



دانشگاه علم و صنعت ایران
دانشکده مهندسی کامپیوتر

عنوان: پروژه سوم درس داده کاوی
خوشه بندی مقالات علمی با استفاده از مجموعه داده covid-19

نام و نام خانوادگی: آیلین نائب زاده

شماره دانشجویی: ۹۹۵۲۲۱۸۵

نیم سال تحصیلی: پائیز ۱۴۰۲

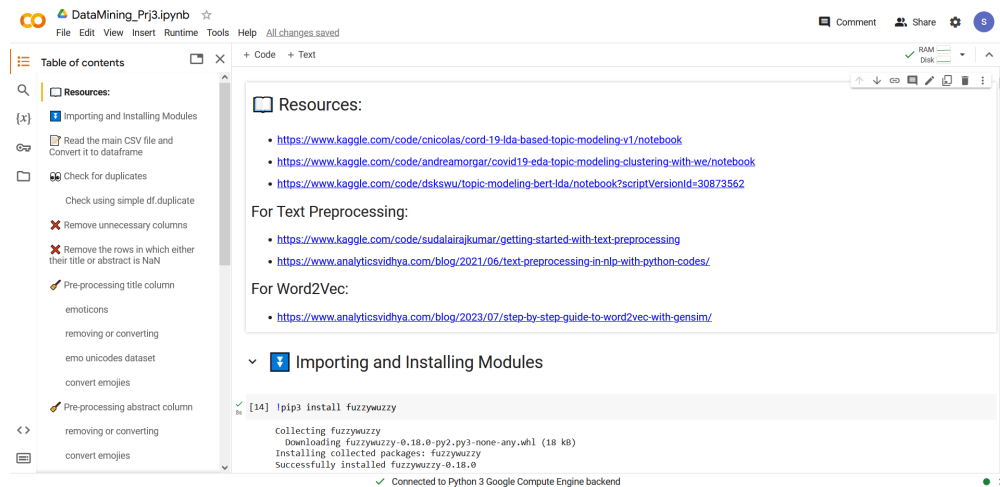
مدرس: دکتر حسین رحمانی

فهرست مطالب

۲	۱ گام اول
۴	۲ گام دوم
۵	۳ گام سوم
۶	۴ گام چهارم
۷	۵ گام پنجم
۹	۶ گام ششم
۱۱	۷ نتیجه گیری
۱۱	۱.۷ کدام یک از روش استخراج ویژگی و مدل خوشه بندی عملکرد بهتری دارند؟
۱۱	۲.۷ دلیل خود برای این انتخاب را ذکر کنید.
۱۲	۳.۷ با بررسی خوشه ها تحلیل کنید که هر خوشه نماینده ی کدام نوع از موضوعات می باشند.
۱۳	۴.۷ نمایش داده ها

