

МГУ им. М.В.Ломоносова Научно-исследовательский вычислительный центр





Университетская информационная система РОССИЯ

Оптимизация параметров алгоритма поиска на основе анализа оценок экспертов

Михаил Агеев, Борис Добров

RCDL'2005, Российский семинар по оценке методов информационного поиска Ярославль, 2005

УИС РОССИЯ (www.cir.ru)

Избранное

Сервис Справка

Ресурсы Университетской информационной системы РОССИЯ

<u>Интегрированная коллекция | Бюджетная система России | Статистика России | Выборы в России | Парламент России |</u>













Новые ресурсы

(сентябрь 2003 г.)

12.11.2003

11.11.2003





Университетская информационная система РОССИЯ

О проекте

Университетская информационная система РОССИЯ (УИС РОССИЯ) создана и поддерживается как база электронных ресурсов для исследований и образования в области экономики, социологии, политологии, международных отношений и других гуманитарных наук и с 2000 года открыта для коллективного доступа университетов, вузов, научных институтов РФ и специалистов.

подробнее...

Полный список коллекций | Академический сервис | Партнеры | Участники | Зеркала | Как к нам пройти

Поиск по ресурсам УИС РОССИЯ

Поиск по источникам: все коллекции (Уровень доступа =

искать 🕨

Расширенный поиск

Изменить уровень доступа (для зарегистрированных Пароль:

пользователей)

Имя:

Об уровнях доступа | Зарегистрироваться | Забыли пароль? | Справка/Практикум

04.11.2003

Краткосрочные экономические показатели Российской Федерации. (август 2003 г.)

Социально-экономическое

положение России (сентябрь 2003 г.)

Краткосрочные экономические

показатели Российской Федерации.

Социально-экономическое положение России. (август 2003 г.)

30.09.2003

Социально-экономическое положение России (июль 2003 г.)

Краткосрочные экономические:

показатели Российской Федерации. (июль 2003 г.)

Соционет/Repec | Ресурсы зарубежных организаций

Бюджетная система России Информационно-аналитический комплекс для изучения бюджетной системы РФ. Формируется открытых первоисточников, предоставляемых органами государственной власти, научными институтами, аналитическими центрами. Содержит бюджетную статистику с 1995 года. Включает тематические публикации научных журналов, центральных СМИ. Представлены материалы профильных учебных курсов Экономического факультета МГУ им. М.В. Ломоносова.

Интегрированная коллекция

Полный список

Документы государственных органов | Издания исследовательских центров | Научные издания | Коллекции зарубежных организаций | Социологические опросы

Статистика России 🗁

(для зарегистрированных пользователей) Интегрированная коллекция статистических и аналитических материалов характеризует социально-экономическое развитие Российской Федерации и регионов в ретроспективе с 1996 года, Формируется на базе первоисточников публикаций Госкомстата России, Минэкономразвития, других государственных организаций, а также изданий независимых аналитических центров.

Реляционная база данных

Соционет / RePEc

Research Papers Economics библиотека библиографических описаний информационных ресурсов, создаваемых специалистами по общественным наукам (экономике, социологии, политологии) во всем мире. Включает архивы электронных публикаций, оглавления онлайновых журналов, каталоги новых поступлений библиотек, планы издательств и др. - более 200 тысяч записей, СОЦИОНЕТ - библиографическое описание материалов по общественным наукам на русском языке.

Выборы в России

Интерактивная карта выборов | Выборы Президента РФ (1996, 2000 год) | Выборы в ГосДуму РФ (1993, 1995, 1999 годы) Выборы в субъектах РФ Административнотерриториальное деление РФ (СОАТО)

Парламент России

ГосДума РФ, Стенограммы пленарных заседаний | 🕼 Гослума РФ. Информационно-Аналитический бюллетень | 🜆

Ресурсы зарубежных организаций

Организация экономического сотрудничества и развития (OECD) - OECD Health Data, 2002, 30 countries,

База данных по системам здравоохранения в 30 странах. Организацией экономического сотрудничества и развития. Включает 1200 показателей. По некоторым показателям временные ряды прослеживаются с 1960 года.

