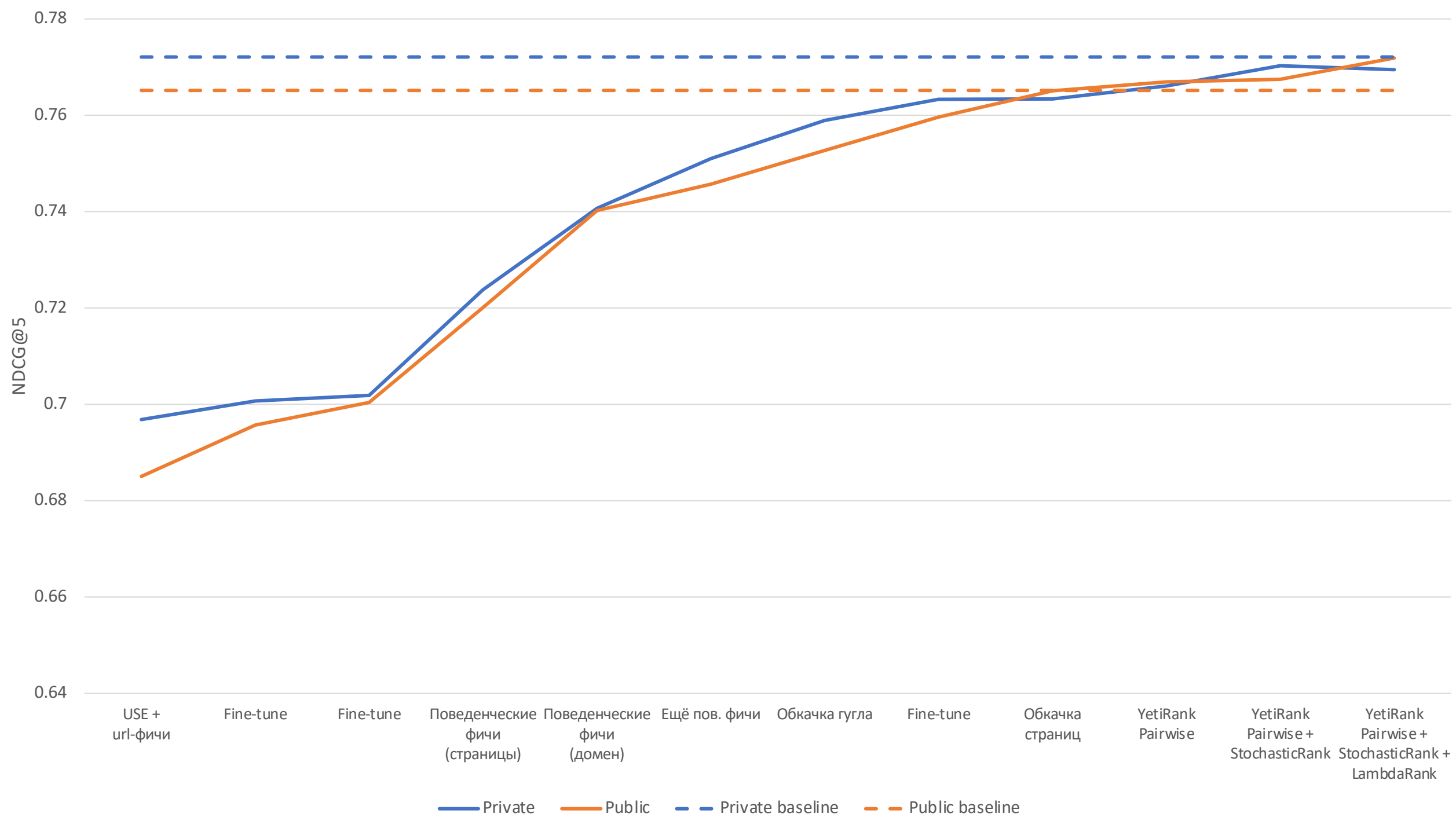


Финальный проект

или как пройти бейзлайн не до конца



Предобработка данных

- Запросы: pyaspeller; pymystem3; удаление стоп-слов nltk
- Заголовки: pymystem3

Дополнительные данные

- SERPy гугла
- Сами страницы с документами

URL и текстовые фиши

- Доменная зона (one-hot)
- Porn / forum / news (простой эмпирический классификатор)
- USE: склярное произведение: запрос (правленный / неправленный) и тайтл, запрос и текст
- USE: разность скоров запрос-тайтл / запрос-текст
- BM25 на тайтлах
- Количество одинаковых слов запрос-тайтл, количество уникальных слов запрос/тайтл
- Год в запросе и документе, совпадают ли года

