

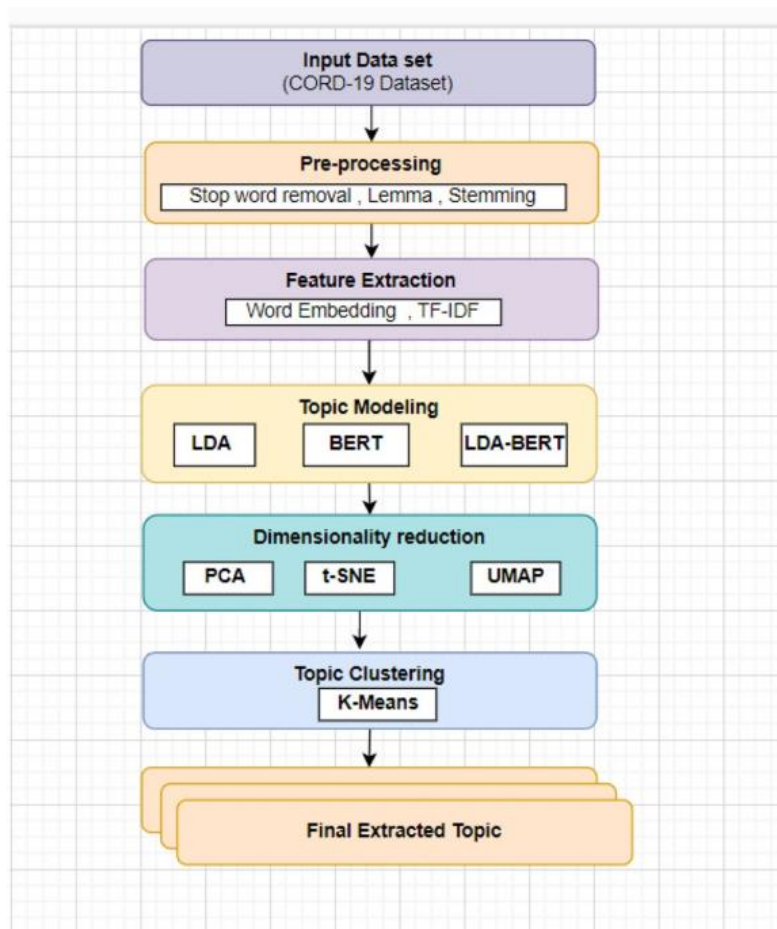
به نام خدا

پروژه ۳ داده کاوی

عاطفه نادری

ابتدا با توجه به مقاله <https://link.springer.com/article/10.1007/s41870-023-01268-w>

مراحل را انجام دادیم:



Pre process:

برای حذف Stop words ابتدا از spacy bio parser رفتیم. این کتابخانه لغات مخصوصی که مرتبط با مقالات علمی حوزه بایو است را جمع آوری کرده است و وابسته به آن stopword هایی دارد. اما در تست ها متوجه شدم خوب نیست و کار درستی را انجام نمیدهد. و به صورت عادی حذف این واژه ها را انجام دادم.

مراحل کار به صورت کلی به شکل زیر است:

- Lower casing
- Removal of punctuations
- Removal of numbers
- Tokenization
- Lemmatization
- stemming
- Removal of stopwords, including a custom list of stopwords

feature Extraction

برای این قسمت، از مدل bow و tf-idf استفاده میکنیم.

topic modeling

برای اهداف topic modeling، از LDA استفاده کردیم. چیزی به عنوان بردار lda وجود ندارد، فقط به عنوان احتمال تعلق به یک موضوع را نشان می دهد.

- اسناد با موضوعات مشابه از گروه های مشابهی از کلمات استفاده می کنند

- موضوعات اسناد، که موضوعات پنهان نامیده می شوند را می توان با جستجوی گروه هایی از کلمات که اغلب با هم در اسناد در سراسر مجموعه وجود دارند، پیدا کرد.

BERT

هم چنین از مدل BERT نیز برای topic modeling استفاده کردم که آن را به دلیل ساعت زیادی که برای fit شدن می گرفت فقط توانستم روی k-means اجرا کنم و نتیجه آن را بر روی هر سه روش کاهش ابعاد تست کردم و در یک جدول قرار دادم.

dimension reduction

برای dimension reduction از سه روش PCA و UMAP و TSNE نیز استفاده نمودم.

PCA : Linear dimension reduction

T-SNE : Non-linear dimensional reduction, preserves local structure in the data.

Umap : Non-linear dimensional reduction, preserves both local and most of the global structure

نتایجی که از umap بدست آمد در برخی از آزمایش ها بهتر از سایرین در ترکیب با tf-idf بود.

ماتریس خروجی tf-idf و bow ماتریس اسپارسی است و ما روی آن کاهش ابعاد میزنیم.

Topic clustering

برای این قسمت از ۴ الگوریتم k-means, mean-shift, DBSCAN, OPTICS استفاده کردیم.

معیارهای ارزیابی

silhouette_score

که این معیار نرخ چسبندگی اعضا به کلاستر خود و جدایی از کلاسترهای دیگر را نشان میدهد. مقدار های نزدیک به یک بهتر است. مقادیر منفی نشان دهنده کلاسترینگ بد است.

برای بدست آوردن k از دو روش elbow و silhouette استفاده شده است. که در silhouette مقداری که بیشترین را داشته است به عنوان k استفاده شده است. چون نتایج elbow قابل اعتماد نیست.

davies_bouldin_score

این شاخص توسط دیویس و بولدین دو دانشمند در رشته برق در سال 1979 معرفی شد و وابسته به تعداد خوشه‌ها و یا الگوریتم خوشه‌بندی نیست. برای محاسبه این شاخص از دو معیار اندازه پراکندگی (Dispersion measure) و عدم شباهت بین خوشه‌ها (Cluster dissimilarity) استفاده می‌کند.

مقدار آن هر چه کمتر باشد بهتر است.

تحلیل:

برای تحلیل با توجه به notebook هایی که قرار داده شده است که همه چیز را تقریباً همان جا توضیح داده ام. ابتدا ما در lda یک عددی را برای استخراج topic ها قرار می‌دهیم. مثلاً در آزمایش آخر عدد ۴۰ را هم برای bow و هم برای tf-idf قرار دادیم. و با استفاده از ویژگی features که روش vectorization برای bow و tf-idf وجود دارد. لغات مرتبط با هر topic را پیدا نمودیم و طبق آن ۴۰ تاپیک یک عنوانی را به هر کدام از آن‌ها اختصاص دادیم. در مراحل بعدی و بعد از کلاسترینگ نیز، لغات پر تکرار هر کلاستر را استخراج می‌کنیم و similarity آن را با هر یک از ۴۰ topic که داشتیم مسنئج و بیشترین شباهت را گزارش می‌کنیم که هر کلاستر به کدام topic مرتبط است. و آن را همانجا چاپ کردیم.

هم چنین با استفاده از بیشترین کلماتی که در هر کلاستر تکرار شده است نیز آن را به صورت ابر کلمات نشان دادیم.

جدول زیر برخی از اجراهایی است که در طول پروژه انجام شده است. و همانطور که قابل مشاهده است هر چقدر در tf-idf مقدار topic ها در lda بزرگتر بود، نتیجه بهتر میشد. شاید به دلیل آن است که ابتدا تمایز خوبی از هر عنوان در هر سند ایجاد میشود و سپس با pca ابعاد به صورت معناداری کاهش می یابد و در نهایت clustering مطلوبی را ایجاد میکند.

در کل طی آزمایشات مختلف و متعدد نتایج حاصل از tf-idf بهتر از bow بوده است. چون که tf-idf اهمیت لغات را بهتر از bow نشان میدهد. بنابراین نتایج بهتری را نیز گرفتیم.

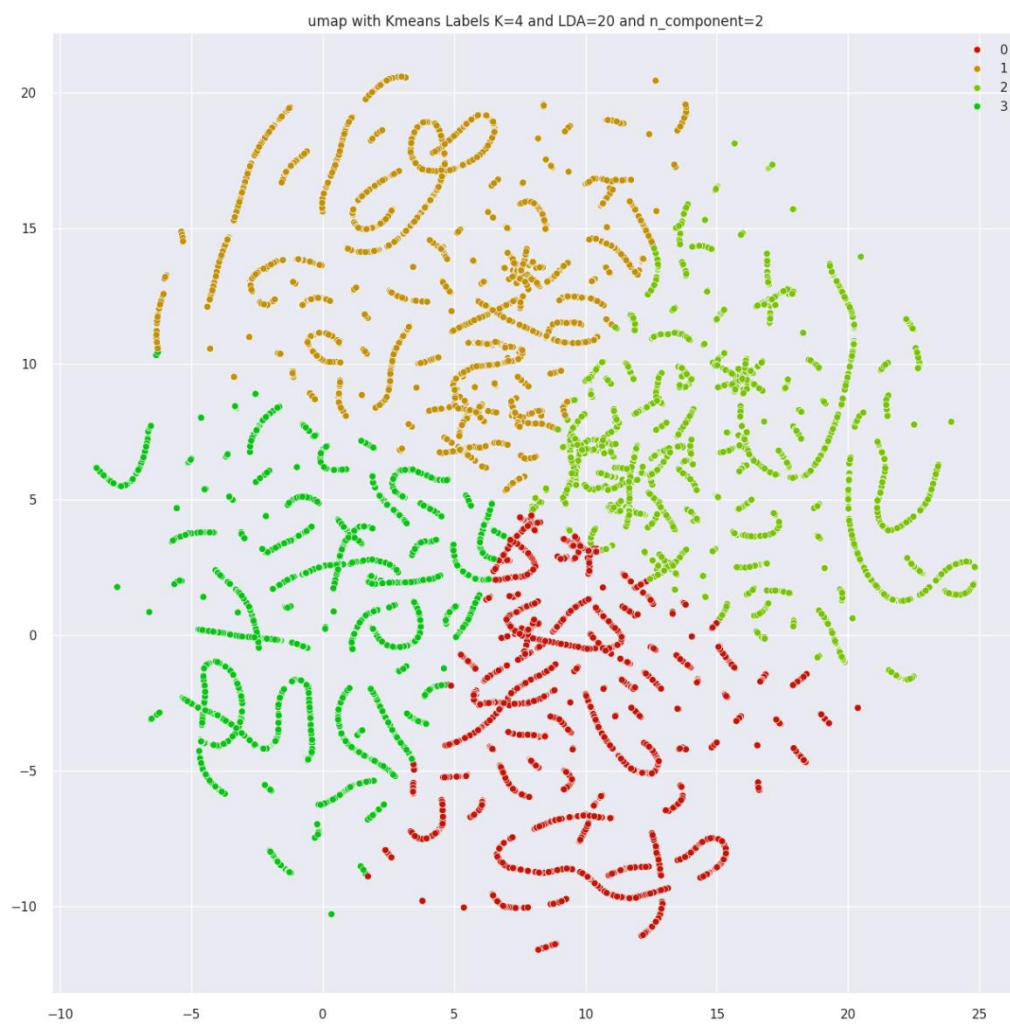
Tf-idf And k-means	K=3, pca=2, lda=20	K=4, pca=2, lda=25	K=3, pac=2, lda=40	K=3, umap, lda =40	K=4 Umpa Lda=2 0	K=3 Umpa Lda=20	
silhouette	0.87705791 2562349	0.874929237 1478099	0.92704557323 77547 0.92932179425 29081 0.90636755661 75831	0.3773 0643	0.352 7427	0.3723 5686	
Davies -			0.199223736 74285676				

