

ЗАДАЧА РАНЖИРОВАНИЯ

» Дано:

- ▶ Объекты x_1, \dots, x_ℓ
- ▶ Порядок на некоторых парах: $\{(i, j) : x_i < x_j\}$

» Найти:

- ▶ Ранжирующую модель $a(x)$, такую что

$$x_i < x_j \Rightarrow a(x_i) < a(x_j)$$

- › Объекты — пары (запрос, документ)
- › Порядок задан только внутри одного запроса:
 $(\text{запрос}, \text{документ}_1) < (\text{запрос}, \text{документ}_2)$
- › Как правило, порядок задают ассессоры

- › Объекты — пары (пользователь, товар)
- › Порядок задан только для пар одного и того же пользователя
- › Порядок определяется оценками или покупками пользователя

ОСОБЕННОСТИ ЗАДАЧИ

- › Объекты не являются независимыми — целевая переменная зависит от пар объектов
- › Сложные метрики качества
- › Много способов сформировать выборку

МЕТРИКИ КАЧЕСТВА РАНЖИРОВАНИЯ

» Дано:

- ▶ Объекты x_1, \dots, x_ℓ
- ▶ Порядок на некоторых парах: $\{(i, j) : x_i < x_j\}$

» Найти:

- ▶ Ранжирующую модель $a(x)$, такую что

$$x_i < x_j \Rightarrow a(x_i) < a(x_j)$$

- › $y(q, d) \in \{0, 1\}$ — релевантность документа d запросу q
- › $a(q, d)$ — оценка релевантности
- › $d_q^{(i)}$ — i -й документ в порядке оценки релевантности
- › Точность среди первых k документов:

$$\text{Precision@}k(q) = \frac{1}{k} \sum_{i=1}^k y(q, d_q^{(i)})$$

- › Точность среди первых k документов:

$$\text{Precision@}k(q) = \frac{1}{k} \sum_{i=1}^k y(q, d_q^{(i)})$$

- › Не учитывает позиции релевантных документов

- › Средняя точность (average precision, AP):

$$AP@k(q) = \frac{\sum_{i=1}^k y(q, d_q^{(i)}) \text{Precision}@i(q)}{\sum_{i=1}^k y(q, d_q^{(i)})}$$

- › Достигает максимума, если все релевантные документы находятся выше всех нерелевантных

› Средняя AP по запросам:

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

- › $y(q, d) \in \mathbb{R}$ — релевантность документа d запросу q
- › $a(q, d)$ — оценка релевантности
- › $d_q^{(i)}$ — i -й документ в порядке оценки релевантности
- › Discounted cumulative gain (DCG):

$$\text{DCG}@k(q) = \sum_{i=1}^k \frac{2^{y(q, d_q^{(i)})} - 1}{\log(i + 1)}$$

- › Метрику DCG принято нормировать:

$$n\text{DCG}@k(q) = \frac{\text{DCG}@k(q)}{\mathbf{max} \text{DCG}@k(q)}$$

- › **max** DCG@k(q) — DCG при идеальном ранжировании

- › Два подхода к оцениванию качества — точность и DCG
- › Точность учитывает долю релевантных документов в топе
- › Модификация: средняя точность
- › DCG учитывает релевантность документа и его позицию

МЕТОДЫ РАНЖИРОВАНИЯ

ПОДХОДЫ К РАНЖИРОВАНИЮ



- Pointwise
- Pairwise
- Listwise

- › $y(q, d) \in \mathbb{R}$ — релевантность документа d запросу q
- › $a(q, d)$ — оценка релевантности
- › Будем предсказывать $y(q, d)$ методами регрессии
- › Например:

$$\sum_{i=1}^{\ell} (a(q, d) - y(q, d))^2 \rightarrow \min$$

- Минимизируем количество дефектных пар:

$$\sum_{x_i < x_j} [a(x_j) - a(x_i) < 0] \rightarrow \mathbf{min}$$

- › Минимизируем количество дефектных пар:

$$\sum_{x_i < x_j} L(a(x_j) - a(x_i)) \rightarrow \mathbf{min}$$

- › $L(M)$ — гладкая функция
- › $L(M) = \log(1 + e^{-M})$ — метод RankNet

- › RankNet, шаг стохастического градиентного спуска для линейной модели:

$$w := w + \eta \frac{1}{1 + \exp(\langle w, x_j - x_i \rangle)} (x_j - x_i)$$

- › Как оптимизировать NDCG вместо $L(M)$?

- › Домножим стохастический градиент по паре (x_i, x_j) на изменение **NDCG** при перестановке x_i и x_j местами:

$$w := w + \eta \frac{1}{1 + \exp(\langle w, x_j - x_i \rangle)} (x_j - x_i) \cdot |\Delta \text{NDCG}_{ij}| (x_j - x_i)$$

- › Эмпирическое наблюдение: такая модификация действительно приводит к оптимизации **NDCG**

- › Три подхода: pointwise, pairwise, listwise
- › Наиболее часто используется попарный подход