

Построение обучаемого алгоритма распознавания научного контента в сети Интернет

Сергей Воронов

Московский физико-технический институт (государственный университет)
Факультет управления и прикладной математики

56-я научная конференция МФТИ

30 ноября 2013 г.

Постановка задачи

- **Цель:** создать систему тематического поиска научного контента в Интернете
- **Задача:** разработать обучаемый алгоритм распознавания научных документов
- **Сложность задачи:**
 - Большой размер коллекции
 - Сильно несбалансированная выборка документов (процент научных документов среди всех не превосходит 2%)
 - Изначально нет системы признаков
 - Изначально нет представительной обучающей выборки

Ненаучный и научный документы

11:04 / ПОНЕДЕЛЬНИК 17 августа 2009 года

Прочитано 2390 раз

Студентов посвятили в чардымовцев

Текст: Альвина ЧАЛОВА , Елизавета ЗИНИНА Фото: Наталия САЛИЙ
В пятницу, 14 августа, в спортивно-оздоровительном лагере ГУ «Чардым» прошло закрытие 4 смены.

По традиции отдыхающих посвятили в «чардымовцев». Для этого студентам необходимо было пройти испытания. Вначале они с закрытыми глазами, держась за руки, ходили по острову, а воспитатели смены рассказывали им чардымские легенды о происхождении лагеря. Студенты услышали историю о красивой девушке и божестве, которые полюбили друг друга. Божество подарило девушке остров «Чардым» и сделало её богиней этого места. Также отдыхающие узнали легенду о чудесном чертополохе, который цветёт разными цветами только один раз в году. После рассказа легенд студентов привели в клуб лагеря. Отдыхающие участвовали в конкурсах на ловкость, скорость и сообразительность. Одним из самых запоминающихся испытаний стало соревнование на лучший танец трёх пар, которые должны были двигаться под музыку с бутылкой, находящейся между ними. Следом за конкурсами отдыхающие стали звать богиню «Чардыма». Её роль исполняла начальник смены Надежда Владимировна Васягина. Она совершила обряд посвящения в «чардымовцев», подарив каждому по веточке чертополоха. Затем самых активных студентов наградили бейсболками и футболками с эмблемой лагеря. Подарок также сделали студентам ФИЖ, обучающимся по специальности «Журналистика», которые каждый день выпускали для отдыхающих стенгазету «В Чардыме». После награждения студенты произнесли клятву, в которой обещали не употреблять алкоголь, делиться ужином с другом и

УДК 159.912.016.77

ПРЕДСТАВЛЕНИЯ О СОЦИАЛЬНО-ПСИХОЛОГИЧЕСКИХ
ХАРАКТЕРИСТИКАХ СУБЪЕКТОВ
ЗАТРУДНЕННОГО ОБЩЕНИЯ

Т.В. Бескова

Институт социального образования (филиал)

Российского государственного социального университета в г. Саратове

E-mail: beskova-t@yandex.ru

В статье анализируются результаты эмпирического исследования представлений студентов о социально-психологических характеристиках субъектов затрудненного общения на примере представителей профессий социомического типа.

Устанавливаются сходства и различия представлений по ряду стимулов (социальный работник, преподаватель, медицинский работник). Выделяются обобщенные портреты трудных партнеров по общению.

Ключевые слова: представления, субъект, общение, трудности, социально-психологические характеристики, ролевые позиции.