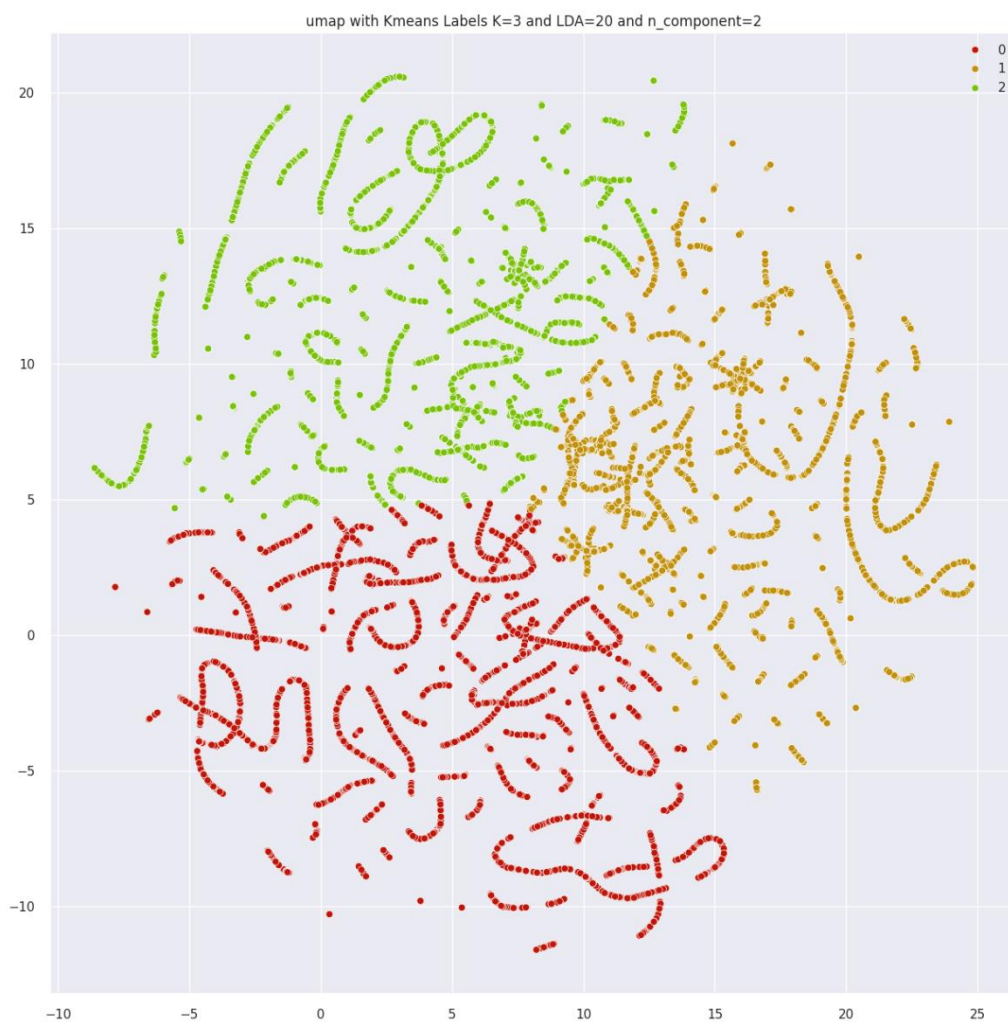
Bouldi n index:							
			شکل ۳				

K=4

Umpa

Lda=20

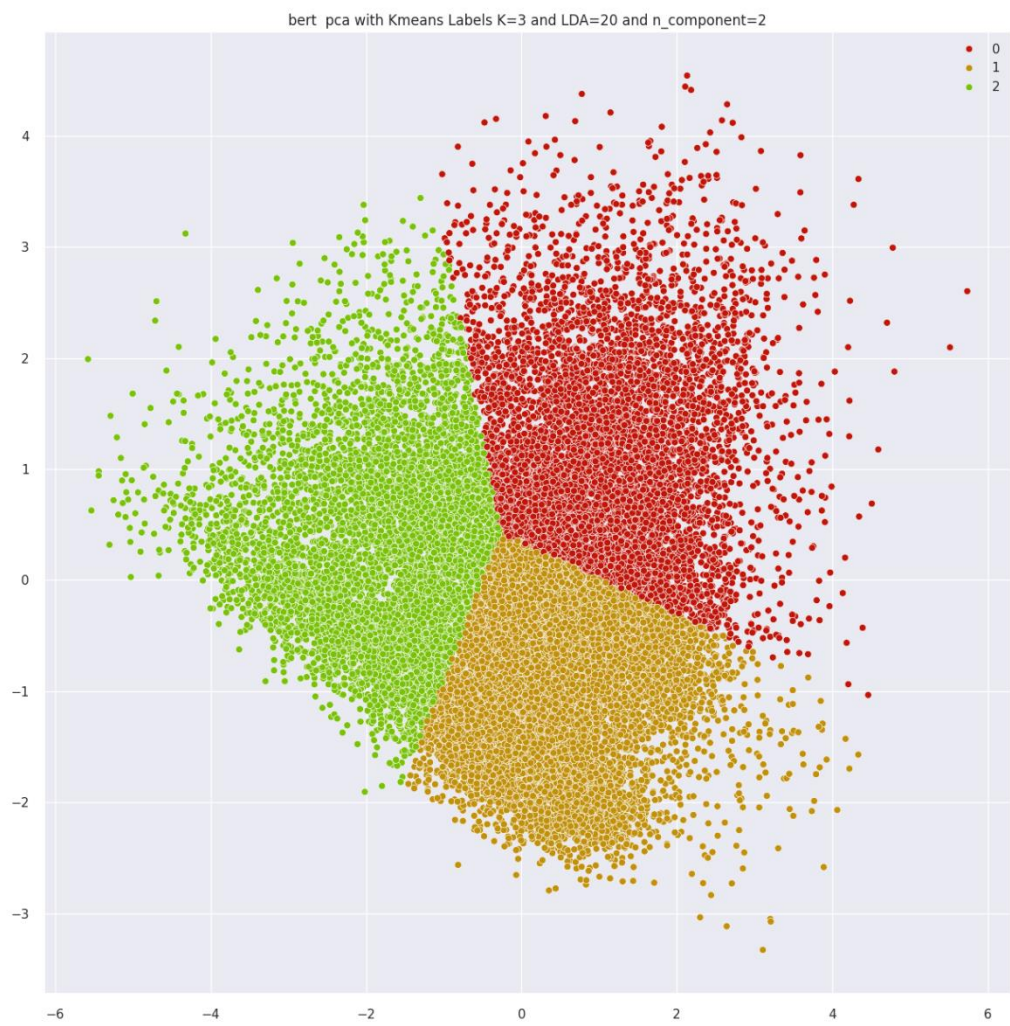




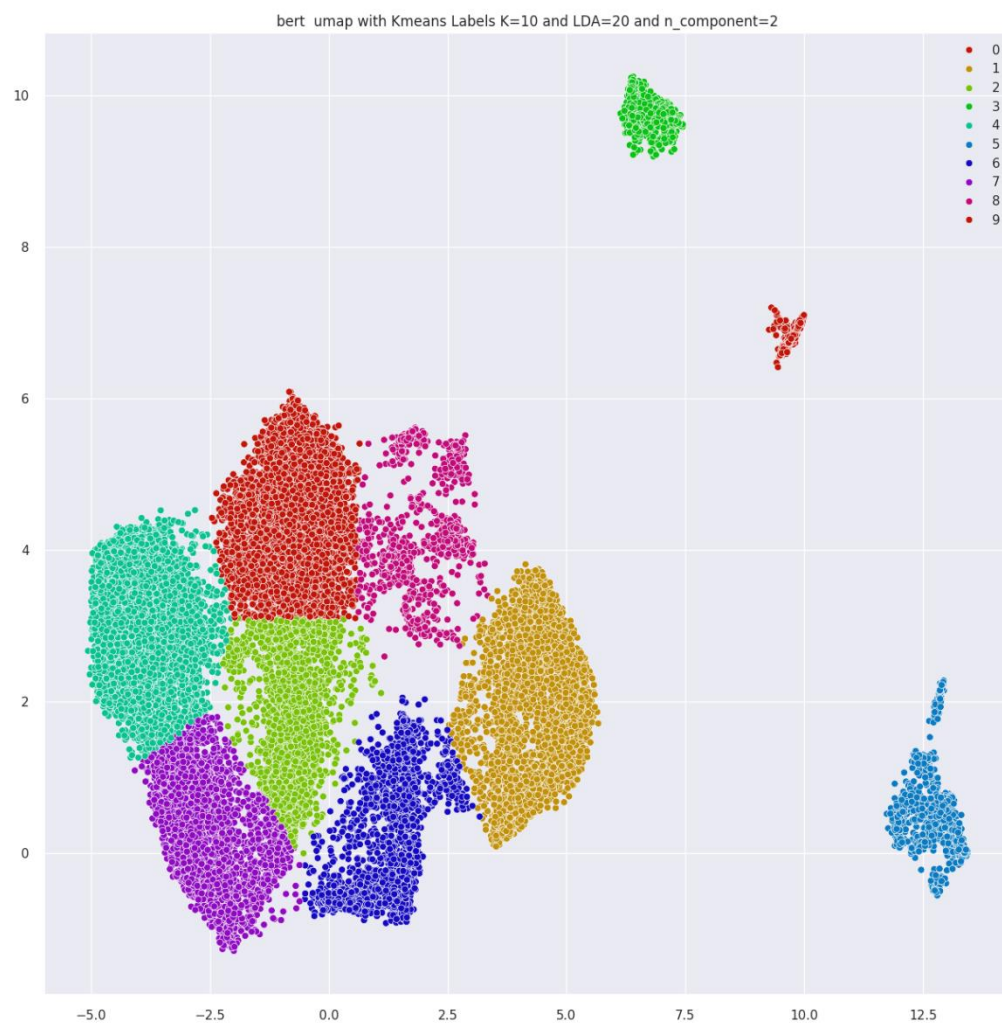
Bert

نمایی از نتایجی که با مدل bert ایجاد شده است با استفاده از هر سه روش کاهش ابعاد تست شده است و مانند نتیجه مقاله با نتیجه بهتری را طبق معیار سیلوئت داشتیم.