Перепись населения в СССР в 1939 году - предоставлены Университетом Торонто, Канада.

Новости

 10.02.2003 10-14 февраля 2003 года в НИВЦ МГУ проведен российскоамериканский семинар "Новые технологии в поддержку государственного управления. Роль университетов в создании национальной информационной инфраструктуры", Программа семинара прилагается.







Информационные ресурсы УИС РОССИЯ 750.000 документов (около 8 Гб)

	Источник	Период	Количество
Нормативные д-ты	НТЦ Система	1993	70,000
Межд. договоры	НТЦ Система	1993	2,300
Стенограммы ГД	Аппарат ГД ФС РФ	1994	144,000
Статистические данные	Госкомстат РФ, Межгос. Стат. Ком. СНГ	1998	54,000
Материалы СМИ	«АиФ» - «Эксперт»,	199(7)	391,000
Аналитические материалы	Минэкономразвития, ЦБ, РЕЦЭП, БЭА, ИПН	1996	22,000
Научные издания	Вестник МГУ, «Соц.исследования»	Экономика 1998	5,300
Социологические опросы	Линейные распределения Мониторинга ВЦИОМ	2000	240

УИС РОССИЯ В РОМИП

2003

- adhoc, коллекция web, TF*IDF
- +эксперимент по учету точной формы слова

2004

- adhoc, коллекции web, legal, TF*IDF
- +эксперимент по учету расстояния между словами
- **■** +эксперимент по повышению веса заголовков
- text categorization, legal

2005

- adhoc, коллекции web, legal, TF*IDF
- +эксперимент по учету расстояния между словами и «мягкий» поиск
- +оптимизация параметров алгоритма на основе оценок РОМИП'2004

Компоненты известных алгоритмов поиска и ранжирования документов

- поиск документов, содержащих все слова запроса
- расширение запроса с учетом морфологической вариации слов
- исключение незначимых слов (стоп-слов) из запроса
- поиск документов, частично удовлетворяющих запросу (т.е. содержащих не все слова запроса)
- расширение запроса пользователя тематически близкими словами/терминами
- учет относительной частоты встречаемости слов запроса в найденном документе (чем чаще слово встречается - тем лучше)
- учет относительной частоты встречаемости слов запроса в документах коллекции (более редкие слова имеют больший вес)
- учет формы слова, которое встретилось в документе (лучше, если слово встретилось в документе в той же форме, что и в запросе)
- учет расстояния между словами запроса в документе (лучше, если слова запроса находятся рядом)
- учет структуры документа (слова в заголовке имеют больший вес)
- учет ссылок между документами (на авторитетный документ часто ссылаются)

«Отправная точка» (TF*IDF)

Вес каждой леммы документа:

$$TFIDF_D(l) = \beta + (1 - \beta) \cdot tf_D(l) \cdot idf_D(l)$$

где "term frequency" – учет частотности леммы в документе:

$$tf_D(l) = \frac{\text{freq}_D(l)}{\text{freq}_D(l) + 0.5 + 1.5 \cdot \frac{dl_D}{\text{avg_dl}}}$$

 $freq_D(l)$ - частотность леммы l в документе, dl_D – мера длины документа, avg_dl – средняя длина документа, $\beta=0.4$

$$id f(l) = \frac{\log \left(\frac{|c| + 0.5}{df(l)}\right)}{\log (|c| + 1)}$$

«Отправная точка» (TF*IDF)

Каждый запрос $Q = w_1 \ w_2 \ w_3 \dots \ w_m$ представлялся в виде формулы

$$L(Q) = L(w_1) \& L(w_2) \& L(w_3) \& ... \& L(w_m)$$
,

 $r\partial e \; L(w) = l_1(w) \; OR \; l_2(w) \; OR \; ...OR \; l_q(w), \; l_k(*) \;$ - леммы морфологического разбора слова.

