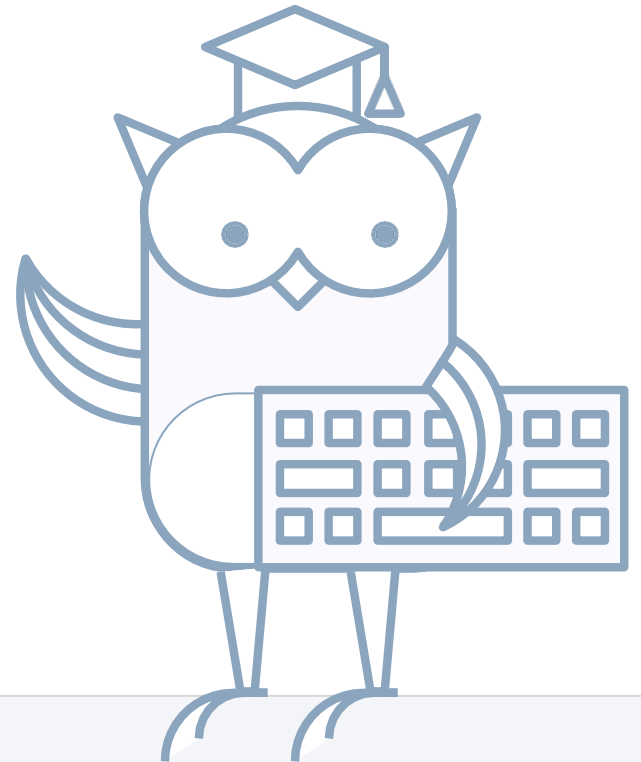


# Learning to rank



- Задача ранжирования
- Метрики
- Подходы и алгоритмы
- Существующие проблемы

Какими качествами должен обладать хороший поиск ?

- Релевантные документы
- Быстро ищет
- Не показывает бред и мусор
- Хочется пользоваться
- Легко найти «мой» сайт
- Приносит деньги

Релевантность – мера того, на сколько документ подходит запросу

Итоговое качество поиска зависит от:

- Оценки качества поиска

«Если вы не можете что-то измерить, то вы не можете это улучшить»

- Способа построения датасета
- Фичей
- Алгоритма

Множество запросов  $Q = \{q_1, q_2, \dots, q_n\}$

Множество документов соответствующих каждому запросу  $q \in Q$

$$q \rightarrow d_1, d_2, \dots$$

Для каждой пары  $(q, d)$  сопоставляется оценка релевантности  $y(q, d)$ , чем выше оценка, тем релевантнее документ  $d$  по запросу  $q$ .

Оценки релевантности сравнимы, только в рамках одного запроса:

$$(q, d_1) < (q, d_2) \Leftrightarrow y(q, d_1) < y(q, d_2)$$

- **Как оценить ранжирование?**
- Сопоставим нашим парам запрос-документы некоторые последовательности чисел:
  - *оценка релевантности*
  - *вероятность, что документ релевантен*
  - *релевантен/не релевантен*
- Оценим последовательности
- Усредним по запросам

- Начнем с простого: случай бинарных ответов  $y_{(i)} \in \{0,1\}$
- Любые стандартные метрики классификации:
  - *точность*
  - *полнота*
  - *F-мера*
  - *ROC-AUC*

- Начнем с простого: случай бинарных ответов  $y_{(i)} \in \{0,1\}$
- Любые стандартные метрики классификации:
  - *точность*
  - *полнота*
  - *F-мера*
  - *ROC-AUC*
- $precision@k(q)$  - точность для первых  $k$  предсказаний модели
- **Недостаток:** не учитывает порядок элементов



- **Average Precision** - равна сумме  $p@i$  по индексам  $i$  от 1 до  $k$  только для релевантных элементов, деленному на  $k$ .

$$AP@k(q) = \frac{1}{k} \sum_{i=1}^k y_{(i)} precision@i(q)$$

- $y_{(i)} \in \{0,1\}$  – бинарная релевантность документа на позиции  $i$
- Учитывает порядок элементов

- **Average Precision** - равна сумме  $p@i$  по индексам  $i$  от 1 до  $k$  только для релевантных элементов, деленному на  $k$ .

$$AP@k(q) = \frac{1}{k} \sum_{i=1}^k y_{(i)} precision@i(q)$$

- $y_{(i)} \in \{0,1\}$  – бинарная релевантность документа на позиции  $i$
- Учитывает порядок элементов

## Пример:

Три элемента, релевантен последний:  $ap@3(q) = \frac{1}{3} (0 + 0 + 1/3) \approx 0,11$ .

- **Average Precision** - равна сумме  $p@i$  по индексам  $i$  от 1 до  $k$  только для релевантных элементов, деленному на  $k$ .

$$AP@k(q) = \frac{1}{k} \sum_{i=1}^k y_{(i)} precision@i(q)$$

- $y_{(i)} \in \{0,1\}$  – бинарная релевантность документа на позиции  $i$
- Учитывает порядок элементов

## Пример:

Три элемента, релевантен последний:  $ap@3(q) = \frac{1}{3} (0 + 0 + 1/3) \approx 0,11$ .

