



# Тематическая классификация WEB-страниц и WEB-сайтов

С.В. Панков, С.П. Шебанин, А.А. Рыбаков

- R&D подразделение компании Ingate
- Специализация – создание массовых сервисов для интернет рекламы
- ROOKEE – первый продукт, SEO

- Опробовать методики предварительной обработки WEB-ресурсов
- Построить и оценить классификатор WEB-ресурсов, приемлемой точности и производительности
- Получить опыт участия в РОМИП

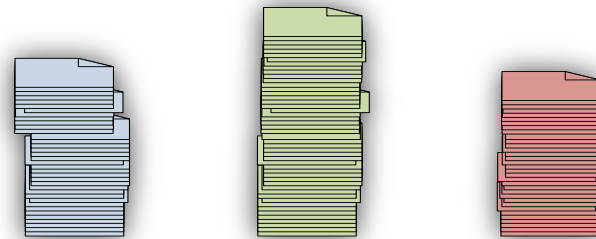
$$D = \{d_1 \dots d_{|D|}\}$$

МНОЖЕСТВО ДОКУМЕНТОВ



$$C = \{c_1 \dots c_{|C|}\}$$

МНОЖЕСТВО ТЕМАТИК



$$\Phi: D \times C \rightarrow \{0, 1\}$$

неизвестная целевая функция

$$\Phi' : D \times C \rightarrow [0, 1]$$

- AOT для работы с морфологией ([aot.ru](http://aot.ru))
- Линейный метод для обучения
- Векторное представление тематик и документов
- Косинус угла для оценки тематической близости

**Но!**

Есть проблемы, связанные со спецификой



Коллекция классифицированных документов (коллекция DMOZ) зашумлена

$$\Omega = \{d_1 \dots d_{|\Omega|}\} = \Omega' \cup \Omega''$$

Коллекция классифицированных документов (коллекция DMOZ) зашумлена

$$\Omega = \{d_1 \dots d_{|\Omega|}\} = \Omega' \cup \Omega''$$

**Не весь текст WEB-страницы информативен**



Главное

В России  
 в СССР  
 В мире  
 Америка  
 Германия

Экономика  
 Финансы  
 Бизнес  
 О регионе

Надежность

Авто  
 Происшествия  
 Масс-медиа  
 О высоком  
 Кино  
 Музыка  
 Спорт

Прогноз  
 Интернет  
 Технологии  
 Игры  
 Судьбы  
 Медицина  
 Из жизни

Психология  
 Опinions  
 Политика  
 Пресс-релизы  
 @lenta.ru  
 RSS

Поиск

08.10.2010, 12:44:00

Версия для печати | RSS-лента



Фото: DAZE

## Хакеры анонсировали инструмент для взлома iPhone 4

Хакеры из группы Citizen Dev Team анонсировали выход инструмента для взлома смартфонов [iPhone 4](#) и других устройств на платформе iOS 4.1, пишет [Redmond Pie](#). Процедура взлома ("джейлбрейк") позволит устанавливать на смартфоны сторонний контент не из официального магазина приложений [Apple App Store](#).

"Джейлбрейк", разработанный для последнего поколения смартфонов iPhone, плееров [Pod touch](#), а также планшетов [iPad](#), будет выпущен 10 октября в 10 часов 10 минут и 10 секунд по британскому времени. Отмечается, что подвергнуть "джейлбрейку" можно будет только устройства на базе процессора A4. Таким образом, он не будет работать на смартфонах и плеерах Apple предыдущих поколений.

Сообщается, что "джейлбрейк" основан на уязвимости SHAzer, [выявленной](#) Citizen Dev Team в начале сентября. При этом "джейлбрейк" будет полным и не потребует подключения устройства к компьютеру при включении аппарата после перезагрузки.

"Джейлбрейк" отличается от процедуры разблокировки SIM-карты (SIM unlock), которая позволяет лишить аппараты привязки к сети определенного мобильного оператора. Отметим, что iPhone с привязкой к оператору продается в некоторых странах, например в США. В России iPhone 4 продается без привязки.