Тогда оценка релевантности документа D для запроса Q вычисляется по формуле:

$$V_D(Q) = \frac{\sum_{i=1}^{N} \sum_{k} (\theta_{ik} \cdot \text{TFIDF}_D(l_{ik}(\mathbf{w}_i)))}{\sum_{i=1}^{N} \sum_{k} |\theta_{ik}|}$$

где $\theta_{ik} = \theta_i = 1.0$ — "вес" леммы в запросе — равен весу, устанавливаемому для соответствующего слова запроса.

Учет дополнительных факторов

$$\operatorname{Rank}_{D}(Q) = \frac{\operatorname{FF}-1}{|Q|} + \frac{1}{|Q|} \cdot \frac{\operatorname{V}_{D}(Q) + \operatorname{Near}_{D}(Q)}{2}$$

- количество слов запроса
- Rank_D(Q) ранг соответствия документа запросу
- FF количество слов запроса, которые содержатся в данном документе

V_D(Q) — ранг, вычисленный по формуле TF*IDF

Near_D
$$(Q) = \frac{1}{\ln(\lambda_D(Q) - |Q| + 4)}$$

$$\lambda_D\left(Q\right)$$
 — длина минимального «куска» документа, в котором содержатся все слова запроса

Оптимизация алгоритма поиска

- 1. Поиск с использованием классической векторной модели TF*IDF
- 2. Расширение алгоритма TF*IDF
 - ранжирование с учетом расстояния между словами
 - поиск документов, содержащих не все слова запроса
- 3. Оптимизация параметров, входящих в предложенный алгоритм:

$$\operatorname{Rank}_{D}(Q) = \operatorname{TFIDF}_{D}(Q) \cdot \frac{\operatorname{FF}}{|Q|} + \beta \cdot \operatorname{Near}_{D}(Q) \cdot 2^{-\alpha \cdot (|Q| - \operatorname{FF})}$$

- β относительный вес кучности слов запроса
- α зависимость веса кучности от количества найденных слов
- 4. Оценка результатов на материалах РОМИП

Оптимизация параметров: перебор на 4 классах функций

Оптимизация параметров производилась для следующих классов функций ранжирования:

1)
$$\operatorname{Rank}_{D}(Q) = \frac{\operatorname{FF}-1}{|Q|} + \frac{1}{|Q|} \cdot \frac{\operatorname{V}_{D}(Q) + \beta \cdot \operatorname{Near}_{D}(Q)}{1 + \beta}$$

— однопараметрическое семейство функций с параметром $\beta > 0$, задающим относительный вес кучности слов запроса;

2)
$$\operatorname{Rank}_{D}(Q) = \frac{\operatorname{FF}-1}{|Q|} + \frac{1}{|Q|} \cdot \frac{1+\beta}{\operatorname{V}_{D}(Q)} + \frac{\beta}{\operatorname{Near}_{D}(Q)}$$

— однопараметрическое семейство функций с параметром $\beta > 0$, задающим относительный вес кучности слов запроса;

