

Санкт-Петербургский государственный университет  
Математико-механический факультет  
Кафедра системного программирования

# **Алгоритм построения оценок весов интентов для многозначных запросов**

Курсовая работа студента 445-ой группы  
Григорьева Артёма Валерьевича

Научный руководитель ..... к. ф.-м. н. Грауэр Л.В.

Санкт-Петербург  
2012

## Оглавление

Введение.....	3
Обзор предыдущих работ.....	5
Алгоритм.....	6
Общая схема алгоритма.....	6
Эталонное множество запросов.....	6
Источник статистики.....	6
Определение набора связанных запросов.....	7
Кластеризация связанных запросов.....	9
Распределение сессий по кластерам.....	12
Веб-инструмент.....	13
Результаты работы алгоритма.....	15
Анализ полноты и точности.....	15
Примеры обработанных результатов.....	16
Заключение.....	17
Список литературы.....	18

## Введение

Проблема поиска в сети Интернет ввиду постоянного роста объёмов общедоступной информации с каждым годом становится всё актуальнее. Основная цель поисковых систем, таких как Яндекс или Google, – отвечать на запросы пользователей релевантными документами.

Для оценки качества поиска используются различные метрики. Метрика в данном смысле – это функция от запроса и списка документов. К примеру,  $r_{found}$  [1] вычисляется на основании оценок релевантности документов по данному запросу. Если из формулировки запроса ясно, какие документы считать релевантными, а какие нет, то выставить оценки по некоторой шкале – несложная задача для человека. Несколько сложнее дело обстоит с классом многозначных запросов. Запрос относится к данному классу, если по нему существует несколько возможных интентов – намерений пользователя. В таблице 1 приведены примеры таких запросов и возможных интентов для них.

Запрос	Некоторые возможные интенты
ягуар	марка автомобиля, животное, алкогольный напиток, фильм
титаник	история корабля, факты о фильме, скачать фильм, где посмотреть в 3D
танки	боевая техника, игра «Танки Онлайн», игра «World of Tanks»

Табл. 1. Примеры интентов для многозначных запросов

Значение метрики качества для многозначного запроса вычисляется по формуле:

$$metrics = \sum_{i=1}^{k+1} w_i \cdot metrics(O_i)$$

как взвешенное среднее значение метрики по каждому интенту  $O_i$  в отдельности с весами  $w_i$  – вероятностями того, что пользователь имел данный интент, задавая данный многозначный запрос [4].

В роли оценки для этих весов выступает доля сессий, в которых пользователь предполагал конкретное значение запроса, по отношению ко всем сессиям, содержащим данный запрос. Здесь и далее под сессией подразумевается последовательность действий пользователя по отношению к поисковой системе (запросы и клики на документы), заключённая в некотором временном отрезке. Задача разграничения отдельных сессий является предметом иных исследований и в данной работе не рассматривается. Для простоты будем считать, что пользователь начал новую сессию после паузы в 10 минут.

Вручную расчёт оценок весов интенгов для конкретного запроса производится по следующему принципу. Из поисковых логов случайным образом выбирается некоторое небольшое количество сессий, содержащих данный запрос, которые затем вместе с предполагаемым набором интенгов отправляются на разметку ассессорам. Задачей последних является понять по действиям пользователя из данной сессии, к какому из интенгов её можно отнести. Итоговые оценки получаются уже описанным выше способом: как доли сессий с данным интенгом среди общего количества сессий в выборке.

Малый размер выборки и ручное составление набора интенгов могут привести к неполноте и неточности полученных оценок. Эти серьёзные недостатки связаны прежде всего с ограниченными возможностями ассессорской службы.

Так возникла тема данной работы. Требуется придумать и реализовать алгоритм, рассчитывающий оценки весов интенгов для многозначных запросов. Для проверки качества его работы будут использоваться уже полученные оценки. Постараемся добиться не меньшей полноты полученных наборов интенгов, а также незначительных различий в оценках при фиксированных наборах.