Три элемента, релевантен первый:  $ap@3(q) = \frac{1}{3} (1/1 + 0 + 0) \approx 0,33$ .

- **Average Precision** - равна сумме  $p@i$  по индексам  $i$  от 1 до  $k$  только для релевантных элементов, деленному на  $k$ .

$$AP@k(q) = \frac{1}{k} \sum_{i=1}^k y_{(i)} precision@i(q)$$

- $y_{(i)} \in \{0,1\}$  – бинарная релевантность документа на позиции  $i$
- Учитывает порядок элементов

## Пример:

Три элемента, релевантен последний:  $ap@3(q) = \frac{1}{3} (0 + 0 + 1/3) \approx 0,11$ .

Три элемента, релевантен первый:  $ap@3(q) = \frac{1}{3} (1/1 + 0 + 0) \approx 0,33$ .

Три элемента, все релевантны:  $ap@3(q) = \frac{1}{3} (1/1 + 2/2 + 3/3) = 1$ .

- **Average Precision** - равна сумме  $p@i$  по индексам  $i$  от 1 до  $k$  только для релевантных элементов, деленному на  $k$ .

$$AP@k(q) = \frac{1}{k} \sum_{i=1}^k y_{(i)} precision@i(q)$$

- $y_{(i)} \in \{0,1\}$  – бинарная релевантность документа на позиции  $i$

- **Average Precision** - равна сумме  $p@i$  по индексам  $i$  от 1 до  $k$  только для релевантных элементов, деленному на  $k$ .

$$AP@k(q) = \frac{1}{k} \sum_{i=1}^k y_{(i)} precision@i(q)$$

- $y_{(i)} \in \{0,1\}$  – бинарная релевантность документа на позиции  $i$
- **Mean Average Precision** – усредняем по всем запросам.

$$MAP@k = \frac{1}{N} \sum_{i=1}^N ap@k(q)$$

- **DCG** (*discounted cumulative gain*)
- Ответы являются вещественными

$$\text{DCG}@k(q) = \sum_{i=1}^k g(y_{(i)})d(i)$$

-

- **DCG** (*discounted cumulative gain*)
- Ответы являются вещественными

$$\text{DCG}@k(q) = \sum_{i=1}^k g(y_{(i)})d(i)$$

- Популярные функции:

$$g(y) = 2^y - 1 \qquad d(i) = \frac{1}{\log(i + 1)}$$

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



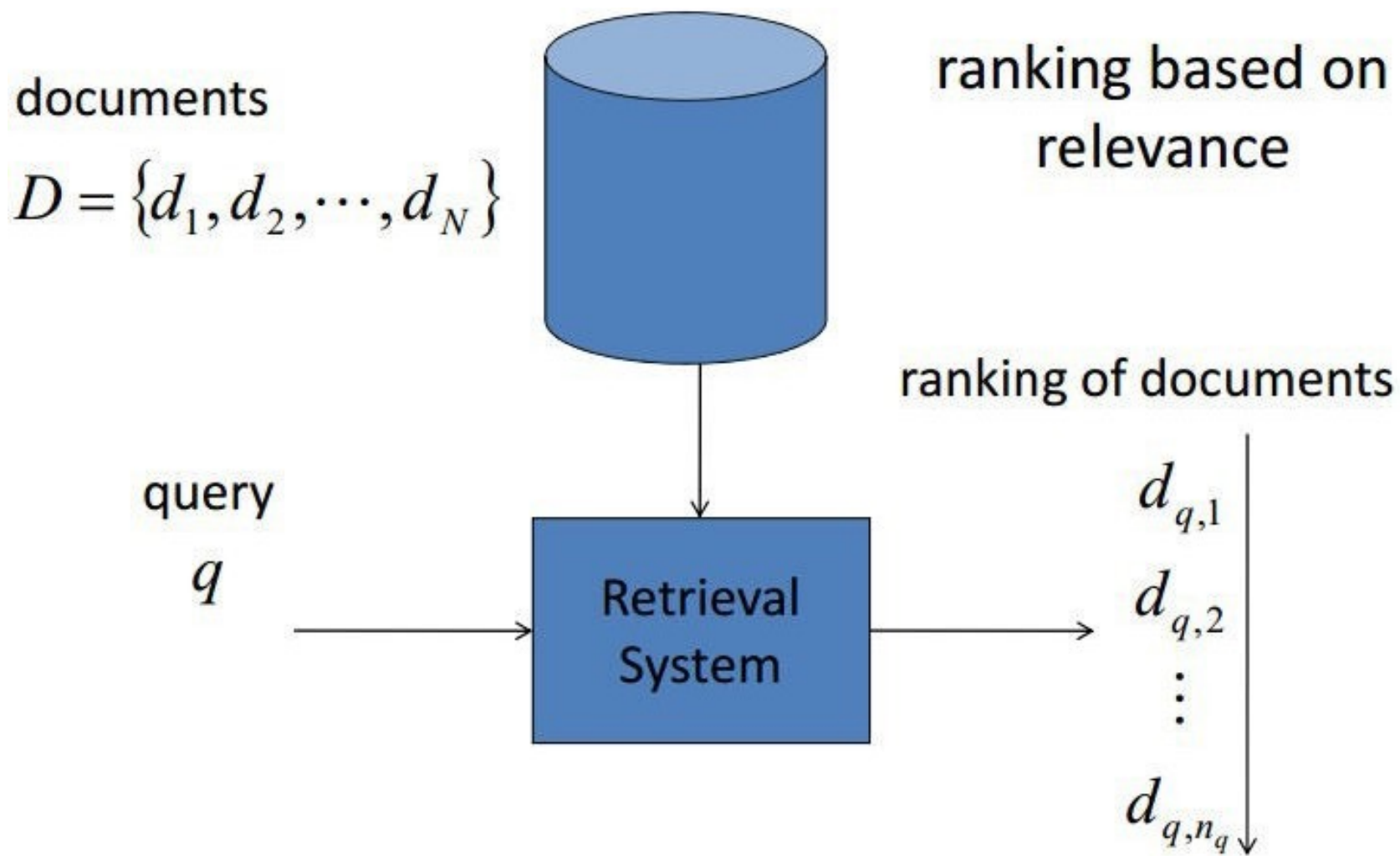
- **DCG** (*discounted cumulative gain*)
- Ответы являются вещественными