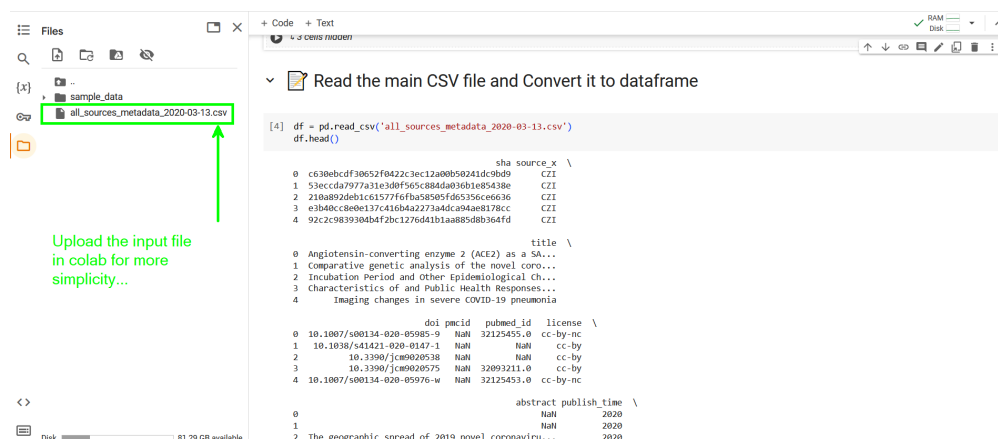
۱ گام اول

در طی این پروژه، هدف تحلیل مقالات مرتبط با بیماری covid-19 می باشد. فایل ورودی شامل ۱۴ ستون می باشد. در اولین قدم با مراجعه به سایت Google Colab و ساخت یک پروژه جدید و فراخوانی بعضی از کتابخانه های معروف زبان Python شروع به مشاهده داده های موجود در هر فایل و انجام تحلیل های ابتدایی می کنیم. محیط کلی پروژه را در تصویر زیر می توانید مشاهده کنید.



شکل ۱: نمای کلی از پروژه notebook در فضای Google Colab

در مرحله اول با استفاده از کتابخانه pandas و تابع read_csv، فایل ورودی را به آبجکتی از نوع dataframe تبدیل می کنیم تا در قدم های بعدی راحت تر بتوانیم تحلیل های مورد نیاز را انجام بدهیم. همچنین با استفاده از توابع head()، info() و describe() اطلاعات بیشتری نسبت به داده های موجود کسب می کنیم.



شکل ۲: خواندن فایل ورودی

همچنین در بخش زیر می‌توانید خروجی مربوط به تابع `info()` را مشاهده کنید. همانطور که مشاهده می‌کنید، در هر ستون تعدادی مقادیر تهی یا به اصطلاح NaN موجود است، که در گام بعدی مورد پردازش قرار می‌گیرند.

```

1 <class 'pandas.core.frame.DataFrame'>
2 RangeIndex: 29500 entries, 0 to 29499
3 Data columns (total 14 columns):
4  #   Column                                Non-Null Count  Dtype
5  ---  ---
6  0    sha                                17420 non-null   object
7  1    source_x                          29500 non-null   object
8  2    title                             29491 non-null   object
9  3    doi                               26357 non-null   object
10  4    pmcid                             27337 non-null   object
11  5    pubmed_id                         16730 non-null   float64
12  6    license                           17692 non-null   object
13  7    abstract                          26909 non-null   object
14  8    publish_time                     18604 non-null   object
15  9    authors                           28903 non-null   object
16  10   journal                           17791 non-null   object
17  11   Microsoft Academic Paper ID      1134 non-null    float64
18  12   WHO #Covidence                    1236 non-null    object
19  13   has_full_text                     17420 non-null    object
20 dtypes: float64(2), object(12)
21 memory usage: 3.2+ MB
22 None

```

۲ گام دوم

در این مرحله که مربوط به پیش‌پردازش داده‌ها می‌باشد، باتوجه به این موضوع که داده‌هایی که در اختیار داریم از نوع متنی می‌باشند، نیاز است بعضی از پیش‌پردازش‌های معمول را اجرا کنیم. بطور خلاصه مراحل زیر بر روی داده‌های انجام شده‌اند:

- حذف ردیف‌های تکراری: براساس جفت عنوان و خلاصه مقاله، داده‌های تکراری شناخته و پیدا شده‌اند و تنها آخرین نمونه تکراری از هر گروه برای ادامه مراحل باقی می‌ماند.

- حذف تمامی ستون‌ها به جز ستون‌های title و abstract

- حذف تمامی ردیف‌هایی که حداقل یکی از مقادیر title یا abstract در آن‌ها تهی می‌باشد.

