

Информационный поиск  
Архитектура поисковых систем  
PageRank

Лекция N 3 курса  
"Алгоритмы для Интернета"

Юрий Лифшиц

ПОМИ РАН - СПбГУ ИТМО

Осень 2006

1 / 29

Авторы алгоритмов ссылочной популярности

В ноябре 1997 при запросе собственного названия только одна из четырех ведущих поисковых систем выдавала себя в первой десятке.

Брин и Пейдж, "Анатомия поисковых систем", 1998

Sergey Brin, Larry Page, Jon Kleinberg:



2 / 29

План лекции

- 1 Модели информационного поиска
  - Булевская модель
  - Векторная модель
  - Вероятностная модель
- 2 Архитектура поисковой системы
- 3 PageRank

3 / 29

Часть I

Формально, что такое документ?

Формально, что такое запрос?

При каком условии мы считаем, что документ соответствует запросу?

4 / 29

Булевская модель

Словарь:  $T = \{t_1, \dots, t_n\}$

Документ:  $D \subset T$ , иначе говоря  $D \in \{0, 1\}^n$

Запрос:  $t_5$  OR  $t_7$  NOT  $t_{12}$

Соответствие:

Формула запроса должна быть выполнена на документе.

Недостатки модели?

5 / 29

Векторная модель

Снова коллекция документов, каждый из которых теперь является **мультимножеством** слов.

Определим матрицу  $M$  по формуле  $M_{ij} = TF_{ij} \cdot IDF_i$ , где:

- Частота термина  $TF_{ij}$  — относительная доля слова  $i$  в тексте  $j$
- Обратная встречаемость в документах  $IDF_i$  — величина, обратная количеству документов, содержащих слово  $i$

Физический смысл  $M_{ij}$  — степень соответствия слова  $i$  тексту  $j$

Запрос:  $t_3$  AND  $t_5$  (разрешаем только AND)

6 / 29

Релевантность в векторной модели

Запишем запрос в виде вектора:

$Q = "t_3 \text{ AND } t_5 - \{0, 0, 1, 0, 1, 0, \dots, 0\}$

Мерой релевантности будет **косинус** между запросом и документом:

$$R(Q, D) = \frac{Q \cdot D}{|D||Q|}$$

7 / 29

Вероятностная модель для чайников

Документ: множество слов (булевский вектор)  $D = \{d_1, \dots, d_n\}$

Запрос:  $Q_k$  — тоже, но храним как множество

Соответствие:

- Зафиксируем запрос  $Q_k$
- Пусть есть распределение вероятностей на всех текстах "быть релевантным запросу  $Q_k$ ": обозначаем  $P(R|Q_k, D)$
- Пусть есть распределение вероятностей на всех текстах "быть нерелевантным запросу  $Q_k$ ": обозначаем  $P(\bar{R}|Q_k, D)$
- Функцией соответствия будет их отношение (или логарифм этой дроби):  $\frac{P(R|Q_k, D)}{P(\bar{R}|Q_k, D)}$

8 / 29

## Вычисляем функцию соответствия

Вспользуемся теоремой Байеса ( $P(a|b) = P(b|a) \frac{P(a)}{P(b)}$ ):

$$\frac{P(R|Q_k, D)}{P(\bar{R}|Q_k, D)} = \frac{P(R|Q_k) P(D|R, Q_k)}{P(\bar{R}|Q_k) P(D|\bar{R}, Q_k)}$$

Первый сомножитель одинаков для всех документов. Предполагая независимость всех слов, второй сомножитель можно представить как произведение:

$$\prod_{i=1}^n \frac{P(x_i = d_i | R, Q_k)}{P(x_i = d_i | \bar{R}, Q_k)}$$

9 / 29

## Вычисляем функцию соответствия II

$$\prod_{i=1}^n \frac{P(x_i = d_i | R, Q_k)}{P(x_i = d_i | \bar{R}, Q_k)}$$

Введем обозначения:  $p_{ik} = P(x_i = 1 | R, Q_k)$  и  $q_{ik} = P(x_i = 1 | \bar{R}, Q_k)$ . Предположим, что для каждого слова  $i$ , не входящего в запрос,

$$p_{ik} = q_{ik}$$

Теперь мы можем переписать нашу дробь:

$$\prod_{i \in Q_k \cap D} \frac{p_{ik}(1 - q_{ik})}{q_{ik}(1 - p_{ik})} \prod_{i \in \bar{Q}_k} \frac{1 - p_{ik}}{1 - q_{ik}}$$