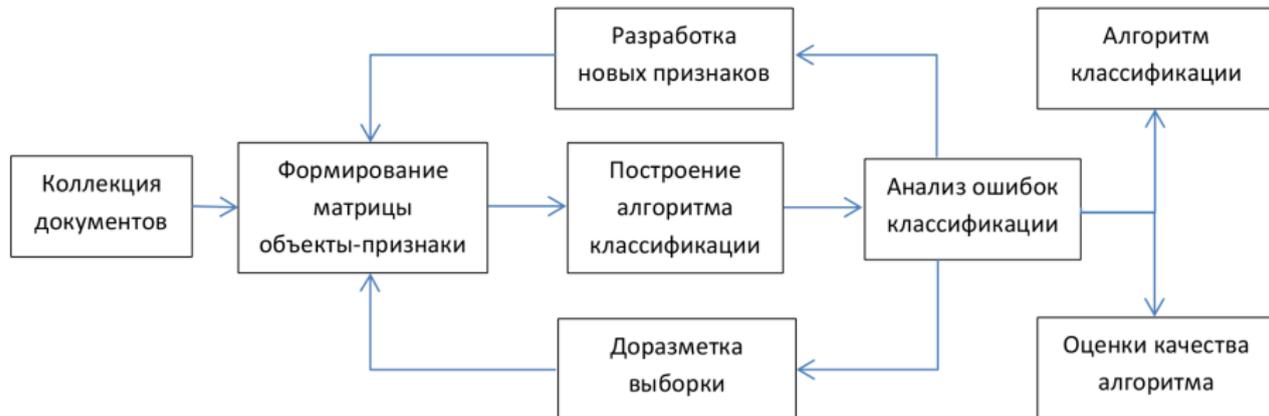
The Imaginations about Social-Psychological Descriptions of the Characters of Difficult Communication

T.V. Beskova

The article analyses the results of empirical research on students' imaginations about socialpsychological descriptions of the characters of difficult communication of the representatives of the social-directional occupations, as the example. It underlines the likenesses and the differences of the

Процесс разработки алгоритма классификации

Автоматизация процесса постепенного наращивания обучающей выборки, разработки признаков и улучшения классификатора.



Коллекция документов

- Размер: $\sim 10^6$ документов
- Источник: Страницы с 2000 сайтов российских ВУЗов
- Доля научных документов не превосходит 2%
- Метрика качества: контроль по 10 блокам (10-fold CV)

Оценка качества алгоритма классификации q -fold CV

Выборка случайным образом разбивается на q непересекающихся блоков почти одинаковой длины k_1, \dots, k_q :

$$X^L = X_1^{k_1} \cup \dots \cup X_q^{k_q},$$

$$k_1 + \dots + k_q = L.$$

Каждый блок по очереди становится контрольной подвыборкой, при этом обучение производится по остальным $q - 1$ блокам. Критерий определяется как средняя ошибка на контрольной подвыборке:

$$CV(\mu, X^L) = \frac{1}{q} \sum_{n=1}^q Q(\mu(X^L \setminus X_n^{k_n}), X_n^{k_n}).$$

Программа для формирования обучающей выборки

GUI v. 0.3

ori	name	state	score list	PMS	PBS	PIW	PGL	NC	NS
1	03/3265452.txt	Unchecked	[0, 0, 1, 0, 0, -1]	0	0	1	0	0	-1
1	20/62399891.txt	Unchecked	[0, 0, 1, 0, 0, 0]	0	0	1	0	0	0
68	F4/59570156.txt	Scientific	[0, 36, 32, 0, 0, 0]	0	36	32	0	0	0
129	B7/62756992.txt	Scientific	[24, 24, 10, 71, 0, 0]	24	24	10	71	0	0
1	FE/62426885.txt	Unchecked	[0, 0, 1, 0, 0, -1]	0	0	1	0	0	-1
1	28/62494441.txt	Unchecked	[0, 0, 1, 0, 0, -1]	0	0	1	0	0	-1
34	D2/59660991.txt	Unknown	[0, 26, 8, 0, 0, 0]	0	26	8	0	0	0
4	DD/62540083.txt	Unchecked	[0, 0, 4, 0, 0, -1]	0	0	4	0	0	-1
10	F0/62178549.txt	Unchecked	[0, 0, 10, 0, 0, 0]	0	0	10	0	0	0
54	62/58693567.txt	Non-scientific	[0, 41, 13, 0, 0, 0]	0	41	13	0	0	0
1	2F/62477065.txt	Non-scientific	[0, 0, 1, 0, 0, -1]	0	0	1	0	0	-1
56	A3/58647152.txt	Scientific	[0, 45, 11, 0, 0, 0]	0	45	11	0	0	0
5	DF/62819465.txt	Non-scientific	[0, 0, 5, 0, -1, -1]	0	0	5	0	-1	-1
1	8E/62424037.txt	Non-scientific	[0, 0, 1, 0, 0, -1]	0	0	1	0	0	-1
4	8E/62665692.txt	Non-scientific	[0, 0, 4, 0, 0, -1]	0	0	4	0	0	-1
10	41/59579186.txt	Non-scientific	[0, 0, 10, 0, 0, -1]	0	0	10	0	0	-1
43	10/62807549.txt	Non-scientific	[0, 42, 1, 0, 0, 0]	0	42	1	0	0	0
11	09/62198822.txt	Unchecked	[0, 0, 11, 0, 0, -1]	0	0	11	0	0	-1
1	83/7128825.txt	Unchecked	[0, 0, 1, 0, 0, 0]	0	0	1	0	0	0
1	E4/62417419.txt	Unchecked	[0, 0, 1, 0, 0, -1]	0	0	1	0	0	-1
1	D8/62357502.txt	Unchecked	[0, 0, 1, 0, 0, 0]	0	0	1	0	0	0
1	45/62455065.txt	Unchecked	[0, 0, 1, 0, 0, -1]	0	0	1	0	0	-1
1	D6/59702493.txt	Unchecked	[0, 0, 1, 0, 0, -1]	0	0	1	0	0	-1
1	48/64728360.txt	Non-scientific	[0, 0, 1, 0, -1, 0]	0	0	1	0	-1	0
1	C9/62350864.txt	Unchecked	[0, 0, 1, 0, 0, -1]	0	0	1	0	0	-1
42	89/62906483.txt	Unchecked	[0, 25, 17, 0, 0, 0]	0	25	17	0	0	0

