



# Обзор доклада

- 1 Проблемы контролируемого обучения
- 2 Новое ранжирование Рамблера
- 3 Метрики качества ранжирования
- 4 Что влияет на долю рынка?



# Обзор доклада

- 1 Проблемы контролируемого обучения
- 2 Новое ранжирование Рамблера
- 3 Метрики качества ранжирования
- 4 Что влияет на долю рынка?



# Проблемы контролируемого обучения

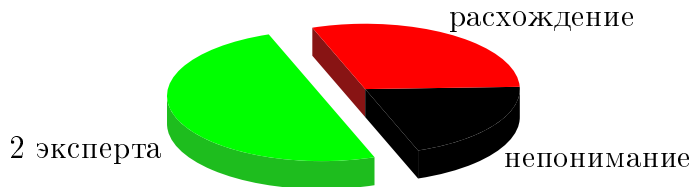
Проблемы **любого** контролируемого обучения

- Дефицит экспертов
- Шумящие факторы

## Проблемы контролируемого обучения

Эксперты: непонимание запросов, расхождение мнений.

Ни один человек не может понимать всего.



# Проблемы контролируемого обучения - Дефицит экспертов

Дефицит экспертов ограничивает качество.

- Возьмем специалиста по взрослым сайтам и обучим на нем MatrixNet
- Возьмем врача и протестируем поиск по лекарствам на обученном алгоритме.

Качество для врача будет низким, так как для лекарств важны другие факторы.

# Проблемы контролируемого обучения

Проблема расхождения мнений между экспертами.

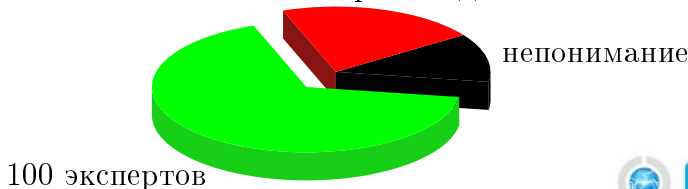
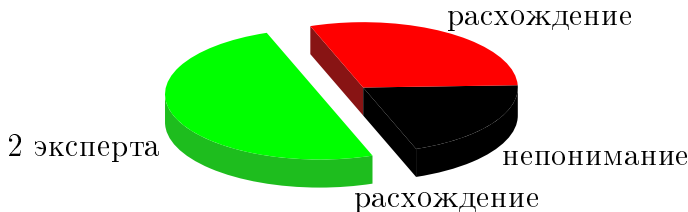
	Эксперт 1	Эксперт 2	...	Эксперт 10
пластиковые окна	сайт1	сайт2	...	сайт10
натяжные потолки	сайт1	сайт2	...	сайт10
межкомнатные двери	сайт1	сайт2	...	сайт10

*Таблица: Каждый эксперт любит свои сайты*

MatrixNet не сможет обучиться на противоречивой информации.  
Расхождение 100%, качество близко к 0.

# Проблемы контролируемого обучения - Дефицит экспертов

Что делать? Увеличивать число и **разнообразие** экспертов.





# Проблемы контролируемого обучения - Дефицит экспертов

Увеличивать число экспертов до каких пор?



транстелеком

По запросу найдено 22 тыс. сайтов, 1 млн. документов

## 1 [Компания ТТК :: Телекоммуникационная компания](#)

ТТК English english Карта сайта Контакты о компании пресс-центр сеть ТТК региональные компании справочная информация карьера операторам связи предприятиям и корпорациям малому и среднему бизнесу частным лицам

[www.ttk.ru/www/nsf/site.nsf](http://www.ttk.ru/www/nsf/site.nsf) - [сохраненный текст](#) - [искать на сайте](#)

### [Компания ТТК :: Контакты](#)

ЗАО "Самара-ТрансТелеКом" 443013, Самара, Московское шоссе, 4а тел.: (846) 973-50-50 факс:(846) 973-50-49 e-mail: [info@samara.ttk.ru](mailto:info@samara.ttk.ru) web: [www.samara-ttk.ru](http://www.samara-ttk.ru).