## Обзор предыдущих работ

Первая работа, о которой хотелось бы упомянуть, стала во многом основополагающей для разработанного алгоритма. Речь о статье E. Sadikov и др. [2] о кластеризации уточняющих запросов по интендам пользователей. В ней рассматривается задача генерации релевантных подсказок по уточнению исходного запроса к поисковой системе, строится марковская цепь, содержащая в качестве состояний запросы и документы, производится случайное блуждание по ней, после чего запросы кластеризуются с использованием предельных вероятностей для документов.

В продолжение этой работы F. Radlinski и другие предлагают при помощи описанной модели выделять поисковые интенды для запроса [5], в дальнейшем изложении будет указано, почему их метод работает далеко не для всех классов запросов (см. метод частотных переформулировок, стр. 7).

Проблему кластеризации запросов к поисковой системе решают в своей работе Ji-Rong Wen и др. [3]. Несмотря на то, что они работают со специфическими данными (поиск по FAQ), из их работы можно почерпнуть сведения о различных используемых в предметной области функциях сходства<sup>1</sup>.

Исследовать сходство запросов, проецируя граф запросов, предлагают I. Bordini и др. [6]. В их модели не рассматриваются документы, описываются различные весовые функции для рёбер, а кластеризация производится благодаря нахождению подграфов.

Более подробно о моделях графа запросов и графа документов рассказывается в работах N. Craswell и M. Szummer [7] и P. Boldi и др. [8] соответственно.

---

<sup>1</sup> Similarity function.

# Алгоритм

## Общая схема алгоритма

Итоговый алгоритм работает по следующей схеме. На вход поступает многозначный запрос  $q$ , затем с использованием различной пользовательской статистики определяется набор связанных с исходным запросом  $\{q_i\}$ . Данный набор кластеризуется, каждый кластер  $C_j$  соответствует интену пользователя (определение этого значения по набору запросов в кластере остаётся за аналитиком). После кластеризации определение искомым оценок весов происходит аналогичным с ручной оценкой образом: выбирается набор случайных сессий, содержащих исходный запрос  $q$ , каждая сессия может с некоторыми весами относиться к различным кластерам. Среднее таких весов по всем сессиям из выборки и есть искомые оценки, получаемые вместе с самими кластерами на выходе.

Таким образом весь процесс можно разделить на три этапа:

- 1) Определение набора связанных запросов;
- 2) Кластеризация связанных запросов;
- 3) Распределение сессий по кластерам.

Далее опишем каждый этап в отдельности.

## Эталонное множество запросов

Во введении были определены два основных требования к результатам работы алгоритма: полнота и точность. В качестве эталонного множества были выбраны 49 многозначных запросов, по которым есть наборы интенгов, размеченные ассессорами сессии и оценки для весов интенгов. На всех шагах построения алгоритма производились расчёты по данному множеству запросов и выбиралось то из альтернативных решений, которые выигрывало у остальных в сравнении с ассессорами. Этим обуславливается размер эталонного множества: все сравнения производились вручную.

## Источник статистики

Источником статистики являются очищенные от роботов поисковые логи «Яндекса» за период в 84 дня с 26 декабря 2011 года по 18 марта 2012 года. Помимо самих логов в работе использовалась агрегированная позапросно статистика по кликам, обратным кликам и переформулировкам.

Запросы нормализуются приведением к нижнему регистру, удалением непечатных символов, заменой подряд идущих пробельных символов на один, удалением конечных

пробелов.

URL-адреса документов нормализуются приведением к нижнему регистру и к главным зеркалам доменов, удалением префикса «www.», хеша, концевых слешей и вопросительных знаков.

## Определение набора связанных запросов

Построение набора связанных запросов влияет на качество результата по двум параметрам. Во-первых, в набор должно входить хотя бы по одному запросу для каждого интента из некоторого идеального набора для данного запроса. Иначе не будет достигнута желаемая полнота. Во-вторых, несколько забегая вперёд, скажем, что при распределении сессий по кластерам происходит матчинг<sup>2</sup> по тексту. Таким образом нужно, чтобы в набор входили разнообразные формулировки одного и того же запроса. Это влияет на количество сессий, которые будут распределены по кластерам, и следовательно точность полученных оценок весов. Например, если в набор входит [андеграунд играть онлайн]<sup>3</sup>, то желательно добавить туда такие запросы, как [онлайн игра андеграунд], [играть андеграунд онлайн] и т.п..

