



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.

Классификация веб-страниц:

машинное обучение vs.
работа экспертов

Машинное обучение: метод ПФА

- Рубрика **135 «Боевые искусства»**

- Recall = 0.52 Precision = 0.88 FMeasure = 0.82

- [Тип = в дереве | имя = БОЕВЫЕ ИСКУССТВА]

- Recall = 0.82 Precision = 0.98 FMeasure = 0.96

- ([Тип = лемма | имя = КАРАТЭ])

- OR ({ [Тип = в тексте | имя = ХОККЕЙНЫЙ КЛУБ]

- OR [Тип = в дереве | имя = ОХРАННОЕ ПРЕДПРИЯТИЕ] }

- AND

- [Тип = в дереве | имя = БЕДСТВИЕ])

- OR ({ [Тип = в тексте | имя = КУЛЬТУРА]

- OR [Тип = в тексте | имя = СЕВЕРО-ЗАПАДНАЯ ЧАСТЬ] }

- AND

- [Тип = в тексте | имя = ОДЕЖДА]

- AND

- [Тип = в дереве | имя = ВЕРОВАТЬ])

- OR ({ [Тип = в тексте | имя = МЕДИЦИНСКОЕ УЧРЕЖДЕНИЕ]

- OR [Тип = в тексте | имя = КРЫЛАТСКОЕ] }

- AND [Тип = в дереве | имя = ВОСТОЧНЫЕ ЕДИНОБОРСТВА])

- OR ([Тип = в тексте | имя = МАСЛЕНИЦА])

- OR ([Тип = лемма | имя = ДЗЭНИН])

- OR ([Тип = в тексте | имя = САМОБОРОНА]

- AND [Тип = в дереве | имя = ИСТОРИЧЕСКИЕ НАУКИ])

Машинное обучение: метод ПФА (2)

- Рубрика 14 «Страны_и_регионы ->Азия»

- Recall = 0.56 Precision = 0.95 FMeasure = 0.88

- { ([Тип = в дереве | имя = ДУНГАНЕ])
- OR ([Тип = в тексте | имя = СРЕДСТВА МАССОВОЙ ИНФОРМАЦИИ]
- AND [Тип = в дереве | имя = ПОЛИТИЧЕСКАЯ ДЕЯТЕЛЬНОСТЬ]
- AND [Тип = в дереве | имя = СРЕДНЯЯ АЗИЯ]) }

- Рубрика 82 «Досуг -> Туризм»

- Recall = 0.92 Precision = 0.99 FMeasure = 0.98

- { ({ [Тип = в дереве | Имя = ТУРИСТИЧЕСКИЙ СЕРВИС]
- OR [Тип = в дереве | Имя = РЕЗЕРВИРОВАТЬ] }
- AND [Тип = в дереве | Имя = ПАССАЖИРСКИЙ ТРАНСПОРТ]
- AND [Тип = в дереве | Имя = АВИАЦИОННЫЕ РАБОТЫ]
- AND [Тип = в дереве | Имя = РАЗРЕШИТЕЛЬНЫЙ ДОКУМЕНТ])
- OR
- ({ [Тип = в дереве | Имя = ГОСТИНИЦА]
- OR [Тип = в дереве | Имя = РУМЫНИЯ] }
- AND [Тип = в дереве | Имя = ПУТЕШЕСТВИЕ]
- AND [Тип = в дереве | Имя = ТУРИСТИЧЕСКИЙ СЕРВИС])
- }

Рабочее место эксперта

Рубрика [Спорт -- Боевые_искусства]

Дизъюнкты	Вес
@БОЕВЫЕ ИСКУССТВА(E)	1

Конъюнкты	Вес
@БОЕВЫЕ ИСКУССТВА(E)	1

Order	Опорные концепты	Знак	Расш.	Вкл.	Подт.
100	БОЕВЫЕ ИСКУССТВА	+	E	True	A--

Приписка концептов

Все остальные концепты	Подт.
АЙКИДО	
БОЕВЫЕ ИСКУССТВА	
ДЖИУ-ДЖИТСУ	
ДЗЮДО	
ДЗЮДОИСТ	+
КАРАТИСТ	+
КАРАТЭ	
САМБИСТ	+
САМБО	