- تبدیل محتوای ستون title و abstract به string

- تبدیل تمامی حروف موجود در ستون‌های title و abstract به حروف کوچک

- حذف تمامی علائم نگارشی از ستون‌های title و abstract

- حذف برخی از حروف پرتکرار مانند for، must و ... از ستون‌های title و abstract

- برگرداندن کلمات موجود در ستون‌های title و abstract به ریشه اصلی آن‌ها با استفاده از عملیات stemming و lemmatization

- حذف emoticons از محتوای ستون‌های title و abstract

- حذف آدرس‌ها از محتوای ستون‌های title و abstract

- حذف تگ‌های html از محتوای ستون‌های title و abstract

- حذف برخی از رشته‌های خلاصه‌شده خاص از محتوای ستون‌های title و abstract

* خروجی تمامی مراحل بالا در دو ستون جدا به نام‌های cleaned_title و cleaned_abstract و در آبجکت df_dropped ذخیره می‌شوند.

۳ گام سوم

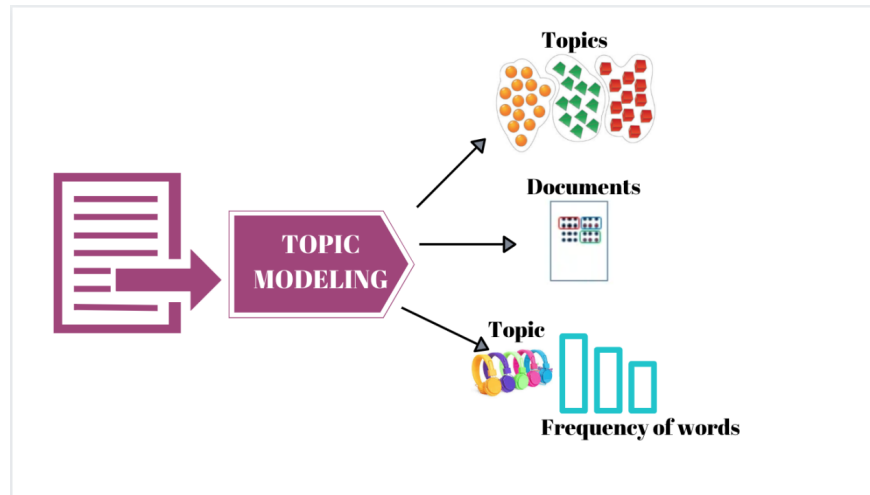
این مرحله مربوط به استخراج ویژگی از داده‌های موجود می‌باشد. در این مرحله باتوجه به داده‌های متنی که در اختیار داریم، نیاز است تا برای پردازش راحت‌تر، کلمات را به مقادیر عددی تبدیل کنیم که برای انجام این کار از دو روش TF-IDF و Bag of Words استفاده شده‌است. همچنین باتوجه به این موضوع که از ما خواسته شده است که تحلیل‌ها را بر روی دو ستون title و abstract انجام دهیم، درواقع مقادیر تجمیع شده این دو ستون را به‌عنوان ورودی به دو تابع TfidfVectorizer و CountVectorizer به‌عنوان ورودی می‌دهیم. این دو تابع در کتابخانه sklearn.feature_extraction.text موجود می‌باشند. و نحوه استفاده از این دو تابع را در کدهای زیر می‌توانید مشاهده کنید:

```

1 tfidf = TfidfVectorizer()
2
3 tf_idf_matrix = tfidf.fit_transform(df_dropped['cleaned_title'] + ' ' +
4     df_dropped['cleaned_abstract'])
5
6 """
7 tf_idf_matrix
8
9 <22819x99910 sparse matrix of type '<class 'numpy.float64'>'
10     with 2082915 stored elements in Compressed Sparse Row format>
11 """
12 bow_vectorizer = CountVectorizer(max_df = 0.90, min_df = 2, max_features =
13     1000, stop_words='english')
14 bow_matrix = bow_vectorizer.fit_transform(df_dropped['cleaned_title'] + ' ' +
15     df_dropped['cleaned_abstract'])
16
17 """
18 bow_matrix
19
20 <22819x1000 sparse matrix of type '<class 'numpy.int64'>'
21     with 1286156 stored elements in Compressed Sparse Row format>
22 """
    
```

۴ گام چهارم

در این مرحله به منظور انجام Topic Modeling از الگوریتم Latent Dirichlet Allocation استفاده می‌کنیم. این الگوریتم را به طور مجزا یک بار با استفاده از ماتریس خروجی از الگوریتم TF-IDF انجام می‌دهیم و بار دیگر نیز الگوریتم را بر روی خروجی Bag of Words اجرا می‌کنیم.