В ходе экспериментов были опробованы различные варианты получения набора связанных запросов. Перечислим их далее в том же порядке, в котором производились эксперименты.

**Частотные переформулировки.** Определим понятие: запрос В является *переформулировкой* запроса А, если он встречается хотя бы в одной сессии после запроса А.

Зачастую пользователь переформулирует свой запрос, если он не удовлетворён поисковыми результатами. Предлагается в качестве исходного множества брать некоторое количество самых частотных переформулировок.

**Общие клики.** Минусом первого способа является необходимость иметь хорошую статистику о переформулировках запроса. В ходе экспериментов было замечено, что по тем запросам, где в поисковой выдаче присутствуют документы по большинству предполагаемых интентов, переформулировки содержат много «мусора». Действительно, пользователю незачем переформулировать запрос, если нужный документ уже найден.

Таким образом встал вопрос об интерпретации сделанных пользователем кликов. Нужно было каким-то образом понять, о чём тот или иной документ. Такую информацию можно извлечь из обратных кликов – статистики, где по URL-адресу документа хранится список запросов, по которым были клики по данному документу как по поисковому

---

<sup>2</sup> Для краткости назовём процесс отнесения действия пользователя к кластеру «матчингом» (от англ. match)

<sup>3</sup> Для обозначения текста поискового запроса внутри статьи будут использоваться квадратные скобки.

результату. Под способом получения связанных запросов «Общие клики» подразумевается следующее: для исходного запроса  $q$  выбираются несколько первых документов среди отсортированных по некоторому параметру (количеству кликов или CTR), для каждого из этих документов выбираются самые частотные запросы из статистики по обратным кликам и добавляются в множество связанных.

**Смешанный способ.** В смешанном способе результаты от двух предыдущих смешиваются в один список, который во избежание «шума» сортируется по относительной частотности. Для переформулировок это  $\frac{ref(q \rightarrow q')}{\max ref(q \rightarrow .)}$ , для общих кликов это

$$\frac{\max_d clicks(q,d)}{\max_d clicks(q',d)} \cdot \frac{clicks(q',d)}{\max_q clicks(q,d)}$$

Из этого списка выбираются несколько (в экспериментах 20) первых запросов.

**Расширенный способ.** Как уже было упомянуто выше, значительную роль при обработке сессий играет текстовое совпадение запросов. В ходе экспериментов было замечено, что оно играет немаловажную роль и что зачастую из-за этого сессия остаётся нераспознанной (см. подробнее про обработку сессий далее на стр. 12). В расширенном способе к уже набранным запросам из предыдущего метода добавляются различные их варианты по следующему алгоритму: просматриваются случайные сессии по исходному запросу (те же, что обрабатываются на третьем этапе), для каждого запроса из них производятся попытки определить, похож ли он на какой-нибудь из уже имеющихся. Для этого оба запроса разбиваются на отдельные слова, затем проверяется, содержит ли запрос из сессий все слова из имеющегося. При чём совпадение проверяется не точное, а при условии, что нормированное расстояние Левенштейна не превосходит порогового значения (в данном случае 0.1).

При таком способе к запросу [снегурочка островский] будут добавлены: [краткое содержание снегурочки], [снегурочка островский смотреть онлайн] и т.п.



В таблице 2 приведены данные о сравнении трёх запусков на эталонном множестве запросов:

	Переформулировки	Смешанный	Расширенный
Количество запросов с полным набором интенгов	28	31	32
Количество запросов с набором не более чем без одного интенга	39	44	46
Средний процент обработанных сессий среди всех запросов	51%	52%	56%
Средний процент обработанных сессий среди запросов с полным набором интенгов	59%	54%	58%

Табл. 2. Сравнение трёх вариантов получения набора связанных запросов

По данным таблицы можно сделать вывод, что расширенный способ получения набора связанных запросов выигрывает у остальных по таким параметрам, как полнота получаемых наборов интенгов и количество обработанных на третьем этапе сессий.