СКУССТВА(E)

Построение дерева

- [-] БОЕВЫЕ ИСКУССТВА
 - [+] выше ВИД СПОРТА
 - [+] ниже АЙКИДО
 - [+] выше БОЕВЫЕ ИСКУССТВА
 - [+] выше(A) ЯПОНСКАЯ БОРЬБА
 - [+] ниже ДЖИУ-ДЖИТСУ
 - [+] выше БОЕВЫЕ ИСКУССТВА
 - [+] ниже ДЗЮДО
 - [+] выше БОЕВЫЕ ИСКУССТВА
 - [+] выше ВОСТОЧНЫЕ ЕДИНОБОРСТВА
 - [+] (#) часть(A) ДЗЮДОИСТ
 - [+] ниже КАРАТЭ
 - [+] выше БОЕВЫЕ ИСКУССТВА
 - [+] выше ВОСТОЧНЫЕ ЕДИНОБОРСТВА
 - [+] (#) часть(A) КАРАТИСТ
 - [+] выше(B) СПОРТСМЕН
 - [+] целое(A) КАРАТЭ
 - [+] ниже САМБО
 - [+] ассоц(1) САМООБОРОНА
 - [+] выше ОБЕРЕГАТЬ, ЗАЩИЩАТЬ
 - [+] целое ЛИЧНАЯ БЕЗОПАСНОСТЬ
 - [+] часть СРЕДСТВО САМООБОРОНЫ

Отображать все концепты для конъюнкта

Добавить с '←'

Расширение

E V

L W

N

Перестроить приписку

Хозяин

Пометка/Снятие подтверждения

Представление смысла рубрики понятиями тезауруса (8 чел.-час.)