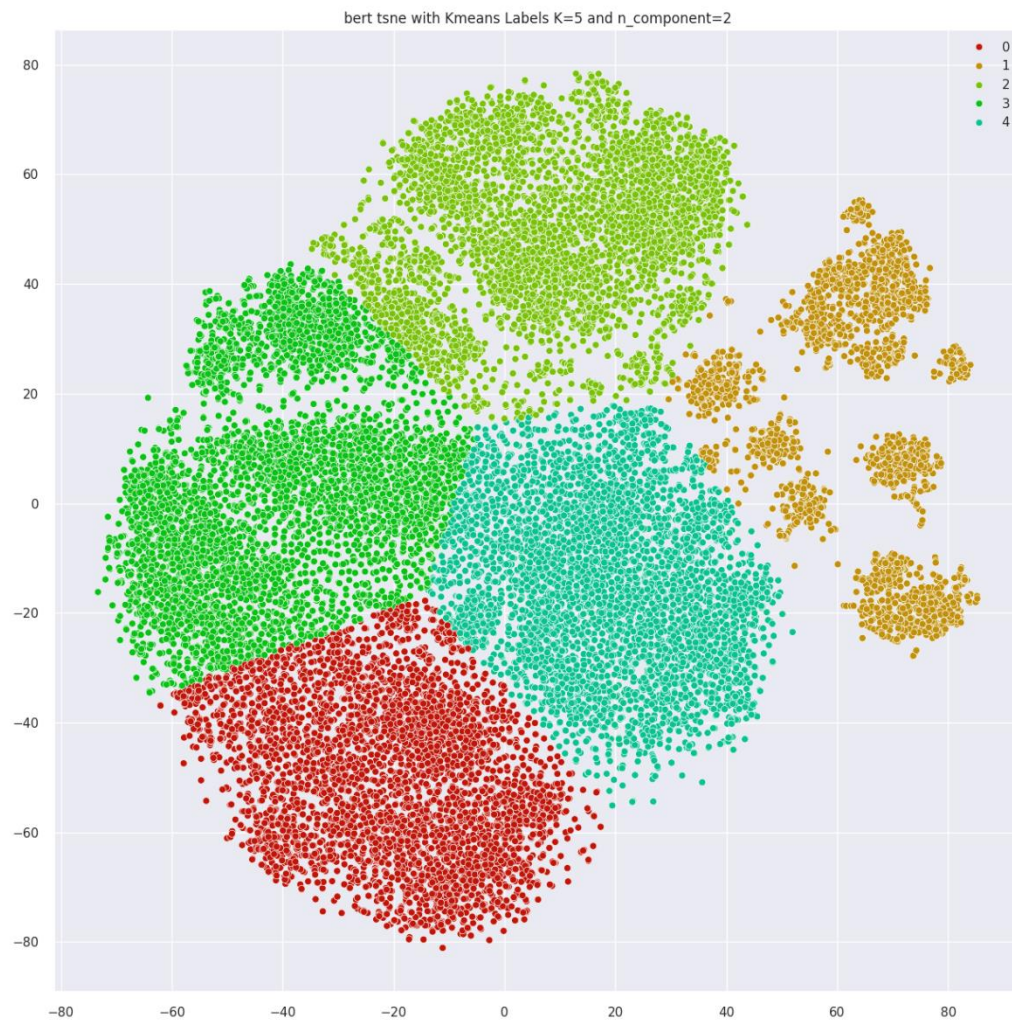
Tf-idf and k-means	Pca, k=3	Umap, k=10	Tsne, k=5	
bert	0.37791735	0.5401364	0.42491844	



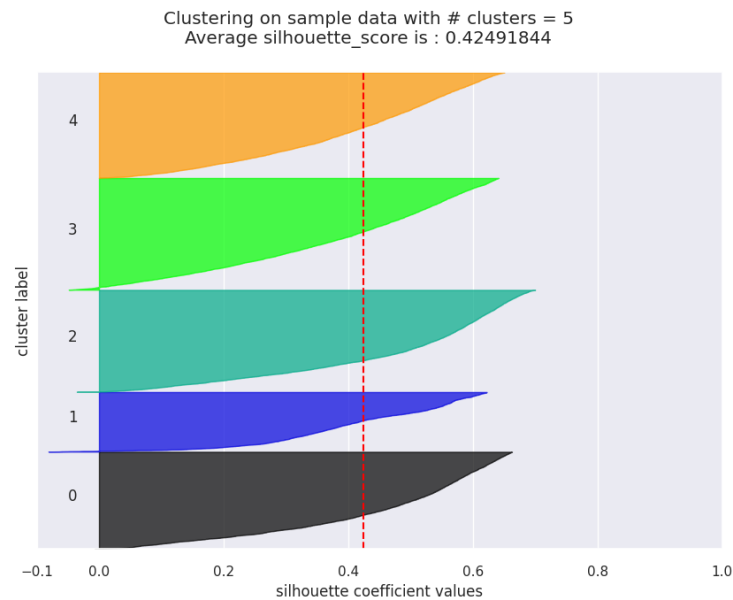
Bert and umap:



Bert & TSNE:



میانگین امتیاز سیلوئتی که برای k بدست آمد:

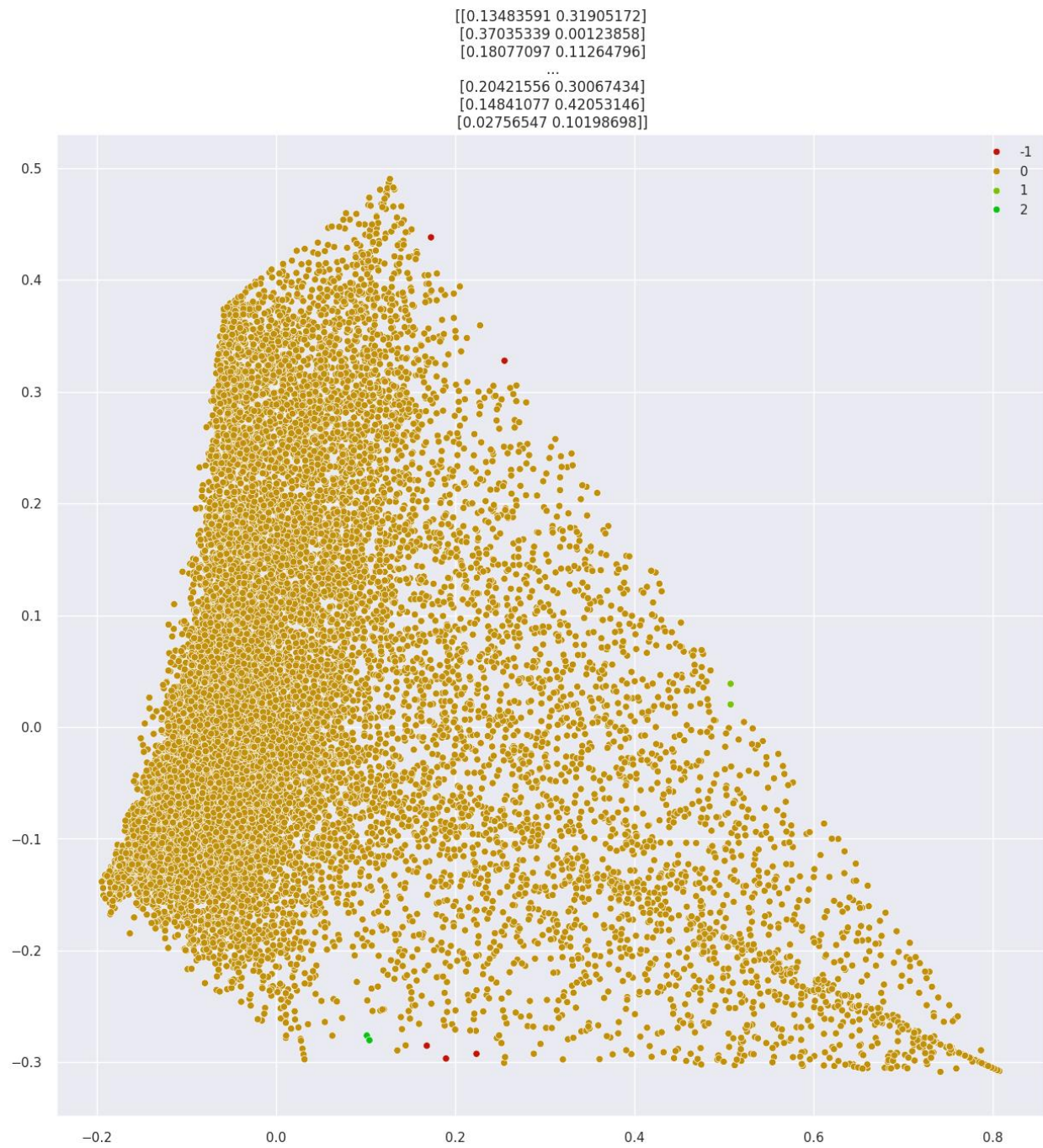


الگوریتم TF-IDF و DBSCAN

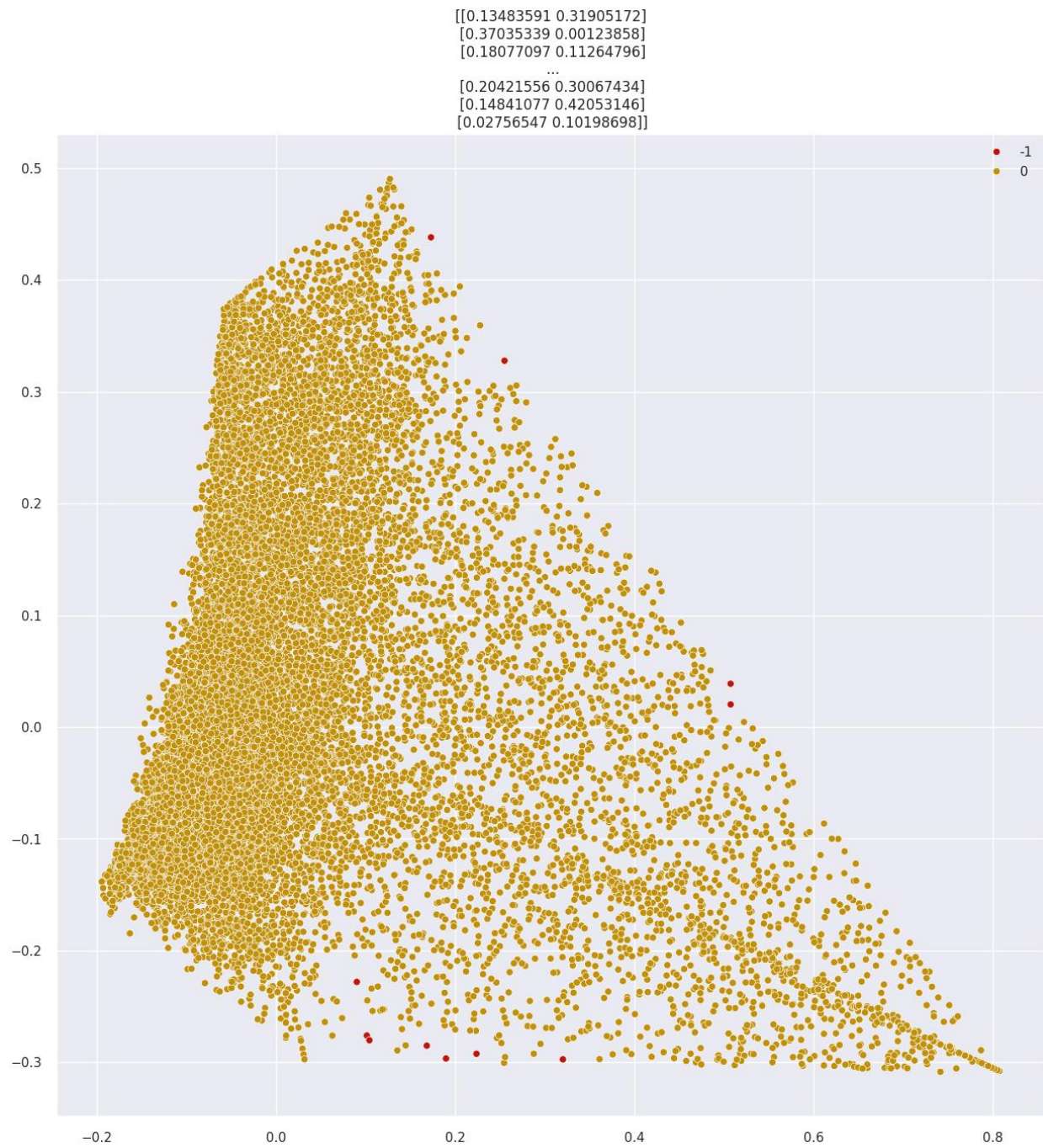
DBSCAN	Pca =2 Lda=30 Cluster=2 Tf-idf eps=0.1, min_samples=2	Pca =2 Lda=30 Cluster=3 Tf-idf eps=0.1, min_samples=5	Pca =2 Lda=30 Cluster=3 Tf-idf eps=0.1, min_samples=5	Pca =2 Lda=40 Cluster=1 5 Tf-idf eps=0.02, min_samples=2	Pca =2 Lda=40 Cluster=6 Tf-idf eps=0.02, min_samples=5	Pca =2 Cluster =4 BoW eps=0.02, min_samples=2	Pca =2 Cluster =2 BoW eps=0.02, min_samples=5
silhouette	0.6788970988952632	0.6577608538557453	0.6761046648666912	0.5009179046544635	0.579432777918427	0.21700325588816796	0.4124630679549918
Davies-Bouldin			0.5476111653961602	2.4550927730647163	4.568920956025978	1.8129676528646153	1.7387124426791212

شماره شکل	1	2	15	14		3	4
--------------	---	---	----	----	--	---	---

تصاویری از اجراهای الگوریتم DBSCAN

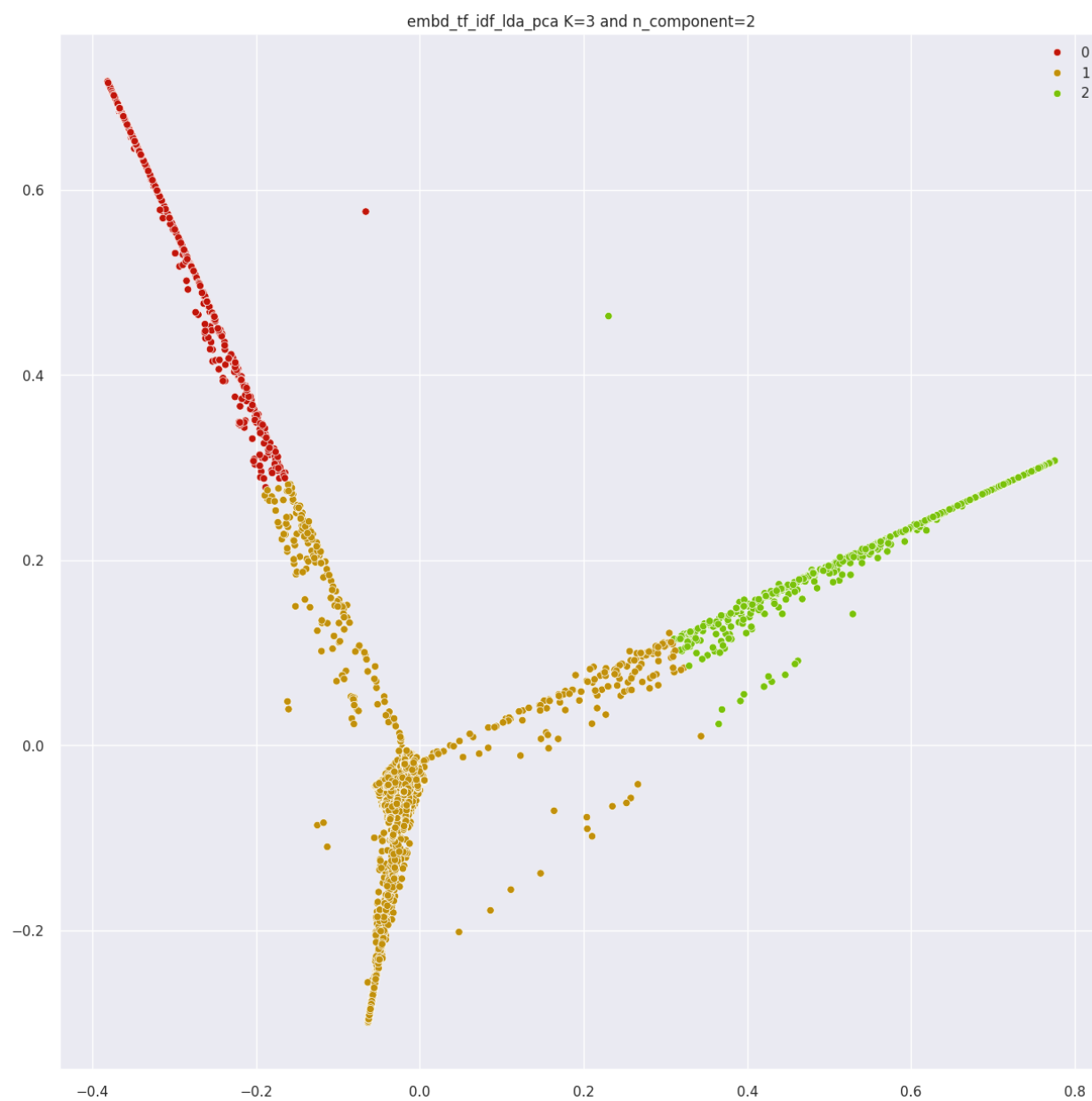


شکل ۳

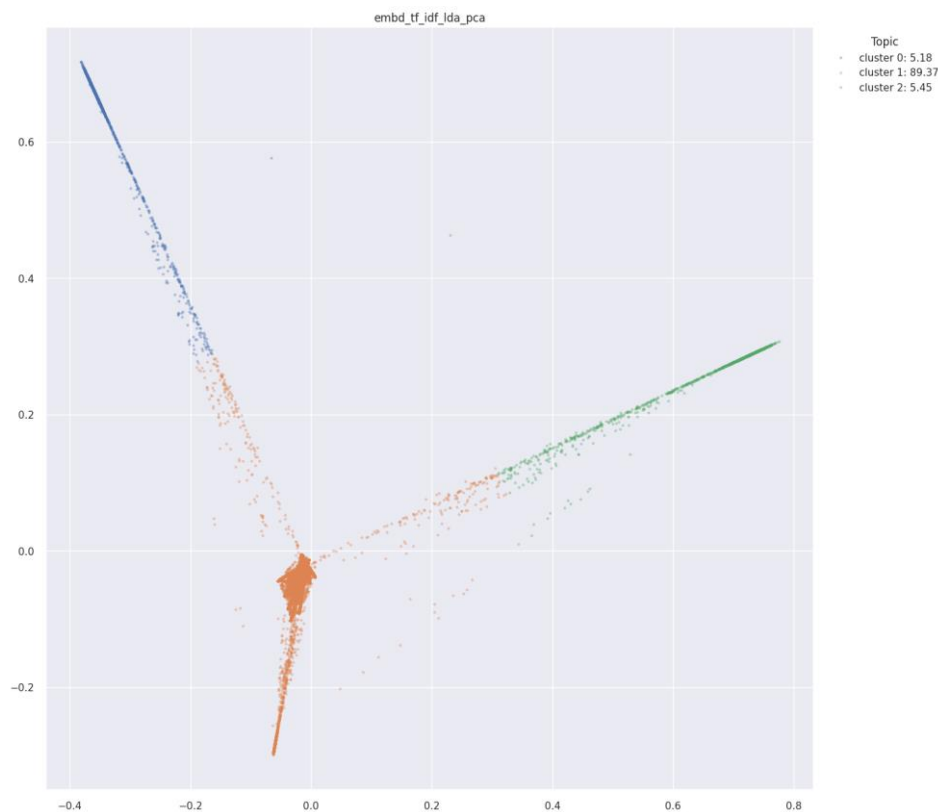


شکل ۴

حالا به توضیح بهترین حالت در **tfidf** میپردازیم. این حالت بهترین امتیاز سیلوئت را داشته است. به تحلیل آن نیز میپردازیم.