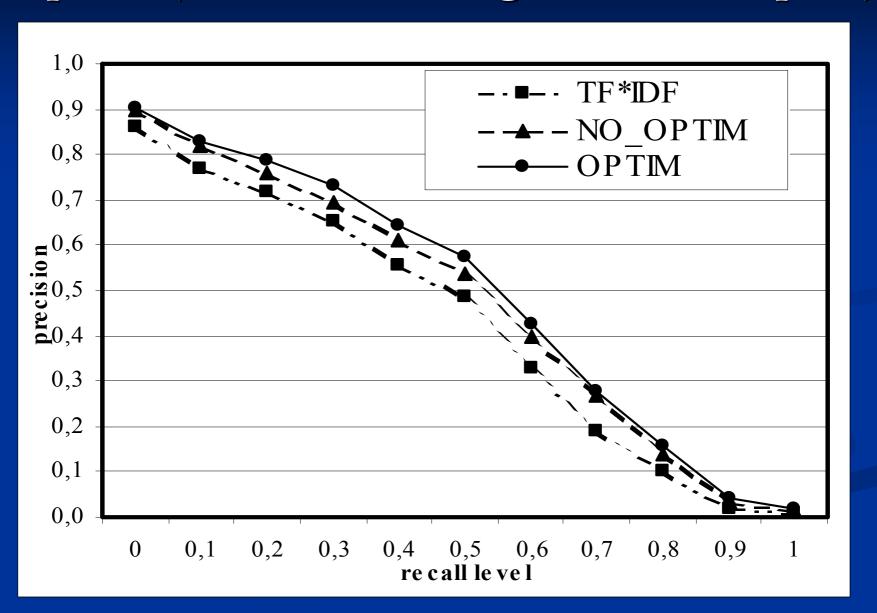
3) Rank_D
$$(Q) = V_D(Q) \cdot \frac{FF}{|Q|} + \beta \cdot Near_D(Q) \cdot 2^{-\alpha \cdot (|Q| - FF)}$$

— двухпараметрическое семейство функций с параметрами $\beta > 0$ (относительный вес кучности слов запроса) и $\alpha > 0$ (зависимость веса кучности от количества найденных слов);

4) Rank_D(Q) =
$$\frac{V_D(Q) + \beta \cdot \frac{1}{\ln^{\gamma} (\lambda_D(Q) - |Q| + 4)} + \alpha \cdot FF}{1 + \beta + \alpha \cdot |Q|}$$

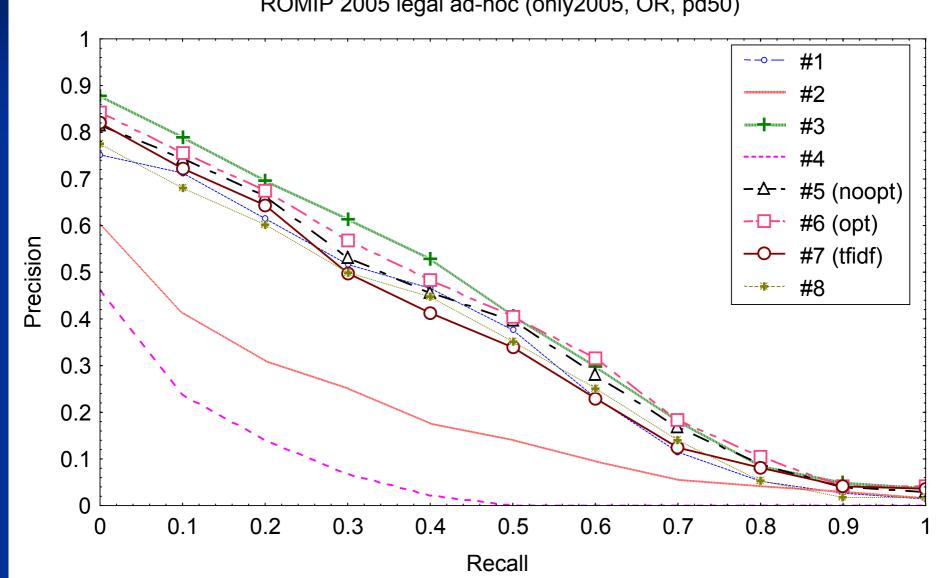
— трёхпараметрическое семейство функций с параметрами $\beta > 0$ (относительный вес кучности слов запроса), $0 < \alpha \le 1$ (вес количества найденных слов запроса) и $\gamma > 0$ (жесткость функции ранжирования по кучности слов).