[www.ttk.ru/www/nsf/site.nsf/docs/contacts.html](http://www.ttk.ru/www/nsf/site.nsf/docs/contacts.html) - [сохраненный текст](#)

## 2 [ОАО "Транстелеком" \ О компании](#)

**Транстелеком**, в настоящее время, действует на территории Москвы и предоставляет клиентам услуги через собственную волоконно-оптическую транспортную сеть...

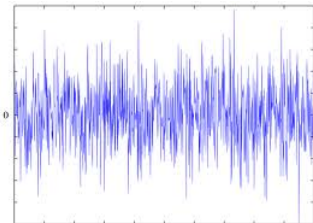
[www.transtelecom.ru/about/](http://www.transtelecom.ru/about/) - [сохраненный текст](#) - [искать на сайте](#)

# Проблемы контролируемого обучения

## Факторы

- Сильные - определяют основное качество ранжирования.
- Слабые - все вместе могут создать силу.
- Шумящие - вредны.

При добавлении факторов качество сначала растет, потом падает из-за “шумящих” факторов.



# Проблемы контролируемого обучения

## Слабые или шумящие?

Проблема в том, что мы не можем точно разделить “шумящие” факторы от “слабых”. Алгоритм обучения с учителем не может справиться с “шумящими” факторами, Сумма шумящих не создает сильный.

P.Long R.Servedio, Random Classification Noise Defeats All Convex Potential Boosters (Google Research, 2010)

# Проблемы контролируемого обучения

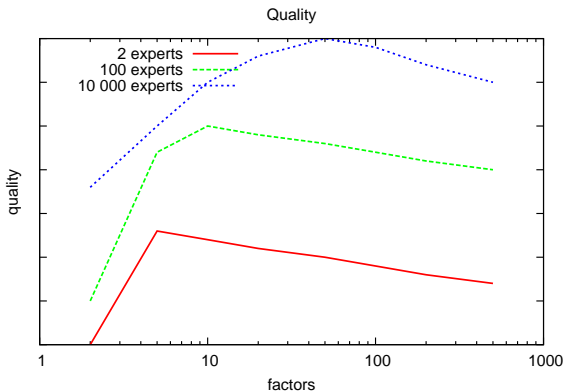
## Надо больше экспертов

Для более точного выявления шумящих факторов надо увеличивать число экспертов.

Чем больше экспертов, тем больше слабых факторов мы можем добавить.

# Проблемы контролируемого обучения

Больше экспертов → Больше понимания → Лучше качество



Качество надо измерять на отдельной группе экспертов.

# Проблемы контролируемого обучения



Обучение экспертами.

- Эксперты смотрят на результаты друг-друга, пользуются поиском как эталоном.
- Быстрое устаревание актуальности оценок.
- **Очень высокое значение цена/качество**

Проблемы приводят к частичному отказу от контролируемого обучения.

# Контролируемое обучение

Где можно использовать контролируемое обучение?

- Фильтрация документов (порноранк)
- **Реинжениринг чужого ранжирования**



# Реинжениринг чужого ранжирования

## Реинжениринг чужого ранжирования

- (+) Легко (дешево) распарсить/разметить 10 млн запросов и обучить сотни факторов
- (-) Не можем подняться выше оригинала по качеству, если у нас нет нужных факторов.
- (-) Возникает **похожесть** результатов



# Похожесть результатов между поисковиками

Похожесть результатов.

	yandex	google	mail	rambler	bing
yandex		46	38	39	37
google	46		40	41	39
mail	38	40		34	32
rambler	39	41	34		34
bing	37	39	32	34	

Таблица: Доля совпадений по первому месту, %, top10k

# Обзор доклада

- 1 Проблемы контролируемого обучения
- 2 Новое ранжирование Рамблера
- 3 Метрики качества ранжирования
- 4 Что влияет на долю рынка?



# Новое ранжирование Рамблера

Новое ранжирование Рамблера - совокупность идей

- **Неконтролируемое** обучение без экспертов
- Вместо мнения экспертов - поведение пользователей
- Стараемся быть уникальными
- Заботимся о первой странице
- **Рандомизация**

# Новое ранжирование Рамблера 2010

## Рандомизация поисковой выдачи

- **Борьба с обратной связью** в поведенческих рангах
- Чем дольше работает, тем лучше качество
- Рандомизируем только плохое, хорошее не рандомизируем
- Затруднение реинжиниринга алгоритмов ранжирования поиска
- Размазывание трафика по первой сотне сайтов.