### Ссылки по теме

- [BREAKING: OpenRadio iOS 4.1 Jailbreak is Confirmed to Release on Sunday, 10th October!](#) - Redmond Pie, 08.10.2010
- [Найден способ взлома последней версии Apple iOS](#) - Lenta.ru, 09.09.2010
- [Хакеры проложили путь к взлому iPhone 4](#) - Lenta.ru, 02.08.2010
- [Библиотека Конгресса США разрешила "джейлбрейк" iPhone](#) - Lenta.ru, 27.07.2010

### Сайты по теме

- [Apple iPhone](#)

[ [Обсудить с другими читателями](#) ]

[ [Сообщить о найденной опечатке](#) ]

[ [Письмо в редакцию](#) ]

URL: <http://lenta.ru/news/2010/10/08/jailbreak/>

### Последние новости

- |  |   |
|--|---|
| 08.10 16:53 Кремль обещал с односторонней <a href="#">возвратом</a> <a href="#">Душкова</a>  | 08.10 17:00 "Врачи без границ" объявили войну <a href="#">американской</a> <a href="#">фармакологической</a>  |
| 08.10 16:36 У ливан "Мерседес S660" <a href="#">появилась</a> <a href="#">проблема</a>   | 08.10 16:29 Китай пригрозил Норвегии <a href="#">улучшением</a> <a href="#">отношений</a> <a href="#">по</a> <a href="#">делам</a> <a href="#">мира</a>             |
| 08.10 17:00 Главная страница "Яндекса" <a href="#">представит</a> <a href="#">изменения</a>  | 08.10 16:05 Ахман Абова резко <a href="#">выразил</a> <a href="#">мнение</a> <a href="#">о</a> <a href="#">ситуации</a> <a href="#">с</a> <a href="#">Microsoft</a> |
| 08.10 17:00 Количество подписчиков <a href="#">World of Warcraft</a> <a href="#">примыслило</a> <a href="#">12</a> <a href="#">миллионов</a> |   |

### Аутсайды

- Mail** [Apple iOS Touchpad](#)  
 Apple не разобралась до конца
- Известно** [Google Phone with integrated air purifier](#)  
 В Японии появился телефон, оснащенный встроенным ионизатором воздуха.
- Новости** [Firefox 4 Beta for Android and Windows is Now Available](#)  
 Вышла первая бета-версия мобильного браузера Firefox 4 для Android.
- Code** [Web Applications port, in Googlecode.com for iPhone 4?](#)  
 Технологичный GdGf рассказал о еще одном недостатке iPhone 4, выявленном инженерами Apple.
- Design** [How E-Cas by Patricio Pineda](#)  
 Британский дизайнер создал концепт телефона, зароняющегося в кармане от тепла тела.

### Технологии

- |  |  |
|--|--|
| 08.10 11:10 OAD отказался от <a href="#">владельца</a> <a href="#">аппарата</a> <a href="#">на</a> <a href="#">BlackBerry</a>  | 07.10 14:47 Motorola <a href="#">пошла</a> <a href="#">в</a> <a href="#">суд</a> <a href="#">за</a> <a href="#">Apple</a> <a href="#">по</a> <a href="#">делу</a> <a href="#">об</a> <a href="#">отказе</a> <a href="#">от</a> <a href="#">iPhone</a> <a href="#">и</a> <a href="#">iPad</a> |
| 07.10 17:02 Samsung <a href="#">представил</a> <a href="#">свой</a> <a href="#">новый</a> <a href="#">смартфон</a> <a href="#">на</a> <a href="#">OS</a> <a href="#">Symbian</a> | 07.10 12:07 В Android Market <a href="#">выпустили</a> <a href="#">аппарат</a> <a href="#">с</a> <a href="#">поддержкой</a> <a href="#">Flash</a>  |
| 07.10 15:19 Sony <a href="#">выпустила</a> <a href="#">16</a> <a href="#">мегапиксельную</a> <a href="#">камеру</a> <a href="#">для</a> <a href="#">телефонов</a>                | 07.10 11:06 В США до конца года <a href="#">запустят</a> <a href="#">крупнейшую</a> <a href="#">в</a> <a href="#">мире</a> <a href="#">LTE</a> <a href="#">сеть</a>  |

### В Отрывке



#### Видеоблоггеры

08.10 15:30  
 О том, как видеообложки постепенно вытеснят "массы"

#### Первое касание

08.10 14:00  
 Обзор нового сенсорного iPad nano.