## Кластеризация связанных запросов

После того, как множество связанных запросов сформировано, следующим этапом работы алгоритма является разбиение этого множества на кластеры, соответствующие пользовательским интенгам. В общем случае задача кластеризации сводится к двум этапам: сопоставлению каждому элементу из кластеризуемого множества вектора в некотором многомерном пространстве, а затем кластеризации этих векторов.

Идея кластеризации запросов, используемая в данной работе, почерпнута из статьи «Clustering Query Refinements by User Intent» [2]. По сравнению с оригинальной была несколько изменена вероятностная модель.

**Построение марковской модели.** Определим взвешенный ориентированный граф, моделирующий поведение пользователей при работе с поисковой системой. Для каждого запроса из множества связанных  $Q = \{q_i, i = 1..N\}$  возьмём первые  $L=100$  самых кликабельных<sup>4</sup> документов, объединение всех документов образует множество документов  $D = \{d_k, k = 1..M\}$ . Вершины в графе будут двух типов: запросы  $q \in Q$  и документы  $d \in D$ .

<sup>4</sup> Т.е. первые L среди всех (отсортированных по количеству кликов) документов, которые когда-либо показывались в качестве поисковых результатов по данному запросу.

Два запроса соединены ребром  $q_1 \rightarrow q_2$ , если конечный запрос является переформулировкой начального. Запрос и документ соединены ребром  $q \rightarrow d$ , если по данному документу были клики как результату в поисковой выдаче по данному запросу.

Определим веса рёбер как статистические оценки вероятностей кликов и переформулировок соответственно:

$$w(u \rightarrow v) = \begin{cases} \varepsilon \cdot ctr \cdot P_{click} \cdot D_{click}, & u \in Q, v \in D \\ (1 - \varepsilon) \cdot P_{reformulation} \cdot D_{reformulation}, & u, v \in Q, \\ 1, & u = v \in D \end{cases}$$

где

$$ctr = \begin{cases} \frac{clicksCount(q, d)}{showsCount(q, d)}, & showsCount(q, d) \neq 0 \\ 0, & showsCount(q, d) = 0 \end{cases};$$

$clicksCount(q, d)$  – количество кликов по документу  $d$  как результату в поисковой выдаче по запросу  $q$ ,

$showsCount(q, d)$  – количество показов документа  $d$  как результата в поисковой выдаче по запросу  $q$ .

$$P_{click} = \frac{clicksCount(q, d)}{\sum_{d' \in D} clicksCount(q, d')};$$

$$P_{reformulation} = \frac{reformulationsCount(q_1, q_2)}{\sum_{q' \in Q} reformulationsCount(q_1, q')};$$

$reformulationsCount(q, q')$  – количество переформулировок запроса  $q$  в запрос  $q'$ .

Для обеспечения бóльшего правдоподобия данных статистики вводятся дисконтирующие множители за малое количество данных:

$$D_{click} = \frac{1}{1 - \exp(10 - requests)}, \text{ где } requests \text{ – общее количество запросов пользователей.}$$

$$D_{reformulation} = \frac{1}{1 - \exp(10 - reformulations)}, \text{ где } reformulations \text{ – общее количество}$$

переформулирований данного запроса.

$\varepsilon$  – нормирующий множитель, в экспериментах брался равным 0.6.

В граф добавлены петли с весом 1 в вершинах с документами.

Таким образом данный граф можно рассматривать как однородную цепь Маркова с поглощающими состояниями – документами. Данная модель не эмулирует поведение

пользователей полноценно: здесь клик на документ происходит после последовательности запросов и означает завершение поиска, однако по мнению автора отвечает задаче определения поискового интента: грубо говоря, пользователь переформулирует запрос, если не находит среди документов нужный, и кликает на один из представленных, если ему кажется, что он отвечает его нуждам.

**Предельные вероятности.** Рассмотрим матрицу смежности  $P$  для построенного выше графа. В ячейке  $P[i,j]$  стоит значение вероятности перехода из состояния  $i$  в состояние  $j$ . Если возвести данную матрицу в степень  $n$ , то в ячейке  $P^n[i,j]$  будет стоять значение вероятности перейти из состояния  $i$  в состояние  $j$  не более чем за  $n$  шагов.

