

Анализ либретто мюзиклов на русском языке

Анастасия Бутакова
ДПО «Компьютерная лингвистика»

старт:

Насколько тексты похожи друг на друга, какие можно выделить тематики?

В итоге:

...лексическое разнообразие?

Сбор корпуса – тексты либретто

1. Из открытых источников — если текст не файлом, то парсинг данных (Beautiful Soup), чтобы сразу было нужное форматирование
 - a. Официальные тексты, выложенные в сеть
 - b. Коллективные расшифровки, сделанные фанатами
2. Запросы к авторам, артистам, критикам, «старожилам» фандома
 - a. Профи игнорируют или не готовы делиться материалами
 - b. Авторы и театры перебрасывают друг на друга ответственность
 - c. У любителей редко что-то есть
3. «Слитые» в сеть материалы

Все тексты оформлены совершенно по-разному и часто содержат доп.информацию, не нужную для анализа

72

текста собрано

30

использовано в проекте

Использованные библиотеки

**os, re, docx2python,
pdfminer.six**

Превращение непонятно
как оформленного текста
непонятно с чем внутри в
рабочий материал

Nltk, pymorphy3

Токенизация,
лемматизация удаление
стоп-слов; определение
частей речи и проч

sklearn

Tf-idf
Косинусная близость

Matplotlib, Seaborn

Визуализация

Pandas, numpy

Работа с
датафреймами и не
только

Stanza, spacy

Попытки работать с
именованными
сущностями (тщетные)

Этапы работы

1. Сбор корпуса
2. Приведение текстов в единообразный вид (форматирование и метки) и сохранение в txt
3. Отсеивание ненужной информации из текста (названия песен, сценическая информация, повествовательные элементы, имена персонажей, которые говорят, технические пометки)
4. Остается только **разговорный текст** (песни, диалоги) => предобработка
5. Различные манипуляции: TF-IDF, косинусная близость, оценка лексического разнообразия

```
# Объединяем редактирование ВСЕХ возможных форматов файлов. Вариант с ОПЦИЯМИ
def v2_complete_formatting_woptions(filename, set_info, speakers, to_remove):
    '''Единая функция для форматирования текстов в формате text, docx/doc и pdf для частных случаев. [C] возможностью ввода специфических паттернов,
    которые позволят обработать конкретно этот текст.
    Настройка параметров:
    1. set_info - информация [O] сценографии: [O] действиях персонажей, технических моментах, характеристиках реплик типа "[C] подозрением", "тихо сме
        - 'none' - явного паттерна нет или все стандартно по нашей схеме
        - ...any custom value here... - паттерн, если есть единый стиль, которым обозначена эта информация
    2. speakers - говорящие/поющие:
        - 'none' - все в порядке, доп. действия здесь не требуется
        - 'absent' - если они не указаны в тексте вообще (поможет избежать ненужных исправлений)
        - ...any custom value here... - паттерн, если [U] спикеров есть явный паттерн, но не тот, что мы взяли по умолчанию
    3. to_remove - специфические пометки, которые захламляют текст и не нужны (например, номера страниц и колонтитулы в пдф, примечания в тексте
        - 'none' - пометки для удаления отсутствуют
        - ...any custom value here... - паттерн(ы) для удаления, задаются в виде списка
```

TF-IDF - провал?

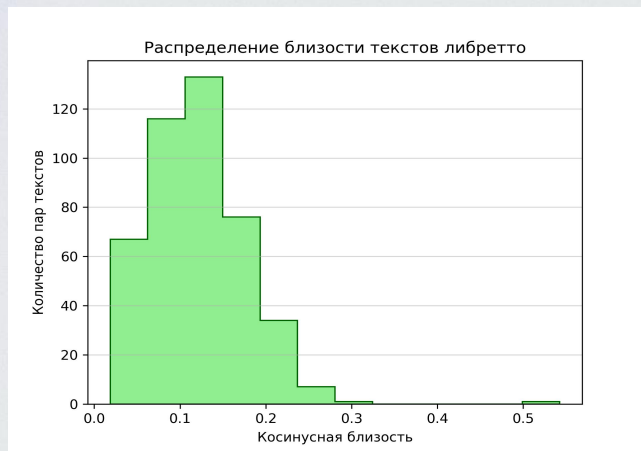
- Проблема именованных сущностей - они становятся ключевыми словами, но их сложно удалить
- Специфика стихотворного текста, которая жестко усугубляет проблему с NER
- Колебания длины текстов (есть немного гораздо более коротких, чем общая масса)