### В Комментариях



#### Двойное касание

08.10 16:02  
 Обзор нового iPod touch

- Главное
- В России
- в СССР
- В мире
- Америка
- Германия
- Экономика
- Эксперты
- Бизнес
- О регионе
- События
- Авто
- Проступки
- Массовые
- О здоровье
- Кино
- Музыка
- Спорт
- Политика
- Инцидент
- Технологии
- Игры
- Судные
- Медицина
- Из жизни
- Пенсии
- Оплата
- Пассажи
- Пресс-релизы
- @lenta.ru
- RSS

Поиск

08.10.2010, 17:44:08

Версия для печати | RSS-лента



Фото: WAPC

## Хакеры анонсировали инструмент для взлома iPhone 4

Хакеры из группы Chronix Dev Team анонсировали выход инструмента для взлома смартфонов iPhone 4 и других устройств на платформе iOS 4.1, пишет [Redmond Pie](#). Процедура взлома ("джейлбрейк") позволит устанавливать на смартфоны сторонний контент не из официального магазина приложений [Apple App Store](#).

"Джейлбрейк", разработанный для последнего поколения смартфонов iPhone, плееров [iPod touch](#), а также планшетов [iPad](#), будет выпущен 10 октября в 10 часов 10 минут и 10 секунд по британскому времени. Отмечается, что подвергнуть "джейлбрейку" можно будет только устройства на базе процессора A4. Таким образом, он не будет работать на смартфонах и плеерах Apple предыдущих поколений.

Сообщается, что "джейлбрейк" основан на уязвимости SHAzer, [выявленной](#) Chronix Dev Team в начале сентября. При этом "джейлбрейк" будет полным и не потребует подключения устройства к компьютеру при включении аппарата после перезагрузки.

"Джейлбрейк" отличается от процедуры разблокировки SIM-карты (SIM unlock), которая позволяет лишить аппараты привязки к сети определенного мобильного оператора. Отметим, что iPhone с привязкой к оператору продается в некоторых странах, например в США. В России iPhone 4 продается без привязки.

### Ссылки по теме

- [BREAKING: OpenPwn iOS 4.1 Jailbreak Confirmed to Release on Sunday, 10th October!](#) - Redmond Pie, 08.10.2010
- [Найден способ взлома последней версии Apple iOS](#) - Lenta.ru, 09.09.2010
- [Хакеры проложили путь к взлому iPhone 4](#) - Lenta.ru, 02.08.2010
- [Информация Интернета США раскрыла "джейлбрейк" iPhone](#) - Lenta.ru, 27.07.2010

### Сайты по теме

- [Apple iPhone](#)

[\[ Обсудить с другими читателями \]](#)

[\[ Сообщить о неточной информации \]](#)

[\[ Письмо в редакцию \]](#)

URL: <http://lenta.ru/news/2010/10/08/jailbreak/>

логин:  пароль:   войти  
 регистрация забыли OpenID

### Последние новости

- |   |   |
|---|---|
| 08.10 16:53 Кремль обеспокоен сценарием возможного прихода Лукашенко        | 08.10 17:00 "Врачи без границ" объявили войну израильской фармацевтике              |
| 08.10 16:06 Упавший "Мерседес S666" попался провозилкам                     | 08.10 16:29 Китай пригрозил Норвегии уступлением отечественной нефти мировому рынку |
| 08.10 17:02 Главная страница "Яндекса" поддается взлому                     | 08.10 16:01 Акции Adobe резко выросли из-за слухов о слиянии с Microsoft            |
| 08.10 17:02 Количество подписчиков World of Warcraft превысило 12 миллионов |   |

### Аутсайды

- Над:** [Nokia N9 Touchfast](#)  
 Nokia N9 разоблачен до деталей
- Известно:** [Dooms iPhone with integrated air purifier](#)  
 В Японии появился телефон, оснащенный встроенным ионизатором воздуха
- Новости:** [Firefox 4 Beta for Android and Windows is Now Available](#)  
 Вышла первая бета-версия мобильного браузера Firefox 4 для Android
- Спойлер:** [Web Applications: part 2, in Glaxo's case for Fox iPhone 4?](#)  
 Технологичный GdF рассказал о еще одном недостатке iPhone 4, выявленном инженерами Apple
- Секреты:** [Nokia E-5a by Patric Hovind](#)  
 Британский дизайнер создал концепт телефона, зароняющегося в кармане от тепла тела

### Технологии

- |  |   |
|--|---|
| 08.10 11:18 OAD отказался от введения запрета на BlackBerry        | 07.10 16:47 Motorola заявила о смене Apple из-за технологий в iPhone и iPad |
| 07.10 17:02 Samsung представил старую модель смартфона на ОС bada  | 07.10 12:57 В Android Market вышел статус с помощью PayPal                  |
| 07.10 15:19 Sony выпустила 16-мегапиксельную камеру для смартфонов | 07.10 11:06 В США до конца года запустят крупнейшую в мире LTE сеть         |

