



# **Ранжирование: от строчки кода до Матрикснета**

**Федор Романенко**

Менеджер отдела качества поиска

Студень, Москва, 1 апреля 2010

# Ранжирование – что это?

- основной алгоритм в поиске
- с помощью **факторов** вычисляет **релевантность** документа в виде числа
- выбирает топ-10
- сложная вещь с простым результатом
- определяет долю на рынке
- самый большой секрет поисковых компаний

# Язык запросов

(Yahoo – 1994 \*)

**Бинарное ранжирование: 0 или 1?**

- формируем правило отбора релевантных документов, используя «язык запросов»
- просматриваем все найденные документы

— *подходит только для специалистов*

— *нельзя использовать, когда найдено много*

\* здесь и далее: первый пример удачного применения в интернет поиске

Я

# Текстовое ранжирование: $tf * idf$

(Altavista – 1995)

**«Близость» текста запроса к тексту документа учитывает:**

- количество слов запроса в документе (term freq.)
- обратную частоту слова в языке (inverted document freq.)
- длину документа и запроса

**Текстовая релевантность также учитывает:**

- слова в заголовках
- близость слов запроса
- совпадение словоформ

# Лингвистика – понимание языка

(Яндекс – 1997)

- морфология: *ребенок шёл = дети идут*
- опечатки: *аднакласники = одноклассники*
- расширения: *МГУ = Московский Государственный Университет*

## Два пути компьютерной лингвистики:

- словарный: морфология, синтаксис
- статистический: языковые модели

# Рейтинги сайтов

(Рамблер Top100 – 1997, Яндекс тИЦ – 1999)

## Использование внешней информации о странице:

- ссылки с тематических сайтов
- переходы из каталогов
- посещаемость по счетчикам

— *тИЦ Яндекса оброс мифами вебмастеров*

# PageRank

(Google - 1998)

**PageRank – глобальная ссылочная авторитетность:**

- на страницу много ссылок
- на страницу есть ссылки с авторитетных страниц

**Модель случайного «блуждателя»:**

1. выбираем случайную страницу
2. с вероятностью  $0.85$  переходим по выходящей ссылке
3. с вероятностью  $0.15$  устаем и переходим к п.1

# Ссылочное ранжирование

(Google - 1998)

**Текст ссылки описывает страницу, на которую ссылается**

- голосование весами релевантных ссылок
- уточнение тематики страницы с помощью сопоставления текстов ссылок