$$R = \prod_i D_i \quad [234] \qquad D_i = \prod_j K_{ij} \quad [265]$$

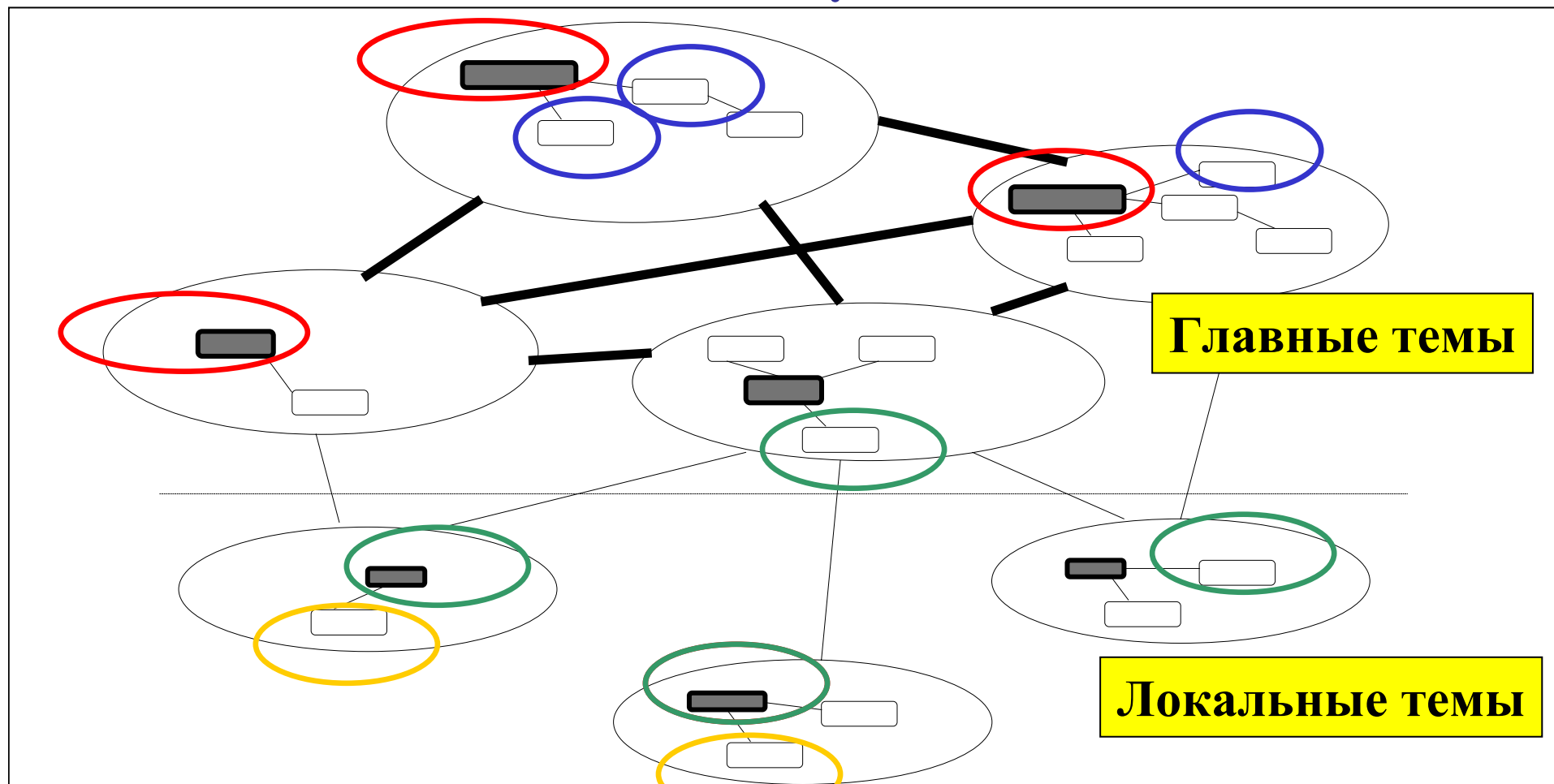
$$K_{ij} = \prod_m f_m(c_{ijm}) \setminus \prod_n f_n(e_{ijn}) = \prod_k d_{ijk}$$

[334] [899]

$$R = \prod_i D_i = \prod_i \left[\prod_j K_{ij} \right] = \prod_i \left[\prod_j \left(\prod_k d_{ijk} \right) \right]$$

[40161]

Вес понятия в тексте: тематическое представления



$$\theta(d) = \alpha \cdot \omega(d; D) + (1 - \alpha) \cdot \frac{\text{freq}(d; D)}{\max_c \text{freq}(c; D)}$$

Расчет веса рубрики: тематическое представление и текстовая связность

$$\theta(D_i) = \frac{\sum_{j=1}^m \theta(K_{ij}) + \sum_{j < k} S(K_{ij}, K_{ik})}{m + C_m^2}$$

$$\theta(K_{ij}) = \min \left\{ 1.0; \max \left(\theta(d_{ijk}), \chi \cdot \theta(p_{ijm}) \right) \right\}$$

$$S(K_{ij}, K_{ik}) = \min \left\{ 1.0; \left(\sum s(c_{ijq} \in K_{ij}, d_{ikw} \in K_{ij}) \right) / \max s(c \in D, d \in D) \right\}$$

Пример простого описания рубрики

- ❖ Рубрика **135 «Боевые искусства»**
(F1-мера [OR] = 0.97, R=0.98, P= 0.96)
- ❖ Опорное булевское выражение состоит из одного понятия

БОЕВЫЕ ИСКУССТВА (E)

с меткой «E» полного расширения по тезаурусу.

- ❖ В состав расширенного булевского выражения входят помимо исходного следующие понятия:

АЙКИДО, ДЖИУ-ДЖИТСУ, ДЗЮДО, КАРАТЭ, САМБО, ДЗЮДОИСТ, КАРАТИСТ, САМБИСТ.