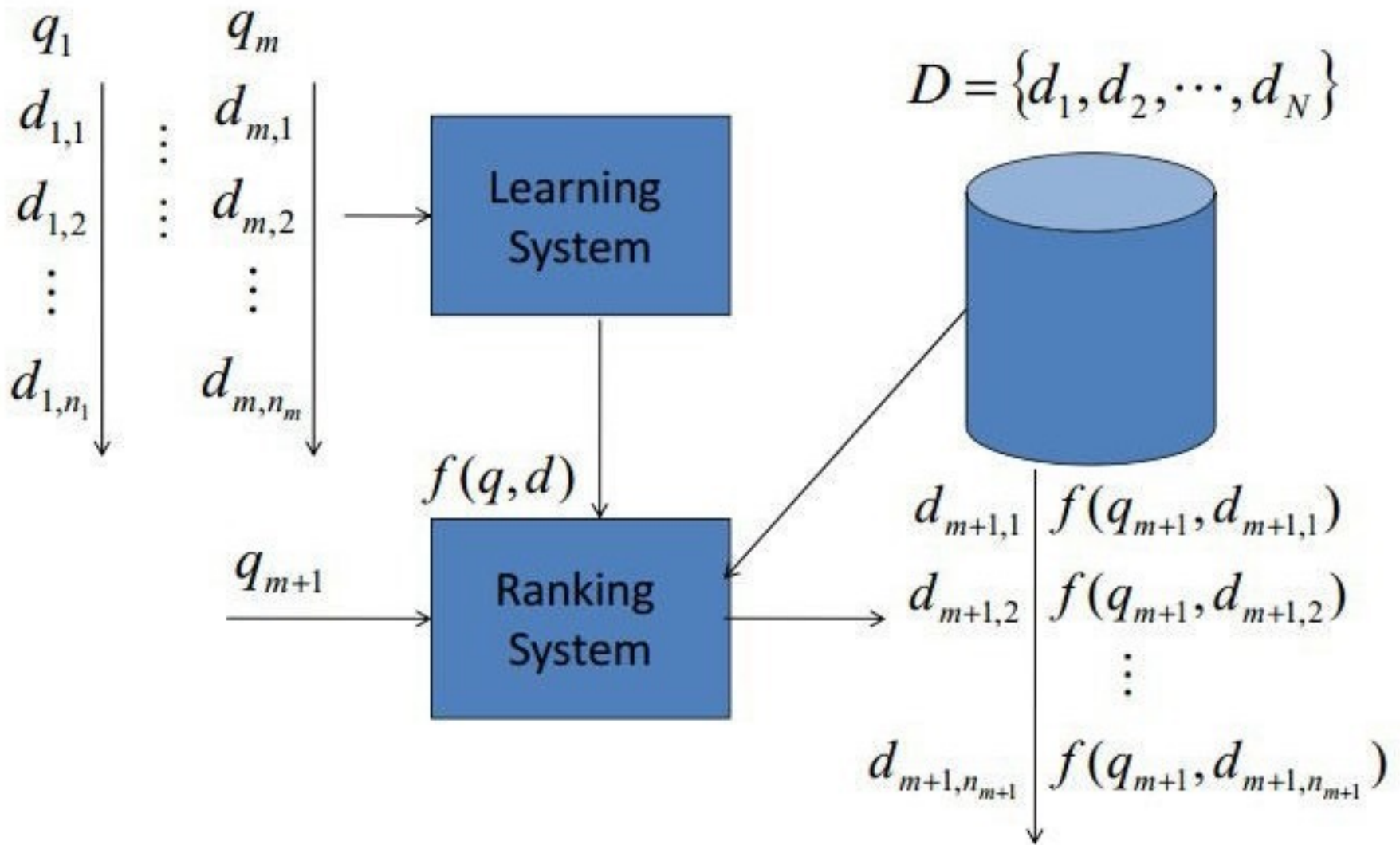
$$\text{DCG}@k(q) = \sum_{i=1}^k \frac{2^{y(i)} - 1}{\log(i + 1)}$$