# Метрика качества поиска

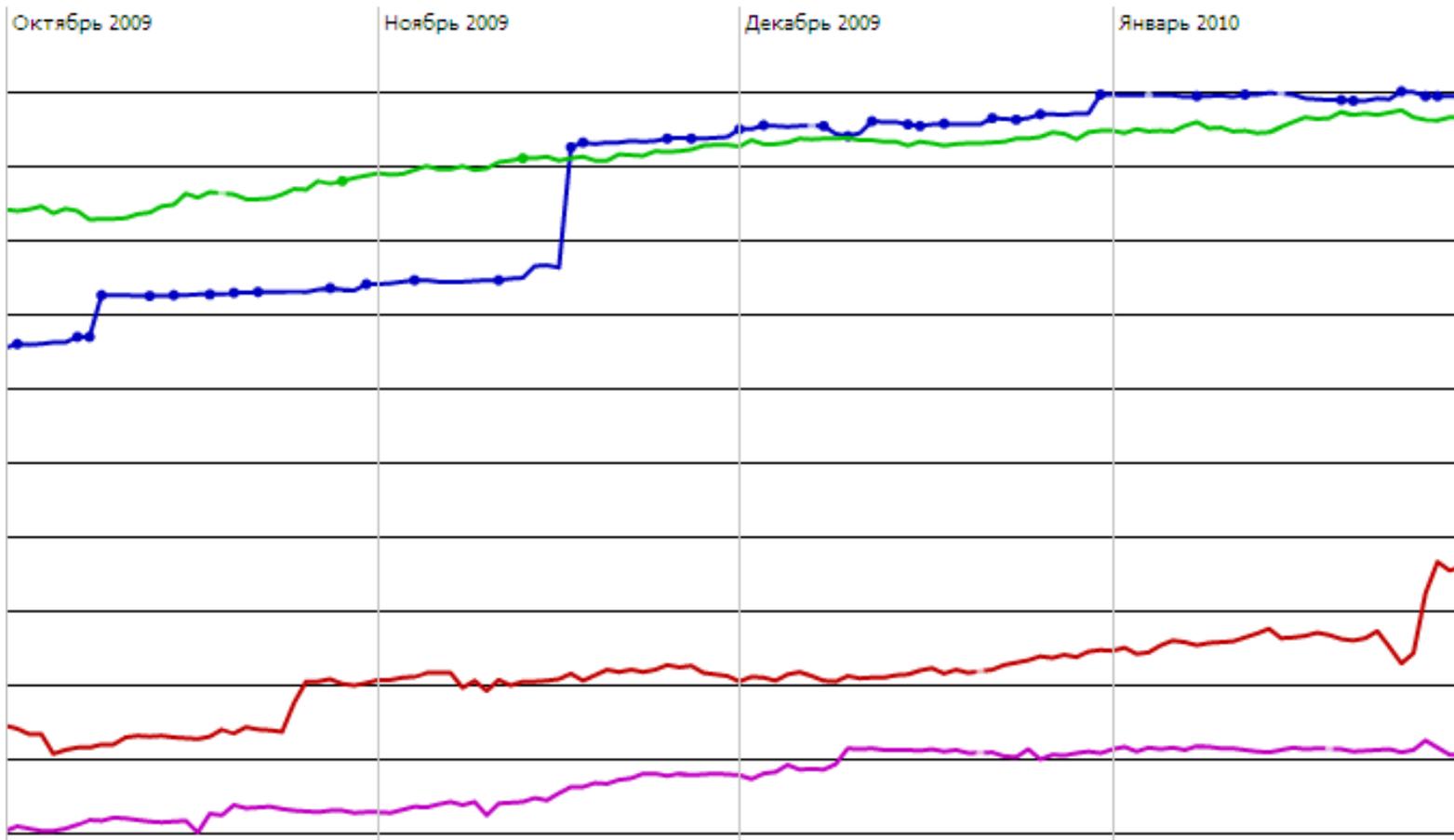
Как измерить удовлетворение пользователей поиском?

- из потока запросов пользователей строим репрезентативную **выборку**
- **ассессоры** (люди) оценивают отдельные результаты поиска
- **оценки** используются в числовой **метрике** качества ранжирования

**Я**ндекс: оценено 4 млн. документов по 100 тыс. запросов



# Метрика качества поиска



# Машинное обучение

(Яндекс, Yahoo, Bing – 200х)

Как учесть в ранжировании **400** параметров документов?

- **модель** (формула на **факторах**) умеет считать релевантность документа
- **машинное обучение** автоматически настраивает модель, максимизируя метрику качества на примерах
- машина учится ранжировать и показывает поведение, не заложенное в нее явно

# Проблемы машинного обучения

**Метрика** – *научить именно тому, что понадобится в бою*  
нужно математически измерить «счастье» пользователей

**Входные данные** – *органы чувств*  
позволяют отличить документы по качеству

**Способность к обобщению** – *аналитические способности*  
учимся на примерах, ранжируем то, что никогда не видели

**Переобучение** – *паранойя*  
мало оценок и умная модель – ложные закономерности

**Кроссвалидация** – *не дадим выучить контрольную*  
нужна не память, а умение

Я

# Матрикснет

(Яндекс – 2009)

Разработка **Яндекса** (релиз «Снежинск», ноябрь 2009)

- автоматически выбирает связанные факторы и диапазоны их значений
  - генерирует тысячи комбинированных факторов
  - очень сложная модель без склонности к переобучению
- *извлекает неиспользованный сигнал из старых данных*
- *позволяет добавлять много данных в ранжирование*