Таким образом, согласно интерпретации построенной модели, после возведения матрицы  $P$  в достаточно большую степень мы получим в ячейках, соответствующих переходам от запроса к документу, вероятности клика по данному документу после цепочки переформулировок данного запроса.

Введём понятие вектора  $V(q)$  – вектора предельных вероятностей документов для запроса  $q$ . В нём на  $k$ -ой позиции стоит значение ячейки предельной матрицы с координатами  $i$  – номером состояния, соответствующего запросу  $q$ , и  $j$  – номером состояния, соответствующего  $k$ -ому документу из предварительно отсортированного множества  $D$ :

$$V(q)_k = P^n[i, j]$$

Теперь, когда для каждого запроса получен вектор, ему соответствующий, можно перейти непосредственно к кластеризации запросов.

**Кластеризация запросов.** Получив для каждого запроса вектор в многомерном пространстве, мы свели задачу кластеризации запросов к задаче кластеризации этих векторов. Существует множество различных функций сходства и алгоритмов кластеризации, вопрос выбора подлежит дальнейшим исследованиям. В данном варианте алгоритма использовалась функция сходства *cosine-similarity* и алгоритм кластеризации *CompleteLink* [2].

$$sim(\bar{v}, \bar{w}) = \frac{\sum v_i w_i}{\sqrt{\sum v_i^2} \sqrt{\sum w_i^2}}$$

Кратко опишем алгоритм кластеризации. Изначально каждый запрос находится в собственном кластере. На каждом шаге перебираем все возможные пары кластеров и объединяем в один ту пару, для которой достигается максимальное значение *completeLink*:

$$completeLink(C_i, C_j) = \min_{v \in C_i, w \in C_j} sim(v, w)$$

Остановка алгоритма может производиться по двум вариантам условия. Во-первых,

можно задать необходимое количество кластеров, однако данный вариант не подходит ввиду того, что заранее это число неизвестно. Во-вторых, можно остановиться, если максимальное значение *completeLink* стало меньше некоторого порогового. Изначально был выбран этот способ со значением порога 0.05. Однако в дальнейшем было замечено, что остаются необъединённые кластеры с одинаковым интентом, в следствие чего порог был уменьшен до 0.01. Более подробное исследование содержимого кластеров показало, что по крайней мере на эталонном множестве запросов дальнейшее снижение порога не приводит к ошибочному объединению кластеров с различными интентами.

## Распределение сессий по кластерам

На выходе этапа кластеризации мы имеем множество кластеров  $C_j$ , каждому из которых соответствует множество запросов  $\{q_{ji}\}$ , в него входящих, и множество документов  $\{d_k\}$  с весами  $u_{jk}$ . Вес документа определяется как сумма предельных вероятностей клика по нему для каждого запроса из кластера:  $u_{jk} = \sum_{q \in C_j} V(q)_k$

Напомним, что сессия пользователя в нашем рассмотрении состоит из последовательности действий: запросов и кликов. В ходе распределения сессий по кластерам последовательно обрабатываются сессии из выборки. Для каждого действия пытаемся получить вектор размерности равной количеству кластеров ( $K$ ). Координаты этого вектора определяются по следующим правилам:

- 1) Для запроса (из списка связанных) полученный вектор содержит 1 на той позиции, номер которой совпадает с номером кластера, в котором данный запрос содержится, и 0 на всех остальных позициях.
- 2) Для клика по документу (из рассматриваемого набора) вектор состоит из его весов  $u_{jk}$  на  $j$ -ой позиции, соответствующей  $j$ -ому кластеру.
- 3) Все остальные действия (клики на документы не принадлежащие рассматриваемому набору и запросы с текстом не принадлежащие списку связанных) не рассматриваются и для них вектор не создаётся.