- **nDCG** (*normalized discounted cumulative gain*) – нормируем на идеально ранжирование

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







- Текстовые
- Ссылочные
- Поведенческие
- Социальные
- Временные

- Запросные:
  - популярность
  - тип
  - число слов и т. п.

- Запросные:
  - популярность
  - тип
  - число слов и т. п.
- Статические (документные) – могут быть посчитаны заранее
  - популярность
  - количество ссылок
  - средний word2vec вектор и т. п.

- Запросные:
  - популярность
  - тип
  - число слов и т. п.
- Статические (документные) – могут быть посчитаны заранее
  - популярность
  - количество ссылок
  - средний word2vec вектор и т. п.
- Динамические (документ – запросные)
  - расстояние между запросом и документом
  - близость тем и т. п.



## *Discounted Cumulative Gain*

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

Как оптимизировать?

Проблемы с дифференцируемостью....

## Discounted Cumulative Gain

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

Как оптимизировать ?

### Три подхода:

- ***pointwise*** (другая целевая функция), обучение на отдельных примерах
- ***pairwise*** (другая целевая функция), обучение на документах в рамках запроса
- ***listwise*** (выучиваем распределение), обучение на отранжированных списках

Будем пытаться решить задачу ранжирования как задачу регрессии (или классификации)

$$L(h) = \sum_q \sum_{(q, d_i)} (y(q, d_i) - h(q, d_i))^2$$

- Работает приемлемо
- Отделяет простые запросы от сложных

## Недостатки

- Нет непосредственной оптимизации порядка документов
- Теряем информацию об упорядоченности
- Нет привязки к запросу
- Количество документов для запроса сильно влияет
- Нельзя влиять на конкретные результаты (YouTube по запросу)

Будем пытаться решить задачу ранжирования как задачу попарной классификации. Переходим к гладкому функционалу ранжирования.

$$\begin{aligned} L(h) &= \sum_q \sum_{(q,d_i) < (q,d_j)} [h(x_j) - h(x_i) < 0] \\ &\leq \sum_q \sum_{(q,d_i) < (q,d_j)} L(h(x_j) - h(x_i)) \rightarrow \min \end{aligned}$$

$h(x)$  – функция ранжирования

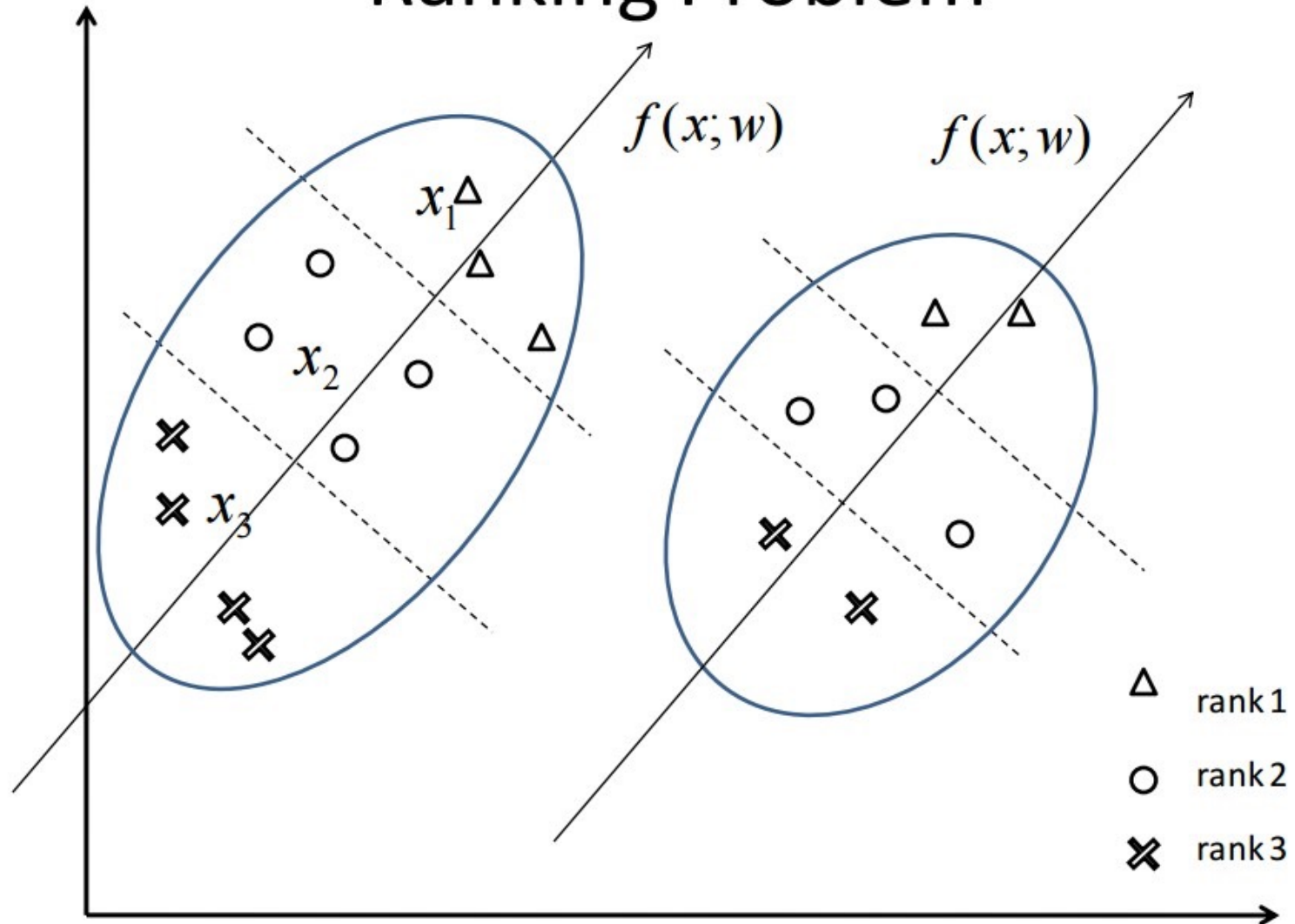
Будем пытаться решить задачу ранжирования как задачу попарной классификации. Переходим к гладкому функционалу ранжирования.