- ❖ Понятия тезауруса, соответствующие людям (***ДЗЮДОИСТ, КАРАТИСТ, САМБИСТ***) входят в рубрику с пометкой подтверждения, поскольку появление соответствующих слов в тексте еще не означает, что текст посвящен боевым искусствам

Более сложное описание рубрики

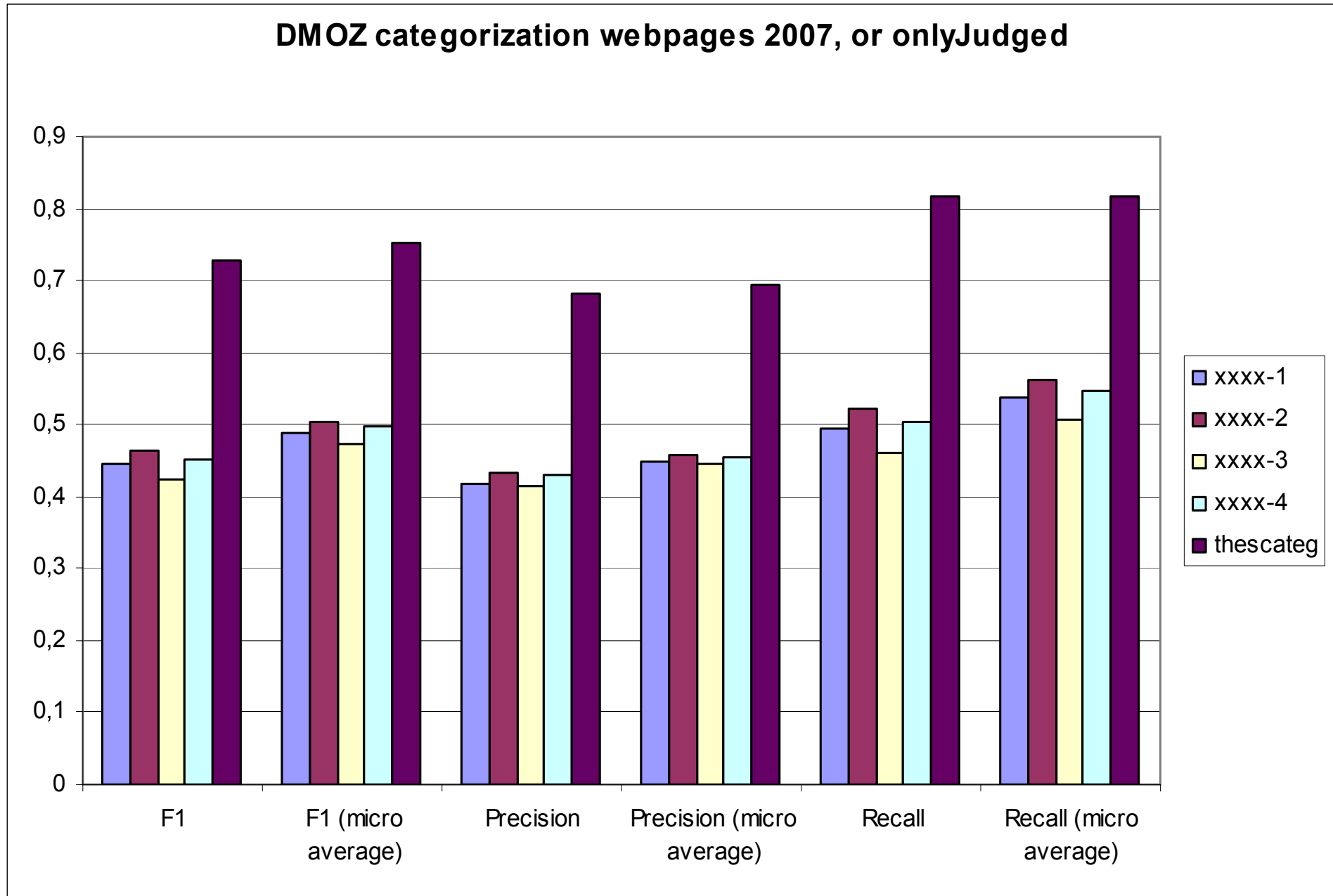
- Рубрика **43** «Домашний ремонт»
F1-мера [OR]= 0.658, R=0.71, P=0.61)