Information

Data loaded in time 101.52 sec

Parameters

Files displayed: 18572
 Documents count: 844206
 Scientific documents: 822
 Non-scientific documents: 1061
 Unknown documents: 45
 Mistakes on sci documents: 76
 Mistakes on non-sci documents: 101

7

Close additional options

Quality control
Quality: 90%

Refresh parameters

Save data
Recheck
Start svm
SVM settings
Display mode: all
Quit

Базовый необучаемый классификатор

Линейный алгоритм классификации:

$$a(x, w) = \text{sign} \left(\sum_{j=1}^n w_j f_j(x) - w_0 \right),$$

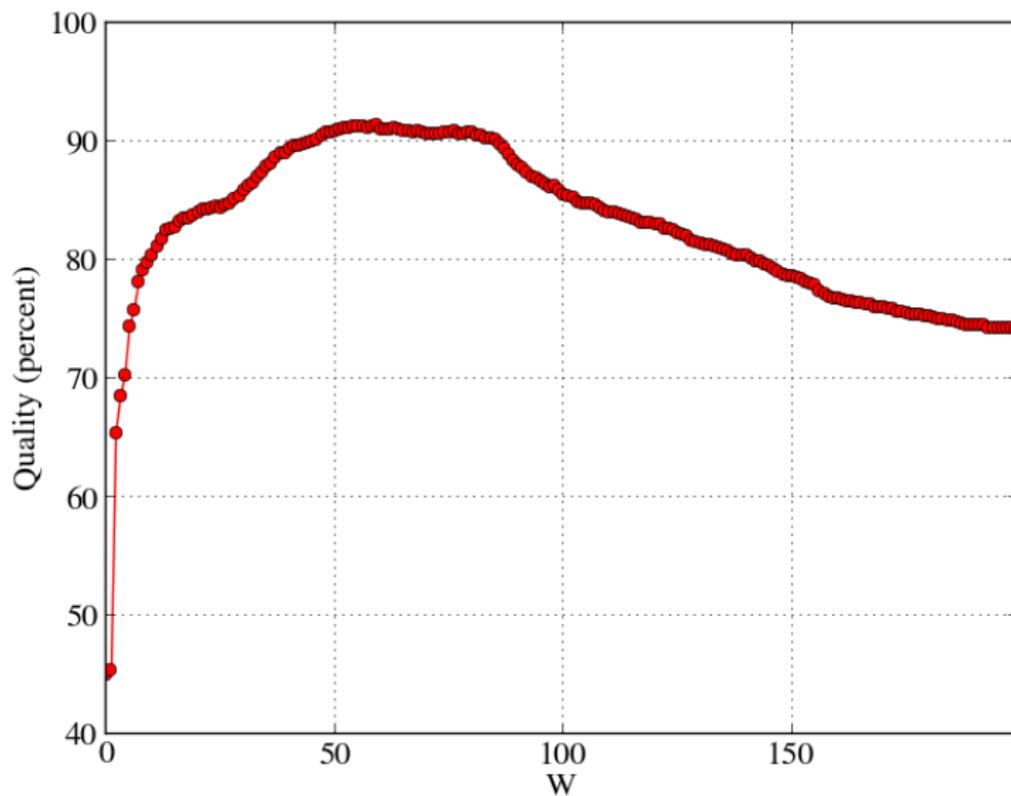
где w_j — вес j -го признака, w_0 — порог принятия решения, $w = (w_1, \dots, w_n)$ — вектор весов признаков.