$$\begin{aligned} L(h) &= \sum_q \sum_{(q,d_i) < (q,d_j)} [h(x_j) - h(x_i) < 0] \\ &\leq \sum_q \sum_{(q,d_i) < (q,d_j)} L(h(x_j) - h(x_i)) \rightarrow \min \end{aligned}$$

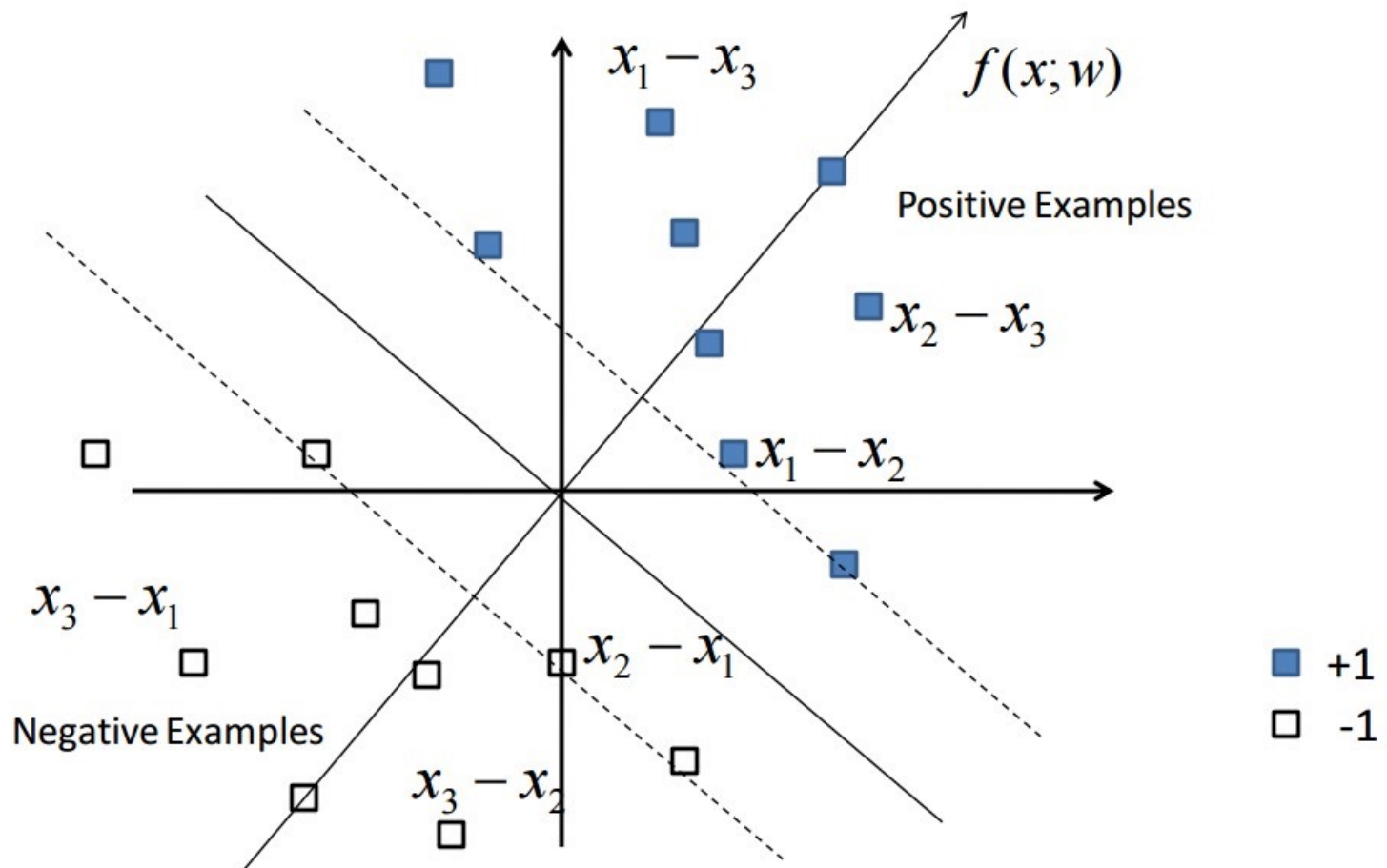
$h(x)$  – функция ранжирования

- $L(m) = (1 - m)_+$  - RankSVM
- $L(m) = \log(1 + e^{-m})$  - RankNet
- $L(m) = e^{-m}$  - RankBoost

# Ranking Problem



# Transformed Pairwise Classification Problem





## Переобучение

- запросы
- документы
- эксперты

## Положительный фидбек

- факторы
- документы

## Шумные данные

- эксперты

## Равномерная выборка из логов

- Несвежие данные
- Шумные по времени
- Скачки при смене набора запросов
- Устаревшие оценки
- Неактуальные запросы

## Равномерная выборка из логов

- Учимся на том, что показываем
- Индекс меняется, часто перестраиваем модель
- Вне топа другое распределение данных
- Могут встречаться аномальные значения

- Пользователи != эксперты
- Асессоры не создают запросы (действуют в вакууме)
- Асессоры не эксперты в некоторых областях, могут читать
- Оценивается по большой сложной инструкции

- Поведенческие факторы (может быть случайным для редкого запроса)
- SEO оптимизация
- Добавляем документы, которые есть у конкурентов, но нет у нас

- Ассессоры тоже люди