Каковы результаты?



# Обзор доклада

- 1 Проблемы контролируемого обучения
- 2 Новое ранжирование Рамблера
- 3 Метрики качества ранжирования
- 4 Что влияет на долю рынка?



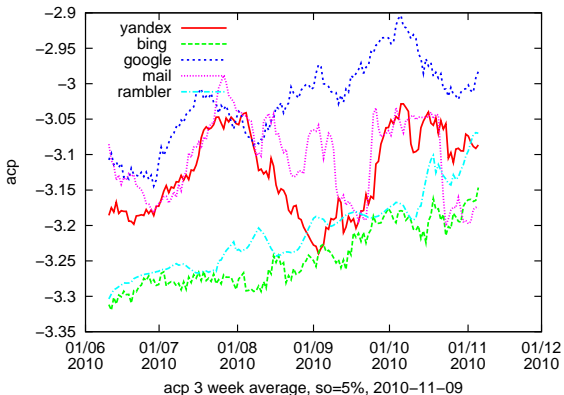
# Графики: Качество поиска p1c

Своих метрик недостаточно.



# Графики: Качество поиска аср конкурентов

*Наилучшее качество не у лидера по доле рынка*



Пользователи Рамблера - не гики.



# Обзор доклада

- 1 Проблемы контролируемого обучения
- 2 Новое ранжирование Рамблера
- 3 Метрики качества ранжирования
- 4 Что влияет на долю рынка?





## Что влияет на долю рынка?

- **Качество** (ранжирование)
- **Маркетинг** (телевидение)
- Новые ноутбуки и плагины к браузерам.
- реклама в результатах поиска **отталкивает**
- Фишки: вертикали, сниппеты, опечаточник, подсказки.

Каково влияние основных двух факторов?

## Оцениваем влияние качества

Падение метрик качества → Падение доли

	дата	Качество %/макс.	Рост доли % в месяц
mail-gogo	2010/01	80	-5
google	2010/08	100	0
yandex	2010/08	90	0
rambler	2010/01	80	-5
rambler-new	2010/11	90	0
mail-gogo	2010/08	90	0
mail-google	2010/10	100	5

Таблица: влияние качества

На Google и Yandex работают дополнительные факторы.

# Оцениваем влияние маркетинга

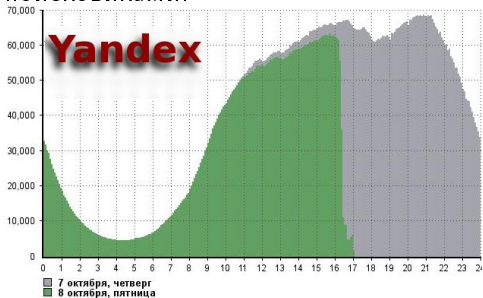


## Маркетинг - сейчас главный фактор

- Аудитория поиска растет на 60% в год
- Основная часть аудитории сформирована в последние два года.
- Влияние качества усилится лишь после того как рост остановится.
- Новые ноутбуки, смартфоны.
- Плагины к браузерам, браузеры.

# Оцениваем влияние маркетинга (по данным LiveInternet)

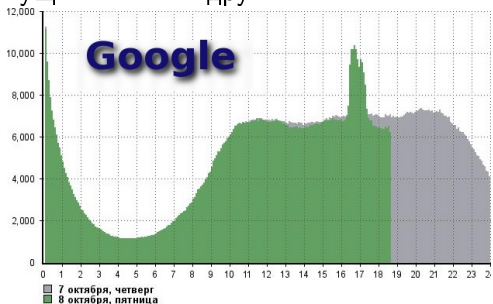
*Yandex: DDoS 2010/10/08* показал, что 2/3 пользователей Яндекса не воспользовались другими поисковиками.



