



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



АНО Центр
информационных
исследований



Университетская информационная система
РОССИЯ

УИС РОССИЯ:

ПОИСК И КЛАССИФИКАЦИЯ В

РОМИТТ 2007

Агеев М.С., Добров Б.В., Красильников П.В.,
Лукашевич Н.В., Сидоров А.В., Штернов С.В.

Ad-hoc поиск документов

- ❖ Поиск по коллекции нормативных документов “Legal2007”
 - ❖ 348410 документов из БД Кодекс
 - ❖ 14797 запросов из лога запросов к правовому разделу портала www.kodeks.ru
- ❖ Поиск по коллекции «Беларусский интернет» “BY.web 2007”
 - ❖ 1525585 документов (выборка из страниц домена .by из индекса yandex.ru)
 - ❖ 19627 запросов (смесь запросов из лога yandex.ru , km.ru, оцененных запросов РОМИП 2003-2006)

УИС РОССИЯ в РОМИП 2003-2005...2007

❖ Технологии:

- ❖ 2003-2004: таблицы Oracle, bag-of-words (*)
- ❖ 2005: (*) + sorted lists: экспериментальные прогоны
- ❖ 2007: новый индексатор и поисковик (sorted lists)

❖ Эксперименты:

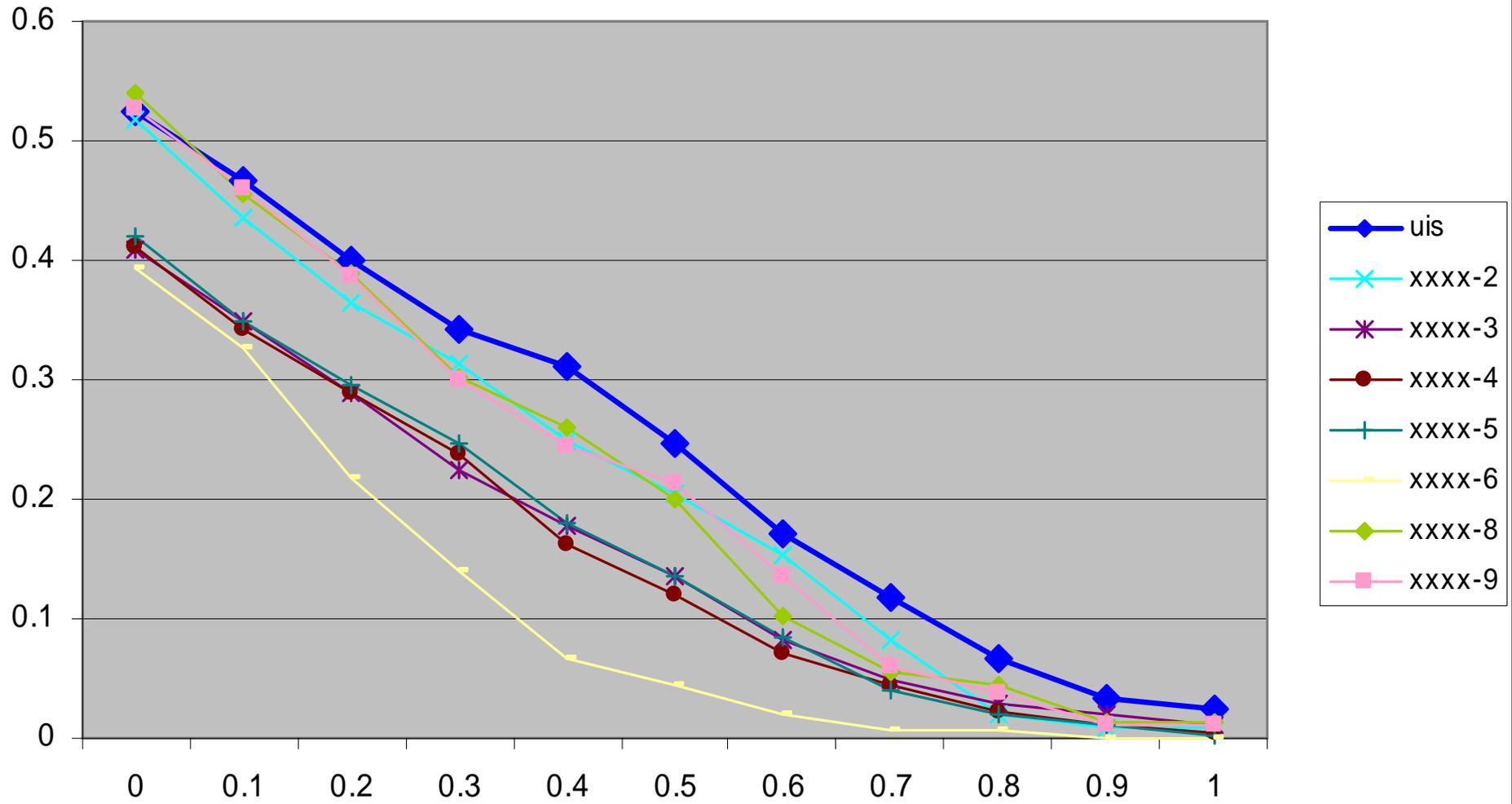
❖ Учет факторов:

- ❖ $TF*IDF$: (2003-2005)
- ❖ минимальное окно: 2005
- ❖ кворум: 2005
- ❖ *близость по парам: 2007*

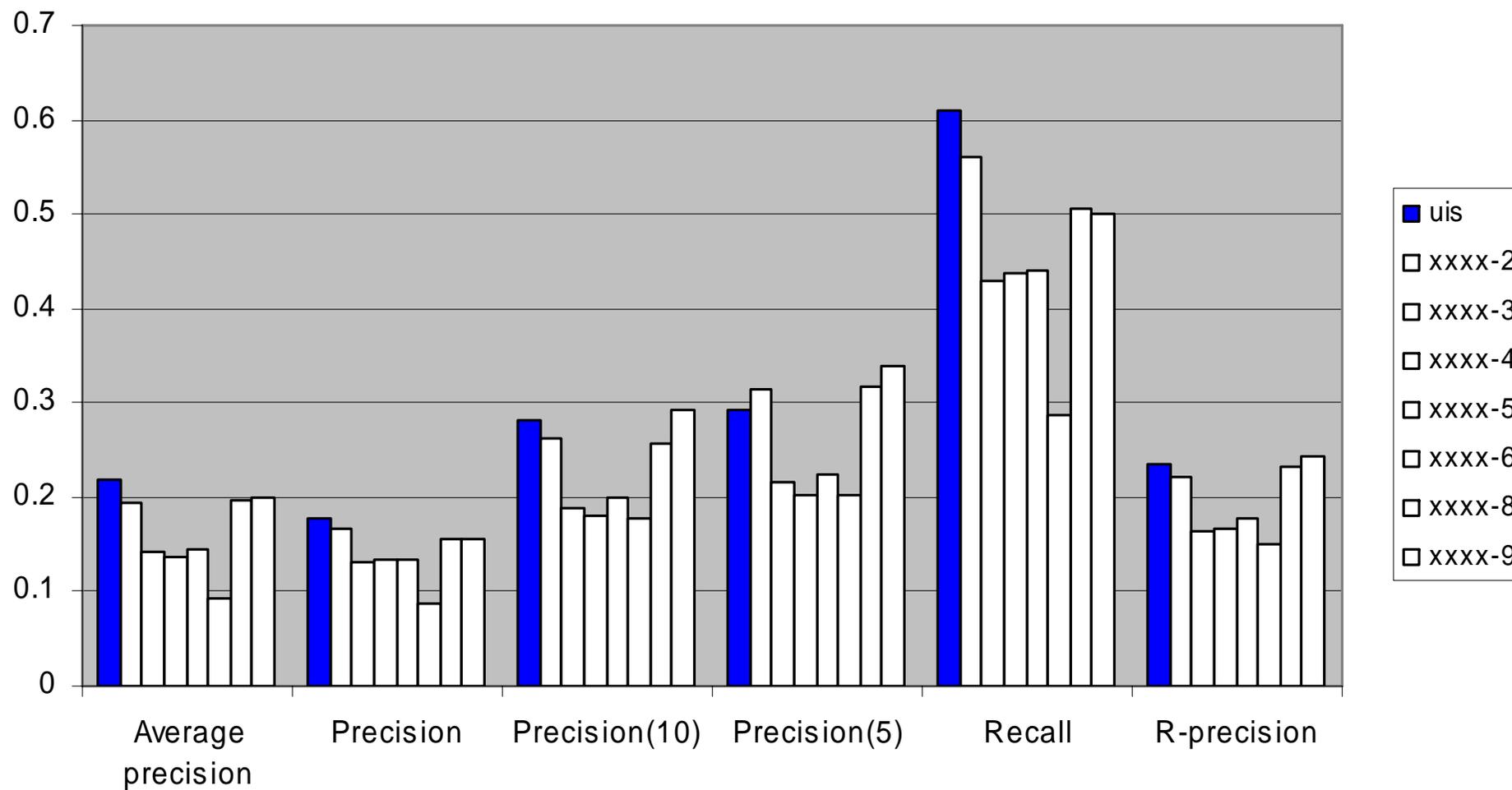
❖ Подбор коэффициентов:

- ❖ *Красильников П. В. «Воспроизведение лучших результатов ad hoc поиска семинара РОМИП»*

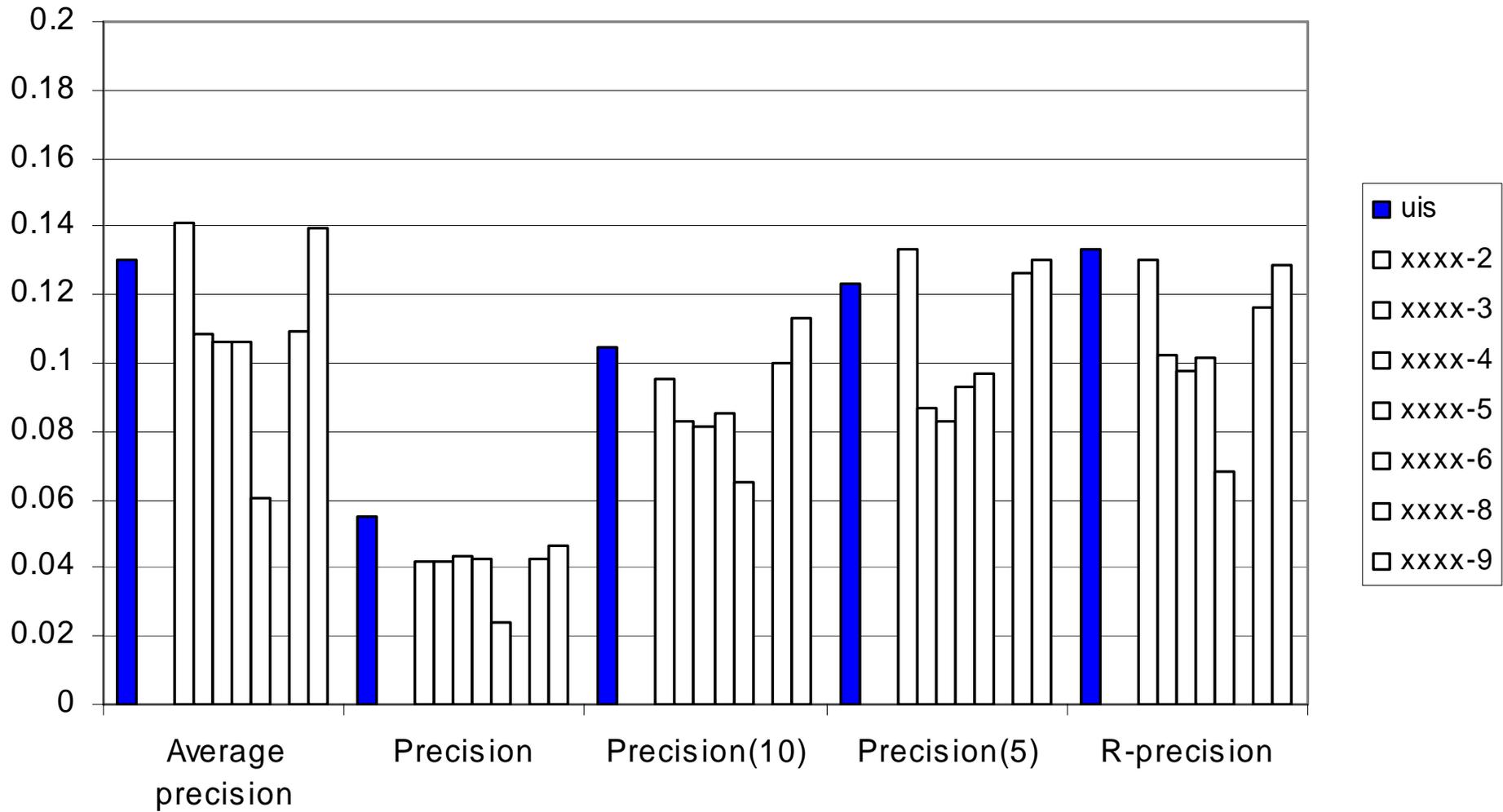
ROMIP 2007, BY.web, or_pd50



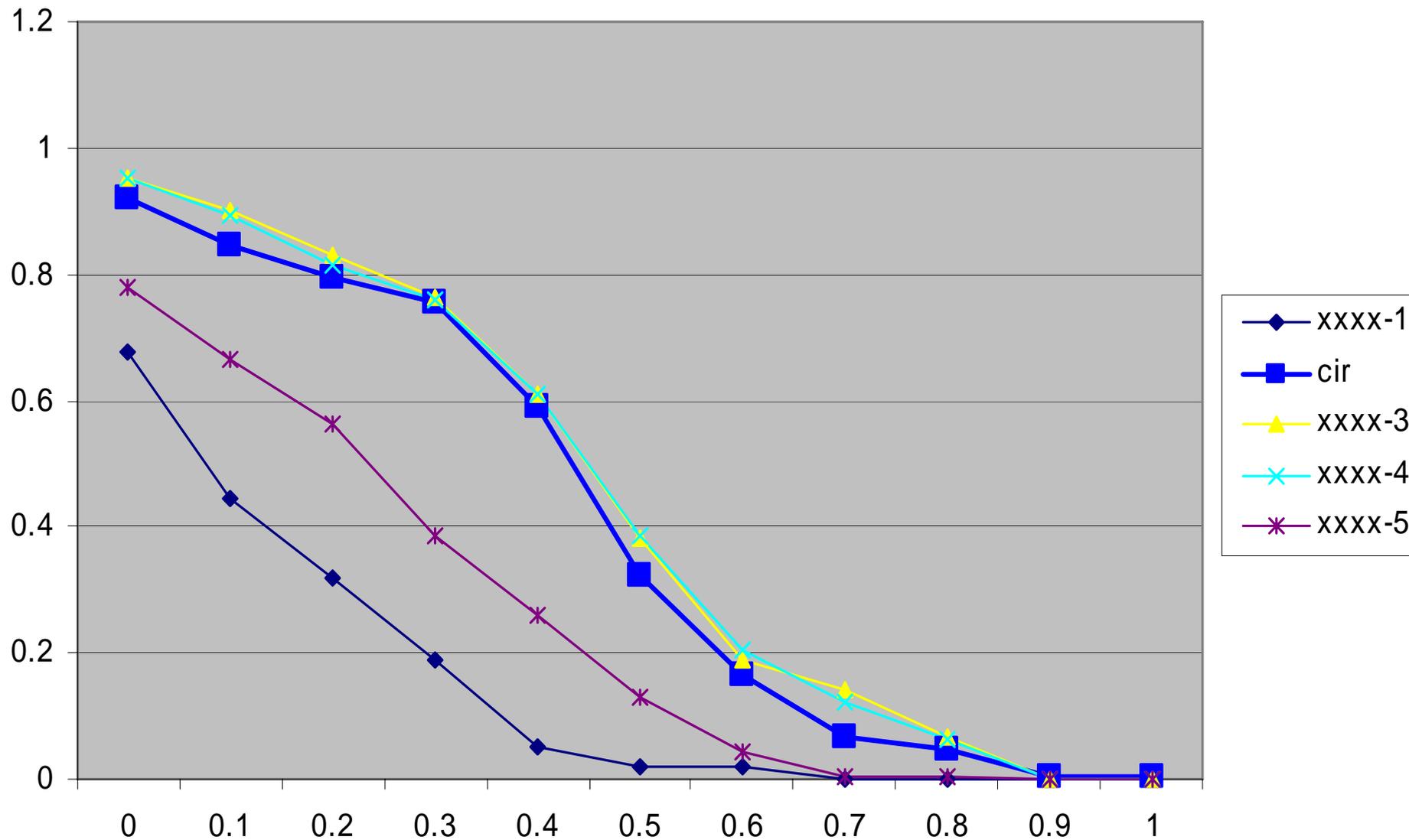
ROMIP 2007, BY.web, or_pd50



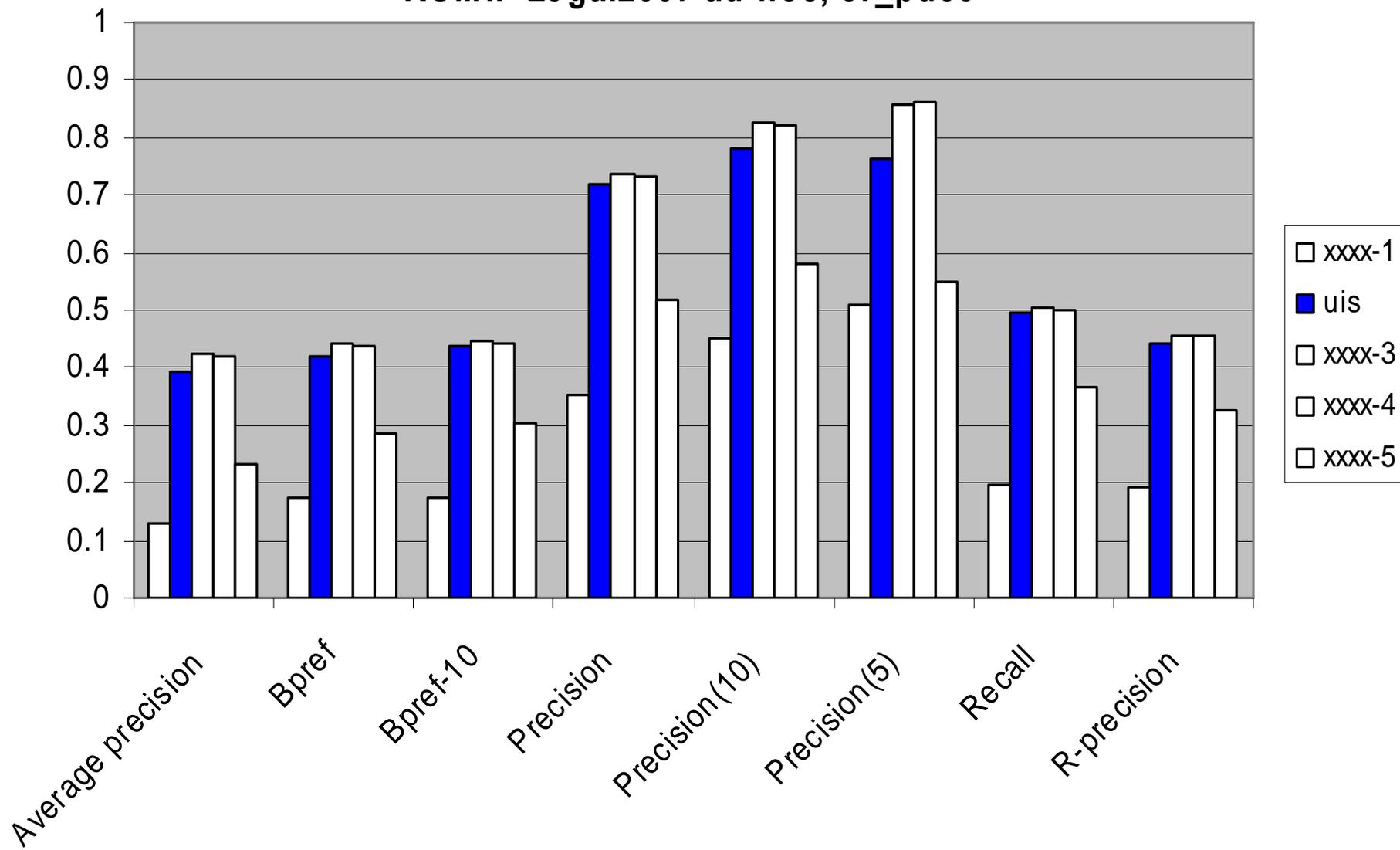