### В Отдыхе



#### Видеоблогеры

08.10 15:30  
 О том, как видеоотомки постепенно входят "в массы"

#### Первое касание

08.10 14:08  
 Обзор нового сенсорного iPad nano

### В Комментариях



#### Двойное касание

08.10 16:52  
 Обзор нового iPod touch

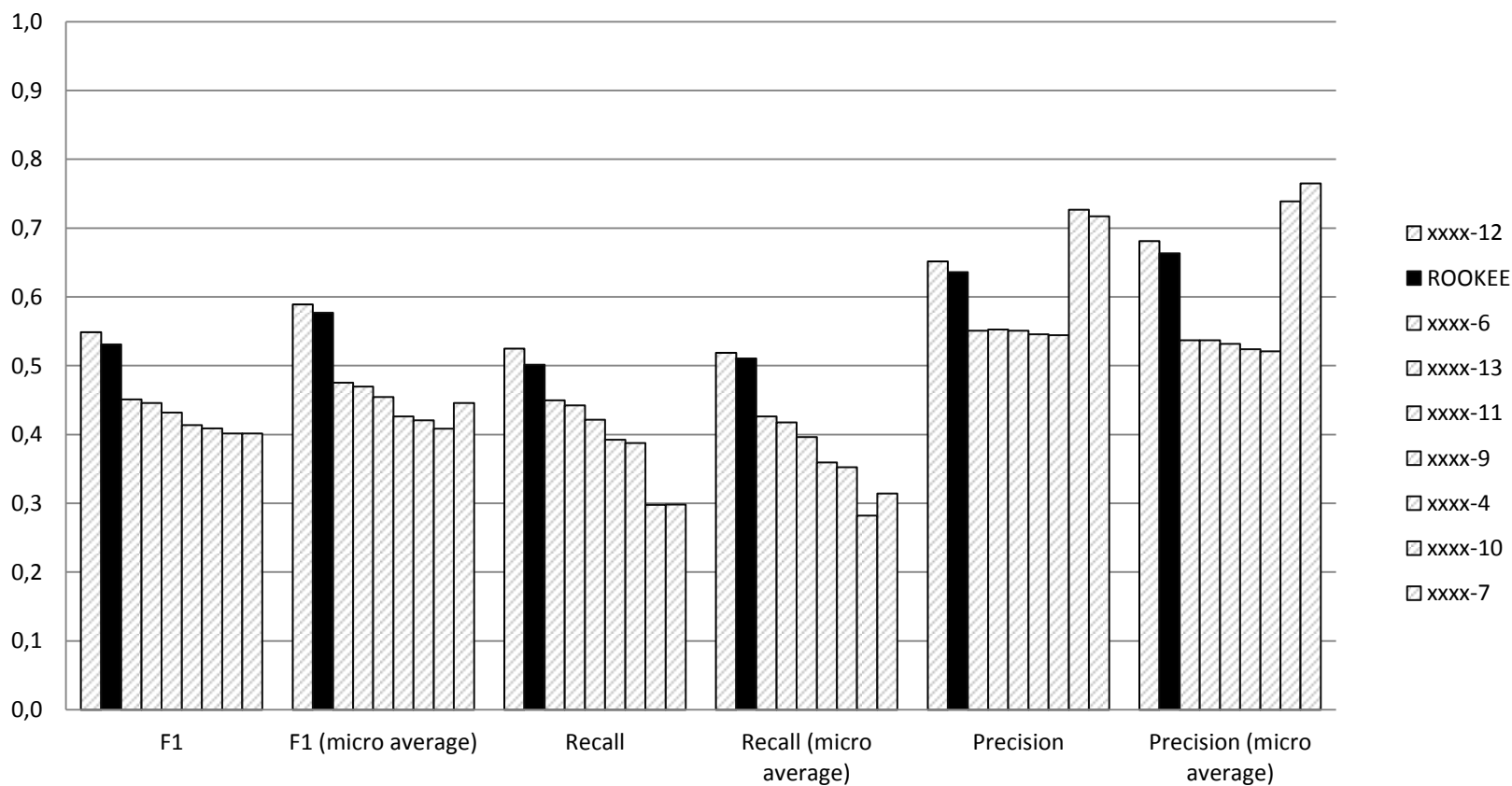
**На сайте присутствуют WEB-страницы,  
не относящиеся к тематике сайта**

