Ранжирование
  (learning to rank)

$Y = \mathbb{R}$

$Y = \{1, \dots K\}$

$X = \{x_1, \dots x_\ell\} \subset \mathbb{X}$

$(i,j) \in R \subset \{1, \dots \ell\}^2 \Rightarrow a(x_i) > a(x_j)$

$(i,j), (j,k), (k,i)$ — может быть

$R$ — вместо целевых переменных

поисковое ранжирование:

$x = (q, d)$

в $R$ — пары $x_i = (q_i, d_i)$ и $x_j = (q_j, d_j)$, где $q_i = q_j$

## Метрики качество ранжирования

Наивный подход:

1) назначить для каждого $x \in X$ $y \in \mathbb{R}$:

   $(i,j) \in R \Rightarrow y_i > y_j$

2) обучаем $a(x)$ на $(x_i, y_i)$
   метрика — MSE

   контрпример: $x_1 > x_2 > x_3$
   
   $y:$ $\quad$ 3 $\quad$ 2 $\quad$ 1
   
   $a:$ $\quad$ 0 $\quad$ -0.05 $\quad$ -10
   
   MSE $\gg$ 0, хотя ранжирование идеальное

Более правильные метрики:
   AUC-ROC, DCG, MAP
                      mean
                      average
                      precision

                                      например:
                                   $g(y) = 2^y - 1$

$$DCG@K(q) = \sum_{i=1} g(y(i)) \cdot d(i)$$

правильный
ответ
документа
с i-й позиции

$$d(i) = \frac{1}{\log(i+1)}$$

$$DCG@K = \frac{1}{|Q|} \sum_{q \in Q} DCG@K(q)$$

nDCG@k - нормируем на DCG для
идеального ранжирования

## Каскадные метрики

ERR

(pFound)

$P_i$ - вероятность дойти до i-й
позиции в выдаче

$$P_1 = 1$$
$$P_{i+1} = P_i (1 - \underset{\uparrow}{y(i)}) (1 - P_{out})$$

вероятность того,          вер-ть того,
что польз. найдет         что пользователь
ответ в i-м документе     заебается

$$pFound@k(q) = \sum_{i=1}^{K} P_i \, y(i)$$

$$pFound@k = \frac{1}{|Q|} \sum_{q \in Q} pFound@k(q)$$

## Признаки в ранжировании

- запросные
- статические (по документам)
- динамические (по запросу и документу)

① BM25

$$BM25(q,d) = \sum_{i=1}^{n} IDF(q_i) \frac{tf(q_i,d)(k_1+1)}{tf(q_i,d) + k_1 \cdot (1-b+b\frac{1}{n}}$$

сумма по
словам запроса
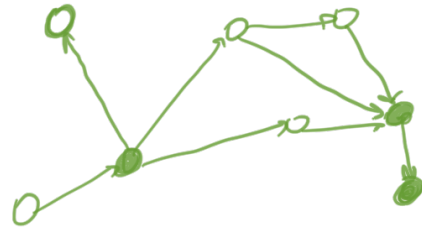
i-слово
запроса

среди
размер
докум.

$к_1, б$ - гиперпараметры

② Page Rank



$$PR(d) = \frac{1-\delta}{|D|} + \delta \sum_{c \in D_d^{in}} \frac{PR(c)}{|D_c^{out}|}$$

документы,
ссылающиеся
на d

документы,
на которые
ссылается c

СЛАУ

$$R = \frac{1-\delta}{|D|} + \delta \cdot A \cdot R$$

„матрица смежности"

$$R = (I - \delta R)^{-1} \frac{1-\delta}{|D|} \cdot \vec{1}$$

инициализируем R, пересчитываем до сходимости

---

Методы ранжирования
- point wise
- pair wise
- list wise

① Point wise

Тот самый наивный подход

$y_i$ - показатели релевантности

$$\frac{1}{l} \sum_{i=1}^{l} L(y_i, a(x_i)) \to \min_a$$

② Pairwise

$$\sum_{(i,j)\in R} [a(x_i) - a(x_j) < 0] \to \min_a \quad ⑤$$

$(i,j) \in R \Rightarrow$
$a(x_i) > a(x_j)$

число дефектных пар

$$[z < 0] \leq \widetilde{L}(z)$$

$$⑤ \leq \sum_{(i,j)\in R} \widetilde{L}(a(x_i) - a(x_j)) \to \min_a$$

если $a$ - дифференц. модель, то можно
обучать SGD, семплируя пары из $R$

1) $\begin{pmatrix} y_i = 100 \\ y_j = 0 \end{pmatrix} = \begin{pmatrix} y_i = 100 \\ y_j = 99 \end{pmatrix}$

чтобы $\overset{\text{все}}{\text{пары}}$ не были равнозначны,
можно семплировать их из $R$
с вероятностями, пропорциональными $|y_i - y_j|$

2) можно сделать модель попарной
$a(x_i, x_j)$

---

Rank Net

$a(x) = \langle w, x \rangle$

$\widetilde{L}(z) = \log(1 + \exp(-\sigma z))$ ← константа

SGD:
$$w := w + h \frac{\sigma}{\exp(\sigma \langle x_j - x_i, w \rangle)} (x_j - x_i)$$

Эмпирическое наблюдение:

$$w := w + h \frac{\sigma}{\exp(\sigma \langle x_j - x_i, w \rangle)} \cdot |\Delta F_{ij}| (x_j - x_i)$$

$|\Delta F_{ij}|$ - изменение метрики ранжирования
(например, nDCG) при обмене
местами $x_i$ и $x_j$

в итоге будет пример но оптимизироваться
метрика ранжирования

Lambda Rank