شکل ۳ از جدول tf-idf-pca



نمایی دیگر از شکل ۳

با توجه به رخداد کلمات در هر کلاستر میتوان گفت برای مثال کلاستر اول مرتبط با بحث تجزیه و تحلیل بیان ژن و پاسخ های ایمنی است. کلاستر دوم مرتبط با بحث پاسخ بافت تنفسی به عفونت آنفلانزا است و کلاستر سوم نیز مرتبط با بحث مدل های اپیدمیولوژیک برای گسترش و کنترل بیماری است.

```
Most Similar List of Words: ['rna', 'mrna', 'replic', 'protein', 'viru',
'genom', 'sequenc', 'transcript', 'viral', 'structur', 'gene', 'site',
'cell', 'region', 'express', 'effici', 'function', 'activ', 'mutat',
'contain']
```

```
Index of Maximum Similarity: 31
```

```
Topic 31: Gene Expression Analysis and Immune Responses
```

```
Most Similar List of Words: ['lung', 'infect', 'cell', 'respiratori',
'mous', 'acut', 'respons', 'tissu', 'diseas', 'sever', 'viru', 'immun',
'express', 'pneumonia', 'influenza', 'viral', 'model', 'activ', 'human',
'role']
```

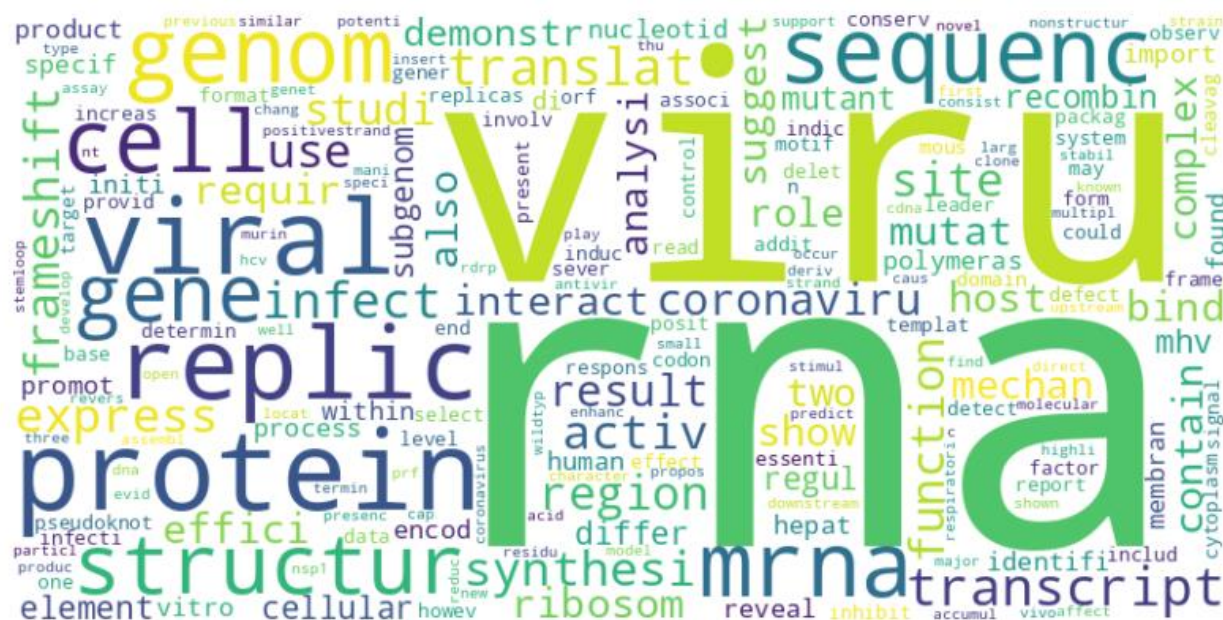
```
Index of Maximum Similarity: 11
```

```
Topic 11: Respiratory Tissue Response to Influenza Infection
```

```
Most Similar List of Words: ['epidem', 'transmiss', 'estim', 'case',
'outbreak', 'model', 'china', 'number', 'diseas', 'spread', 'data',
```

Topic 7: Epidemiological Models for Disease Spread and Control

کلاستر اول که مرتبط با مباحث ریز تر مثل سلول، ویروس و rna است.



کلاستر دوم: که میتوان در مورد آلودگی و بحث درمان و شیوع بیماری صحبت کرد.

Most Similar List of Words: ['infect', 'cell', 'viru', 'viral', 'host', 'cultur', 'human', 'replic', 'infecti', 'primari', 'system', 'use', 'studi', 'tissu', 'caus', 'establish', 'howev', 'spread', 'also', 'mechan']

Index of Maximum Similarity: 8

Topic 8: Protein Structure, Binding, and Functional Domains

Most Similar List of Words: ['health', 'public', 'surveil', 'global', 'emerg', 'system', 'countri', 'diseas', 'care', 'respons', 'inform', 'report', 'china', 'develop', 'research', 'need', 'outbreak', 'provid', 'data', 'region']

Maximum Similarity: system

Index of Maximum Similarity: 22

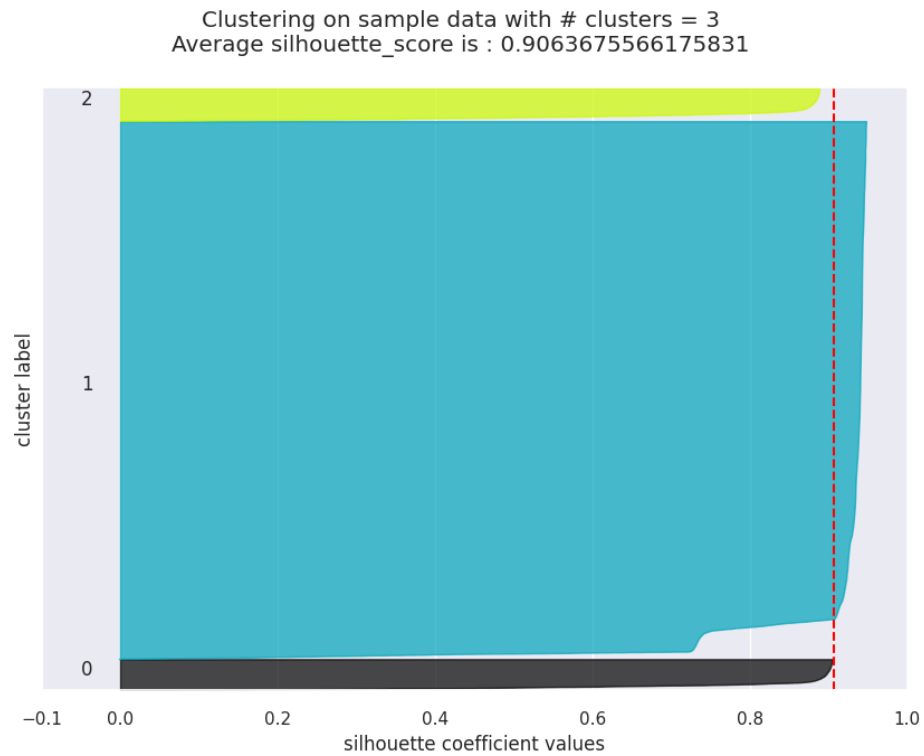
Topic 22: Viral Receptor Binding, Mutations, and Entry

برای هر یک از کلاستر ها بیشترین لغاتی که تکرار شده است نیز جمع آوری شده است. و تاپیکی که برای کلاستر اول و دوم سوم انتخاب شده است ۸ و ۲۲ و ۳۱ است. و برای مثال نتیجه کلاستر اول مرتبط با ساختار پروتئین، اتصال، و دامنه های عملکردی است.

بطور مثال این ابر کلمه ای است که برای کلاستر اول ایجاد کردیم.