**⇒ чтобы сравнивать содержание векторизованных текстов ,
нужна еще более глубокая предобработка**

```
{('АББА', 383),  
 ('Алёхин', 4317129024397789502),  
 ('Анатолий', 4317129024397789502),  
 ('Анатолий Сергиевский', 4317129024397789502),  
 ('Анотолия', 4317129024397789502),  
 ('Англии', 385),  
 ('Андерсен', 385),  
 ('Ареной идеологической борьбы', 383),  
 ('Ах', 4317129024397789502),  
 ('Бангкок', 385),  
 ('Бангкоке', 385),  
 ('Безумствуй', 4317129024397789502),  
 ('Ботвинник', 4317129024397789502),  
 ('Будапешт', 385),  
 ('Будапеште', 385),  
 ('Будь', 4317129024397789502),  
 ('Быть', 4317129024397789502),  
 ('Бъём', 4317129024397789502),  
 ('Важно', 4317129024397789502),  
 ('Ваш', 4317129024397789502),  
 ('Вдруг', 4317129024397789502),  
 ('Ведь', 4317129024397789502),
```

```
  "text": "Врёшь Поздно Лжец Поздно Поздно Сумел смогла одно я Есть ли кто-нибудь",  
  "type": "MISC",  
  "start_char": 30562,  
  "end_char": 30632  
},  
{  
  "text": "Мой триумф",  
  "type": "MISC",  
  "start_char": 30683,  
  "end_char": 30693  
},  
{  
  "text": "Дамы",  
  "type": "PER",  
  "start_char": 30715,  
  "end_char": 30719  
},  
{  
  "text": "Анатолий Сергиевский",  
  "type": "PER",  
  "start char": 30730,
```


Косинусная близость текстов либретто



- 0.54257 - косинусная близость РОМЕО И ДЖУЛЬЕТТА и РОМЕО VS. ДЖУЛЬЕТТА: XX ЛЕТ СПУСТЯ
- 0.30709 - косинусная близость ЭЛОЯ и МАСТЕР И МАРГАРИТА
- 0.27332 - косинусная близость ПЁТР ПЕРВЫЙ и БЕЛЫЙ. ПЕТЕРБУРГ
- 0.26868 - косинусная близость ЭЛОЯ и НОТР-ДАМ ДЕ ПАРИ
- 0.24766 - косинусная близость МЕТРО и ПРАЙМТАЙМ
- 0.24691 - косинусная близость ЭЛОЯ и ПЕТЯ И ФОЛК
- 0.24638 - косинусная близость ЭЛОЯ и МАММА МΙΑ!
- 0.24461 - косинусная близость ЭЛОЯ и РОМЕО VS. ДЖУЛЬЕТТА: XX ЛЕТ СПУСТЯ
- 0.24274 - косинусная близость МАСТЕР И МАРГАРИТА и НОТР-ДАМ ДЕ ПАРИ

