

2025-03-10

Тематическое моделирование

Тематическое моделирование (моделирование топики, topic modelling) — вид машинного обучения без учителя, основанный на статистическом моделировании.

Применяется для выявления тем в коллекции документов. Принцип тематического моделирования предполагает автоматическую кластеризацию слов, которые часто встречаются в документах в одном контексте (т.е. вместе), с целью выявления групп слов, характеризующих отдельные темы в коллекции документов. Конечная цель тематического моделирования — выявление K основных тем (топиков) в корпусе текстовых данных.

Основные сферы применения тематического моделирования

1. Кластеризация и классификацию документов: распределение документов по группам на основе их содержания (либо выявление этих групп, либо, когда топики уже есть, отнесение документов к какому-либо из них)
2. Информационный поиск: помощь поисковым системам в подборе наиболее релевантных документов
3. Квазиреферирование текста: сжатие больших фрагментов текста в более короткие рефераты с основной информацией
4. Сегментация клиентов: группировка клиентов на основе их отзывов или обзоров
5. Анализ тональности: определение положительного, отрицательного, нейтрального или иного заданного тона для большой коллекции текстов
6. Исследовательский анализ слабоструктурированных данных: обнаружение скрытых закономерностей и тем в большом корпусе текстовых данных

Основные этапы

1. Загрузить текстовые данные
2. Предварительная обработка данных: удаление стоп-слов, знаков препинания, другой нерелевантной информации; лемматизация.

3. Формирование матрицы «документ-термин»: предварительно обработанные текстовые данные преобразуются в матрицу количества слов, где каждая строка — документ, каждый столбец — уникальная лемма.
4. Запуск модели: наиболее популярные подходы включают латентное размещение (аллокацию) Дирихле (Latent Dirichlet Allocation, LDA), неотрицательное матричное разложение (non-negative matrix factorization, NMF), а также методы вроде [структурного тематического моделирования \(Structural Topic Models, STM\)](#). На вход модели подается сформированная ранее матрица.
5. В ходе моделирования формируется матрица «термин-топик»:

	hobby	pursue	ability	giving	...
Topic 4	0.5	0.3	0.0	0.0	
Topic 5	0.0	0.0	0.0	0.8	
...					
Topic 88	0.0	0.2	0.5	0.0	
...					

6. Далее формируется набор наиболее релевантных терминов для каждого топика (обычно по 10):

Topic 4: hobby, rewarding, pursue, pursue hobby, leisure time, leisure, passion, family hobby, time hobby, friend, social
Topic 5: giving, woman, grateful, charitable, giving community, nonprofit, philosophy
Topic 88: ability, contribute, ability travel, time, ability work, pursue, time pursue, family, community, contribute community, contribute society

7. Полученные топики можно далее использовать для визуализации, анализа и решения иных задач на корпусе текстов.

Например, для группировки текстов по темам достаточно составить матрицу «документ-топик», опираясь на статистику встречаемости ключевых терминов топика в каждом конкретном документе (т.е. через определение функции принадлежности) — можно сформировать нечеткое множество для нечеткого отнесения текстов в множестве тем.

Document	Topic 4	Topic 5	...	Topic 88	...
I have a number of great hobbies and plenty of time to pursue them	0.7	0.3		0.0	
Life is most rewarding when I'm giving back to the community	0.5	0.2		0.8	
I have a great job that gives me the ability to pursue meaningful goals at work	0.3	0.7		0.0	
...					

Лабораторная работа №5: Тематическое моделирование

1. Дана коллекция текстовых документов `2021_SPORT`
2. Необходимо провести тематическое моделирование коллекции методами латентного размещения (аллокации) Дирихле и неотрицательного матричного разложения и визуализировать темы с разным количеством топ-слов, реализовав функцию `plot_top_words`.
 1. Для тематического моделирования понадобится `matplotlib`, `sklearn.feature_extraction.text` (`TfidfVectorizer`, `CountVectorizer`), `sklearn.decomposition` (`NMF`, `LatentDirichletAllocation`)
 2. Для препроцессинга нужно воспользоваться материалами прошлой лабораторной (в части стоп-слов, токенизации и лемматизации)
 3. Затем построить через `CountVectorizer`, `TfidfVectorizer` (попробуйте оба) векторное представление каждого текста (для визуализации полезно также сделать что-то вроде `tfidf_feature_names = tfidf_vectorizer.get_feature_names()`)
 4. Полученные вектора передавать уже в модель для тематического моделирования
 5. Количество признаков — начать с 1000, топиков — 10, топ-слов — 20