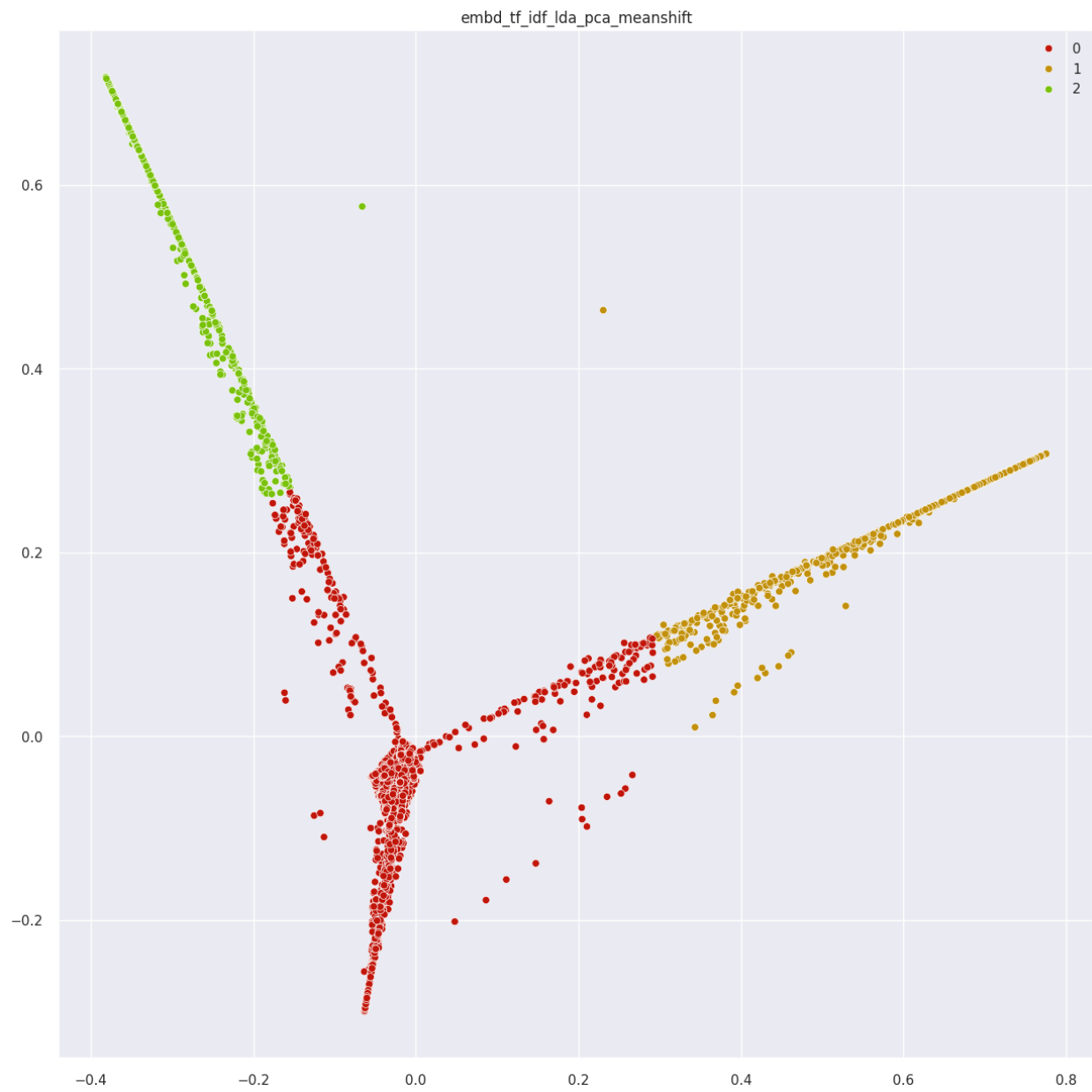
به صورت میانگین ضریب سیلویت برای k ها مختلف در بازه ۰.۸ تا ۱ قرار گرفته است.



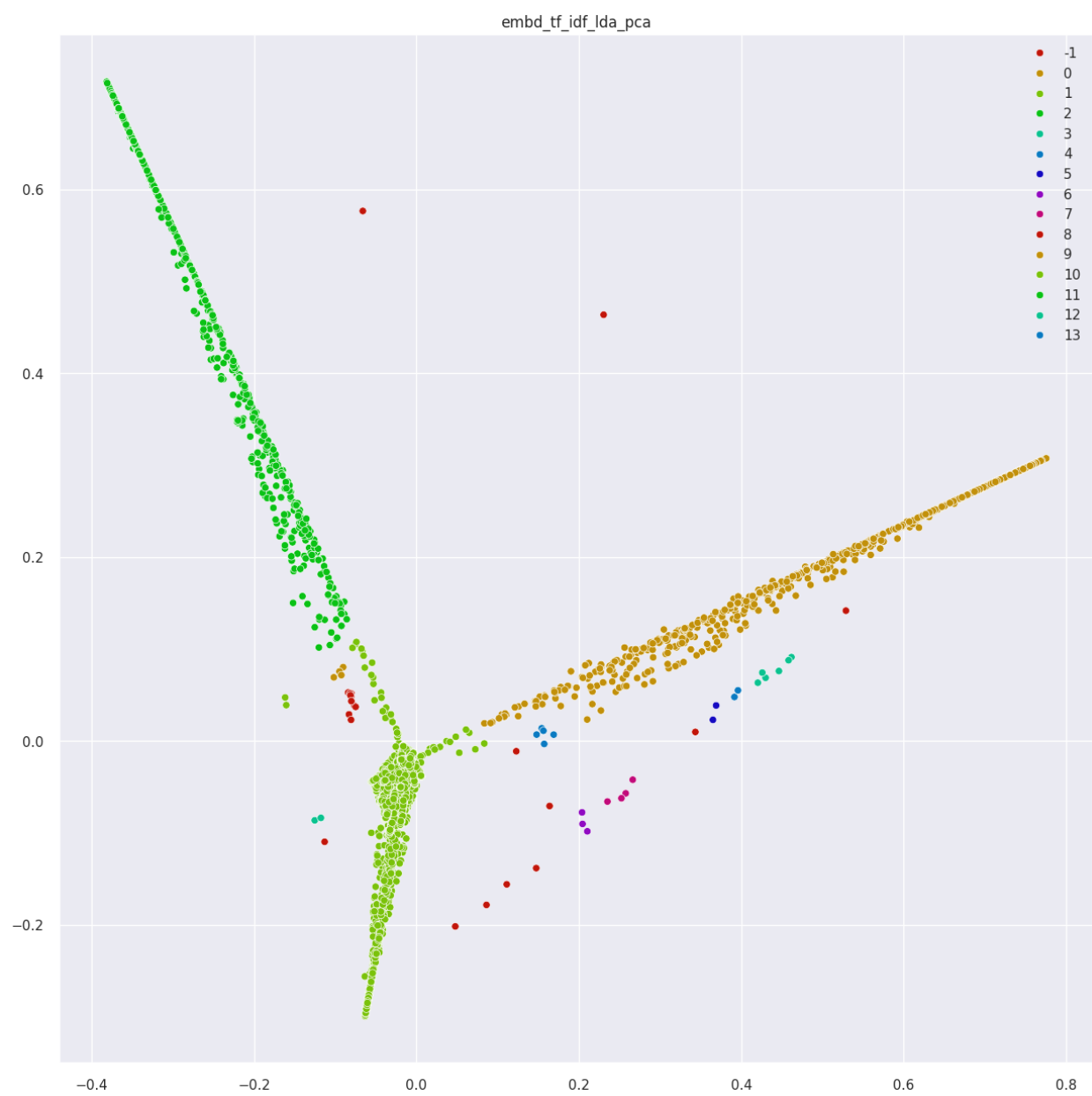
نمایی از اجراهای مختلف Mean_shift که نسبت به دیگر الگوریتم های استفاده شده هم در bow و هم در tf-idf نتایج خوبی را نشان داد.

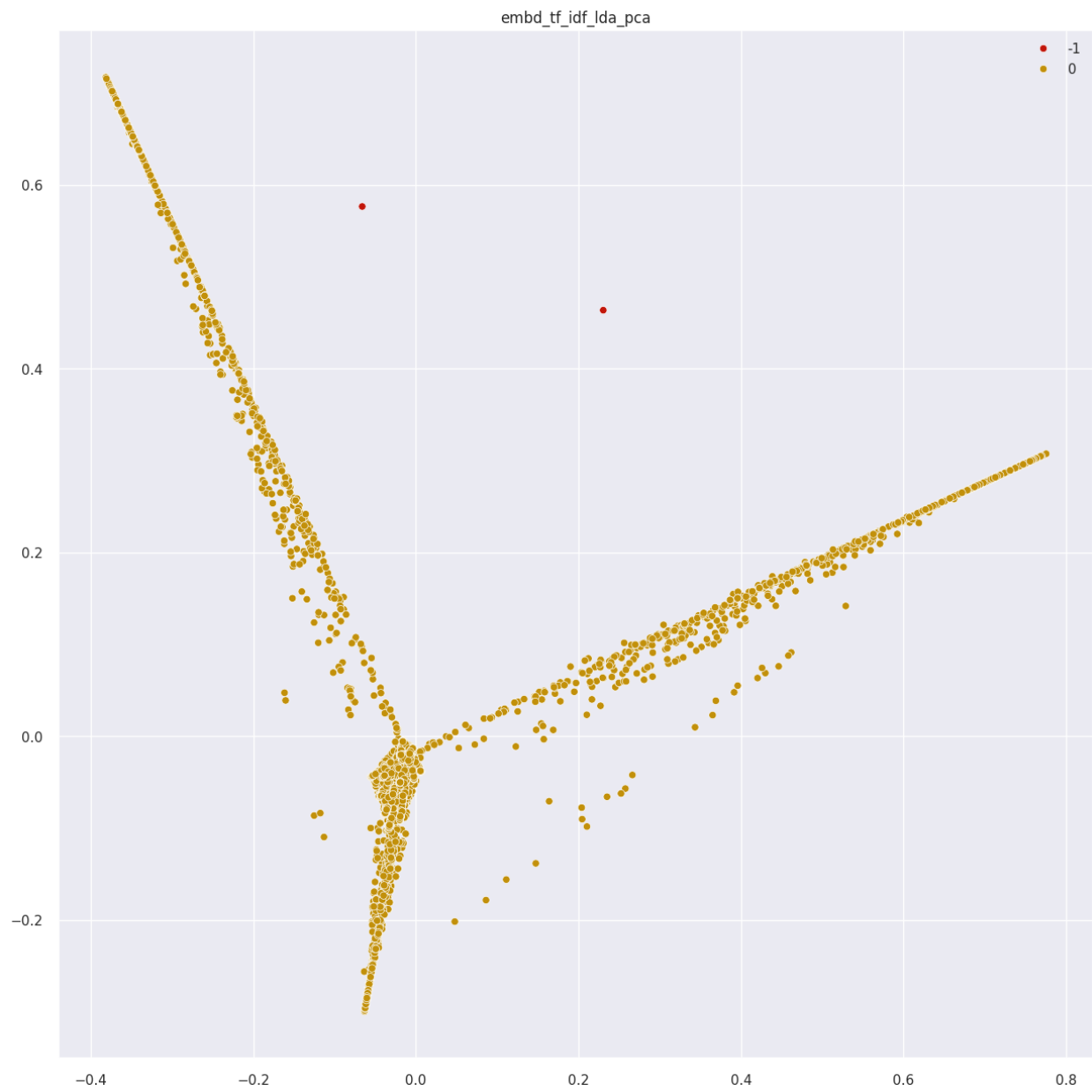
0.9058838612252527

Mean_shift	K=3, pca=2, lda=20	K=4, pca=2, lda=25	K=3, pac=2, lda=40	K=3, umap, lda =40	K=4 Umpa Lda=2 0	K=3 Umpa Lda=20	
silhouette	0.87705791 2562349	0.874929237 1478099	0.90588386122 52527	0.3773 0643	0.352 7427	0.3723 5686	
Davies - Bouldi n index:			0.211458311 75495547				
			13				