(*РЕМОНТ*
OR *КАПИТАЛЬНЫЙ РЕМОНТ*
OR *ТЕКУЩИЙ РЕМОНТ*
OR *РЕМОНТНО-СТРОИТЕЛЬНЫЕ РАБОТЫ*)

AND

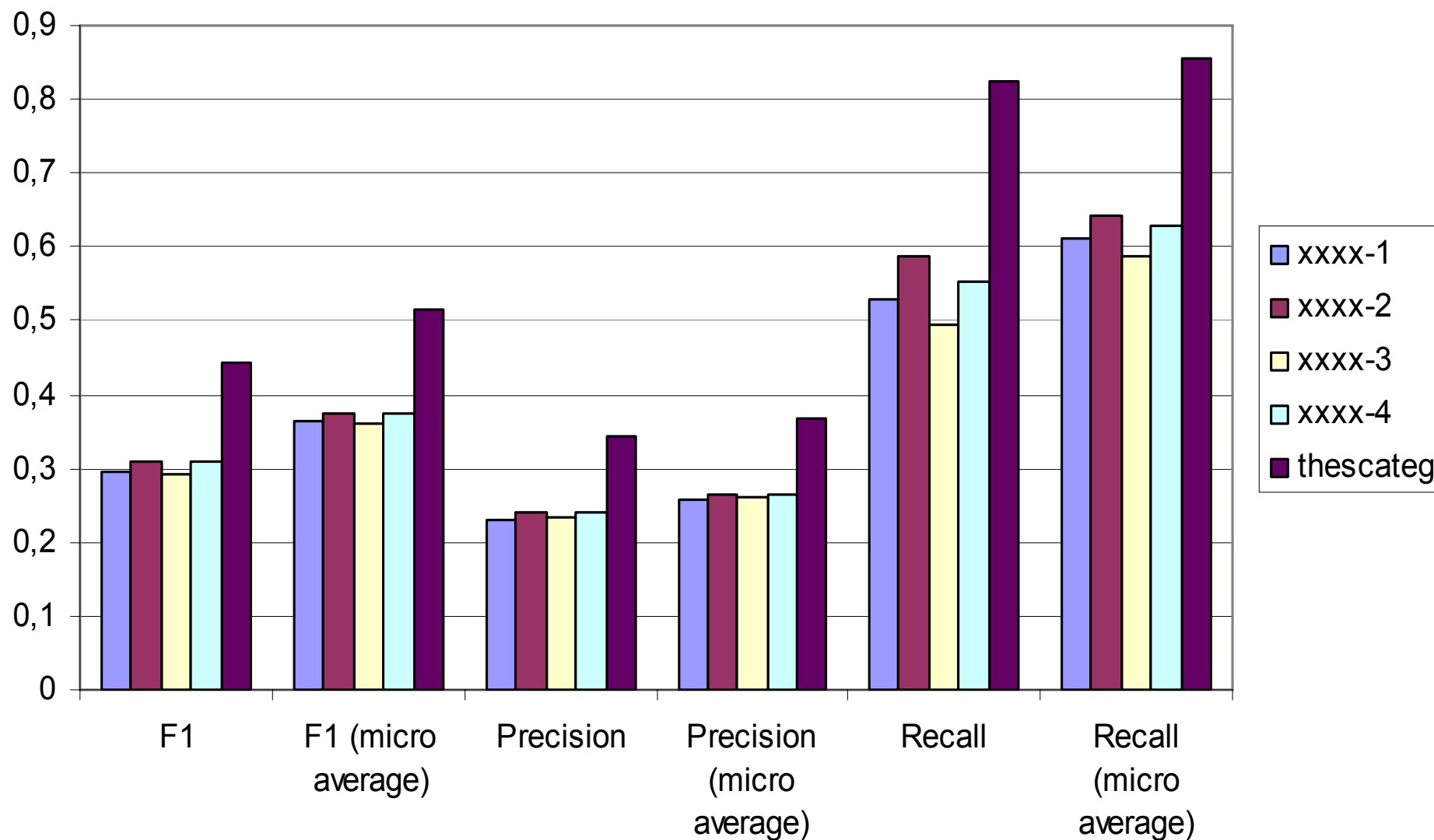
(*ЖИЛОЕ ЗДАНИЕ*_L
OR *ЖИЛОЕ ПОМЕЩЕНИЕ*_L
OR *КВАРТИРА*_L)