Результаты: оценка на «тренировочных» запросах (РОМИП'2004, legal-adhoc, OR-pd50)

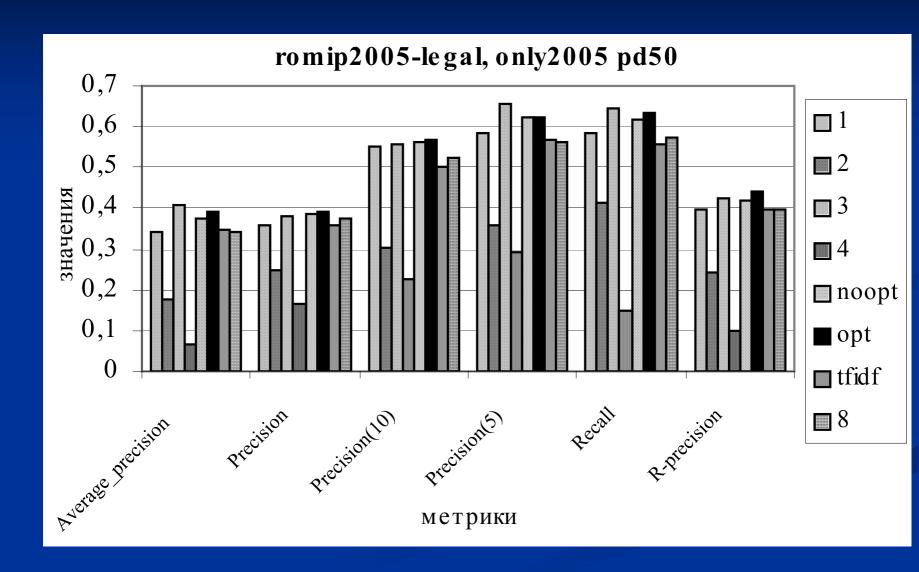


Результаты: оценка на новых запросах (POMI/Π'2005, legal-adhoc, OR-pd50)

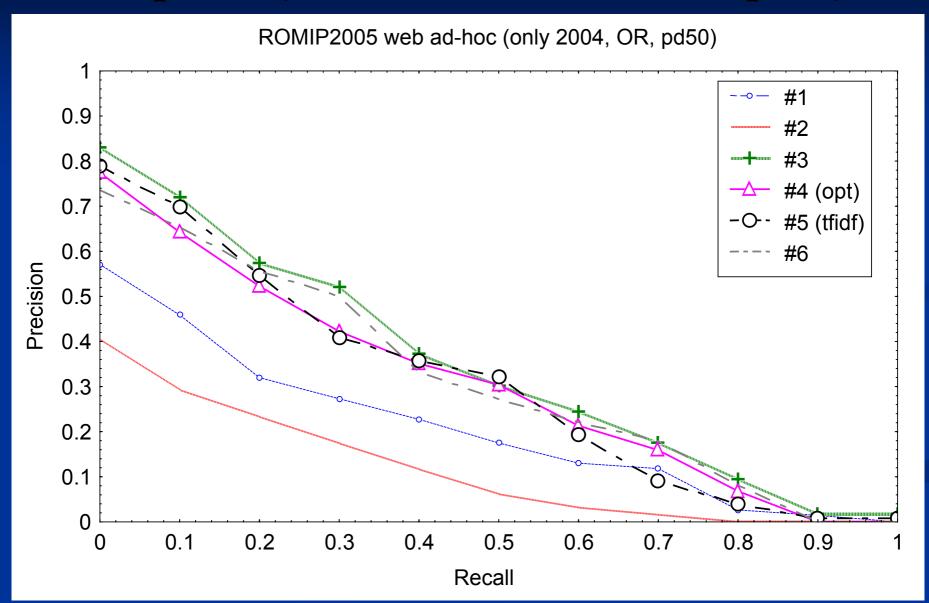
ROMIP 2005 legal ad-hoc (only2005, OR, pd50)