как это может выглядеть (без реализации логики функций)

```
n_features = 1000
n_components = 10
n_top_words = 20
```

```
# тут все, собственно, происходит
```

```
if __name__ == '__main__':
    tfidf_vectorizer, tfidf = create_vectors_tf_idf('./sports')

    nmf = NMF(n_components=n_components, random_state=1, alpha=.1,
l1_ratio=.5).fit(tfidf)
    tfidf_feature_names = tfidf_vectorizer.get_feature_names()
    plot_top_words(nmf, tfidf_feature_names, n_top_words, 'Topics in NMF
model (Frobenius norm)')

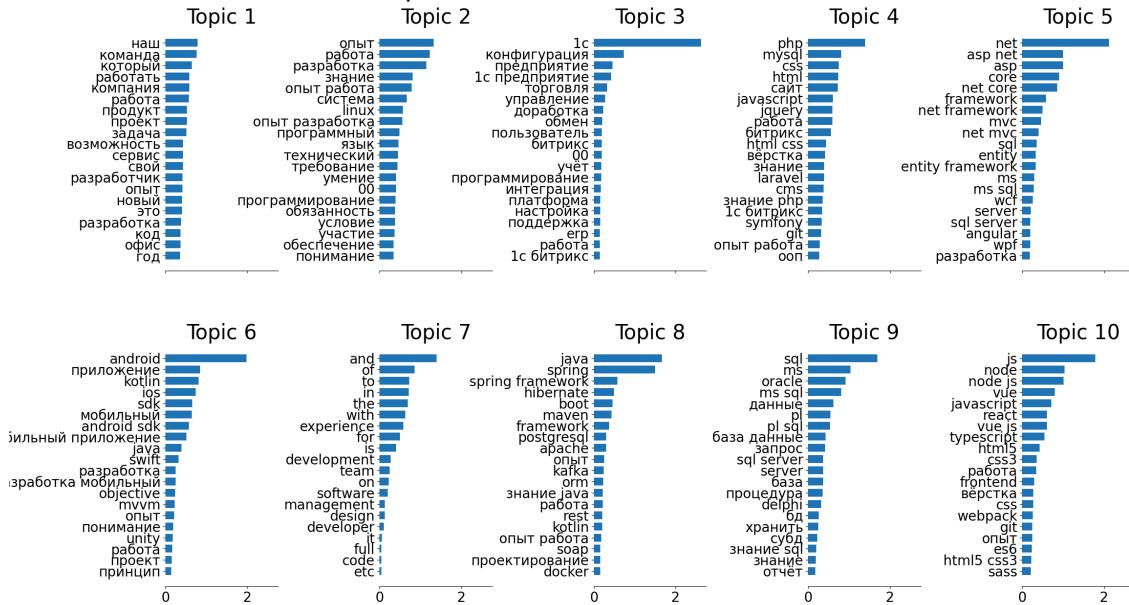
    nmf_k = NMF(n_components=n_components, random_state=1,
        beta_loss='kullback-leibler', solver='mu', max_iter=1000,
alpha=.1,
        l1_ratio=.5).fit(tfidf)

    plot_top_words(nmf_k, tfidf_feature_names, n_top_words, 'Topics in
NMF model (generalized Kullback-Leibler divergence)')

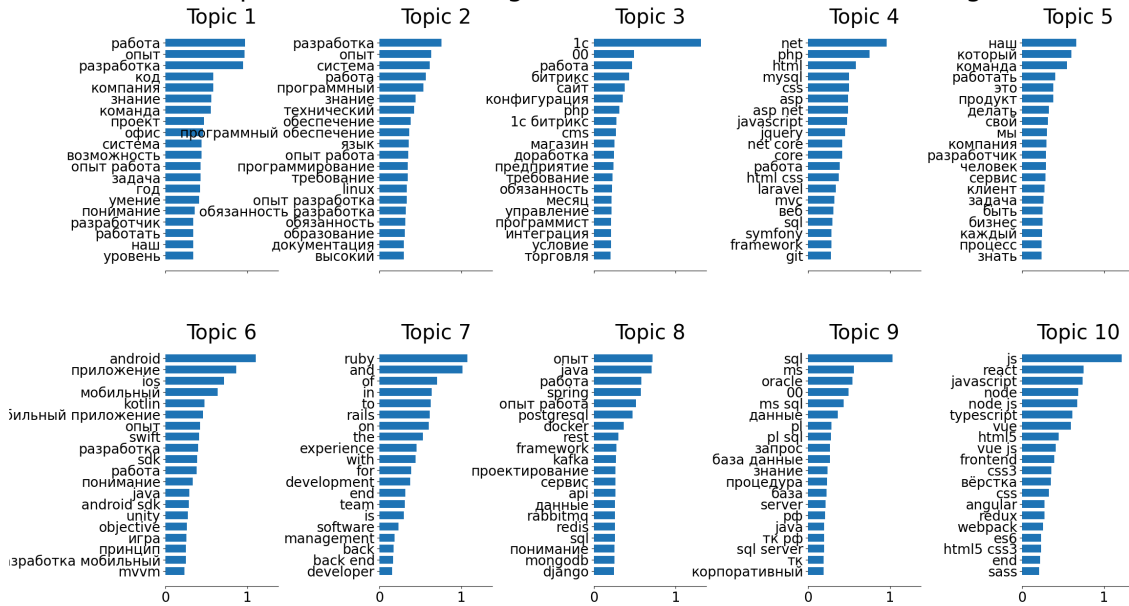
    tf_vectorizer, tf = create_vectors_count('./sports')
    lda = LatentDirichletAllocation(n_components=n_components,
max_iter=5,
                                learning_method='online',
                                learning_offset=50.,
                                random_state=0)

    lda.fit(tf)
    tf_feature_names = tf_vectorizer.get_feature_names()
    plot_top_words(lda, tf_feature_names, n_top_words, 'Topics in LDA
model')
```