ROMIP 2007, BY.web, and_pd50

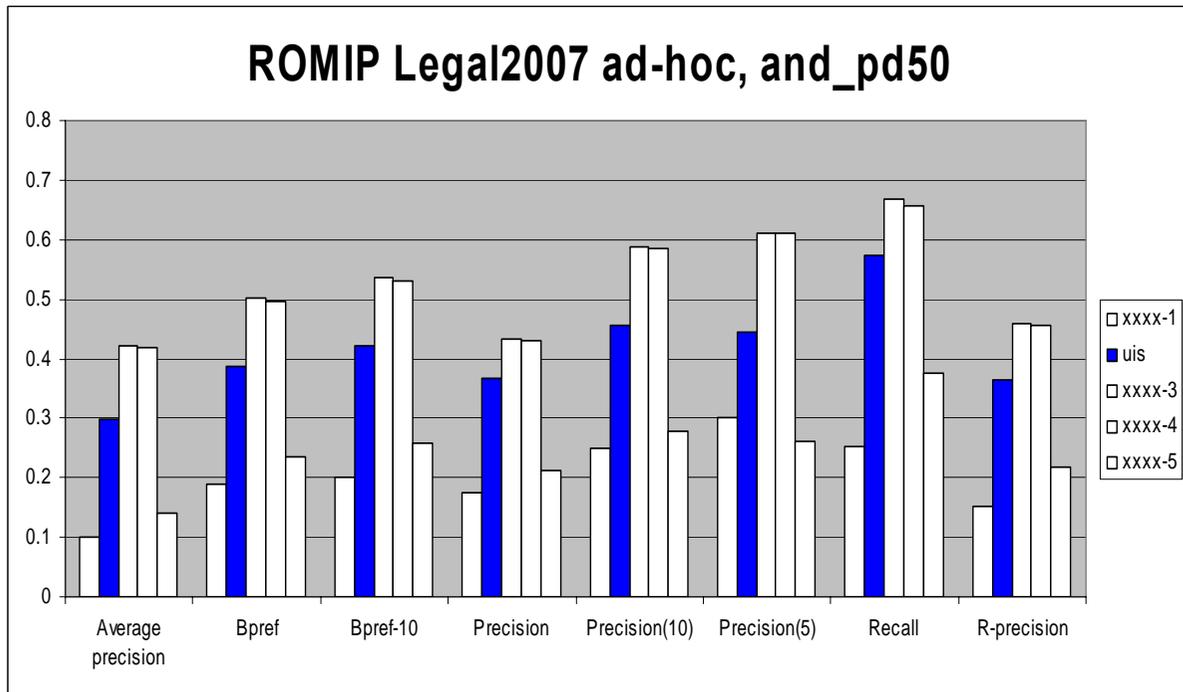
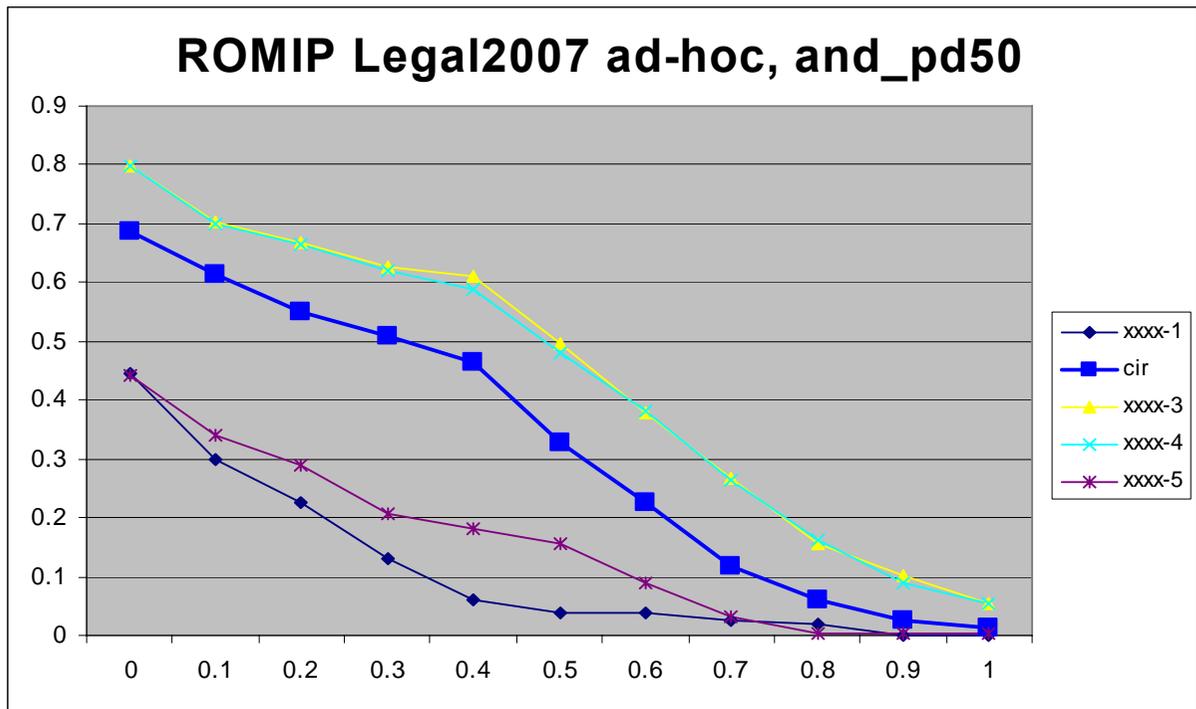


ROMIP Legal2007 ad-hoc, or_pd50



ROMIP Legal2007 ad-hoc, or_pd50





Анализ ad-hoc

- ❖ **Лучшие результаты получены по метрике Average Precision, or_pd50**
 - ❖ По этой же метрике и обучали алгоритм
- ❖ **Результаты по P(10) и/или and_pd50 заметно хуже**
 - ❖ Обучаться по этим метрикам?
 - ❖ Проанализировать найденные документы!
 - ❖ Можно ли улучшить одновременно по всем метрикам?
- ❖ **Замешивание оценок по КМ и смешанной коллекции – как интерпретировать результаты?**

Ad-hoc: выводы

- ❖ **Участие в РОМИП'2007 позволило отладить новые технологические модули и проверить эффективность формулы ранжирования**
- ❖ **Предложенная в статье «Воспроизведение лучших результатов...» формула ранжирования показывает хорошие результаты поиска**
- ❖ **Дальнейшие планы:**
 - ❖ **Учет новых факторов: relevance feedback, структура (заголовки, ссылки)**
 - ❖ **Анализ отдельных запросов – поиск «равномерного» улучшения**