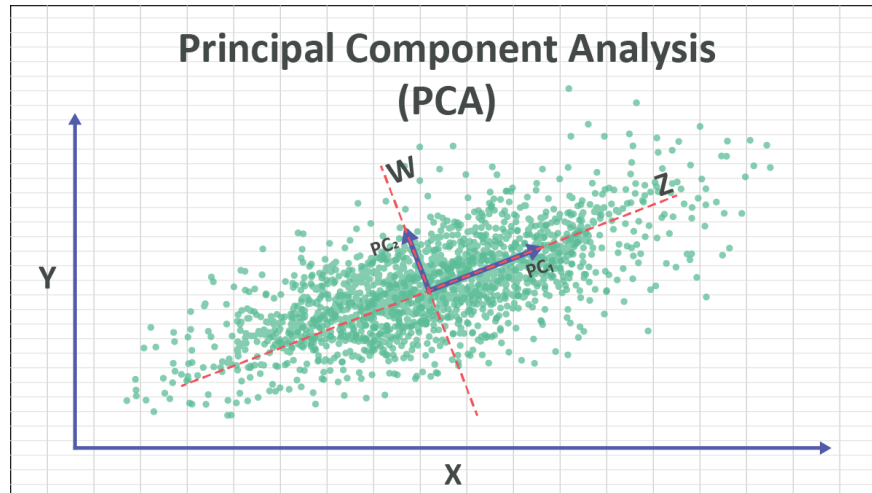
شکل ۳: الگوریتم LDA

کدهای مربوط به این بخش را در بخش زیر می‌توانید مشاهده کنید:

```
1 # Create an LDA object with 20 topics
2 lda_tf_idf = LatentDirichletAllocation(n_components=20)
3 lda_bow = LatentDirichletAllocation(n_components=20)
4
5 tf_idf_topics = lda_tf_idf.fit_transform(tf_idf_matrix)
6 print(tf_idf_topics.shape) # Output will be (22819, 20)
7
8 bow_topics = lda_bow.fit_transform(bow_matrix)
9 print(bow_topics.shape) # Output will be (22819, 20)
```

۵ گام پنجم

در این مرحله از ما خواسته شده است تا ابعاد ماتریس‌هایی که در مرحله قبلی محاسبه شده بودند را کاهش دهیم. به منظور انجام این کار از الگوریتم Principle Component Analysis استفاده می‌کنیم. در هنگام استفاده از این الگوریتم همانگونه که در تصویر زیر مشاهده می‌کنیم، تعداد بخش‌های اصلی را با ۲ مقداردهی می‌کنیم. همانند مرحله پیشین نیاز است یک بار الگوریتم را با استفاده از ماتریس خروجی TF-IDF و بار دیگر با استفاده از ماتریس خروجی Bag of Words اجرا نمائیم.

