

# Решение конкурса по ранжированию

Шилов Валентин

Public score: 0.78486

Private score: 0.78966

# Подготовка данных

- ▶ Коррекция запросов через Яндекс спеллер
- ▶ Нормализация заголовков, запросов и доков
  - ▶ Rymorphy2
- ▶ Разделение файлов с контентом доков на файлы по запросам
- ▶ Подсчёт поведенческих признаков

# Поведенческие факторы

- ▶ Усреднённые фичи по документам и доменам
  - ▶ Фичи для документов по каждому запросу
  - ▶ Фичи для документов для ближайшего запроса
- 
- ▶ CTR
  - ▶ Среднее время, проведённое на странице
  - ▶ Сколько раз был кликнут последним
  - ▶ ~~PBA~~
    - ▶ Бустинг на CTR по позициям, а так же на долях кликов (сколько раз кликнут первым, вторым, ... / число кликов)
  - ▶ Минимальное расстояние от запроса из логов до интересующего запроса
    - ▶ (расстояние Левенштейна, число общих слов)
  - ▶ Среднее число кликов до
  - ▶ Среднее число кликов после
  - ▶ Cascade model

# Классические текстовые факторы

- ▶ Tfidf (слова, словарные нграммы, буквенные нграммы)
  - ▶ Косинусное расстояние между заголовком документа и запросом
  - ▶ Косинусное расстояние между заголовком документа и заголовком ранее найденного топового документа
  - ▶ Косинусное расстояние между запросом и содержимым документа
  - ▶ Косинусное расстояние между содержимым документа и содержимым ранее найденного топового документа
  - ▶ Добавление синонимов
- ▶ BM25F

# Семантические факторы

- ▶ FastText
  - ▶ заголовок документа <-> запрос
- ▶ Universal Sentence Encoder
  - ▶ (universal-sentence-encoder-multilingual-large, universal-sentence-encoder-multilingual-qa)
- ▶ Bert
  - ▶ Slavic BERT
  - ▶ RuBERT
  - ▶ Multilingual BERT
- ▶ Word2vec

# CATBOOST

YetiRank , YetiRankPairwise  
learning\_rate: 0.3  
iterations 1500  
depth: 5

# LightGbm

YetiRank  
n\_estimators=1000  
learning\_rate=0.05  
max\_depth=60

Public score: 0.78486  
Private score: 0.78966

# Feature Importance

USE\_query\_answer 0.0014121655472196792  
content\_tfidf\_prev\_top 0.0012800984544819594  
host\_[nafter] 0.0008331748547898421  
host\_[nlast] 0.0007681707342075406  
host\_[nbefore] 0.0006412293517727985  
bm\_features@3 0.0006239362354457656  
tf\_features@0 0.0004994392099945122  
doc\_click[nlast] 0.0004630720782545783  
host\_pos\_done\_clicks@8 0.0004137773966830238  
host\_pos\_done\_clicks@14 0.00039347154205482404  
tf\_ngram\_features\_prev@2 0.00038750350570704306  
content\_tfidf@3 0.0003874106351922091  
fasttext\_norm 0.0003819986551691912  
host\_max\_pos\_show 0.0003690063306374114  
fq\_dom\_features@23 0.0003412027989799782  
doc\_log(1.0+self.nclicks) 0.0003331871723959434  
content\_tf@0 0.0003241621496963276

Спасибо за внимание!