شکل ۱۳



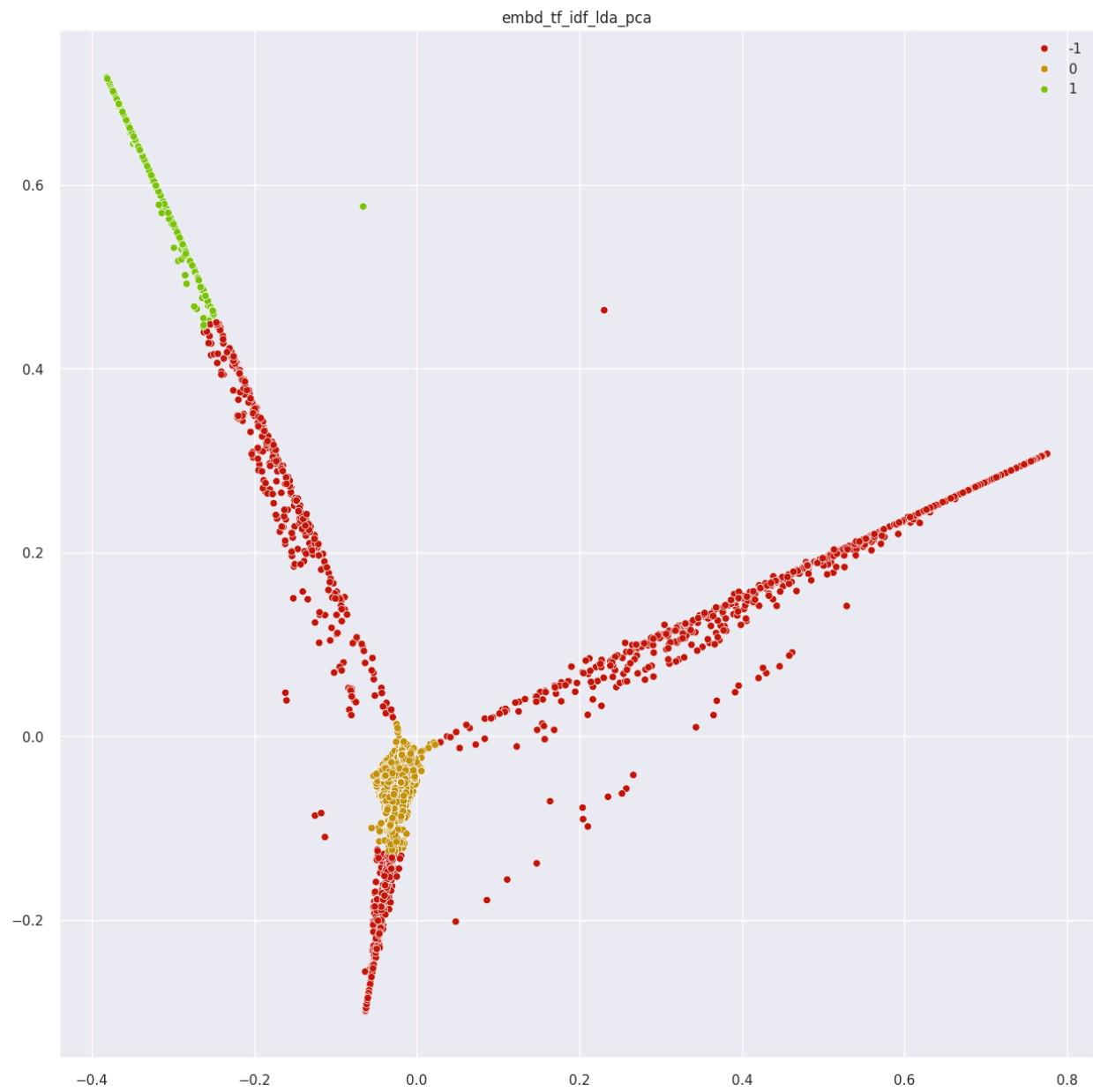


Optics and tf-idf

```
min_samples=1000  
number of cluster = 3
```

```
silhouette_score  
0.7916178385946345
```

```
Davies-Bouldin index  
1.2208399694536198
```