# Формула Матрикснет

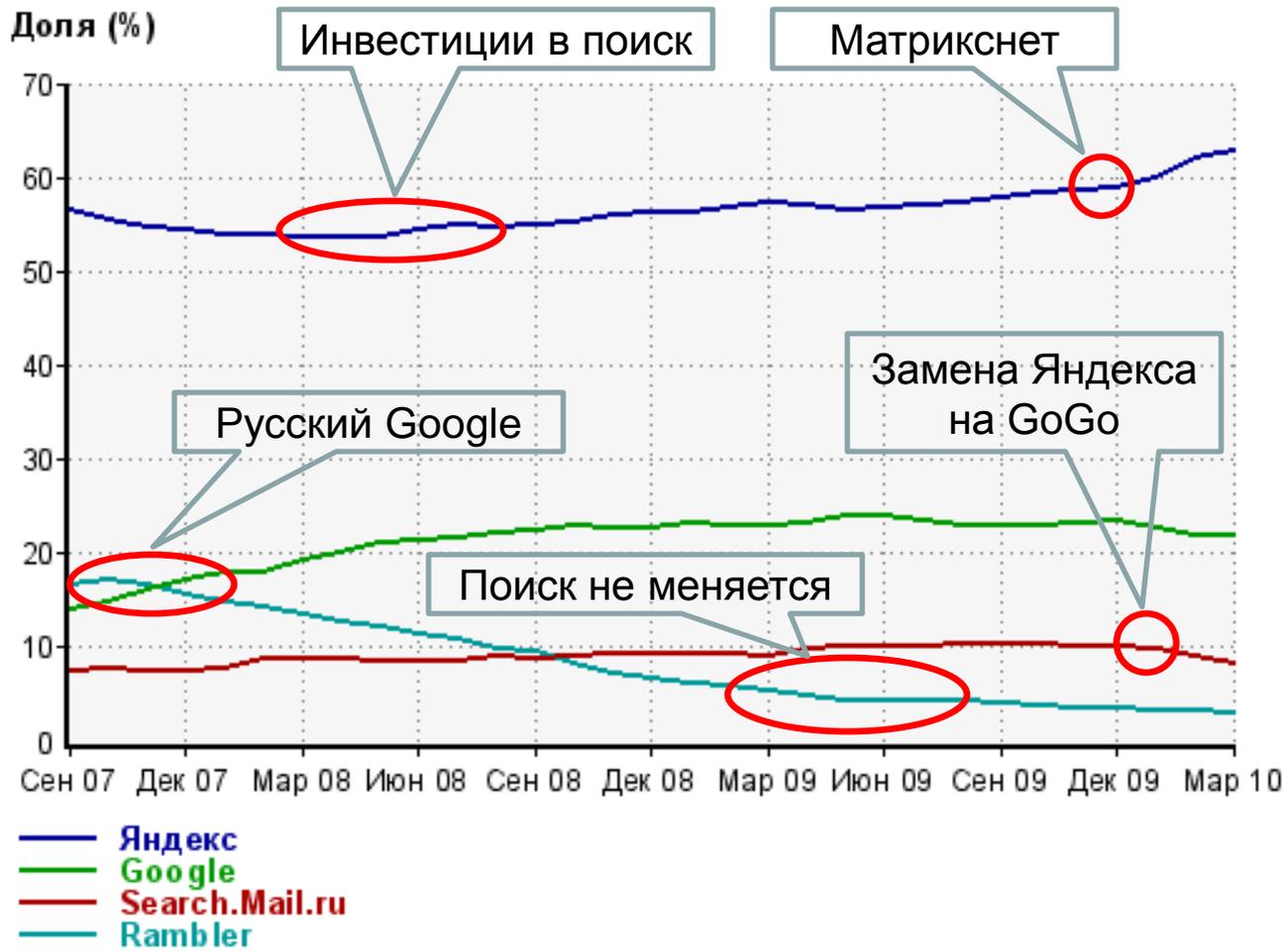
## Одна из формул

(вычисляет релевантность по факторам найденного документа)

```
...  
  
-4025627,483990,-36437960,39979596,  
-92842056,-50086892,-100233272,243162224,  
-22282850,57163664,-24991620,-9889194,  
  
...  
  
vars[5699] = fFactorInt[376] > 1060970280 ? 1 : 0; // 0.738757  
vars[5700] = fFactorInt[376] > 1061923687 ? 1 : 0; // 0.795584  
vars[5701] = fFactorInt[376] > 1049721454 ? 1 : 0; // 0.284137  
vars[5702] = fFactorInt[376] > 948291011 ? 1 : 0; // 6.37877e-05  
  
...  
  
, {376, .0f}  
, {376, 0.6251018047f}  
, {-1, .0}  
, {376, 0.05682743713f}  
, {376, 0.4546194971f}  
  
...
```

- 20 тыс. строк на C++
- объем 3МБ
- 10 тыс. коэффициентов

# Liveinternet: доля рынка, Россия



# Независимые метрики

У поиска много «качеств»!