РОМИТ2007: классификация веб-страниц [or]



РОМИТТ2007: классификация веб-страниц [and]

DMOZ categorization webpages 2007, and onlyJudged



Проблемы при описании рубрик

- Наиболее низкие по качеству результаты были показаны для рубрики **033 «Сад и огород»** мера F1 по метрике OR составила всего **0.32**
- Это связано с тем, что в состав булевского выражения рубрики были введены понятия ***ДАЧНЫЙ УЧАСТОК*** и ***ЗАГОРОДНАЯ ДАЧА***, без дополнительных условий конъюнкции.

Это описание недостаточно точно, поскольку появление в тексте соответствующих слов еще не гарантирует то, что текст обсуждает проблемы сада и огорода

В ы в о д ы

по задаче web-классификации РОМИТТ 2007

- ❖ **Существуют задачи классификации текстов, когда нет достаточно качественной обучающей коллекции:**
 - ❖ Нет достаточного множества обучающих примеров или ручная классификация проведена недостаточно последовательно
 - ❖ В таких условиях применения методов машинного обучения очень проблематично
- ❖ **При машинном обучении системы извлекают некоторые знания о языке и мире, которые можно условно подразделить на:**
 - ❖ Общие знания о языке и мире, необходимые для работы различных приложений в разнообразном круге предметных областей, и
 - ❖ «Текущие знания», характерные именно для текущей задачи, текущей коллекции, данного типа пользователей и т.п.
 - ❖ Значимую часть знаний о современной жизни общества и современном языке деловой прозы нам удалось упорядочить в рамках понятийных структур Общественно-политического тезауруса

В ы в о д ы (2)

по задаче web-классификации РОМИТТ 2007

- ❖
- ❖ **В текущем эксперименте у нас не было возможности сделать предварительный прогон, оценить и исправить ошибки и неточности описания рубрик**
 - ❖ **В обычной практике проводится несколько итераций, консультаций с экспертами**
 - ❖ **Имеются средства анализа расхождения между системой и экспертами, расхождения также описываются через понятия тезауруса**
 - ❖ **Поэтому имеются определенные возможности улучшения полученных результатов на основе тезаурусных знаний**

Классификация веб-сайтов

Расчет веса рубрики для сайта (только вес страниц)

$$\text{Rank } R_i(S_j) = \text{Avg } R_i(S_j) \cdot \frac{\text{cnt80 } R_i(S_j) + 1}{\text{cnt } R_i(S_j) + 1} \cdot \min\left\{\frac{\text{cnt } R_i(S_j)}{N_S}, 1.0\right\} \cdot \min\left\{\frac{\text{cnt } R_i(S_j)}{\text{cnt } S_j \cdot \frac{PS}{100}}, 1.0\right\}$$