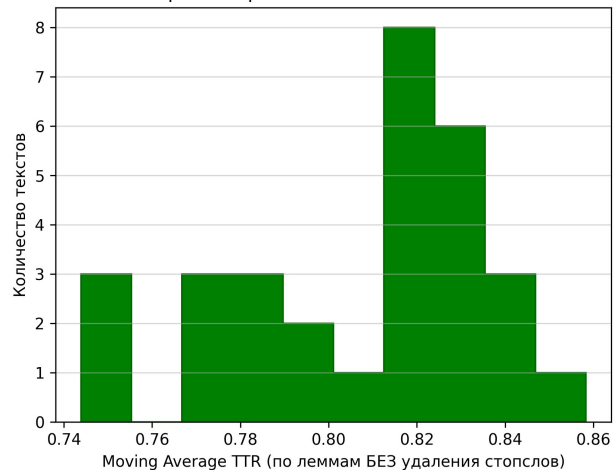
Косинусная близость либретто музыков																															
Айседель	1	0.09	0.0760	0.091	0.09	0.12	0.0340	0.0660	0.075	0.1	0.11	0.1	0.0230	0.0780	0.0770	0.086	0.1	0.12	0.11	0.1	0.08	0.12	0.09	0.078	0.11	0.0760	0.0650	0.058	0.13	0.11	
Бел Валентин	0.09	1	0.09	0.15	0.13	0.23	0.0540	0.0940	0.17	0.17	0.19	0.23	0.034	0.13	0.13	0.19	0.21	0.16	0.098	0.17	0.13	0.12	0.16	0.16	0.18	0.17	0.11	0.11	0.2	0.14	
Кошки	0.076	0.09	1	0.0810	0.0690	0.0840	0.0270	0.0570	0.0650	0.0830	0.092	0.11	0.0210	0.0540	0.0780	0.0720	0.0890	0.0930	0.092	0.08	0.0590	0.0690	0.0760	0.0540	0.0860	0.0720	0.0580	0.056	0.12	0.072	
Шахматы	0.091	0.15	0.081	1	0.16	0.2	0.052	0.11	0.16	0.17	0.22	0.16	0.032	0.13	0.13	0.14	0.14	0.17	0.1	0.17	0.11	0.17	0.17	0.13	0.19	0.15	0.07	0.15	0.2	0.17	
Принцесса цирка	0.09	0.13	0.069	0.16	1	0.13	0.0420	0.0950	0.094	0.15	0.14	0.12	0.028	0.1	0.1	0.13	0.12	0.15	0.13	0.097	0.1	0.15	0.11	0.09	0.15	0.11	0.0670	0.089	0.2	0.12	
Элоя	0.12	0.23	0.084	0.2	0.13	1	0.057	0.11	0.21	0.23	0.25	0.31	0.034	0.13	0.13	0.22	0.27	0.18	0.11	0.25	0.14	0.15	0.19	0.23	0.24	0.22	0.14	0.16	0.19	0.18	
Волосы	0.0340	0.0540	0.0270	0.0520	0.0420	0.057	1	0.0410	0.0390	0.0510	0.0670	0.0560	0.0260	0.026	0.06	0.04	0.05	0.0440	0.0340	0.057	0.04	0.0580	0.0430	0.0330	0.0490	0.056	0.03	0.0530	0.054	0.04	
Джейн Эйр	0.0660	0.0940	0.057	0.11	0.095	0.11	0.041	1	0.078	0.1	0.11	0.1	0.024	0.08	0.0820	0.0890	0.097	0.1	0.0880	0.0890	0.0730	0.0910	0.0870	0.083	0.13	0.0790	0.0560	0.063	0.12	0.076	
Христос Суперзвезда (Моссвет)	0.075	0.17	0.065	0.16	0.094	0.21	0.0390	0.078	1	0.15	0.14	0.22	0.0260	0.0990	0.096	0.14	0.18	0.14	0.077	0.15	0.13	0.12	0.2	0.15	0.17	0.16	0.11	0.098	0.14	0.15	
Анна Каренина	0.1	0.17	0.083	0.17	0.15	0.23	0.051	0.1	0.15	1	0.19	0.21	0.03	0.12	0.12	0.19	0.21	0.17	0.13	0.18	0.17	0.17	0.18	0.15	0.11	0.1	0.22	0.19	0.1		
МАММА МΙΑ!	0.11	0.19	0.092	0.22	0.14	0.25	0.067	0.11	0.14	0.19	1	0.22	0.036	0.11	0.16	0.16	0.22	0.17	0.1	0.19	0.12	0.17	0.15	0.14	0.2	0.16	0.096	0.13	0.22	0.12	
Мастер и Маргарита	0.1	0.23	0.11	0.16	0.12	0.31	0.056	0.1	0.22	0.21	0.22	1	0.038	0.13	0.13	0.18	0.24	0.18	0.094	0.21	0.16	0.13	0.18	0.19	0.22	0.18	0.13	0.12	0.2	0.16	
Наугли	0.0230	0.0340	0.0210	0.0320	0.0280	0.0340	0.0260	0.0240	0.026	0.03	0.0360	0.03	1	0.0190	0.0480	0.0290	0.0280	0.0280	0.0220	0.0360	0.0320	0.037	0.03	0.0210	0.0290	0.034	0.02	0.0270	0.0570	0.032	
Мертвые души	0.078	0.13	0.054	0.13	0.1	0.13	0.026	0.08	0.099	0.12	0.11	0.13	0.015	1	0.075	0.12	0.12	0.13	0.1	0.097	0.13	0.12	0.14	0.1	0.13	0.0910	0.0980	0.058	0.14	0.2	
Метро	0.077	0.13	0.078	0.13	0.1	0.13	0.06	0.0820	0.096	0.12	0.16	0.13	0.0480	0.075	1	0.1	0.13	0.13	0.067	0.13	0.1	0.25	0.1	0.081	0.11	0.0820	0.0660	0.088	0.14	0.1	
Монте-Кристо	0.086	0.19	0.072	0.14	0.13	0.22	0.04	0.089	0.14	0.19	0.16	0.18	0.029	0.12	0.1	1	0.21	0.16	0.11	0.13	0.12	0.11	0.15	0.18	0.19	0.13	0.0920	0.091	0.19	0.14	
Нотр-Дам де Пари	0.1	0.21	0.089	0.14	0.12	0.27	0.05	0.097	0.18	0.21	0.22	0.24	0.028	0.12	0.13	0.21	1	0.16	0.1	0.18	0.12	0.12	0.16	0.23	0.2	0.16	0.13	0.099	0.21	0.14	
