

Анализ поисковых запросов

Павел Браславский

Санкт Петербург,
ноябрь-декабрь 2010

План на ноябрь

- Введение
- Характеристики потока запросов
- Данные для анализа
- Сегментация
- Тематическая классификация
- Близкие запросы
- Практические задания

ВВЕДЕНИЕ

Предварительные замечания

- Логи запросов – «опыт» и богатство МП
- Современный поиск: «меньше информации, больше контекста»
- Мало открытых данных (→ проблема для академических исследований)
- Проблемы с персональными данными (приватность)
- Очень короткие тексты – сложность анализа
- Недостаток информации компенсируется большими объемами данных (веб)

Статистика запросов к Яндексу

День ▲	Посетители	Хиты	Хитов с посетителя
11.11.2010	14 248 496	103 547 538	7,27
10.11.2010	15 506 565	122 113 251	7,87
09.11.2010	15 164 865	122 395 405	8,07
08.11.2010	15 198 131	134 154 973	8,83
07.11.2010	12 173 980	106 138 000	8,72
06.11.2010	11 434 498	95 857 719	8,38
05.11.2010	11 517 047	95 532 682	8,29
04.11.2010	11 647 878	95 698 951	8,22

Прямой эфир

20 последних запросов пользователей Яндекса: - Mozilla Firefox

Файл Правка Вид Журнал Закладки Инструменты Справка

http://stat.yandex.ru/queries/last20.xml

Самые популярные Начальная страница Лента новостей

20 последних запросов польз...

Поиск Почта Карты Маркет Новости Словари Блоги Видео Картинки ещё ▾

Яндекс

статистика

[Посещаемость Яндекса](#)

20 последних запросов пользователей Яндекса:

бабирусса (8469)	стеклосетка (210083)
потенциал (11844146)	анжелина джоли (2765825)
форум нож белый медведь (132672)	папарацци (5530596)
обои с знаменитостями (4)	la2base.ru (49064)
на улице под юбкой (2886289)	IE 4.01 (194447)
ионизирующая радиация (446035)	www.dom2.ru (54372)
отправить сообщение (199672215)	фильмы с хилари дафф (394511)
китайская книга перемен (1491732)	владивосток авторынок (899125)
сновидения (4481337)	что такое snapshots (18875)
ghptthdfnbds (22)	Формат cdr (1384912)

[Следующие 20 запросов →](#)

Личная история запросов

Web History for pb@yandex-team.ru

All History

- [Web](#)
- [Images](#)
- [News](#)
- [Products](#)
- [Sponsored Links](#)
- [Video](#)
- [Maps](#)
- [Blogs](#)
- [Books](#)

- [Pause](#)
- [Remove items](#)

Trends

Bookmarks

Today

- 2:24pm Searched for [weka](#) - [Viewed 2 results](#)
- 2:24pm Searched for [rapid miner](#) - [Viewed 1 result](#)

Yesterday

- 8:08pm Searched for [google zeitgeist](#) - [Viewed 1 result](#)
- 6:46pm Searched for [малахит сандей](#) - [Viewed 1 result](#)
- 5:13pm Searched for [YetiRank: Everybody Lies](#) - [Viewed 4 results](#)

Jul 28, 2010

- 4:19pm [☆](#) Searched for [57.002445 60.326752](#)
- 4:16pm [☆](#) Searched for [гать](#)
- 4:15pm [☆](#) Searched for [56.941744 60.347053](#)
- 4:15pm [☆](#) Searched for [undefined](#)
- 4:15pm [☆](#) Searched for [57.002445 60.326752](#)
- 4:12pm [☆](#) Searched for [56.941744 60.347053](#)
- 4:10pm [☆](#) Searched for [N54 04.531 E036 02.141](#)