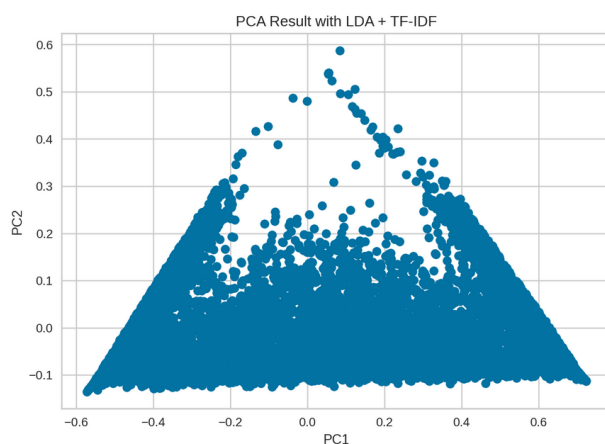


شکل ۴: الگوریتم PCA

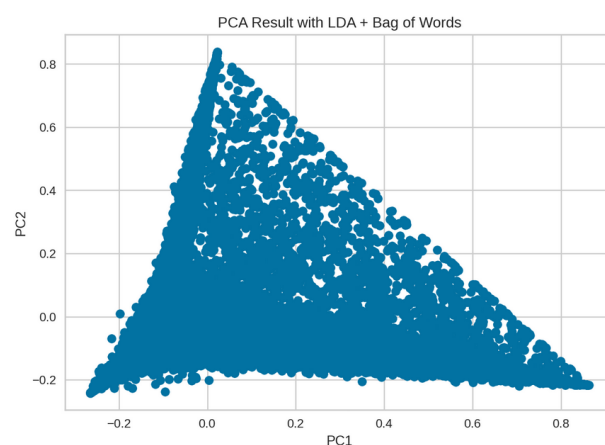
کدهای مربوط به این بخش را در بخش زیر می‌توانید مشاهده کنید:

```
1 # Create a PCA object with 2 components
2 pca_tf_idf = PCA(n_components=2)
3 pca_bow = PCA(n_components=2)
4
5 tf_idf_reduced = pca_tf_idf.fit_transform(tf_idf_topics)
6 print(tf_idf_reduced.shape) # Output will be (22819, 2)
7
8 bow_reduced = pca_bow.fit_transform(bow_topics)
9 print(bow_reduced.shape) # Output will be (22819, 2)
```


همچنین تصاویر مربوط به خروجی این الگوریتم را بر روی دو ماتریس TF-IDF و Bag of Words می‌توانید مشاهده کنید:



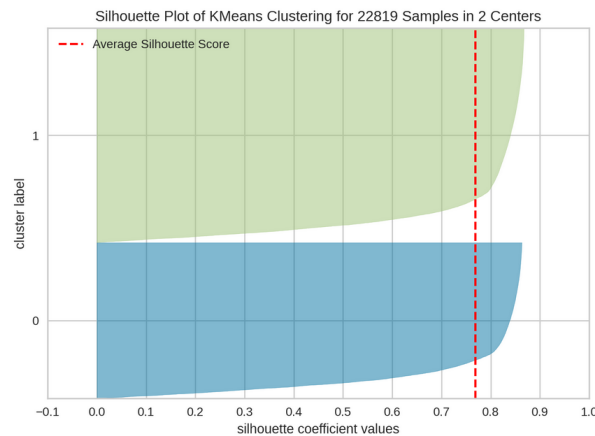
شکل ۵: خروجی الگوریتم PCA برای ماتریس ورودی TF-IDF



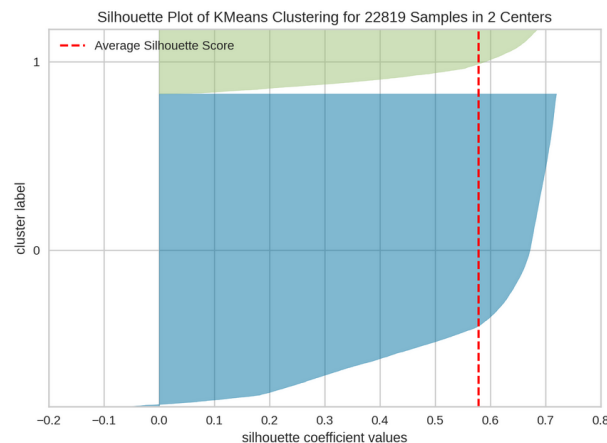
شکل ۶: خروجی الگوریتم PCA برای ماتریس ورودی BoW

