



ТЕХНОСФЕРА

Лекция 1

Лингвистика в поиске Основы обработки текстов

Владимир Гулин

13 сентября 2016 г.

Что мы знаем на текущий момент?

- ▶ Как качать интернет
- ▶ Что класть в индекс
- ▶ Фильтрация нежелательного контента (антиспам/антипорн)
- ▶ Исправление опечаток и саджесты
- ▶ Модель булева поиска
- ▶ Как формировать сниппеты
- ▶ Что еще?

Структура курса

1. Лингвистика в поиске. Основы обработки текстов.
2. Коллокации, n-граммы, скрытые марковские цепи ^H
3. Текстовое ранжирование. Базовые модели ^H
4. Оценка качества поиска. ^H
5. Оценка качества поиска. Интерливинг.
6. Ссылочное ранжирование. ^H
7. Поведенческое ранжирование.
8. Машинное обучение в ранжировании 1 ^H
9. Машинное обучение в ранжировании 2
10. Текстовое ранжирование. Хитрые модели
11. Разбор ключевых ДЗ
12. Мультимедия поиск
13. Технологии поиска по лицам ^H
14. Контентные рекомендательные системы

Лингвистика в поиске

Этапы ранжирования



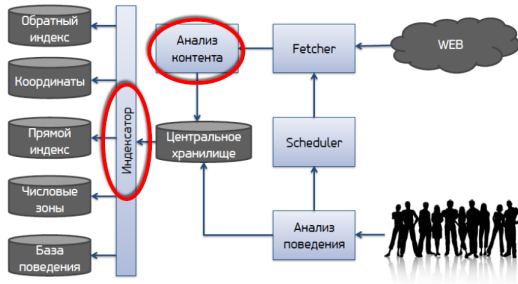
Терминология

- ▶ *Токен* - экземпляр последовательности символов в документе, объединенных в семантическую единицу для обработки
- ▶ *Термин* - “нормализованный” токен (регистр, морфология, исправление ошибок и т.п.)

Нормализация

- ▶ Необходимо “нормализовывать” термины как в индексируемом тексте, так и в запросе
- ▶ Например: Желательно считать одинаковыми термины U.S.A. и USA
- ▶ Обычно термины объединяются в классы эквивалентности
- ▶ Можно поступить наоборот, расширять:
 - ▶ window → window, windows
 - ▶ windows → Windows, windows
 - ▶ Windows (нет расширения)
- ▶ Такой подход более гибкий, но более ресурсоемкий

Обработка документа



Документы

Делаем неявное предположение:

- ▶ Мы знаем, что такое документ
- ▶ Каждый документ доступен для автоматического разбора

Вопрос:

- ▶ Какие тут есть проблемы?

Лингвистика при обработке документов

- ▶ Определение формата документа (pdf, word, html и т.д.)
- ▶ Определение кодировки документа
- ▶ Определение языка документа
- ▶ Токенизация и сегментация
- ▶ Нормализация и лемматизация
- ▶ Выделение объектов и зон
- ▶ Вычисление текстовых факторов

Нормализация

Нормализация зависит от языка документа

- ▶ PETER WILL NICHT MIT. → MIT = mit
- ▶ He got his PhD from MIT. → MIT \neq mit

Нормализация

Ударения и диакритика

- résumé vs. resume
- Умуляуты: Universität vs. Universitaet
(заменяем на специальную
последовательность «ae» или даже «æ»)
- Самый важный вопрос: как пользователи предпочитают писать запросы с этими словами?

Нормализация

Классы эквивалентности

- Soundex
 - фонетическая эквивалентность, Muller = Mueller
- Тезаурус
 - семантическая эквивалентность, car = automobile

Нормализация

Регистр

- Понизить регистр всех букв.
- Возможны исключения, например, для капитализированных слов внутри предложения.
 - MIT и mit
 - Fed и fed
 - КОТ и кот (Калининградская областная таможня)
- NB: немецкий → существительные с большой буквы
- Часто лучше понижать всё, потому что пользователи не заботятся о капитализации в запросах.