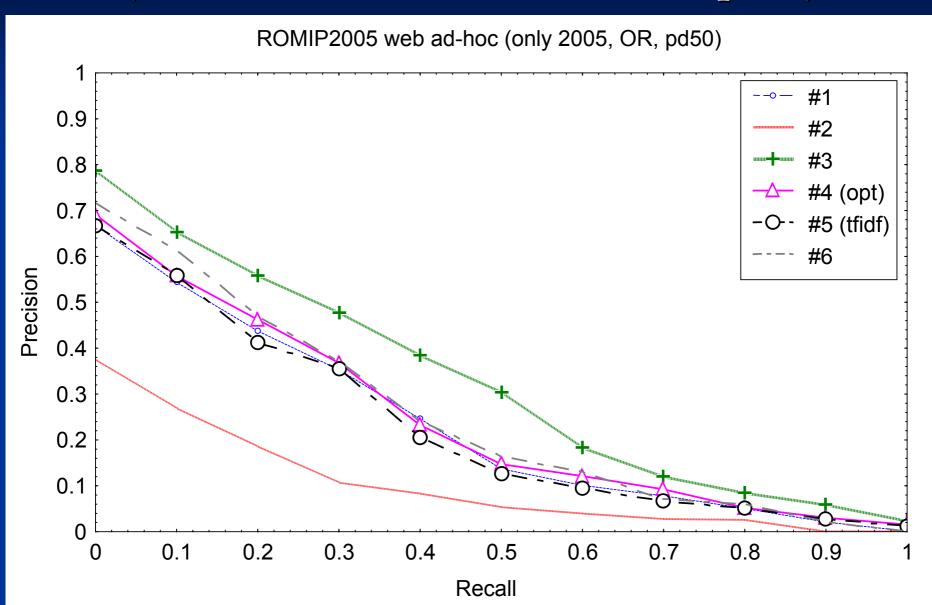
Результаты: оценка на новых запросах (РОМИП'2005, legal-adhoc, OR-pd50)



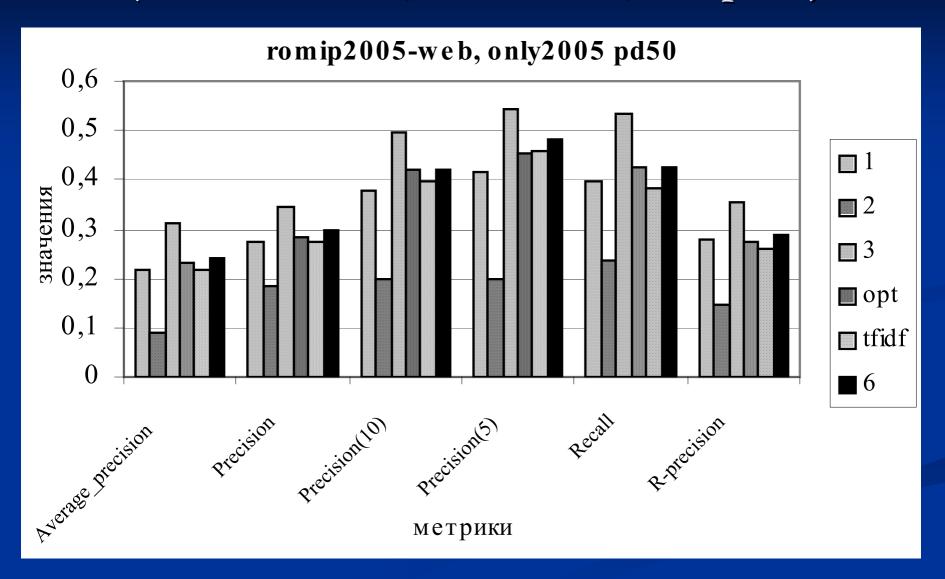
Результаты: оценка на тренировочных запросах (РОМИП'2004, web, OR-pd50)



Результаты: оценка на новых запросах (РОМИП'2005, web-adhoc, OR-pd50)



Результаты: оценка на новых запросах (РОМИП'2005, web-adhoc, OR-pd50)



Заключение

Предложенный метод автоматической настройки параметров алгоритма поиска может быть использован для повышения эффективности поисковых машин

■ В дальнейшем мы планируем развивать данный подход в направлении расширения набора параметров, влияющих на работу поискового алгоритма