- Пытаться учесть ошибки в модели

- Оптимизировать целевую метрику качества вряд ли получится из-за дискретности
- But we will try to do our best!
- Введем некоторое вероятностное распределений
- Свели задачу к задаче восстановления распределения
- ListNet – списочный метод, учитывает порядок документов

- Вместо максимизации NDCG будем минимизировать расстояние между истинным ранжированием и ранжированием, порождаемым ранжирующей функцией
- Метки релевантности или значения ранжирующей функции на документах будут порождать вероятностное распределение на перестановках
- Максимизируем близость распределения, порождаемого значениями функции и истинными метками



- Как это работает?
- Для каждой перестановки документов  $\pi$  вводим вероятность этой перестановки при ранжировании нашим алгоритмом  $f$ :

$$P_f(\pi)$$

- Получаем два вероятностных распределения: для алгоритма ( $p$ ) и истинное ( $q$ ) для оценок ассессоров

- Как будем оптимизировать?
- Дивергенция Кульбаха-Лейблера!

$$\text{KL}(p\|q) = \int p(x) \log \frac{p(x)}{q(x)} dx.$$

- Функционал имеет смысл «расстояния между распределениями»

$$\text{KL}(p\|q) \geq 0, \quad \forall p, q;$$

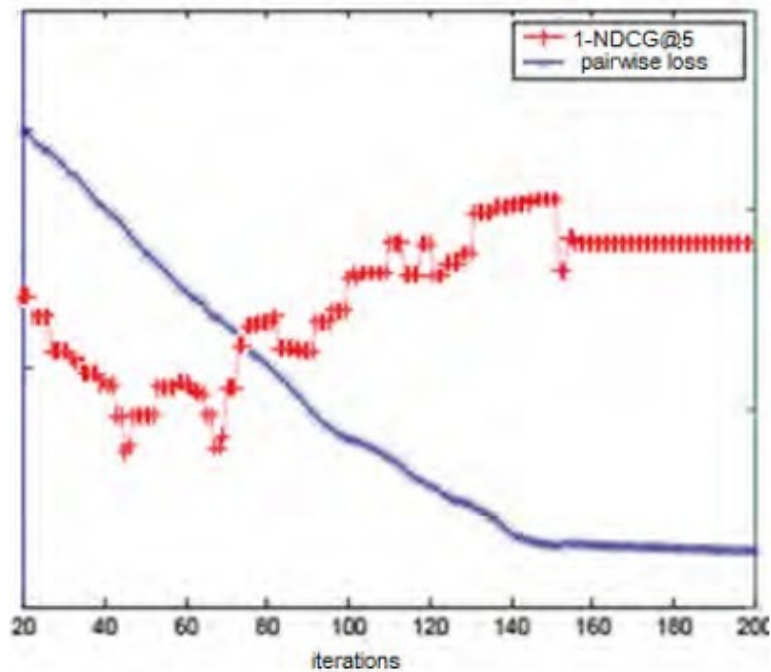
$$\text{KL}(p\|q) = 0 \iff p = q.$$

- Является ли он метрикой?