Норд-Ост	0.12	0.16	0.093	0.17	0.15	0.18	0.044	0.1	0.14	0.17	0.17	0.18	0.028	0.13	0.13	0.16	0.16	1	0.14	0.15	0.16	0.16	0.16	0.12	0.17	0.14	0.1	0.09	0.19	0.2	
Остров сокровищ	0.11	0.0980	0.092	0.1	0.13	0.11	0.0340	0.0860	0.077	0.11	0.1	0.0940	0.022	0.1	0.067	0.11	0.1	0.14	1	0.11	0.082	0.11	0.12	0.069	0.11	0.09	0.0640	0.061	0.13	0.11	
ПЕТА И ФОЛК	0.1	0.17	0.08	0.17	0.097	0.25	0.0570	0.089	0.15	0.18	0.19	0.21	0.0360	0.097	0.13	0.13	0.18	0.15	0.11	1	0.13	0.14	0.17	0.16	0.17	0.11	0.12	0.12	0.19	0.14	
Преступление и наказание	0.08	0.13	0.059	0.11	0.1	0.14	0.04	0.073	0.13	0.13	0.12	0.16	0.032	0.13	0.1	0.12	0.12	0.16	0.082	0.13	1	0.13	0.14	0.12	0.14	0.11	0.12	0.077	0.14	0.15	
Праймтайм	0.12	0.12	0.069	0.17	0.15	0.15	0.0580	0.091	0.12	0.18	0.17	0.13	0.037	0.12	0.25	0.11	0.12	0.16	0.11	0.14	0.13	1	0.13	0.092	0.14	0.12	0.0830	0.096	0.2	0.15	
ПЁТР ПЕРВЫЙ	0.09	0.16	0.076	0.17	0.11	0.19	0.0430	0.087	0.2	0.17	0.15	0.18	0.03	0.14	0.1	0.15	0.16	0.16	0.12	0.17	0.14	0.13	1	0.14	0.16	0.15	0.12	0.098	0.17	0.27	
Ромео и Джульетта	0.078	0.16	0.054	0.13	0.09	0.23	0.0330	0.083	0.15	0.17	0.14	0.19	0.021	0.1	0.081	0.18	0.23	0.12	0.069	0.14	0.12	0.092	0.14	1	0.54	0.13	0.0860	0.072	0.17	0.11	
VS. ДЖУЛЬЕТТА: XX ЛЕТ СПУСТЯ	0.11	0.18	0.086	0.19	0.15	0.24	0.049	0.13	0.17	0.18	0.2	0.22	0.029	0.13	0.11	0.19	0.2	0.17	0.11	0.16	0.14	0.14	0.16	0.54	1	0.16	0.1	0.11	0.2	0.14	
Сон у Красной горы	0.076	0.17	0.072	0.15	0.11	0.22	0.0560	0.079	0.16	0.15	0.16	0.18	0.0340	0.0910	0.082	0.13	0.16	0.14	0.09	0.17	0.11	0.12	0.15	0.13	0.16	1	0.0950	0.098	0.14	0.13	
Вайолет	0.065	0.11	0.058	0.07	0.067	0.14	0.03	0.056	0.11	0.11	0.096	0.13	0.02	0.0980	0.0660	0.092	0.13	0.1	0.064	0.11	0.12	0.083	0.12	0.086	0.1	0.095	1	0.05	0.1	0.11	
Волух 2112	0.058	0.11	0.056	0.15	0.089	0.16	0.0530	0.0630	0.098	0.1	0.13	0.12	0.0270	0.0580	0.0880	0.0910	0.099	0.09	0.061	0.12	0.0770	0.0960	0.0980	0.072	0.11	0.098	0.05	1	0.13	0.088	
Всё о Золушке	0.13	0.2	0.12	0.2	0.2	0.19	0.054	0.12	0.14	0.22	0.22	0.2	0.057	0.14	0.14	0.19	0.21	0.19	0.13	0.19	0.14	0.2	0.17	0.17	0.2	0.14	0.1	0.13	1	0.18	
Белый. Петербург	0.11	0.14	0.072	0.17	0.12	0.18	0.04	0.076	0.15	0.19	0.12	0.16	0.032	0.2	0.1	0.14	0.14	0.2	0.11	0.14	0.15	0.15	0.27	0.11	0.14	0.13	0.11	0.088	0.18	1	
Айседель																															
Бел Валентин																															
Кошки																															
Шахматы																															
Принцесса цирка																															
Элоя																															
Волосы																															
Джейн Эйр																															
Христос Суперзвезда (Моссвет)																															
Анна Каренина																															
МАММА МΙΑ!																															
Мастер и Маргарита																															
Наугли																															
Мертвые души																															
Метро																															
Монте-Кристо																															
Нотр-Дам де Пари																															
Норд-Ост																															
Остров сокровищ																</															