10 / 29

## Вычисляем функцию соответствия III

$$\prod_{i \in Q_k \cap D} \frac{p_{ik}(1 - q_{ik})}{q_{ik}(1 - p_{ik})} \prod_{i \in \bar{Q}_k} \frac{1 - p_{ik}}{1 - q_{ik}}$$

Второй сомножитель одинаков для всех документов. Забудем про него и возьмем логарифм от первого:

$$\sum_{i \in Q_k \cap D} c_{ik}, \quad \text{где } c_{ik} = \log \frac{p_{ik}(1 - q_{ik})}{q_{ik}(1 - p_{ik})}$$

11 / 29

## Подбор параметров

$$\sum_{i \in Q_k \cap D} c_{ik}, \quad \text{где } c_{ik} = \log \frac{p_{ik}(1 - q_{ik})}{q_{ik}(1 - p_{ik})}$$

Для использования полученной формулы нужно знать  $p_{ik}$  и  $q_{ik}$ .

**Рецепт:** пусть у нас уже есть некий набор текстов, про которые мы знаем, релевантны ли запросу  $Q_k$  или нет. Тогда мы можем использовать формулы:

$$p_{ik} = \frac{r_i}{r} \quad \text{и} \quad q_{ik} = \frac{f_i - r_i}{f - r},$$

Угадываете смысл обозначений?

12 / 29

## Подбор параметров II

$$p_{ik} = \frac{r_i}{r} \quad \text{и} \quad q_{ik} = \frac{f_i - r_i}{f - r},$$

Тут  $f$  — общее число документов,  $r$  — число релевантных документов,  $r_i$  число релевантных документов, содержащих слово  $i$ , а  $f_i$  — общее число документов со словом  $i$ .

13 / 29

## Часть II

В каком формате запоминать интернет-страницы?

В какой структуре данных их хранить?

Как обрабатывать запрос пользователя?

14 / 29

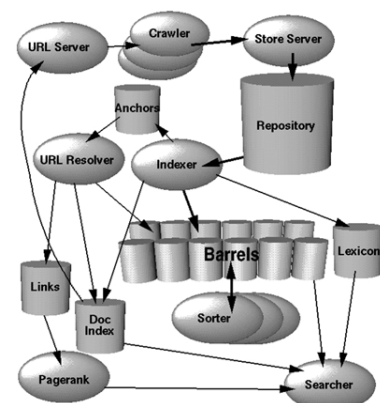
## Анатомия поисковой системы

Любая поисковая система содержит три базовые части:

- Робот (он же краулер, спайдер или индекатор)
- Базы данных
- Клиент (обработка запросов)

15 / 29

## Схема из [Brin,Page, 1998]



16 / 29

## Прямой и обратный индекс

### Прямой индекс — записи отсортированы по документам

- Номер документа
- Отсортированный список слов
- Для каждого слова: первые несколько вхождений, частота вхождений, формат вхождений

### Обратный индекс — записи отсортированы по словам

- Номер слова
- Отсортированный список документов
- Для каждого документа: информация о вхождении

17 / 29

## Релевантность

Характеристики, влияющие на позицию в списке ответов?

- Наличие слов на сайте
- Частота слов
- Форматирование
- Близость слов друг к другу
- Количество ссылок с других страниц на данную
- Качество ссылок
- Соответствие тематик сайта и запроса
- Регистрация в каталоге, связанном с поисковой системой

18 / 29

## Как работает клиент?

- Разбирает запрос на слова
- Переводит слова в их идентификаторы
- Для каждого слова находит в обратном индексе список документов, его содержащих
- Одновременно бежит по этим спискам, ища общий документ
- Для каждого найденного документа вычисляет степень релевантности
- Сортирует образовавшийся список по релевантности

19 / 29

## Качество поиска

Как оценить качество поиска?

- **Полнота:** отношение количества найденных релевантных документов к общему количеству релевантных документов
- **Точность:** доля релевантных документов в общем количестве найденных документов
- **Benchmarks:** показатели системы на контрольных запросах и специальных коллекциях документов
- **Оценка экспертов**

Не пропустите, 23 ноября — приглашенная лекция Игоря Некрестьянова "Оценка качества интернет-поиска"

20 / 29

## Часть III

Как определить ссылочную популярность страницы (PageRank)?

Как быстро вычислить приближение PageRank?



21 / 29

## PageRank: постановка задачи

Хотим для каждой страницы сосчитать показатель ее "качества".