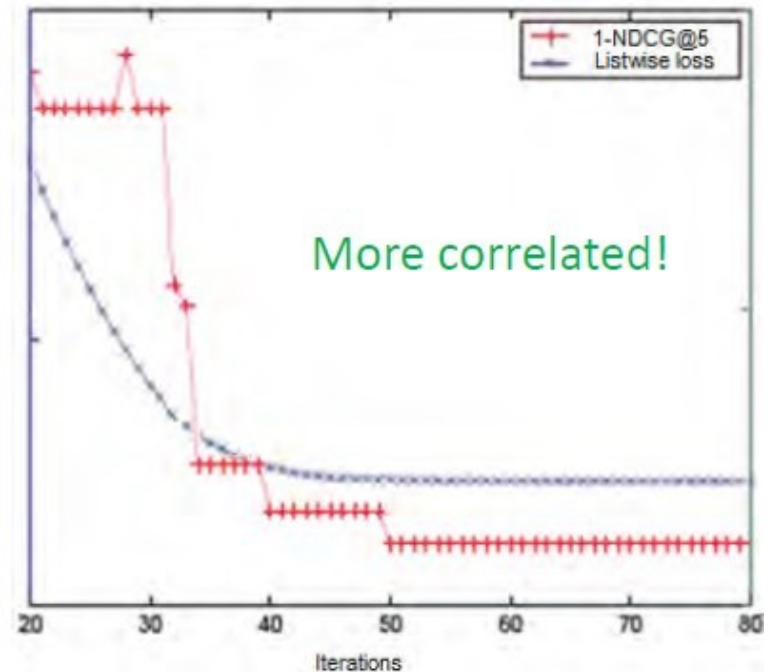
- Используем нейронки
- Обучаем градиентным спуском
- Функция потерь – KL дивергенция между распределениями

$$D_{\text{KL}}(P \parallel Q) = \int_X \log \frac{dP}{dQ} dP,$$

$$L(h) = - \sum_{q \in Q} \sum_{G(j_1, \dots, j_n)} \left( \prod_{t=1}^n \frac{e^{\text{rel}_{j_t}}}{\sum_{u=t}^n e^{\text{rel}_{j_u}}} \right) \log \left( \prod_{t=1}^n \frac{e^{h(x_{j_t})}}{\sum_{u=t}^n e^{h(x_{j_u})}} \right)$$



Pairwise (RankNet)



Listwise (ListNet)

Training Performance on TREC Dataset

- Статья на хабр про метрики ранжирования:  
<https://habr.com/ru/company/econtenta/blog/303458/>
- Конспект лекции Соколова ВШЭ: <https://github.com/esokolov/ml-course-hse/blob/master/2019-spring/lecture-notes/lecture24-ranking.pdf>
- Видео лекции Соколова: <https://youtu.be/cqI3bVrFEOQ>
- Лекция Воронцова  
<http://www.machinelearning.ru/wiki/images/8/89/Voron-ML-Ranking-slides.pdf>
- Статья про RankNet <https://medium.com/@nikhilbd/intuitive-explanation-of-learning-to-rank-and-ranknet-lambdarank-and-lambdamart-fe1e17fac418>