Bow

Lda = 40

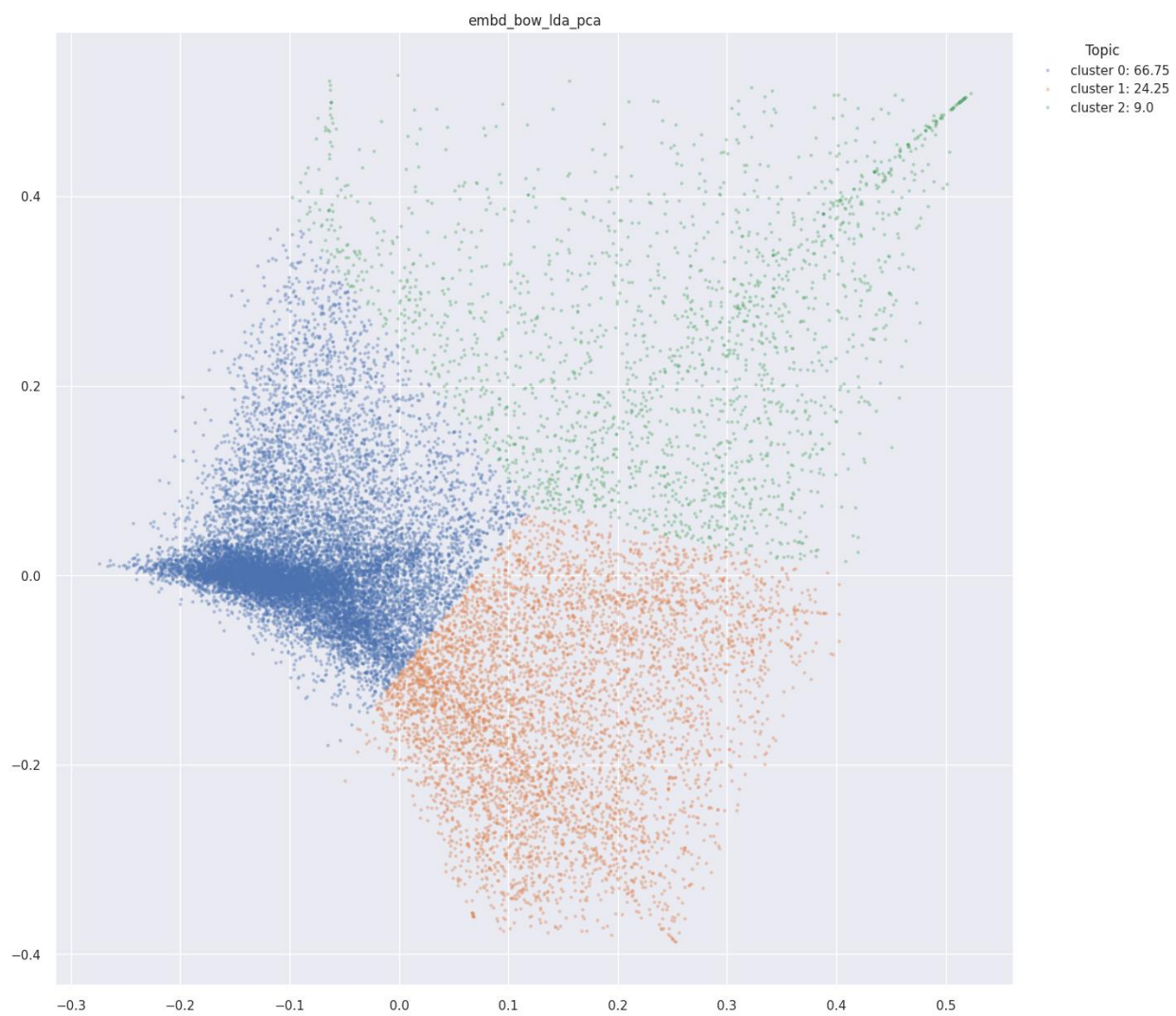
Pca =2

K=3

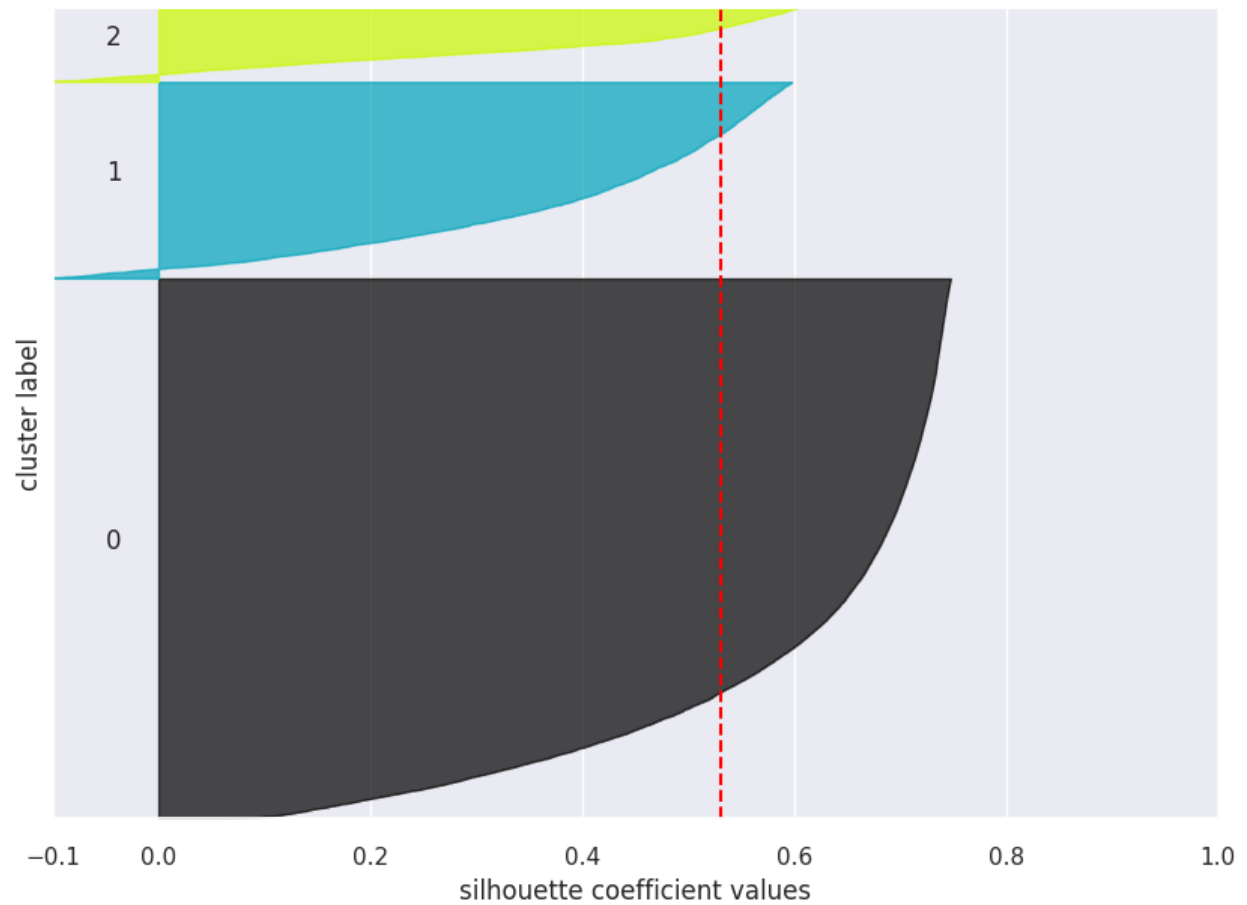
silhouette_score 0.5303344261095964

Davies-Bouldin index: 0.7447208458599722



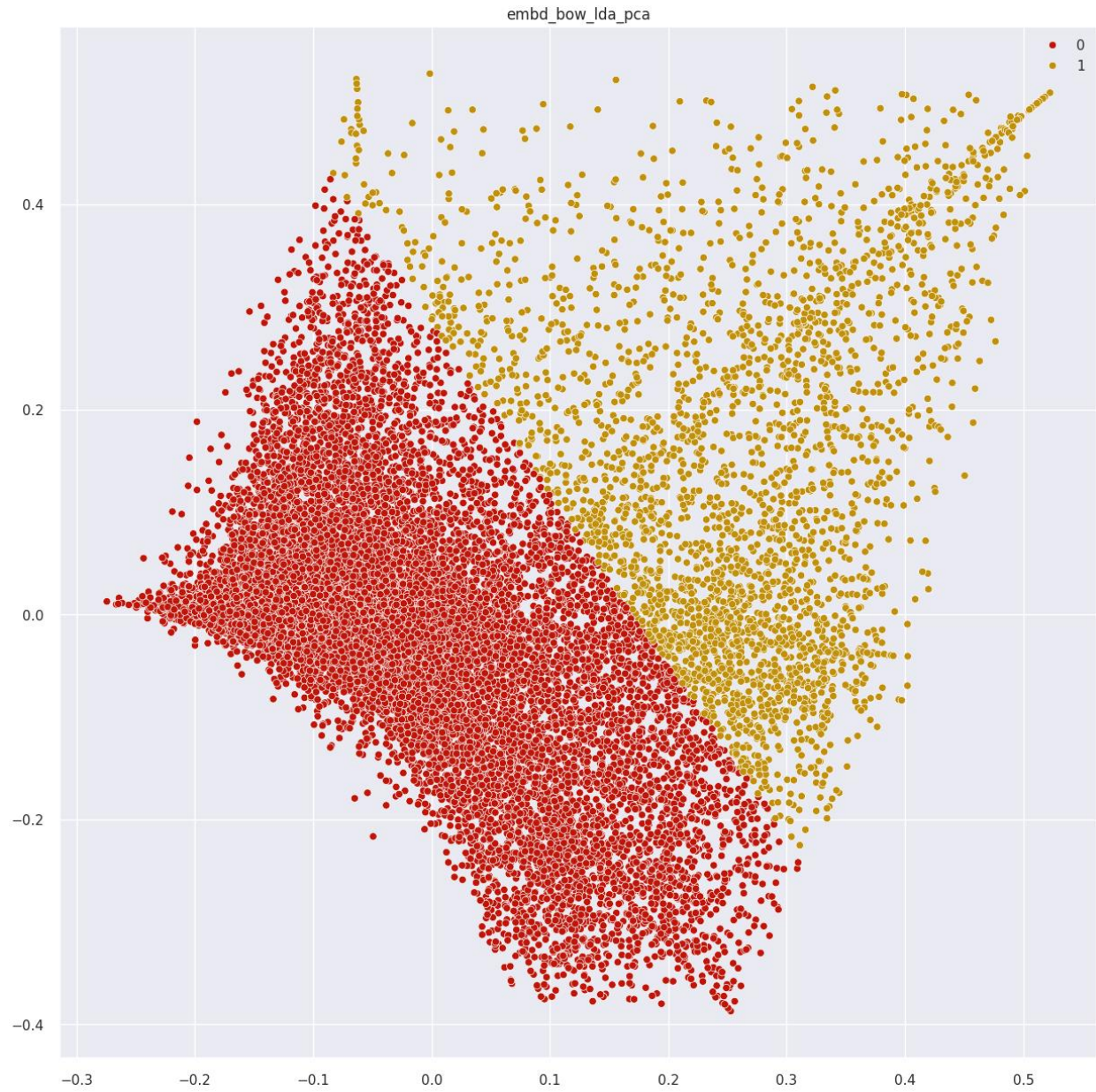


Clustering on sample data with # clusters = 3
Average silhouette_score is : 0.5303344261095964



Mean- shift

```
bandwidth=0.1  
silhouette_score  
0.5226578911593829  
Davies-Bouldin index: 0.9092823906391729
```

Dbscan

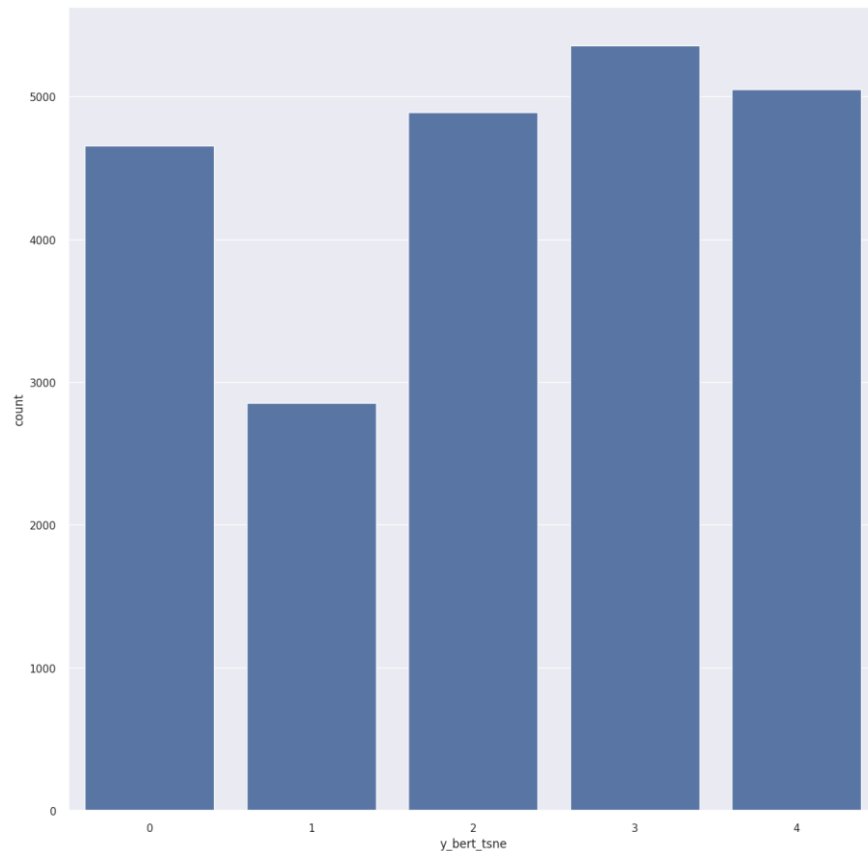
`eps=0.03, min_samples=5`

`silhouette_score`

0.4867697732251089

Davies-Bouldin index: 1.650833686858342





در آخر هم نموداری داشتیم که میتوانست مقدار و بالانس بودن هر یک از کلاسترها را نشان دهد. که به دلیل جلوگیری از زاید شدن مطالب صرفاً در انتها آورده شد که در صورت لزوم بتوان از آن استفاده نمود.

نتیجه گیری:

بنظرم از میان همه روش های vectorization بین tf-idf و bow، نتیجه بهتری داشته است. چون بردارهای آن بر اساس وزن کلماتی است که تاثیر بیشتری دارند. و بنابراین به نسبت عملکرد بهتری داشته است.

در مورد الگوریتم ها وابسته به مقادیر داده شده، دو الگوریتم k-means و mean shift در هر دو حالت bow و tf-idf عملکرد بهتری را با توجه به معیار سیلوئت داشته اند. اما به نسبت باز نتایج در tf-idf بهتر شده است.

در مورد دیگر الگوریتم های optics و DBscan مجدداً نتایج در tf-idf همانطور که به دلیل آن بیان شد، نتایج بهتری داشته است.

اما نتایج k-means و mean-shift بهتر بود آن ممکن است به خاطر جنس داده های متنی و اثری که وکتوری شدن روی آن ها گذاشته است باشد. احتمال دارد که ساختار گروهی شکل تاثیر بهتری را گذاشته باشد.

در کل از روش DB SCAN راضی نبودم و اینجا خوب جواب نداد و همچنین optics. شاید در این جور مسایل از روش های چگالی محور استفاده نشود بهتر باشد.