# Практическая работа №1.

#### Тямгин Иван

# Латентно-семантический анализ.

Для обучения поисковой системы были использованы 15 тем по 50-75 документов. Всего 950 документов. Названия документов для каждой темы были получены с помощью поискового запроса в википедии вида:

http://ru.wikipedia.org/w/index.php?title=Служебная:Поиск&limit=50&offset=0&profile=default&search=Таврия

Затем скачаны по полученным ссылкам средствами языка python.

Были выбраны следующие темы:

- Тригонометрия
- Грипп
- Автомобиль
- Хлеб
- Москва
- Техника
- Таврия
- Биржа
- Депутат
- Газ
- Яндекс
- Водка
- Русский
- Видео
- Собака

Первым шагом были вырезаны из текстов все символы, кроме русских букв. Таким образом слово — это непрерывная последовательность русских букв. Каждое слово поддалось воздействию алгоритма Портера, чтобы одинаковые слова с разными окончаниями не считались как разные.

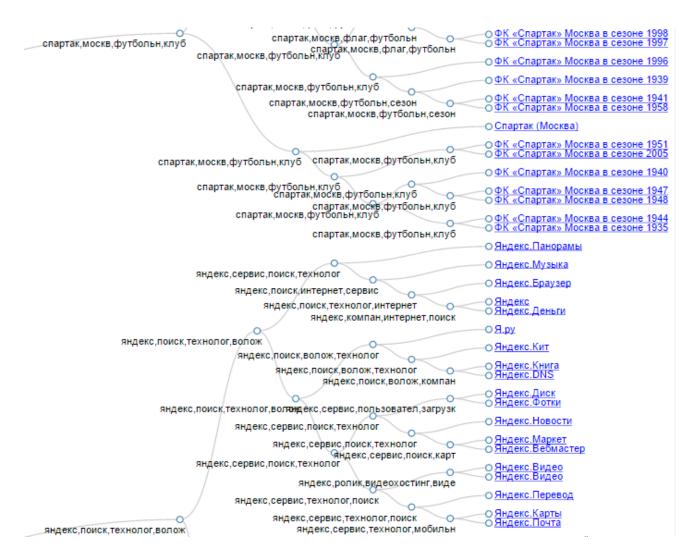
Была построена матрица 31238x950, которая означала частоту і-го слова в ј-м документе. Матрица была нормализована с помощью TF-IDF, чтобы ключевые слова получили большой вес, а предлоги и союзы не значительный.

После чего матрица была разложена с помощью svd по 10-ми признакам. (разложение по 2-3 признакам дали существенно худший результат).

# Дендрограмма

Далее построена иерархическая кластеризация документов. Расстоянием между двумя документами являлся угол между 10-мерными векторами.

Так как дендрограмма очень большая, приведена её часть:



Также попадались довольно абсурдные объединения документов на первом уровне:



Дерево было построено с помощью стандартной функции **hclust** в R. И нарисовано с помощью библиотеки **d3js** в *javascript*.

#### Поиск

Для каждого слова был построен чемпионский список – список документов и расстояние до этого документа. Чем меньше расстояние - тем он релевантнее для данного слова.

Поисковой запрос – неупорядоченный набор слов. Ответ на запрос – упорядоченный список документов по релевантности. Релевантность документа D - это

$$\sum_{\text{сдово S из запроса}} \frac{\text{расстояние от S до D}}{1 + \text{количество слов из запроса, присутствующих в D}}$$

### Примеры запросов:

### {Москва, вокзал}

- 1. <a href="https://ru.wikipedia.org/wiki/Казанский вокзал">https://ru.wikipedia.org/wiki/Казанский вокзал</a> <sup>91.13%</sup>
- 2. <a href="https://ru.wikipedia.org/wiki/Ярославский вокзал">https://ru.wikipedia.org/wiki/Ярославский вокзал</a> <sup>90.95%</sup>
- 3. <a href="https://ru.wikipedia.org/wiki/Ленинградский вокзал">https://ru.wikipedia.org/wiki/Ленинградский вокзал</a> 90.84%
- 4. <a href="https://ru.wikipedia.org/wiki/Белорусский вокзал">https://ru.wikipedia.org/wiki/Белорусский вокзал</a> 90.17%
- 5. <a href="https://ru.wikipedia.org/wiki/Киевский вокзал">https://ru.wikipedia.org/wiki/Киевский вокзал</a> 90.02%
- 6. <a href="https://ru.wikipedia.org/wiki/Курский вокзал">https://ru.wikipedia.org/wiki/Курский вокзал</a> 89.61%
- 7. https://ru.wikipedia.org/wiki/Москва-Сити 89.37%
- 8. <a href="https://ru.wikipedia.org/wiki/Mocква-Каланчёвская">https://ru.wikipedia.org/wiki/Mocква-Каланчёвская</a> 89.19%
- 9. <a href="https://ru.wikipedia.org/wiki/Mocква">https://ru.wikipedia.org/wiki/Mocква</a> (гостиница в Москве) 88.91%
- 10. https://ru.wikipedia.org/wiki/Памятник Кутузову (Москва)

# {Яндекс}

- 1. https://ru.wikipedia.org/wiki/Яндекс Нано <sup>99.88%</sup>
- 2. <a href="https://ru.wikipedia.org/wiki/Яндекс.Календарь">https://ru.wikipedia.org/wiki/Яндекс.Календарь</a>
- 3. <a href="https://ru.wikipedia.org/wiki/Элементы Яндекса">https://ru.wikipedia.org/wiki/Элементы Яндекса</a> 99.87%
- 4. <a href="https://ru.wikipedia.org/wiki/Яндекс.Услуги">https://ru.wikipedia.org/wiki/Яндекс.Услуги</a> <a href="https://ru.wikipedia.org/wiki/Яндекс.Услуги">99.85%</a>
- 5. <a href="https://ru.wikipedia.org/wiki/Яндекс.Открытки">https://ru.wikipedia.org/wiki/Яндекс.Открытки</a>
- 6. https://ru.wikipedia.org/wiki/Яндекс.Shell 99.84%
- 7. https://ru.wikipedia.org/wiki/Yandex.SpeechKit 99.84%
- 8. https://ru.wikipedia.org/wiki/Яндекс.XML 99.83%
- 9. <a href="https://ru.wikipedia.org/wiki/Яндекс.Недвижимость">https://ru.wikipedia.org/wiki/Яндекс.Недвижимость</a> 99.83%
- 10. <a href="https://ru.wikipedia.org/wiki/Яндекс.Навигатор">https://ru.wikipedia.org/wiki/Яндекс.Навигатор</a> 99.83

Более простая версия формулы:

$$\sum_{\text{слово S из запроса}}$$
 расстояние от S до D

давала на вид более худший результат.