Web Activity

« Jun Jul 2010

S	M	T	W	T	F	S
27	28	29	30	1	2	3
4	5	6	7	8	9	10
11	12	13	14	15	16	17
18	19	20	21	22	23	24
25	26	27	28	29	30	31

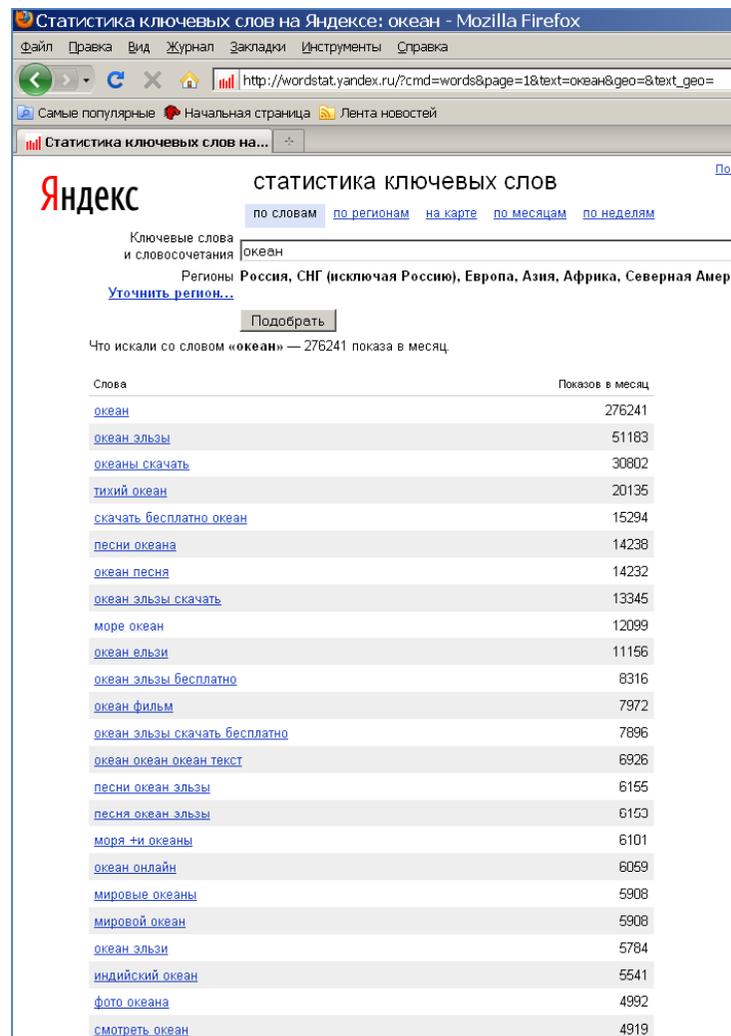
Today, Jul 30

1 - 5	6 - 10	11 - 20	21+
-------	--------	---------	-----

Total Google searches: 6222



Статистика слов запросов



Статистика ключевых слов на Яндексе: океан - Mozilla Firefox

Файл Правка Вид Журнал Закладки Инструменты Справка

http://wordstat.yandex.ru/?cmd=words&page=1&text=океан&geo=&text_geo=

Самые популярные Начальная страница Лента новостей

Статистика ключевых слов на...

Яндекс статистика ключевых слов [Поиск](#)

[по словам](#) [по регионам](#) [на карте](#) [по месяцам](#) [по неделям](#)

Ключевые слова и словосочетания

Регионы **Россия, СНГ (исключая Россию), Европа, Азия, Африка, Северная Америка**

[Уточнить регион...](#)

Что искали со словом «океан» — 276241 показа в месяц

Слова	Показов в месяц
океан	276241
океан эльзы	51183
океаны скачать	30802
тихий океан	20135
скачать бесплатно океан	15294
песни океана	14238
океан песня	14232
океан эльзы скачать	13345
море океан	12099
океан эльзи	11156
океан эльзы бесплатно	8316
океан фильм	7972
океан эльзы скачать бесплатно	7896
океан океан океан текст	6926
песни океан эльзы	6155
песня океан эльзы	6153
моря +и океаны	6101
океан онлайн	6059
мировые океаны	5908
мировой океан	5908
океан эльзи	5784
индийский океан	5541
фото океана	4992
смотреть океан	4919

Find Trends

Discover

Time



Popular Searches

▼ Today

- patrick swayze
- bond market
- kate middleton chelsy davy
- julianne hough
- derek hough

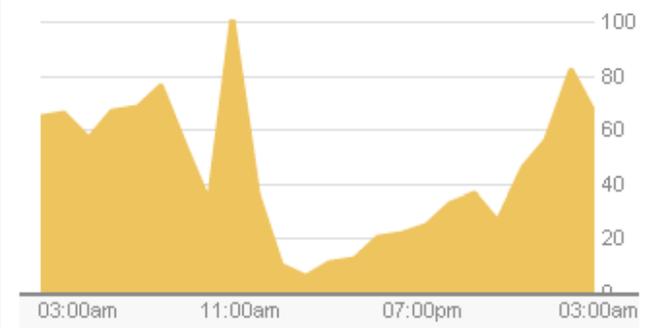
► Past 7 Days

► Past 30 Days

emma watson

Expand

SEARCHES OVER TIME



BY DEMOGRAPHIC

Age | Gender | Both

Women



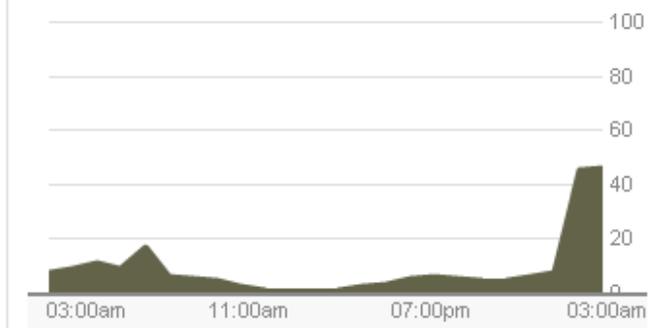
Men



daniel radcliff...