Окей, посмотрим что-то еще

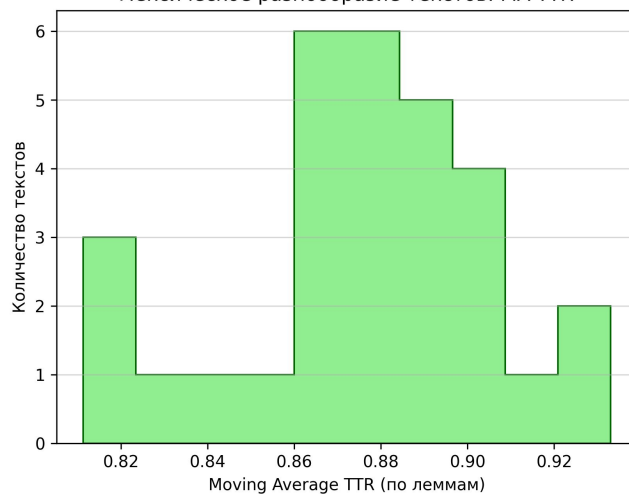
	Название	N токенов	N уникальных лемм	TTR	MA TTR	MA TTR (леммы)	MA TTR (леммы со стопсловами)
0	Айсвилль	5616	1654	0.294516	0.863413	0.877921	0.819641
1	Бал Вампиров	6231	1874	0.300754	0.863986	0.901561	0.820450
2	Кошки	4434	1636	0.368967	0.856078	0.899166	0.826878
3	Шахматы	5777	1534	0.265536	0.874654	0.922497	0.826201
4	Принцесса цирка	10282	2070	0.201323	0.870605	0.880896	0.824300
5	Эля	6003	1475	0.245710	0.827608	0.871410	0.775724
6	Волосы	1407	709	0.503909	0.870692	0.874671	0.835228
7	Джейн Эйр	3698	1219	0.329638	0.871044	0.902343	0.822532
8	Иисус Христос Суперзвезда (Моссовет)	4056	1169	0.288215	0.831270	0.875627	0.782511
9	Анна Каренина	3931	1157	0.294327	0.786425	0.814414	0.743854
10	МАММА MIA!	2429	740	0.304652	0.839529	0.866393	0.771655
11	Мастер и Маргарита	4034	1481	0.367129	0.897882	0.933068	0.858334
12	Маугли	1207	570	0.472245	0.869585	0.862733	0.826511
13	Мёртвые души	5556	1964	0.353492	0.861812	0.887866	0.831589

Лексическое е разнообраз ие текстов

Лексическое разнообразие текстов: МА TTR со стопсловами



Лексическое разнообразие текстов: МА TTR



Лексическое разнообразие текстов: TTR