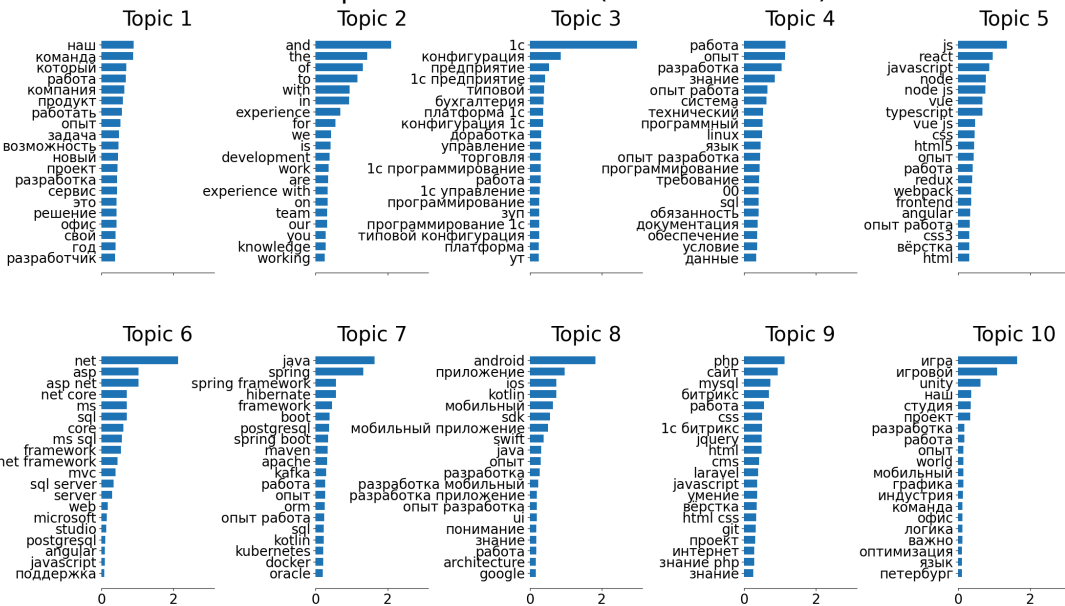
Topics in NMF model (Frobenius norm)