Для каждой сессии значения векторов, соответствующих действиям внутри неё, усредняются. Таким образом для каждой сессии мы получаем оценки вероятности того, что пользователь её совершивший придерживался интента, соответствующего кластеру. После чего усредняются результаты по всем сессиям. Сессии, в которых не содержалось ни одного действия, соответствующего пунктам 1 или 2, в усреднении не участвуют.

Полученный вектор согласно общей схеме алгоритма, описанной выше, содержит искомые оценки весов интентов.

## Веб-инструмент

Для удобства проведения экспериментов, а также для дальнейшего использования аналитиками, на языке Java был разработан инструмент с веб-интерфейсом, доступный во внутренней сети компании «Яндекс».

Инструмент включает в себя форму для запуска расчёта оценок весов для запроса с указанием региона, при желании можно указать дополнительные настройки: способ построения набора связанных запросов, вероятностную модель, алгоритм кластеризации и его параметры (см. рис. 1).

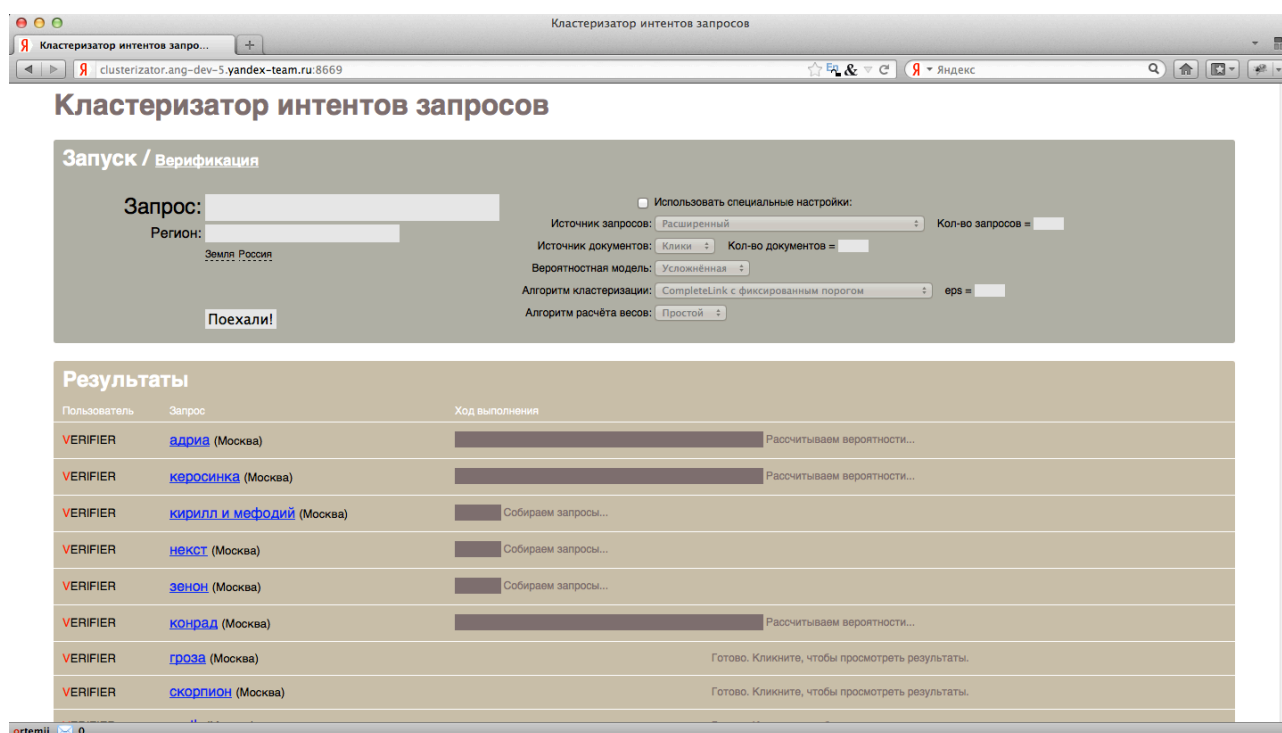


Рис. 1. Интерфейс веб-инструмента «Кластеризатор интенгов запросов»

Внутри инструмента производится запуск описанного выше алгоритма расчёта весов, параллельно обрабатывается несколько запросов. Необходимые данные пользовательской статистики и сессии загружаются непосредственно с кластера MapReduce.

