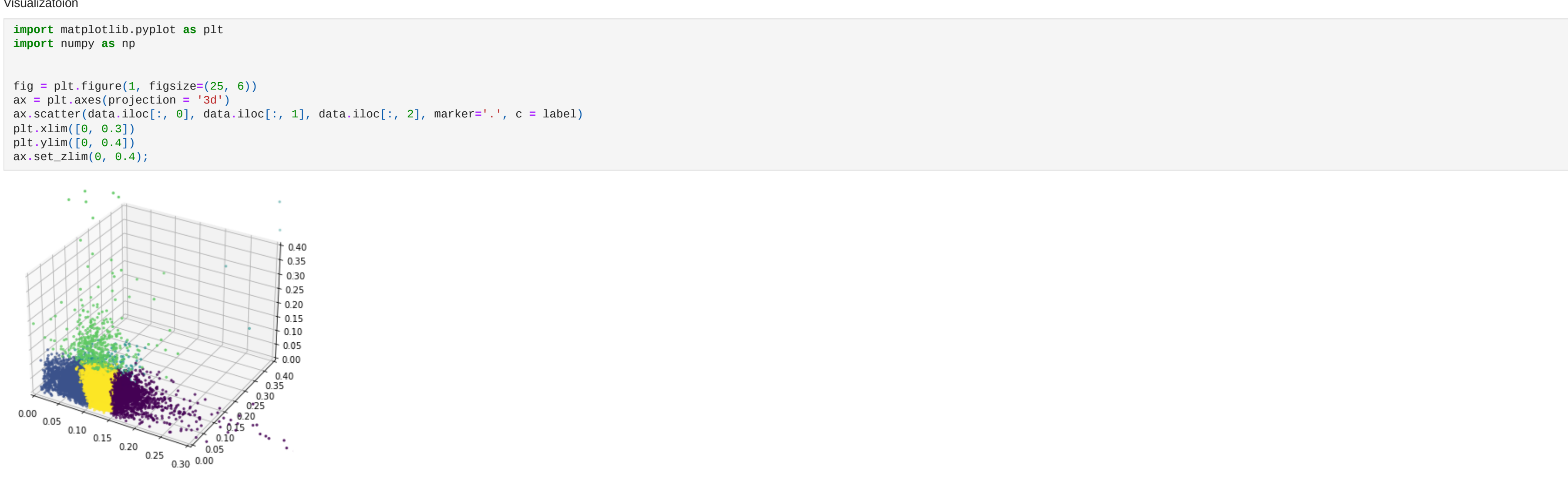


Final Clustering K= 5

```
In [15]: kmeans = KMeans(n_clusters=5, random_state=0)

#predict the labels of clusters.
label = kmeans.fit_predict(train[['comp_title', 'comp_body', 'comp_abstract']])

Visualizaition
```



Looks like a good separation

Analysis

In order to see if the clusters actually presenting similar papers we will take the example from part A and check if they are in the same cluster.

```
In [17]: example_paper = '88b80f02e54d2ed8c9f05f422b45f4818a6a405b'
example_paper_raw_data = train[train['paper_id'] == example_paper]
predicted_cluster_example_paper = kmeans.predict(example_paper_raw_data[['comp_title', 'comp_body', 'comp_abstract']])[0]
counter = 0

for similar_paper in similarity_df['paper_id'].values[1:]:
    paper_raw_data = train[train['paper_id'] == similar_paper]
    paper_raw_data = paper_raw_data[['comp_title', 'comp_body', 'comp_abstract']]
    paper_prediction = kmeans.predict(paper_raw_data)
    if paper_prediction == predicted_cluster_example_paper:
        counter += 1

counter
```

Out[17]: 7

Wow! Seven out of ten were in the same cluster. In addition, we saw in Part A an example paper in the field of machine learning. therefore, This cluster seems to represent a good collection of papers in the field of machine learning.

Further Work
Based only on the compression method and using only text data in this work, we got pretty good results. For future work, we would like to add some more methods, such as word embedding (DOC2VEC), and create a model that incorporates both methods and more metadata features. Using this approach, we believe that we will achieve better performance and similarity.