Expand

SEARCHES OVER TIME



BY DEMOGRAPHIC

Age | Gender | Both

Women



Men



Find Trends

snowboard

ski

Discover

Time



Today | Past 7 Days | Past 30 Days

Popular Searches

▼ Today

patrick swayze

bond market

kate middleton chelsy davy

julianne hough

derek hough

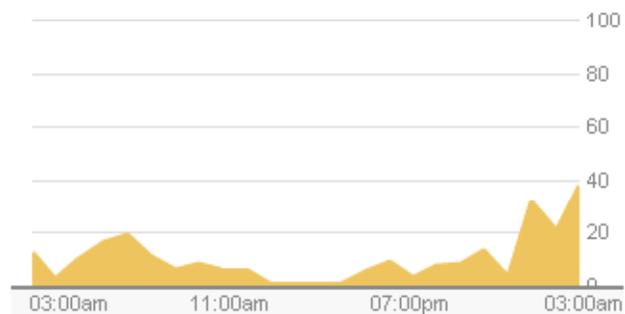
► Past 7 Days

► Past 30 Days

snowboard

Expand

SEARCHES OVER TIME



BY DEMOGRAPHIC

Age | Gender | Both

Age

Below 24 21%

25 to 34 56%

35 to 44 13%

45 to 54 8%

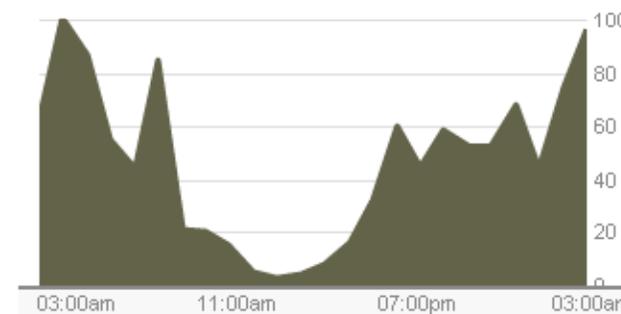
55 to 64 2%

Above 65 0%

ski