# Оцениваем влияние маркетинга (по данным LiveInternet)

*Yandex: DDoS 2010/10/08*

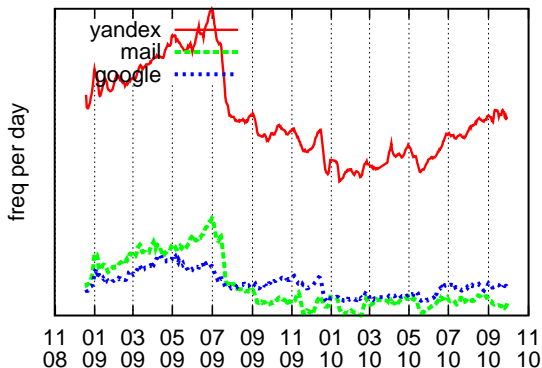
Возможно, что 2/3 пользователей ничего не знают о существовании других поисковиков.



# Маркетинг: сила бренда (по данным Рамблер)

Графики частотности запросов “яндекс”, “гугл”, “майл” в поиске Рамблера

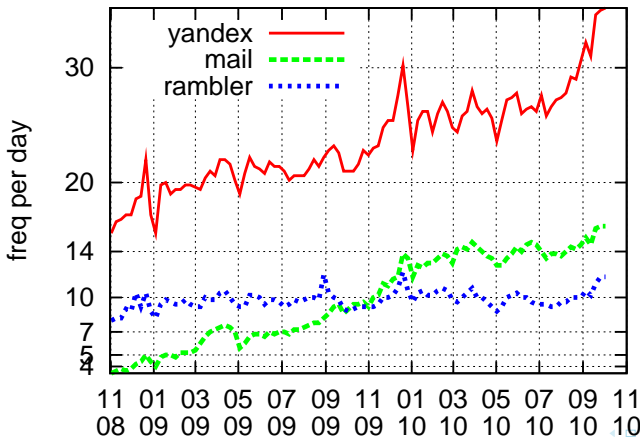
Яндекс в 2 раза более на слуху.



## Маркетинг: сила бренда (Google Trends)

Графики частотности запросов “рамблер”, “яндекс”, “майл” в поиске Google

Яндекс в 3 раза более на слуху.



# Маркетинг: доля новых пользователей (Rambler Top 100)

## *Rambler Top 100 cookie creation age*

Кука Рамблера (например '05F9DAFE4BED91CF000156714BA4DF01') внутри себя содержит дату создания.

	возраст куки, дней	доля новых, проценты
yandex	0.0051	11.1
google	7.85	4.0
rambler	34.9	2.61
bing	27.1	0.35

Таблица: Средний возраст куки

Яндекс иногда откуда-то берет толпы новых пользователей.  
Реклама на ТВ?



# Качество и маркетинг

## *Качество и маркетинг*

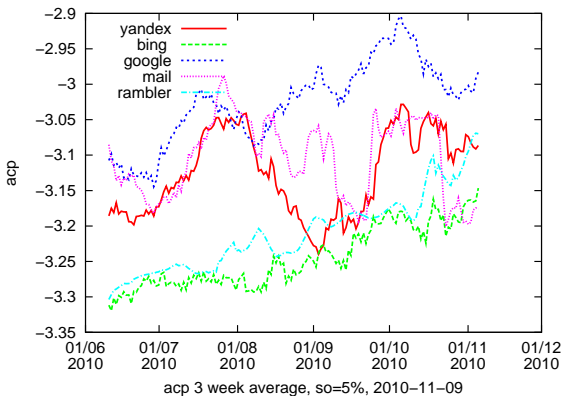
	качество место	маркетинг место	прогноз 1 год
yandex	II	I	рост
google	I	II	падение

Таблица: *Качество и маркетинг*

# Новое ранжирование Рамблера

- Выкатка: октябрь 2010

Будет ли эффект в росте доли рынка?



## Новое ранжирование: доля рынка

Доля рынка LiveInternet, как оценка поисковой системы.

- Зависит от чужого ранжирования (LiveInternet как фактор)
- Зависит от активности антивирусных роботов поисковиков.



## Новое ранжирование: доля рынка

Стабилизация по доле рынка означает рост на 60% в год.  
 Чем дольше работает - тем лучше качество.