В базовой необучаемой версии w_i задаются экспертом.

Начальные признаки:

- Длина текста (длина текста больше заданного порога)
- Число греческих букв в тексте
- Число математических символов в тексте

Зависимость качества базового классификатора от w_0



Классификатор

Итоговый классификатор: SVM.

$$a(x, w) = \text{sign} \left(\sum_{j=1}^n w_j f_j(x) - w_0 \right),$$

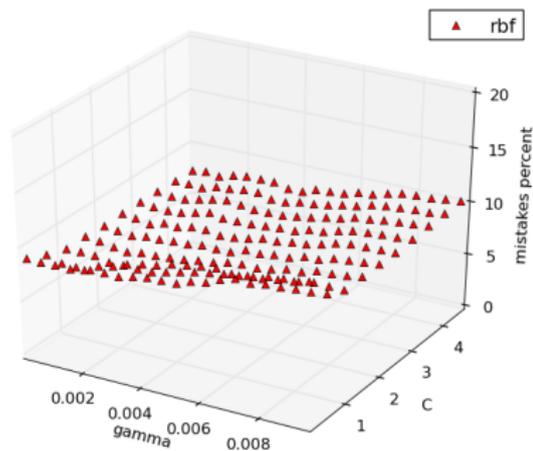
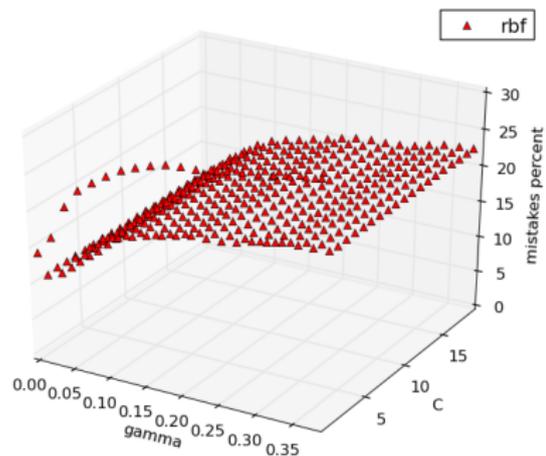
w_j определяются из условия:

$$\begin{cases} \frac{1}{2} \langle w, w \rangle + C \sum_{i=1}^n \xi_i \rightarrow \min_{w, w_0, \xi} \\ y_i (\langle w; x_i \rangle - w_0) \geq 1 - \xi_i, \quad i = 1 \dots n \\ \xi_i \geq 0 \end{cases}$$

Выбор оптимальных параметров:

- ядро rbf (радиальные базисные функции, $e^{-\gamma(x-x_0)^2}$)
- γ и C (штраф за неверную классификацию) — по скользящему контролю