۶ گام ششم

در این مرحله که مربوط به استفاده از الگوریتم‌های خوشه‌بندی می‌باشد، همانطور که در تصویر زیر مشاهده می‌کنید از دو الگوریتم KMeans و DBSCAN استفاده می‌کنیم. البته هر یک از این دو الگوریتم را بصورت جداگانه با استفاده از دو ماتریس TF-IDF و Bag of Words مورد آموزش قرار می‌دهیم. همچنین هنگام استفاده از الگوریتم KMeans با استفاده از الگوریتم Elbow و Silhouette سعی می‌کنیم تا بهترین مقدار K یا به عبارتی بهینه‌ترین تعداد خوشه‌ها را انتخاب کنیم. مقدار بازه تعداد خوشه‌ها را از ۲ تا ۲۰ تعریف می‌کنیم و همانطور که در خروجی مشخص است، در هر دو حالت تعداد خوشه دو بهترین پاسخ را باتوجه به معیار Silhouette می‌دهد. چراکه همانطور که می‌دانیم هر چه مقدار این معیار به ۱ نزدیک‌تر باشد، یعنی خوشه‌بندی بهتر انجام شده‌است.

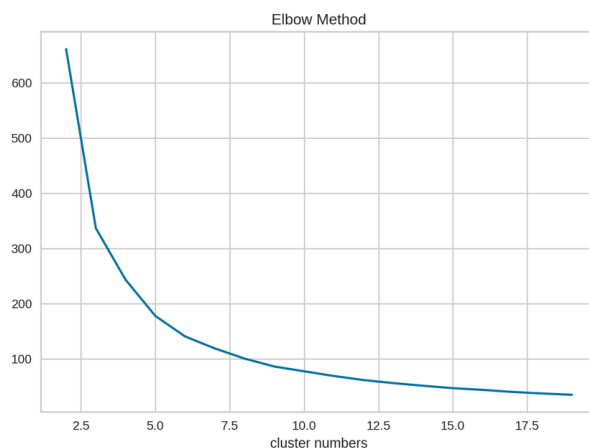


شکل ۷: خروجی تابع Silhouette برای الگوریتم KMeans در حالت استفاده از الگوریتم TF-IDF

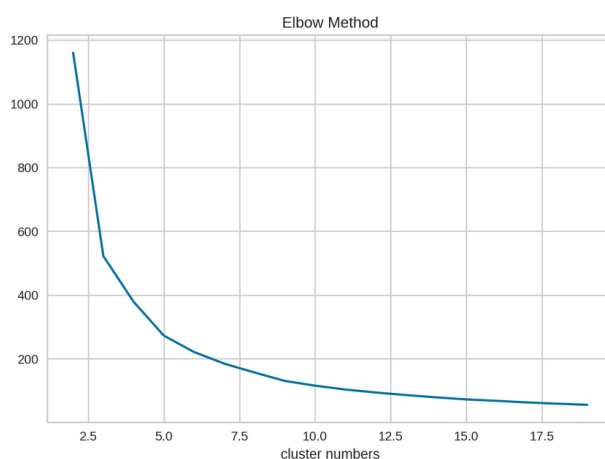


شکل ۸: خروجی تابع Silhouette برای الگوریتم KMeans در حالت استفاده از الگوریتم Bag of Words

نمودارهای Elbow مربوط به الگوریتم‌های KMeans را نیز می‌توانید در تصاویر زیر مشاهده کنید:



شکل ۹: Elbow for KMeans TF-IDF



شکل ۱۰: Elbow for KMeans BoW

*تمامی خروجی‌های مربوط به این بخش در خروجی کدها موجود است.

۷ نتیجه گیری

۱.۷ کدام یک از روش استخراج ویژگی و مدل خوشه بندی عملکرد بهتری دارند؟
در حالت استفاده از الگوریتم KMeans و ماتریس نتایج TF-IDF به همراه $K=2$ به بهترین عملکرد می‌رسیم.