По всем когда-либо запущенным запросам можно посмотреть параметры запуска, результаты, сравнить их со значениями ассессоров, если таковые имеются (см. рис. 2).

[СОЛНЦЕСТОЯНИЕ](#) (Москва)

**Параметры запуска**

eps: 0.990000  
Кол-во шагов цепи Маркова: 20  
Кол-во запросов: 20  
Кол-во документов: 100  
Источник запросов: Расширенный  
Источник документов: Клики  
Вероятностная модель: Усложнённая  
Алгоритм кластеризации: CompleteLink с фиксированным порогом  
Алгоритм расчёта весов: Простой

**Информация о результате**

Начало: 2012-05-08 19:26  
Конец: 2012-05-08 19:26  
Время: 10 сек  
[Сессий](#) обработано: 52, подошло: 43 (83%)  
[Схожесть](#) запросов (cosine similarity)  
[Схожесть](#) кластеров (complete link)

0.484718	<a href="#">22 декабря день зимнего солнцестояния, день зимнего солнцестояния, зимнее солнцестояние, солнцестояние 2012, солнцестояние зимнее</a>
0.243669	<a href="#">солнцестояние онлайн, солнцестояние фильм, смотреть солнцестояние, смотреть фильм солнцестояние, фильм солнцестояние смотреть онлайн, солнцестояние смотреть онлайн бесплатно в хорошем качестве, фильм солнцестояние, смотреть онлайн солнцестояние, солнцестояние смотреть онлайн</a>
0.150166	<a href="#">день весеннего равноденствия, равноденствие, что такое равноденствие, равноденствия, весеннее равноденствие, осеннее равноденствие, день весеннего равнодействия</a>
0.081912	<a href="#">солнцестояние трейлер</a>
0.000000	<a href="#">солнцестояние скачать торрент</a>
0.000000	<a href="#">колдоврат</a>
0.000000	<a href="#">долгота дня</a>
0.039535	other

0.660000 природное явление

0.320000 фильм 2008 года

0.020000 other

Результаты ассессоров по запросу [солнцестояние](#) (Москва)

[Свернуть](#)

*Рис. 2. Пример отображения результатов запуска*

Инструмент позволяет посмотреть случайные сессии по которым рассчитывались веса, а также матрицу со значениями функций сходства запросов и полученных кластеров.

# Результаты работы алгоритма

## Анализ полноты и точности

При постановке задачи упоминалось о двух критериях, по которым будут проверяться результаты работы алгоритма.

Первый критерий – это полнота. Для каждого запроса из эталонного множества вручную производилась проверка, для какого количества интенгов из известного набора присутствуют кластеры в наборе, полученном автоматически. В результате для 65% эталонных запросов был получен полный набор интенгов, для 94% недоставало максимум одного.

Второй критерий – точность. Оценивалось максимальное различие между посчитанными ассессорами весами и весами, полученными в результате работы алгоритма. В случае, если один интенг из «ручного» набора был представлен несколькими в «автоматическом» (к примеру, интенг «богиня» для запроса [афина] алгоритмом разделяется на более мелкие «богиня», «священное животное Афины» и «другие греческие богини»), значения весов для разделённых наборов суммировались.

Проверка точности вызвала трудности, связанные с тем, что данные о весах, полученные ассессорами, устарели по сравнению с тем периодом пользовательской статистики, по которому производились расчёты. В ближайшее время планируется организовать такую проверку.

Однако, если использовать устаревшие данные о весах, максимальное различие, усреднённое по всем запросам с полным набором интенгов, составляет 0.17.

## Примеры обработанных результатов

В таблице 3 приведены примеры обработанных результатов работы алгоритма по нескольким запросам в сравнении с результатами, полученными после работы ассессоров. Обработка заключалась в сопоставлении кластеру интента запроса.