Expand

SEARCHES OVER TIME



BY DEMOGRAPHIC

Age | Gender | Both

Age

Below 24 8%

25 to 34 25%

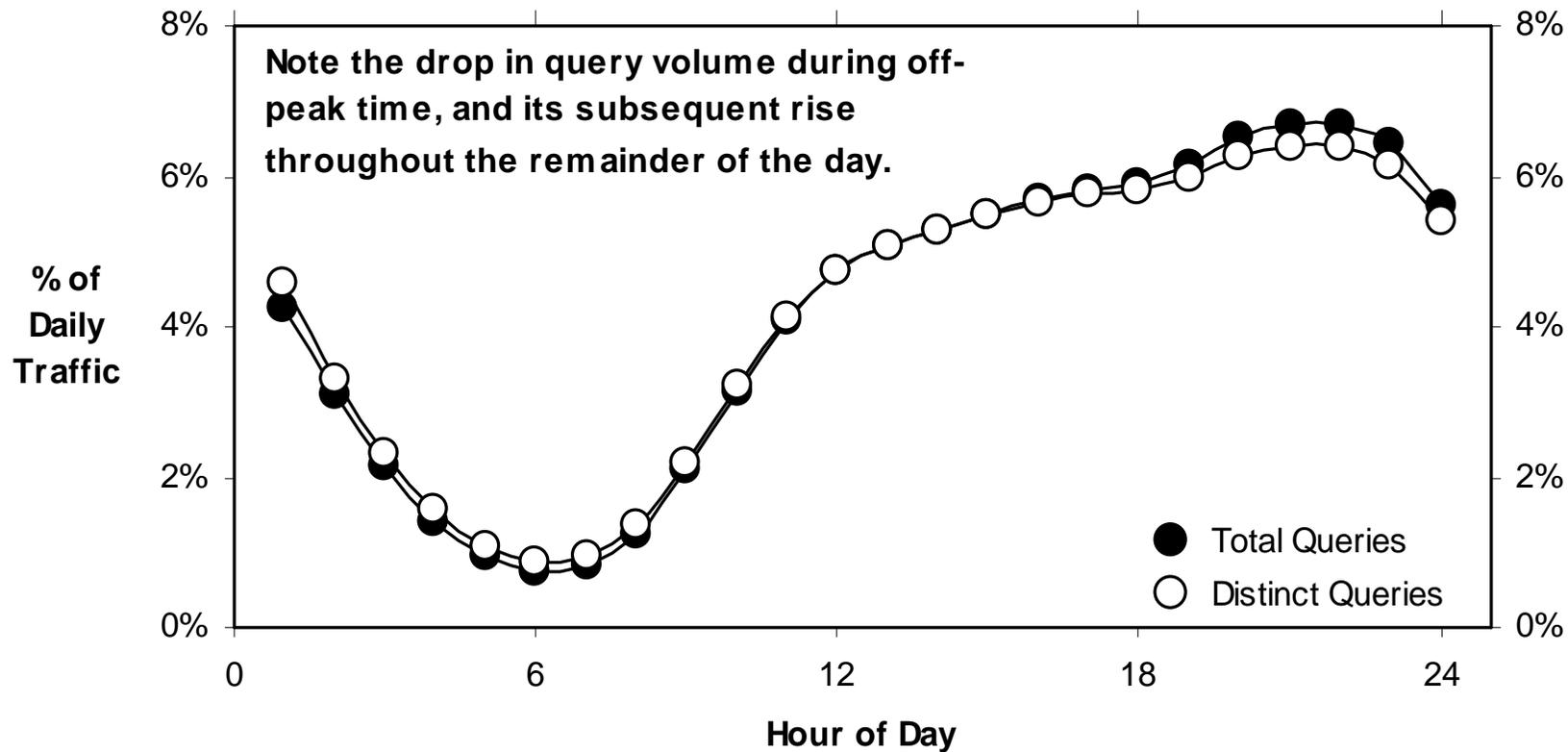
35 to 44 23%

45 to 54 24%

55 to 64 13%

Above 65 5%

Traffic Volume Over a Day

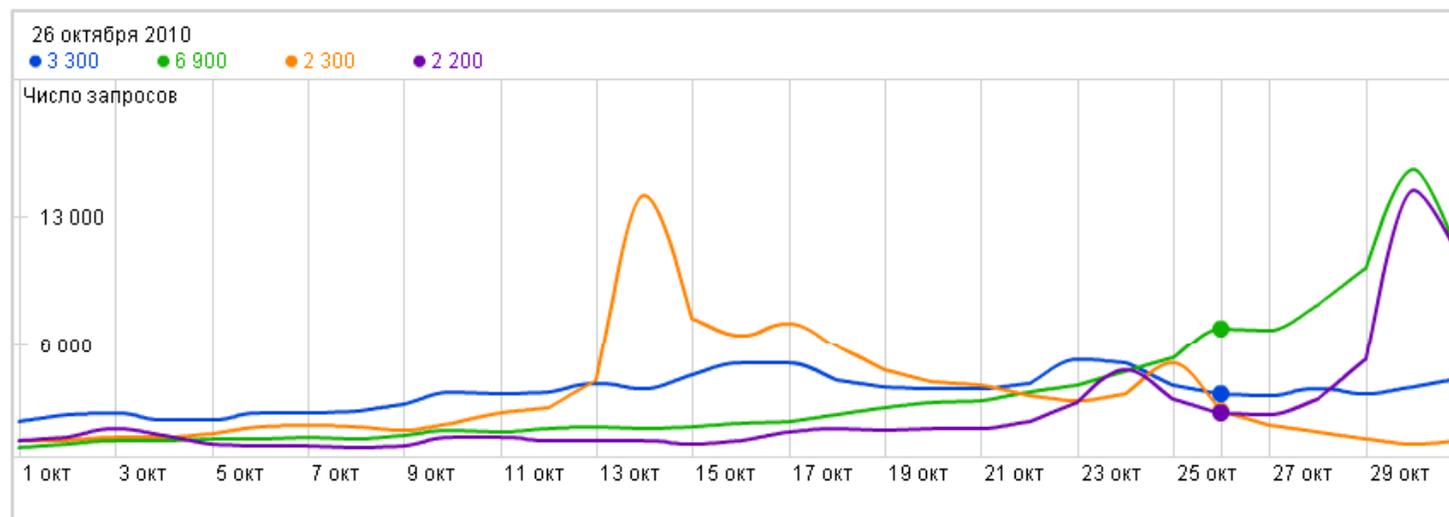


Динамика запросов

Развитие интересов

В регионе [Северо-Запад](#)

- Зимняя одежда
- Хэллоуин
- Перепись населения — 2010
- Фильм «Рэд»
- Переход на зимнее время
- Фильм «Паранормальное явлени...
- Ноябрьские праздники
- Сериал «Глухарь» 3 сезон
- Игра «Fallout: New Vegas»
- Сериал «Сплетница»