**Спасибо  
за внимание!**



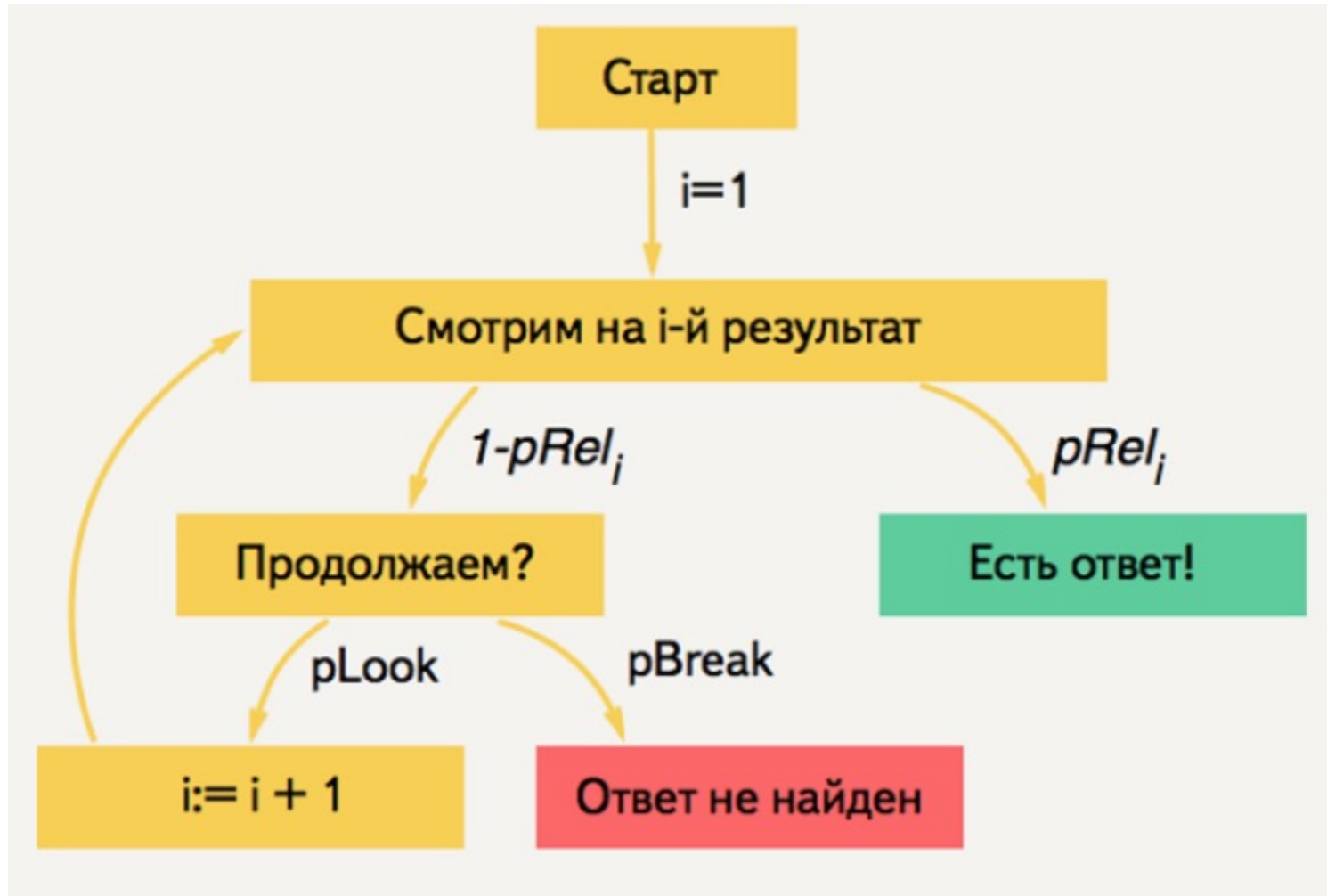
- $p_{Found}$
- Ответы на отрезке  $[0, 1]$ , выражают вероятность найти ответ в документе

$$p_{i+1} = p_i(1 - y_{(i)})(1 - p_{out})$$

- $p_1 = 1$
- $p_{out}$  - вероятность того, что пользователь уйдет, не найдя ответ

$$p_{Found@k}(q) = \sum_{i=1}^k p_i y_{(i)}$$

- Является каскадной метрикой





$$\begin{aligned}
 L(h) &= \sum_q \sum_{(q,d_i) < (q,d_j)} [h(x_j) - h(x_i) < 0] \\
 &\leq \sum_q \sum_{(q,d_i) < (q,d_j)} L(h(x_j) - h(x_i)) \rightarrow \min
 \end{aligned}$$

Берем в качестве оценки сверху

$$L(m) = (1 - m)_+$$

Отступ

$$m_{i,j} = w^T(x_j - x_i)$$

В итоге получили следующую задачу минимизации:

$$L(h) = \sum_q \sum_{(q,d_i) < (q,d_j)} (1 - m_{i,j})_+ + \frac{1}{2C} \|w\|^2 \rightarrow \min_w$$

- Probability of permutation  $\pi$  is defined as

$$P_s(\pi) = \prod_{j=1}^n \frac{\varphi(s_{\pi(j)})}{\sum_{k=j}^n \varphi(s_{\pi(k)})}$$

- Example:

$$P_f(ABC) = \frac{\varphi(f(A))}{\varphi(f(A)) + \varphi(f(B)) + \varphi(f(C))} \cdot \frac{\varphi(f(B))}{\varphi(f(B)) + \varphi(f(C))} \cdot \frac{\varphi(f(C))}{\varphi(f(C))}$$

$P(\text{A ranked No.1})$

$P(\text{B ranked No.2} \mid \text{A ranked No.1})$   
 $= P(\text{B ranked No.1}) / (1 - P(\text{A ranked No.1}))$

$P(\text{C ranked No.3} \mid \text{A ranked No.1, B ranked No.2})$

- Two queries in total
- Same error in terms of pairwise classification  $780/790 = 98.73\%$ .
- Different errors in terms of query level evaluation **99%** vs. **50%**.

		Case 1	Case 2
Document pairs of $q_1$	correctly ranked	770	780
	wrongly ranked	10	0
	Accuracy	98.72%	100%
Document pairs of $q_2$	correctly ranked	10	0
	wrongly ranked	0	10
	Accuracy	100%	0%
overall accuracy	document level	98.73%	98.73%
	query level	99.36%	50%