Запрос Регион	Алгоритм		Ассессоры	
	Вес	Интент	Вес	Интент
Солнцестояние (Москва)	0.48	Солнцестояние	0.66	Природное явление
	0.27	Фильм	0.32	Фильм
	0.10	Смотреть фильм онлайн	0.02	Другое
	0.09	Равноденствие		
	0.07	Другое		
Афина (Москва)	0.43	Богиня Афина	0.69	Богиня
	0.20	Священное животное Афины	0.07	Певица
	0.20	Другие греческие богини	0.07	Салон красоты
	0.10	Певица	0.17	Другое
	0.04	Салон красоты		
	0.03	Другое		
Понедельник (Москва)	0.19	Светотехническая компания	0.81	День недели
	0.10	Афоризмы про понедельник	0.17	Светотехническая компания
	0.09	Картинки про понедельник	0.02	Другое
	0.04	Виктор Понедельник		
	0.04	Понедельник по-английски		
	0.03	Понедельник начинается в субботу		
	0.53	Другое		
Наполеон (Земля)	0.73	Наполеон I Бонапарт	0.85	Наполеон I Бонапарт
	0.14	Торт	0.05	Игра "Napoleon: Total War"
	0.05	Фильм 2002 года	0.03	Торт
	0.03	Игра "Napoleon: Total War"	0.07	Другое
	0.01	Фильм 2011 года		
	0.04	Другое		
Вертикаль	0.57	Группа компаний	0.58	Группа компаний
	0.07	Понятие	0.08	Центр образования
	0.06	Магазин	0.03	Фильм
	0.05	Программа	0.01	Фестиваль туристических фильмов
	0.04	Роллы в Екатеринбурге	0.30	Другое
	0.02	Турфирма		
	0.02	Смотреть фильм онлайн		
	0.07	Другое		

Табл. 3. Примеры результатов работы алгоритма



## Заключение

Был построен алгоритм, позволяющий получать наборы интенгов для многозначных запросов и оценки их весов и показывающий лучшую полноту по сравнению с результатами ассессоров на эталонном множестве запросов.

В завершении работы по проверке точности полученные данные будут использоваться для расчёта метрик в службе оценки качества поиска компании «Яндекс», обновления наборов запросов, по которым производится оценка качества поиска.

В ходе работы кроме самого алгоритма:

- 1) Был создан инструмент на языке Java, имеющий веб-интерфейс, удобный для использования аналитиками.
- 2) Написана утилита на языке C++ для расчёта статистики по переформулировкам запросов на кластере MapReduce.
- 3) Написаны скрипты для сбора по поисковым логам сессий пользователей, используемых в расчётах.

В дальнейшем предполагается работа в следующих направлениях:

- 1) Создание полуавтоматической системы проверки точности и полноты;
- 2) Фильтрация не относящихся к данному интенгу данных;
- 3) Улучшение матчинга сессий с использованием лингвистических данных;
- 4) Исследования по определению интенга из коротких, малоинформативных сессий;
- 5) Другие алгоритмы кластеризации и функции сходства.

## Список литературы

- [1] Оптимизация алгоритмов ранжирования методами машинного обучения, А. Гулин, П. Карпович, Д. Расковалов, И. Сегалович // РОМИП'09.
- [2] Clustering Query Refinements by User Intent, E. Sadikov, J. Madhavan, Lu Wang, A. Halevy // WWW'10.
- [3] Clustering User Queries of a Search Engine, Ji-Rong Wen, Jian-Yun Nie, Hong-Jiang Zhang // WWW'01.
- [4] Diversifying Search Results, R. Agrawal, S. Gollapudi, A. Halverson, S. Jeong // WSDM'09.
- [5] Inferring Query Intent from Reformulations and Clicks, F. Radlinski, M. Szummer, N. Craswell // WWW'10.
- [6] Query Similarity by Projecting the Query-Flow Graph, I. Bordino, C. Castillo, D. Donato, A. Gionis // SIGIR'10.
- [7] Random Walks on the Click Graph, N. Craswell, M. Szummer // SIGIR'07.
- [8] The Query-flow Graph: Model and Applications, P. Boldi, F. Bonchi, C. Castillo, D. Donato, A. Gionis, S. Vigna // SIKM'08.