

Министерство науки и высшего образования Российской Федерации
федеральное государственное автономное образовательное учреждение высшего
образования
«НАЦИОНАЛЬНЫЙ ИССЛЕДОВАТЕЛЬСКИЙ УНИВЕРСИТЕТ ИТМО»

Отчет

по заданию №9
«Кластеризация»

по дисциплине «Компьютерная лингвистика»

Автор: Лакиза Александр Николаевич

Факультет: ИКТ

Группа: К3242

Преподаватель: Чернышева Анастасия Владимировна



УНИВЕРСИТЕТ ИТМО

Санкт-Петербург 2021

Цель работы: произвести кластеризацию своего корпуса по именованным сущностям

Ход работы:

Весь использованный функционал библиотеки `patasha` был взят мной из файла [NER_LDA](#)

Сначала подготовил свои данные

```
1. with open("../data/corpus_as_dict.json") as f:
2.     docs = json.load(f)
3. titles = list(docs.keys())
4. texts = list(docs.values())
5. df = pd.DataFrame.from_dict({'title':titles, 'text':texts})
```

Далее воспользовался функцией `get_ner`

```
1. def get_ner(transcript):
2.     script = Doc(re.sub(r'\((.*?)\)', "", transcript))
3.     script.segment(segmenter)
4.     script.tag_morph(morph_tagger)
5.     for token in script.tokens:
6.         token.lemmatize(morph_vocab)
7.     script.tag_ner(ner_tagger)
8.     for span in script.spans:
9.         span.normalize(morph_vocab)
10.    named_ents = [(i.text, i.type, i.normal) for i in script.spans]
11.    normed_ents = []
12.    for word, tag, norm in named_ents:
13.        if len(word.split()) == 1 and tag == "LOC":
14.            for gram in range(len(analyzer.parse(word))):
15.                if "Geox" in analyzer.parse(word)[gram].tag:
16.                    normed_ents.append((analyzer.parse(word)[gram].normal_form))
17.                    break
18.            elif gram == len(analyzer.parse(word)) - 1:
19.                normed_ents.append((norm.lower().strip(". , ! ? ; -")))
20.        else:
21.            normed_ents.append((norm.lower().strip(". , ! ? ; -")))
22.    return sorted(normed_ents)
```

Далее я посчитал, что из 363 документов лишь в 223 есть именованные сущности, именно с этими документами дальше и работал

```
1. vocabulary = sorted(set(ner_voc))
2. corpus = df_ner.named_entities.apply(str).tolist()
3.
```

```

4. pipe = Pipeline([('count', CountVectorizer(vocabulary=vocabulary)),
5.                  ('tfidf', TfidfTransformer())]).fit(corpus)
6. X = pipe.fit_transform(corpus)
7. km = KMeans(n_clusters=18, init='k-means++', max_iter=600,
8.             algorithm="full", precompute_distances=True)
9.
10. km.fit(X)

```

Затем оценил свою кластеризацию

```

1. print(metrics.silhouette_score(X, km.labels_, sample_size=1000))
2. print(metrics.davies_bouldin_score(X.toarray(), km.labels_))

```

Результат:

```

0.0868968053447224
1.3695426552227496

```

Вообще, мне кажется странным судить по данным метрикам, так как первый показатель всё время колеблется у одинакового значения, а второй становится больше (что плохо) при уменьшении кол-ва кластеров и наоборот уменьшается при увеличении числа кластеров

После предсказания получилось следующее распределение по кластерам:

2	170
13	10
14	8
3	7
12	4
15	4
17	3
6	3
8	3
5	2
11	2
16	1
0	1
10	1
1	1
7	1
4	1
9	1

(В первом столбце номер кластера, во втором – кол-во документов, относящихся к нему)

Посмотрим на результат кластеризации

Cluster 0: солонина японское море ехидну заморская территория залив белфас
т зал слава закавказье жюль мишле жюль ардуэн-мансар жоау фернандеша лэвра
дур

Cluster 1: итчен тест портсмут англия великобритания европа европейский со
юз евразия заморская территория залив белфаст

Cluster 2: канада уайт сша европа титаник великобритания белфаст франция м
орзе кунард

Cluster 3: сша америка nautilus вмс штат биштекс парадайз кларк стрип нев
ад

Cluster 4: англия шотландия нидерланды оранский японское море ж. картые за морская территория залив белфаст зал слава закавказье
Cluster 5: атлантика северная азовское евразия шпицберген черное америка а тлантический жоау фернандеша лэврадур жан кальвин
Cluster 6: гренландия америка северная земля атлантический южная австралия европа исландия антарктида
Cluster 7: непотопляемость трюм японское море ж. б. ламарком запад заморская территория залив белфаст зал слава закавказье жюль мишле
Cluster 8: испания мадрид европа африка португалия андорра марокко канарские гibraltar нато
Cluster 9: массачусетс бостон содружество англия атлантический сша жан калъвин женева жоау фернандеша лэврадур японское море
Cluster 10: эллипс западное полушарие западная европа запад заморская территория залив белфаст зал слава закавказье жюль мишле жюль ардуэн-мансар
Cluster 11: конгресс сша японское море ехидну заморская территория залив белфаст зал слава закавказье жюль мишле жюль ардуэн-мансар
Cluster 12: россия российская англия сша эстония ленинград нева снг кронштадт эрмитаж
Cluster 13: великобритания ирландия лондон англия северная шотландия уэльс георг уайтхолл европа
Cluster 14: земля луна российская солнце марс тейи венера арктика меркурий жан кальвин
Cluster 15: европа пари́ж шотландия франция рিশелье сена азия океания зевсом евразия
Cluster 16: грирсон японское море западное полушарие запад заморская территория залив белфаст зал слава закавказье жюль мишле жюль ардуэн-мансар
Cluster 17: франция пари́ж республика сена нато паризии сите евросоюз версаль оон

Видно, что в отличие от тематического моделирования из Задания 10, тут можно в каждом кластере примерно четко определить тему, что говорит нам о хорошем результате. Число кластеров – 18 взято, можно сказать, случайно. При любом кол-ве кластеров, тематика каждого кластера хорошо прослеживается

Ссылка на исходный код и json файл с корпусом:

<https://github.com/alexanderlakiza/cs224/tree/main/task9>