Avg $R_i(S_j)$

– средний вес страниц, отнесенных к рубрике

cnt $R_i(S_j)$

– кол-во страниц сайта, отнесенных к рубрике

cnt80 $R_i(S_j)$

– кол-во страниц сайта с весом 80+,
отнесенных к рубрике

cnt S_j

– общее кол-во страниц сайта

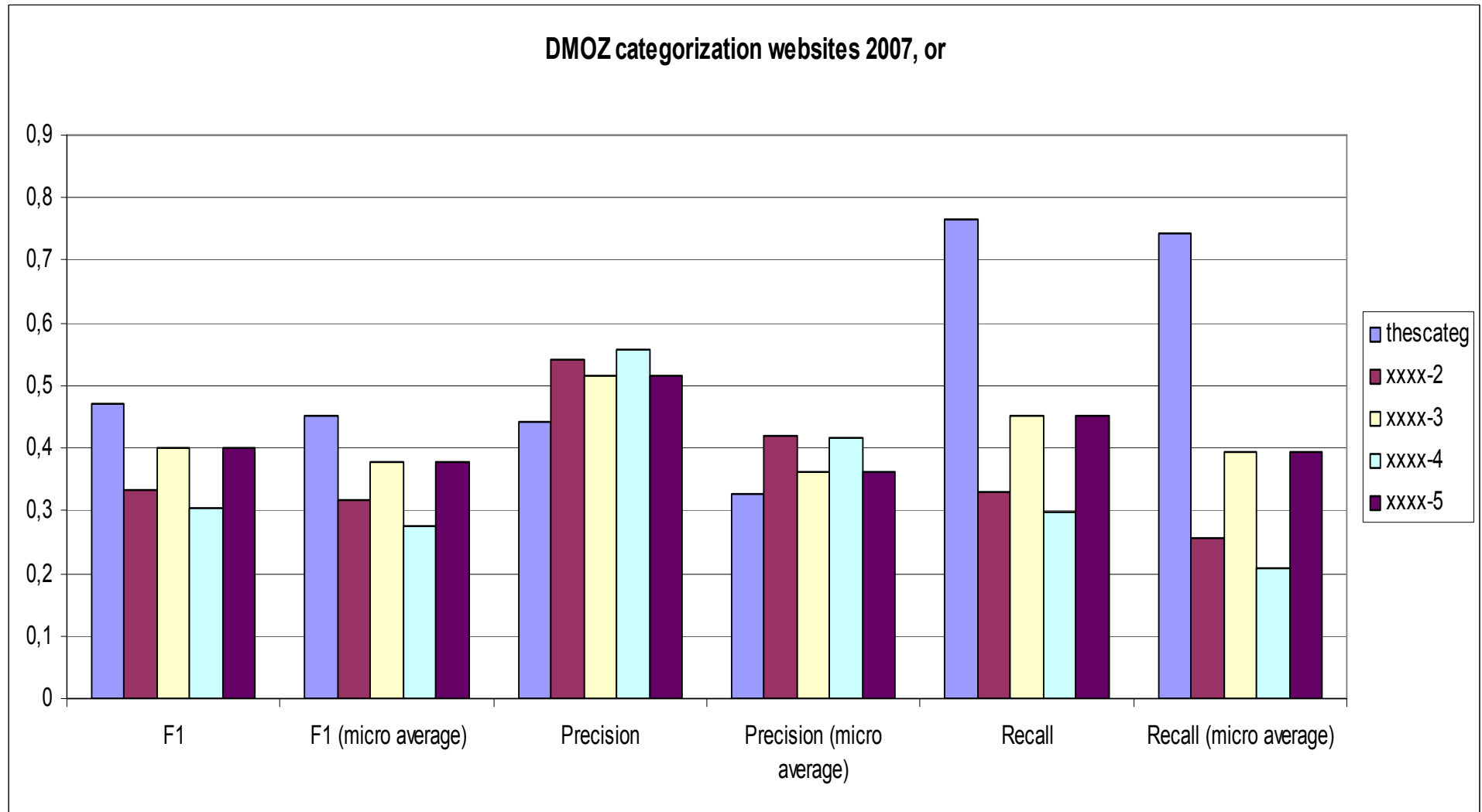
PS

– мин. процент страниц, посвященных рубрике

NS

– мин. кол-во страниц, посвященных рубрике

РОМИПТ2007: классификация веб-сайтов [or]



РОМИПТ2007: классификация веб-сайтов [and]