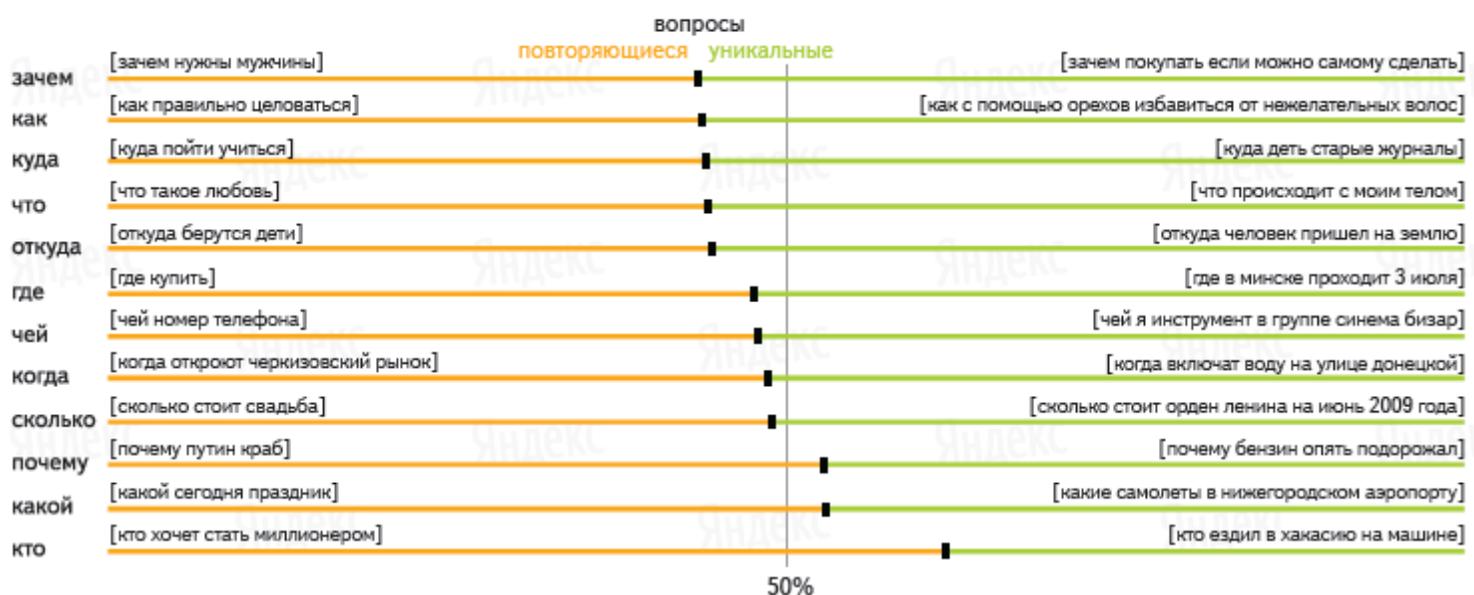
<http://interes.yandex.ru/>



Несколько цифр

- Длина запроса 2-3 слова
- Поисковая сессия в среднем 3 запроса
- 2-3% сформулированы как вопрос
- 12-15% запросов содержат опечатки

Запросы – вопросы



http://company.yandex.ru/facts/researches/ya_search_2009.xml

Классификация запросов / информационных потребностей (Broder, 2002)

- **Informational** – want **to learn** about something (~40% / 65%)

Low hemoglobin

- **Navigational** – want **to go** to that page (~25% / 15%)

Warsaw Airport

- **Transactional** – want **to do something** (web-mediated) (~35% / 20%)

- Access a service

Barcelona weather

- Downloads

Mars surface images

- Shop

Canon S410

- Gray areas

- Find a good hub

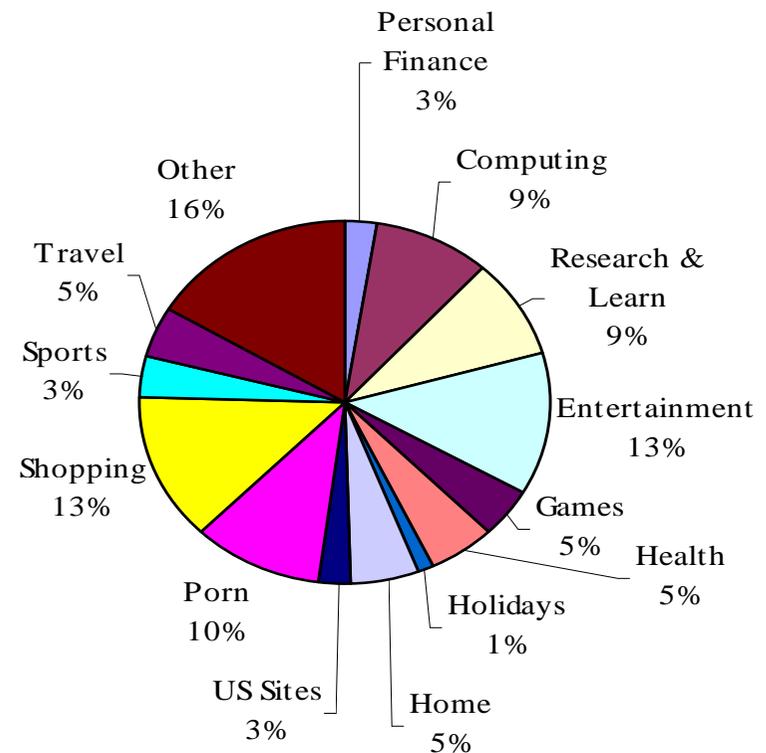
Car rental Poland

- Exploratory search “see what’s there”