Навигационный поиск	Тематический поиск	Подсказки	Опечатки	Цитатный поиск	Оригиналы	Синонимы
Я 98.0 ↑3.1	Я 31.6 ↑1.1	@ 96.1 ↑0.6	b 81.8 ↑1.1	Я 66.0 ↑1.0	Я 55.0	Я 63.8 ↑0.4
Я 97.0 ↑1.1	Я 29.3 ↑3.6	Я 95.5 ↓0.3	Я 60.8 ↓0.1	Я 60.0	Я 42.0	Я 62.4 ↓0.2
Я 94.9 ↑2.0	Я 27.2 ↑2.1	Я 93.4 ↓0.9	@ 55.2 ↑0.3	Я 42.0 ↑2.0	Я 42.0	Я 48.3 ↓0.3
Я 94.9	Я 18.1 ↓1.6	Я 91.4 ↑4.0	Я 46.1 ↓3.1	Я 36.0 ↑1.0	Я 38.0 ↓2.0	@ 47.7 ↓0.4
Я 91.9 ↓2.0	@ 17.5	Я 91.0 ↓2.0	Я 22.5 ↑0.3	Я 35.0 ↓1.0	Я 29.0	Я 47.3 ↓0.1
Я 90.9 ↑3.1	Я 17.5 ↓0.6	Я 85.5 ↑0.7	Я 4.4 ↓1.2	@ 35.0 ↓1.0	Я 26.0 ↑2.0	Я 47.3 ↓0.1
@ 89.9 ↓4.0	Я 15.6 ↑0.9	Я 77.8 ↑0.6	Я 1.4 ↓0.4	Я 33.0	@ 26.0 ↑1.0	Я 37.4 ↑0.8
Я 78.8 ↓3.9	Я 2.7 ↓0.3	Я N/A	Я 0.3	Я 23.0 ↑4.0	Я 8.0 ↓4.0	Я 36.0 ↓0.4

Поисковый спам	SEO-прессинг	Порнография	Полнота индекса	Апдейты	Переходы
Я 0.8 ↑0.3	Я 27.7 ↑0.9	Я 6.3 ↑0.1	Я 75.0 ↑0.8	Я 0.2 ↑0.1	Я 63.2 ↑0.8
Я 3.3 ↑0.4	Я 34.2 ↑0.4	Я 9.6 ↓0.8	Я 64.5 ↑1.7	Я 4.3 ↓0.6	Я 22.2 ↓1.8
Я 4.4 ↓0.2	@ 34.7 ↑0.2	Я 14.1 ↑0.2	Я 59.2 ↓5.3	Я 4.6 ↑1.8	@ 8.9 ↑0.4
Я 4.6	Я 40.8 ↓0.7	Я 15.1	Я 48.5 ↑1.5	Я 4.9 ↑0.3	Я 3.2 ↑0.4
@ 5.0 ↓0.4	Я 42.4 ↑0.1	Я 15.8 ↓1.0	@ 42.7	Я 8.1 ↓2.8	Я 1.0 ↑0.2
Я 5.0 ↓0.1	Я 43.3 ↓0.5	Я 16.8 ↑0.2	Я 15.9 ↑3.3	Я 15.5 ↓0.8	Я 0.1
Я 6.3 ↓1.3	Я 45.5 ↓0.2	Я 20.3 ↓0.7	Я 7.8 ↓0.5	@ 15.9 ↓1.6	Я 0.0
Я 6.5 ↓0.3	Я 48.7 ↓0.8	@ 20.4 ↓0.8	Я 5.2 ↓0.1	Я 36.0 ↓35.4	Я 0.0

<http://www.analyzethis.ru>



# Поиск - вечная задача

- поиск нельзя «доделать»: тот, кто думает, что все сделал, проигрывает
- размер интернета и количество пользователей растут, люди решают в интернете новые задачи
- много направлений: локальность, глобальность, свежесть, разнообразие, скорость, представление, ...
- сильные математики и программисты имеют дело с уникальными задачами на огромных объемах данных
- результаты их работы сразу видят миллионы людей



Федор Романенко

Я