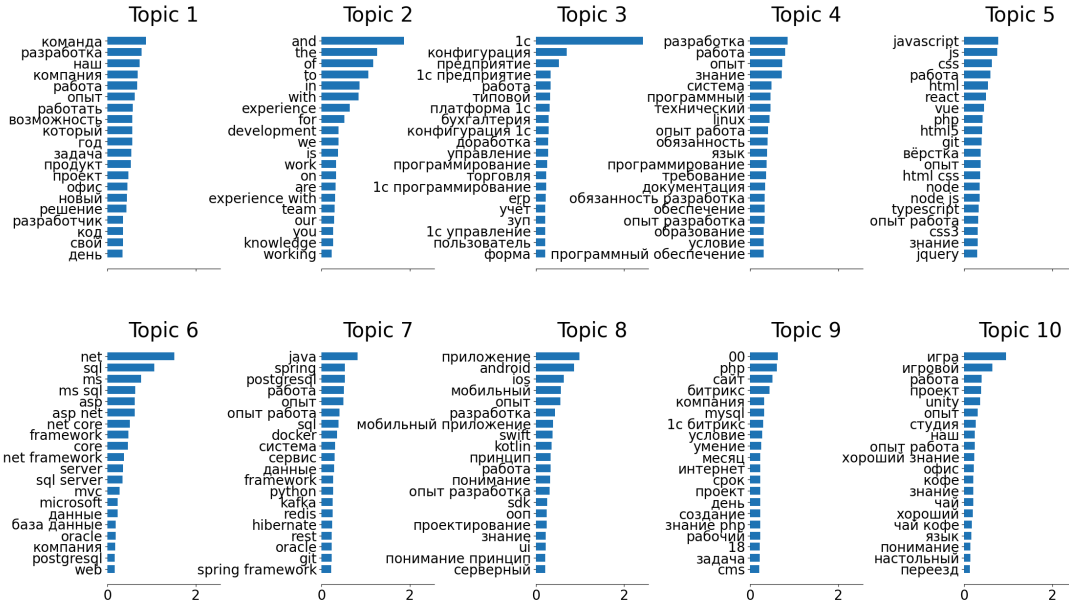
Topics in NMF model (generalized Kullback-Leibler divergence)



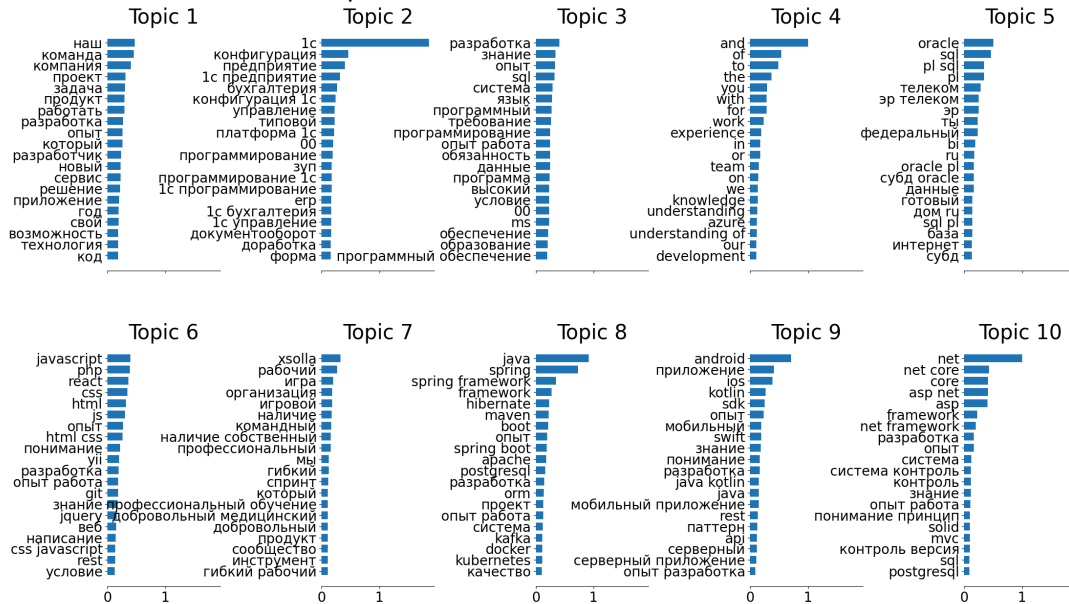
Topics in NMF model (Frobenius norm)



Topics in NMF model (generalized Kullback-Leibler divergence)



Topics in NMF model (Frobenius norm)



Topics in NMF model (generalized Kullback-Leibler divergence)