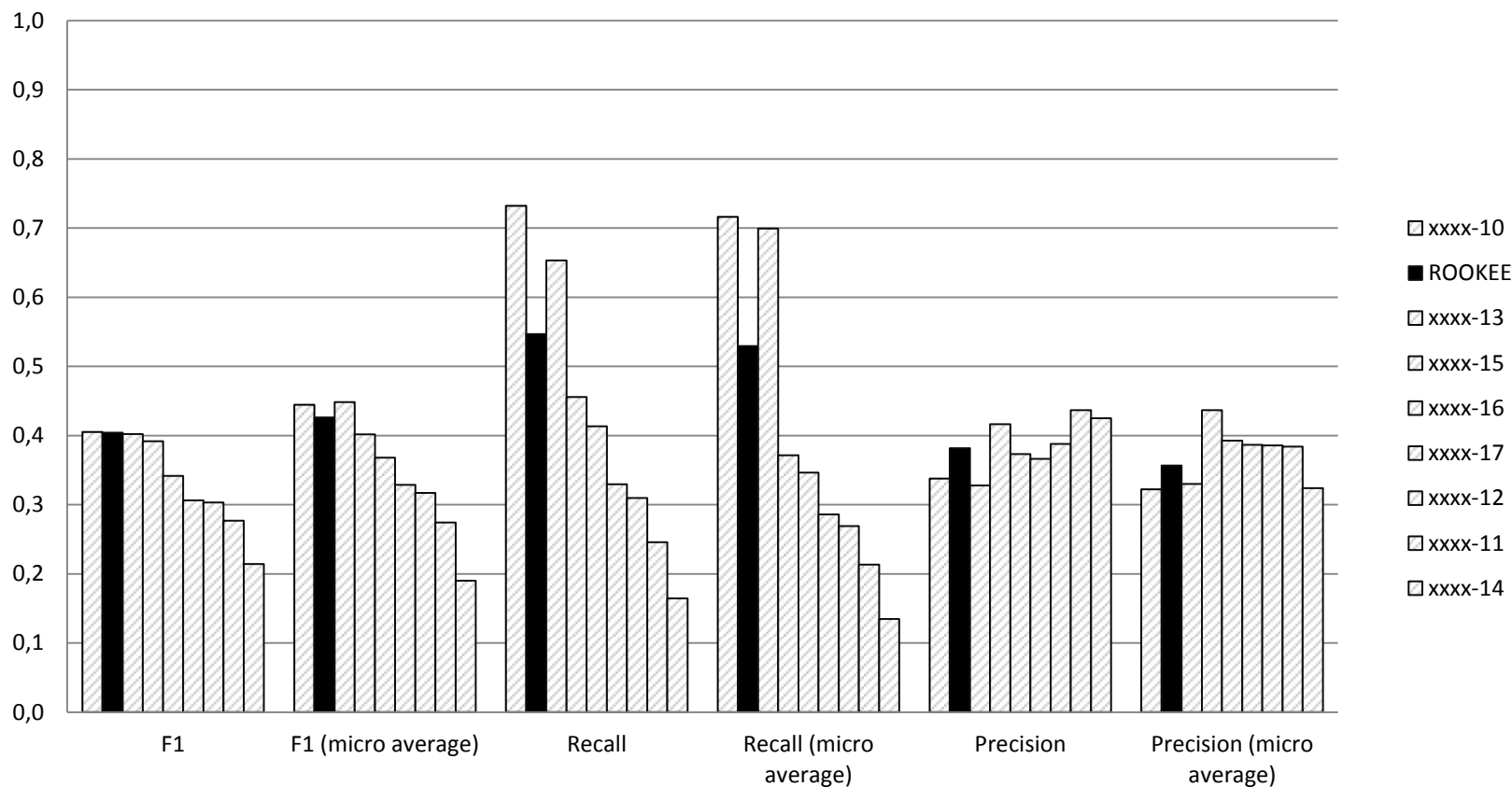


### **Гипотеза:**

Большая часть страниц соответствует тематике сайта

Отсекаем девиантные отклонения





- Использованная методика позволяет получать неплохие результаты рубрикации
- Методика выделения значимой информации на странице может получить широкое применение



- Учитывать позицию текста на странице (заголовки, html-теги и т.д.)
- Корректно обрабатывать случаи мультитематичных сайтов
- Протестировать использование прочих способов для оценки тематической близости
- ...

**Спасибо за внимание!**

**Вопросы?**