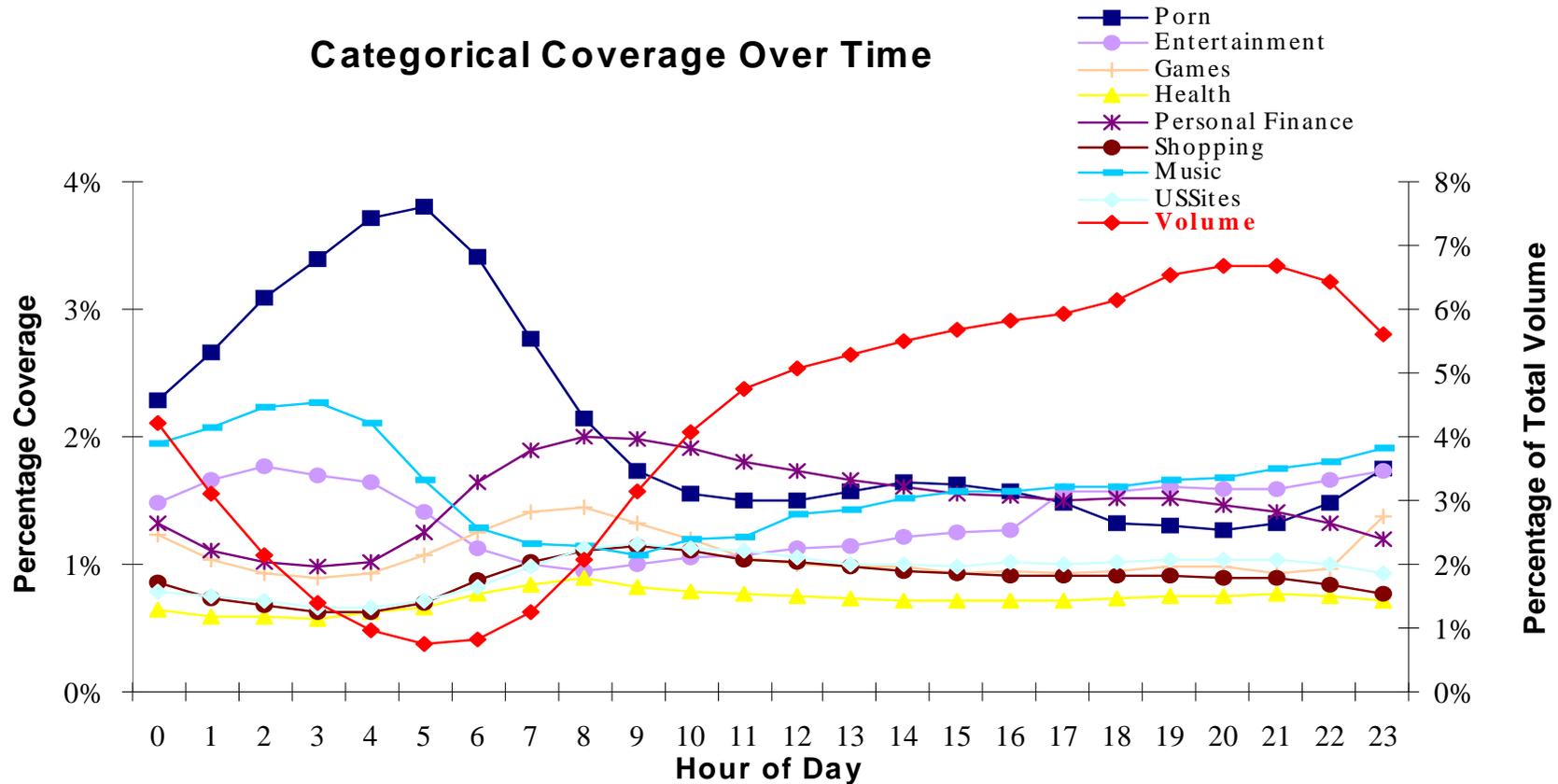
Category Breakdown

- Query lists for each category formed by a team of human editors
- Query stream classified by exactly matching each query to category lists

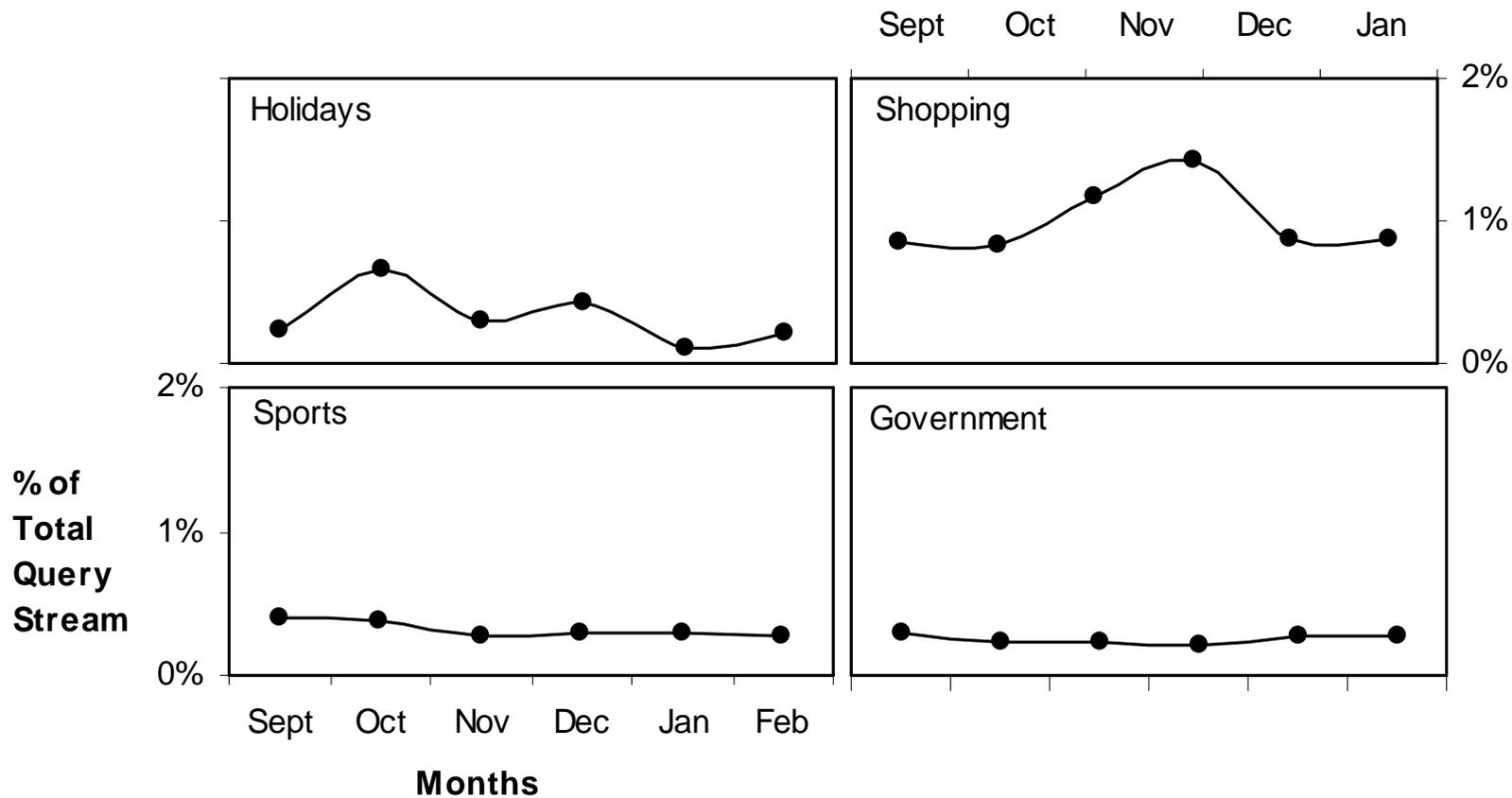
Sampled Categorized Query Stream Breakdown