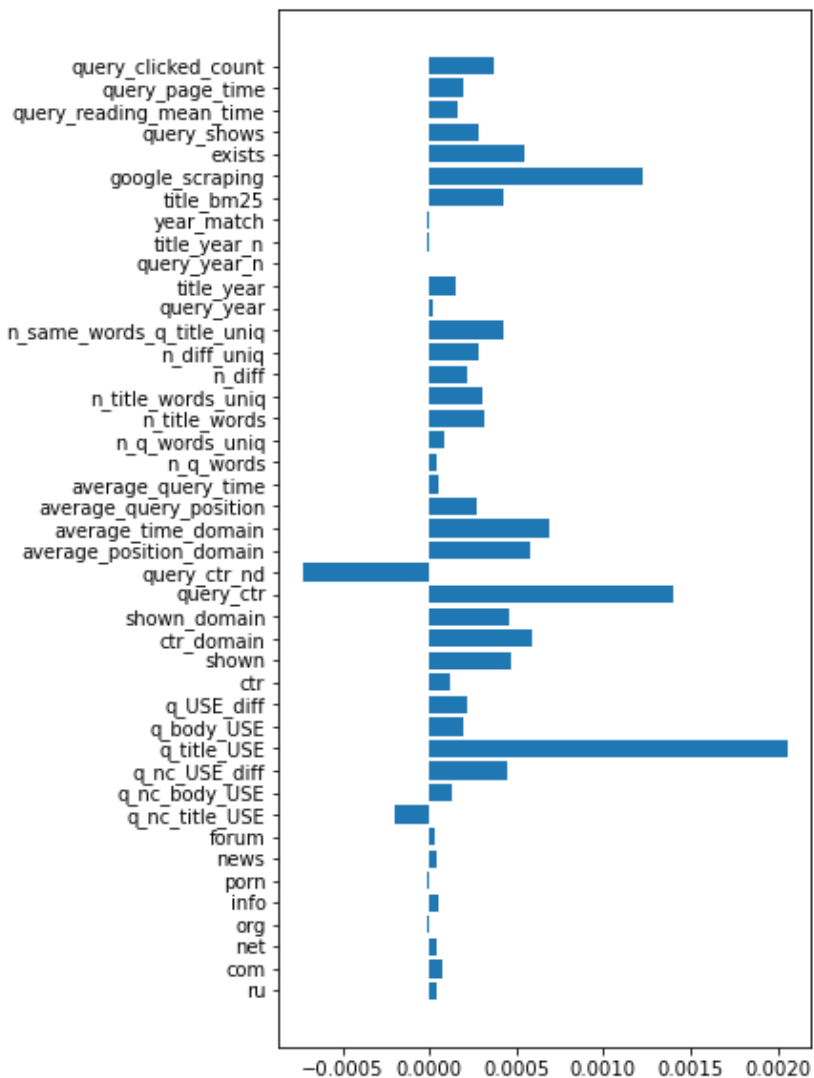
Поведенческие фичи

- По запросу, документу и домену:
 - Количество показов / кликов
 - Количество показов / кликов | запрос
 - Среднее время на SERPe
 - Среднее время перед кликом
 - Средняя позиция документа в SERPe

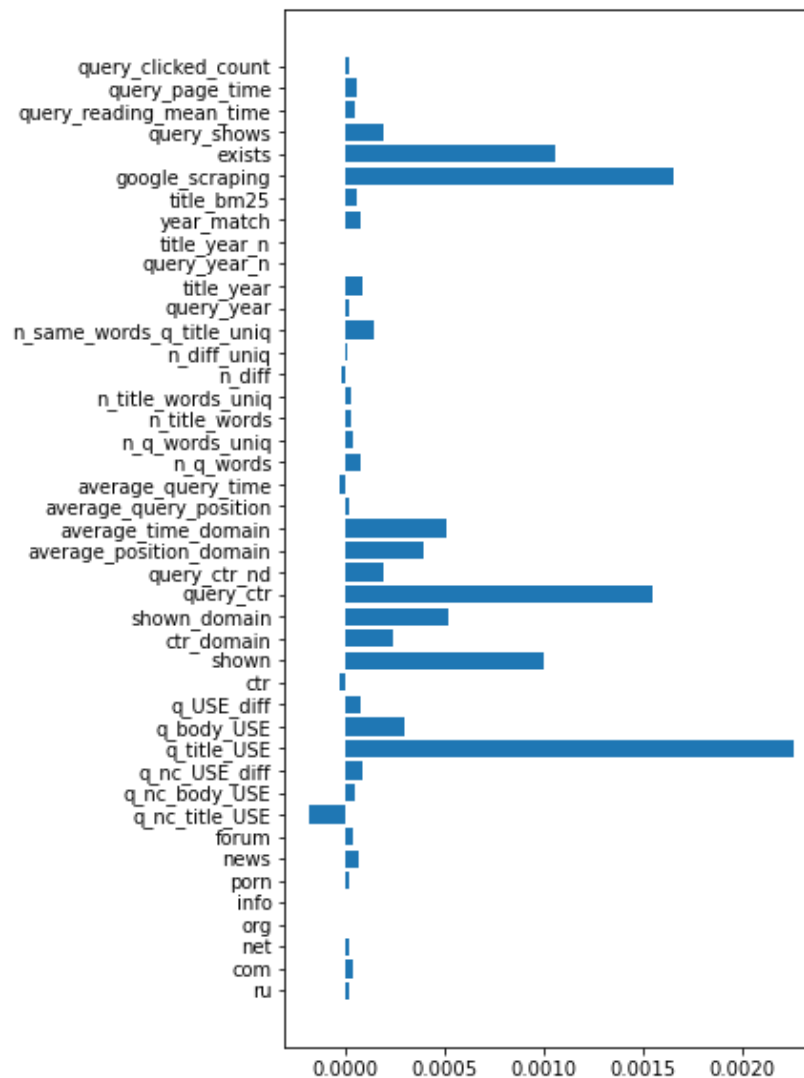
Хитрые фичи

- Наличие документа в выдаче гугла по запросу и его позиция
- Доступность страницы с документом и её размер в байтах

YetiRankPairwise



StochasticRank



LambdaRank

