

# Конкурс по текстовой релевантности

Ира Букреева

- обработка документов
- обработка запросов
- ранжирование

# обработка документов

- лемматизация - батчами по 10 текстов стеммеру Mystem
- удалила стоп-слова из текстов документов: словарь nltk + еще 112 добавленных руками, всего 264
- сначала брала просто тексты и заголовки, потом еще заголовки различных уровней
- 3 разных словаря

# обработка запросов

много поисковых расширений

- исправление опечаток
- синонимы
- транслиты
- ошибочные написания
- более сложные синонимы: почему -> причины
- класс запроса
- разные части речи
- в аббревиатуры и обратно
-

# ранжирование

- BM25

итоговая формула:  $\text{score} = 3.0 * \text{titlescore} + 1.2 * \text{textscore} + 1.0 * \text{headersscore}$

- `gensim.models.doc2vec`

обучение на тайтлах: 37971 тайтл

лучший скор: переранжирование документов на основе косинусного расстояния между эмбедингом запроса(без поисковых расширений) и заголовка документа с учетом сора BM25

у 143 документов нет заголовков