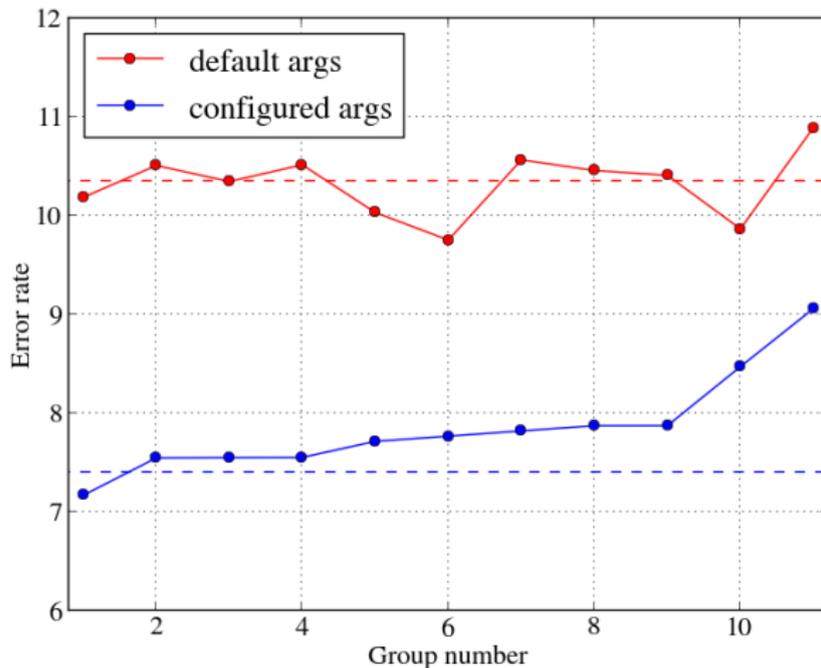
Оптимизация параметров SVM



Признаковое описание документов

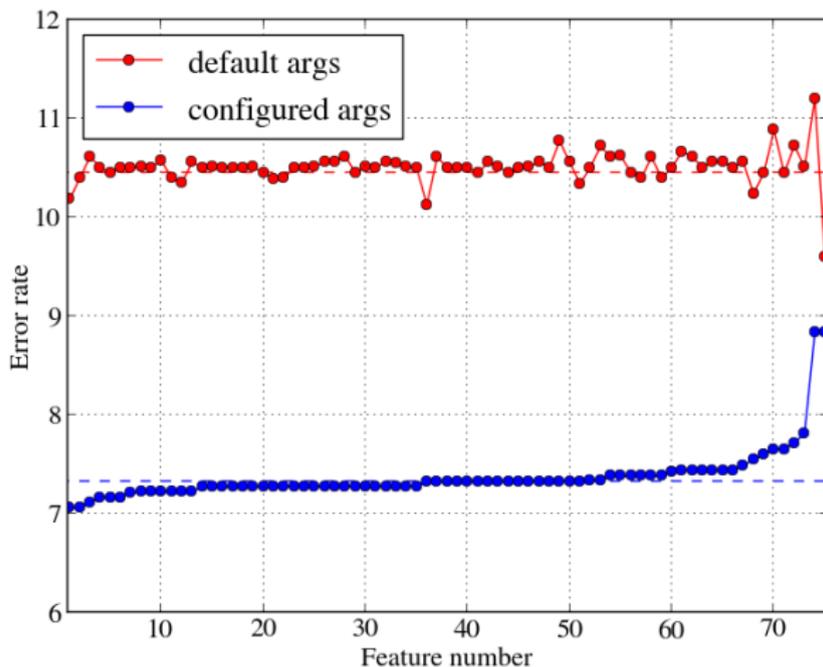
- 1 Количество ключевых слов в тесте (все слова разбиты на группы)
- 2 Наличие цифр в документе (их количество более 5 или нет)
- 3 Длина текста (длина текста больше заданного порога)
- 4 Минус-слова, присутствие которых понижает “научность” (количество часов, профком, экзамен)
- 5 Остальные слова, не попавшие в предыдущие группы
- 6 Слова, относящиеся к оформлению (редактор, редколлегия)
- 7 Слова, которые должны встречаться не более трех раз (Труды)
- 8 Число математических символов в тексте
- 9 Научные термины
- 10 Размер текста (логарифм размера текста в символах)
- 11 Число греческих букв в тексте

Отбор наиболее значимых признаков



Результаты проверки групп признаков на значимость
(удаляется по одной группе)

Отбор наиболее значимых признаков



Результаты проверки всего набора признаков на значимость
(удаляется по одному признаку)

Результаты

Стадия алгоритма	Ошибка 1 рода	Ошибка 2 рода
Базовый классификатор	4,9%	6,5%
SVM	3.2%	7.0%
Настройка параметров	4.2%	3.8%
Слова из оформления	4.2%	3.7%
Научные термины	4.2%	3.3%
Стоп слова	4.2%	3.2%

- Ошибка первого рода: научный классифицирован как ненаучный.
- Ошибка второго рода: ненаучный классифицирован как научный.