۲.۷ دلیل خود برای این انتخاب را ذکر کنید.
سادگی پیاده سازی - قابلیت فهم بیشتر الگوریتم و نحوه عملکرد - کمبود زمان:

۳.۷ با بررسی خوشه‌ها تحلیل کنید که هر خوشه نماینده‌ی کدام نوع از موضوعات می‌باشند.

همانطور که در بخش زیر می‌توانید مشاهده کنید، ده عنوان پرتکرار در حالتی که تنها دو خوشه داشته باشیم، به‌صورت زیر می‌باشند:

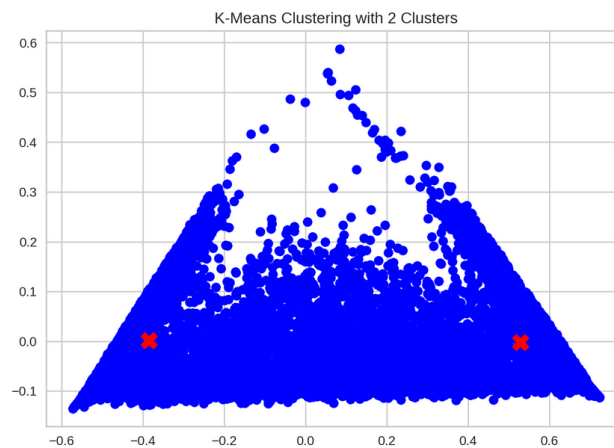
```

1
2 Cluster 0:
3 Title + Abstract:
4 [('vte', 1.678757792866809), ('cistran', 1.6842943889715567), ('rp',
   1.7129443327427645), ('hbz', 1.9769338882057215), ('de',
   1.9923199119796313), ('ppias', 2.2181387904324343), ('nw',
   2.2671036443696733), ('stat3', 4.269071400018623), ('cypa',
   5.930938308329624), ('cyclophilin', 7.042773950588994)]
5 Cluster 1:
6 Title + Abstract:
7 [('csa', 3.3479223037548627), ('nendou', 3.463791236165143), ('omtas',
   3.5563027857379423), ('la', 3.58139521730169), ('rc', 3.7748316966902937),
   ('nsp14', 4.299709957546419), ('nsp16', 4.828983563894078), ('nsp11',
   5.2362523138395884), ('de', 5.585708033407686), ('cchfv',
   8.213499671505836)]
8 .....
9
10 Cluster 0:
11 Title + Abstract:
12 [('recent', 1470.153988758314), ('research', 1594.534216970629), ('global',
   1638.6642640546884), ('new', 1987.7613982795474), ('drug',
   2023.8265750613432), ('emerg', 2292.504527338654), ('review',
   2448.7748359447), ('health', 3010.6509597508707), ('develop',
   3332.1425793052204), ('diseas', 4115.477387357077)]
13 Cluster 1:
14 Title + Abstract:
15 [('anim', 1172.0006459933547), ('syndrom', 1454.5165365477444), ('virus',
   1591.727393376288), ('infect', 1689.7183916321692), ('respiratori',
   1699.9138172302232), ('east', 1796.7737427966777), ('middl',
   1877.3738146026203), ('bat', 3053.279585291326), ('human',
   3479.644308713773), ('merscov', 3559.049999997403)]

```

۴.۷ نمایش داده‌ها

تصویر مربوط به خروجی الگوریتم KMeans را با دو خوشه می‌توانید در تصویر زیر مشاهده کنید. با توجه به اینکه پیش از استفاده از خوشه‌بندی، از الگوریتم کاهش ابعاد استفاده کردیم، دقت مدل متاسفانه کاهش یافته است و به این دلیل که داده‌ها را به دو بعد کاهش داده‌ایم، امکان دارد که بعد دوم در زیر این بعد قرار داشته باشد. * سایر نمایش‌های مربوط به خوشه‌ها در خروجی برنامه قرار دارند.



شکل ۱۱: Clustering Visualize