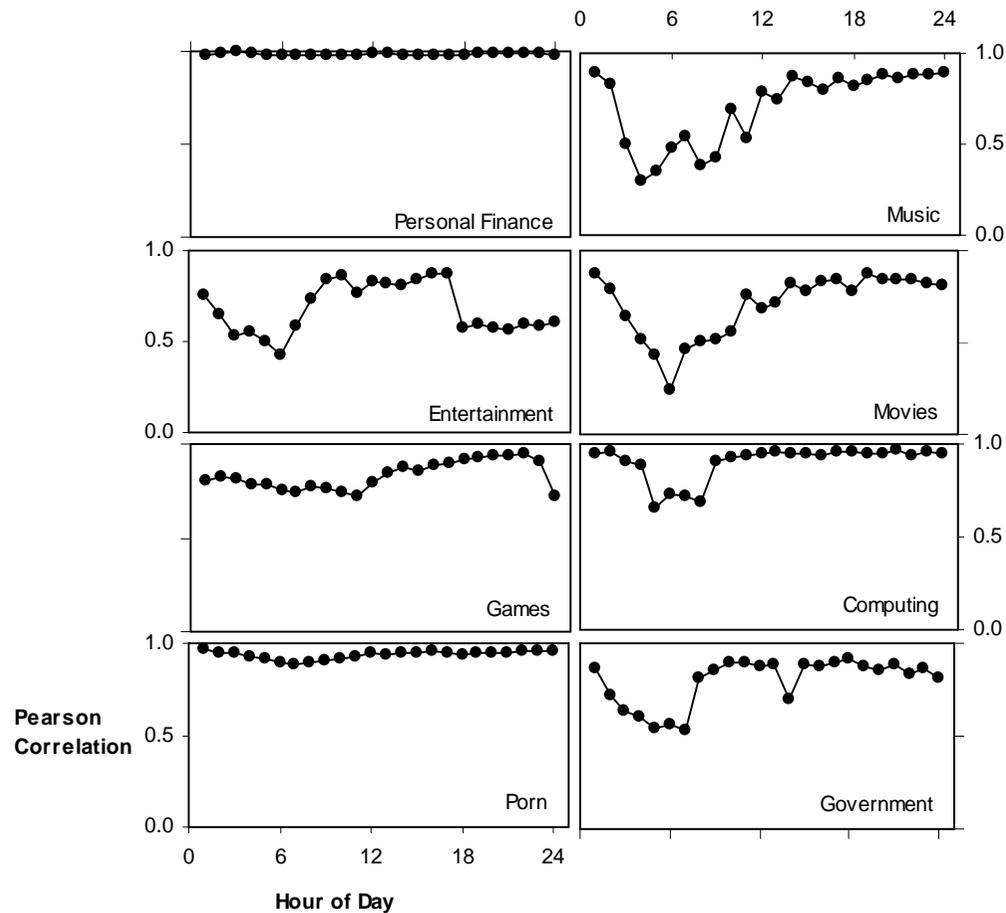
Category Popularity Over a Day



Category Popularity Over Six Months



Pearson Correlations for Selected Categories Over A Day



Источники и типы данных



Jiang et al. 2010

Какая информация у нас есть?

- текст запроса
- время
- IP → география
- Cookie → (уникальный) пользователь
- клики на результатах поиска
- (персональные данные, соцдем)

Данные (→ методы)

- отдельные запросы
- список запросов
- + время
- + сниппеты/документы
- + клики
- ...

Яндекс: ~6 Кбайт/запрос + ~0,5 Кб/клик

Доступные данные

- Excite 1997, 1999, 2001
- AOL 2006
- ИМАТ 2004
- MSN Search query Log excerpt (RFP 2006 dataset)
- ...

Excite 1997

- Запросы за один день (16 сентября 1997 г.)
- userId, timeStamp, query
- ~1М запросов (много повторов)

0C6B5395895CD808	970916125351	henri rousseau
0C6B5395895CD808	970916125511	henri rousseau+tiger
949946B881F137F0	970916115517	"pharmacy"
949946B881F137F0	970916115550	prescriptions
91A98BC9BEDCF053	970916075435	australian+chat+victoria
61305D2ADC74BC78	970916095742	dailyplanet
61305D2ADC74BC78	970916095846	dailyplanet
33D1A0D49E8DB2AB	970916144916	maizehighschool
33D1A0D49E8DB2AB	970916144951	maize high school
FCBB8401805D783F	970916212508	warez strata studio pro
FCBB8401805D783F	970916212541	warez mac
12FE04344578F249	970916202819	"midwife conference"
12FE04344578F249	970916202924	midwifery
477CC4190EF76EB4	970916165602	nrwmac
477CC4190EF76EB4	970916172706	npac
477CC4190EF76EB4	970916175242	nrwmac

ИМАТ 2004

- 7 дней * 10% от 5-10 миллионов запросов в день → 3,5 Гб

<UID1>

<запрос1> <время> <найдено документов> <номер страницы>

<URL1> <время выбора>

<URL2> <время выбора>

...