**Идея [Брин, 1998]:** Определить рейтинг страницы через количество ведущих на нее ссылок и рейтинг ссылающихся страниц

### Другие методы:

- Учет частоты обновляемости страницы
- Учет посещаемости
- Учет регистрации в каталоге-спутнике поисковой системы

22 / 29

## Модель случайного блуждания

### Сеть:

- Вершины
- Ориентированные ребра (ссылки)

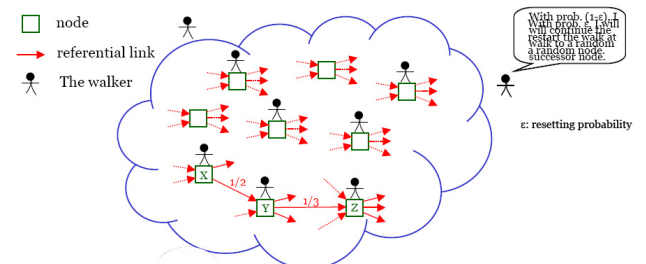
### Передвижение пользователей по сети

- Стартуем в случайной вершине
- С вероятностью  $\epsilon$  переходим в случайную вершину
- С вероятностью  $1 - \epsilon$  переходим по случайному исходящему ребру

### Предельные вероятности

- Для каждого  $k$  можно определить  $PR_k(i)$  как вероятность оказаться в вершине  $i$  через  $k$  шагов
- Факт:  $\lim_{k \rightarrow \infty} PR_k(i) = PR(i)$ , то есть для каждой вершины есть предельная вероятность находится именно в ней

23 / 29



24 / 29

## Основное уравнение PageRank

Пусть  $T_1, \dots, T_n$  — вершины, из которых идут ребра в  $i$ ,  $C(X)$  — обозначение для исходящей степени вершины  $X$ .

Утверждение:  $PR(i) = \varepsilon/N + (1 - \varepsilon) \sum_{i=1}^n \frac{PR(T_i)}{C(T_i)}$

Кто может доказать?

По определению  $PR_k(i)$  верно следующее:

$$PR_0(i) = 1/N$$

$$PR_k(i) = \varepsilon/N + (1 - \varepsilon) \sum_{i=1}^n \frac{PR_{k-1}(T_i)}{C(T_i)}$$

Нужно просто перейти к пределу!

**Практическое решение:** вместо  $PR(i)$  используют  $PR_{50}(i)$ , вычисленное по итеративной формуле.

25 / 29

## Задачи


Докажите, что по расстояние между векторами  $PR_k(i)$ ,  $PR(i)$  экспоненциально быстро (по  $k$ ) стремится к нулю


27 / 29

## Источники

Страница курса <http://logic.pdmi.ras.ru/~yura/internet.html>

Использованные материалы:

 [Sergey Brin and Larry page  
The Anatomy of Search Engine  
http://www-db.stanford.edu/pub/papers/google.pdf](http://www-db.stanford.edu/pub/papers/google.pdf)

 [Илья Сегалович  
Как работают поисковые системы  
http://company.yandex.ru/articles/article10.html](http://company.yandex.ru/articles/article10.html)

 [Langville and Meyer  
Deeper Inside PageRank  
http://meyer.math.ncsu.edu/Meyer/PS\\_Files/DeeperInsidePR.pdf](http://meyer.math.ncsu.edu/Meyer/PS_Files/DeeperInsidePR.pdf)

 [Norbert Fuhr  
Probabilistic Models in Information Retrieval  
http://www.is.informatik.uni-duisburg.de/bib/fulltext/ir/Fuhr:92.pdf](http://www.is.informatik.uni-duisburg.de/bib/fulltext/ir/Fuhr:92.pdf)

29 / 29

## PageRank как собственный вектор

Определим матрицу  $L$ :

Если нет ребра из  $i$  в  $j$ , то  $l_{ij} := \varepsilon/N$

Если ребро есть, то  $l_{ij} := \varepsilon/N + (1 - \varepsilon) \cdot \frac{1}{C(j)}$

**Введем обозначения:**

$$\overline{PR}_k = (PR_k(1), \dots, PR_k(N))$$

$$\overline{PR} = (PR(1), \dots, PR(N))$$

**Получаются соотношения:**

$$PR_k = L^k PR_0$$

$$PR = L PR$$



26 / 29

## Главные моменты

**Сегодня мы узнали:**

- Модели информационного поиска: булевская, векторная, вероятностная
- Поисковая система (1) скачивает и анализирует интернет-страницы, (2) записывает в базу и сортирует ее, (3) обрабатывает запросы, выводя лучшие страницы по функции релевантности
- PageRank — это предельная вероятность оказаться на web-сайте в результате случайного блуждания по ссылкам.

Вопросы?

28 / 29