<запрос2> <время> <найдено документов> <номер страницы>

<URL1> <время выбора>

<URL2> <время выбора>

...

...

<UID2>

...

http://company.yandex.ru/academic/grant/datasets_description.xml

AOL 2006

- Большой скандал!!!
- ~20M web queries from ~650k users over three months

AnonID	Query	QueryTime	ItemRank	ClickURL
993	myspace.co	01.03.2006 12:13		
993	myspace.com	01.03.2006 12:13		
993	googl	01.03.2006 15:03		
993	chasebadkids.net	03.03.2006 16:55	1	http://www.chasebadkids.net
1268	ozark horse blankets	01.03.2006 17:39	8	http://www.blanketsnmore.com
1268	www.ghostrockranch.com	04.03.2006 13:58		
1268	openrangeht.zachsairforce.com	09.03.2006 22:38		
1268	sstack.com	11.03.2006 0:17		
1268	www.mecab.org	12.03.2006 18:59		
1268	www.raindanceexpress.com	18.03.2006 20:13		
1268	www.victoriacostumiere.com	19.03.2006 0:26		
1268	osteen-schatzberg.com	21.03.2006 17:55		
1268	osteen-schatzberg.com	21.03.2006 17:55	1	http://www.osteen-schatzberg.com
1268	osteen-schatzberg.com	21.03.2006 17:55	2	http://www.osteen-schatzberg.com

MSN Search query Log excerpt

- 15 million queries
- Sampled over one month
- Queries from the US site (mostly English)

Per query attributes included:

- Session ID
- Time-stamp
- Query string
- Number of results on results page
- Results page number

Data per query for each result clicked:

- URL
- Associated query
- Position on results page
- Time-stamp

СЕГМЕНТАЦИЯ ЗАПРОСОВ

Сегментация запросов

Сегментация:

1. поиск

2. дальнейшая обработка запросов

международный почтампт | москва
молодежный отдых | в турции
официальный сайт | автоваз
Купить | кроссовки | Nike Zoom BB
магазин | рыбачьте с нами
ирина круг | пусть сейчас я плачу | слушать
сбербанк россии | в алтайском крае
знак зодиака | близнецы
нино катамадзе | билеты
Смысл названия | рассказа | матренин двор
Сергей Тармашев | Корпорация | скачать | бесплатно

вечерняя москва vs пицца москва

банк москвы vs банки москвы

Сегментация: подходы

- похоже на выделение устойчивых словосочетаний
- + микросинтаксис
- лог vs корпус текстов
- ML (больше признаков, более богатое описание)
- внешние ресурсы (Wikipedia)
- эвристики ([дима билан] ← [димабилан])

Сегментация на основе лога

$$\text{conn}(S) = \text{freq}(S) * I(w_1 \dots w_{n-1}; w_2 \dots w_n)$$

$$I(X; Y) = \sum_{y \in Y} \sum_{x \in X} p(x, y) \log \frac{p(x, y)}{p(x)p(y)}, \quad SI(x, y) = \log \frac{p(x, y)}{p(x)p(y)}.$$

```
msdn library visual studio
```

```
34259: (msdn library)[5110] (visual studio)[29149]
```

```
29149: msdn[47658] library[209682] (visual studio)[29149]
```

```
5110: (msdn library)[5110] visual[23873] studio[53622]
```

```
41: (msdn library visual studio)[41]
```

```
7: msdn[47658] (library visual studio)[7]
```

```
0: msdn[47658] library[209682] visual[23873] studio[53622]
```

Risvik et al. WWW2003

Сегментация на основе ML

нино | катамадзе | билеты



(0, 1, 0, 1, 1, 0)

(1, 1, 0, 1, 0, 1)

Table 1: Indicator features.

Name	Description
is-the	token $x = \text{“the”}$
is-free	token $x = \text{“free”}$
POS-tags	Part-of-speech tags of pair $x_{LO} x_{RO}$
fwd-pos	position from beginning, i
rev-pos	position from end $N - i$

Table 2: Statistical features.

Name	Description
web-count	count of “ x ” on the web
pair-count	web count “ $w x$ ”
definite	web count “the $w x$ ”
collapsed	web count “ wx ” (one word)
and-count	web count “ w and x ”
genitive	web count “ w ’s x ”
Qcount-1	Counts of “ x ” in query database
Qcounts-2	Counts of “ $w x$ ” in database

Bergsma and Wang, 2007

Сегментация: «наивный подход»

На основе веб-корпуса n-грамм

$$score(S) = \sum_{s \in S, |s| \geq 2} |s|^{|s|} \cdot count(s)$$

Hagen et al. SIGIR2010