Токенизация

Проблемы токенизации

- ▶ Hewlett-Packard
- ▶ State-of-the-art
- ▶ co-education
- ▶ San Francisco
- ▶ York University vs. New York University

Проблемы токенизации

Числа

- ▶ 3/20/91
- ▶ 20/3/91
- ▶ Mar 20, 1991
- ▶ B-52
- ▶ 100.2.86.144
- ▶ (800) 234-2333
- ▶ 800.234.2333

Обработка запроса



Обработка запроса

Запросы задают не по-русски

- ▶ Распознавание языка
- ▶ Исправление опечаток
- ▶ Токенизация
- ▶ Нормализация и лемматизация
- ▶ Кореференция (расширение запроса)
- ▶ Переформулировки запросов
- ▶ Сегментация запроса
- ▶ Извлечение объектов

Проблемы токенизации

В китайском нет пробелов

李克强说，当前国际和地区形势发生复杂深刻变化，中越都处于发展的关键阶段，双方要从战略高度和长远角度出发，在发展中越关系十六字方针和“四好”精神指引下，坚定不移推进中越友好。中方愿同越方保持高层战略沟通，加强治国理政经验交流，坚持经济优先、民生优先，深化务实合作，推动中越全面战略合作伙伴关系迈上新台阶。

(c) news.xinhuanet.com

Другие случаи отсутствия пробелов

Компаунды в датском, немецком, шведском

- ▶ Computerlinguistik → Computer + Linguistik
- ▶ Lebensversicherungsgesellschaftsangestellter → leben + versicherung + gesellschaft + angestellter
- ▶ Kallistuksenvaimennusjärjestelmä - система, предотвращающая крен (в погрузчиках).

Льезоны в романских языках

- ▶ em os → nos
- ▶ por a → pela

Эскимосы: tusaatsiarunnanngittualuujunga (Я не очень хорошо слышу)

Японский

- 4 разных алфавита.
- ローマ字 Romaji
- 平仮名 Hiragana
- 片仮名 Katakana
- 漢字 Kanji

Запрос может быть задан в любом из них

Named Entity Recognition

Извлечение объектов (группа слов/токенов в запросе, которые означают одно понятие)

- ▶ ФИО
- ▶ Телефоны
- ▶ Адреса
- ▶ Даты
- ▶ Названия песен, фильмов, книг и т.д.

Кодировки

- ▶ ASCII (ISO 646) - 7-битный стандарт
- ▶ ISO 8859
 - ▶ 8859-1 (ISO Latin-1)
 - ▶ ISO 8859-5
- ▶ Русские кодировки
 - ▶ CP1251 (windows)
 - ▶ 866 (dos)
 - ▶ KOI8-R (unix)
- ▶ Unicode
 - ▶ UTF-8
 - ▶ UTF-16
 - ▶ UTF-32

Кодировки

ASCII

American Standard Code for Information Interchange (1967)

ASCII Code Chart																
	0	1	2	3	4	5	6	7	8	9	A	B	C	D	E	F
0	NUL	SOH	STX	ETX	EOT	ENQ	ACK	BEL	BS	HT	LF	VT	FF	CR	SO	SI
1	DLE	DC1	DC2	DC3	DC4	NAK	SYN	ETB	CAN	EM	SUB	ESC	FS	GS	RS	US
2		!	"	#	\$	%	&	'	()	*	+	,	-	.	/
3	0	1	2	3	4	5	6	7	8	9	:	;	<	=	>	?
4	@	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O
5	P	Q	R	S	T	U	V	W	X	Y	Z	[\]	^	_
6	`	a	b	c	d	e	f	g	h	i	j	k	l	m	n	o
7	p	q	r	s	t	u	v	w	x	y	z	{		}	~	DEL

Кодировки

koi8-r

	0	1	2	3	4	5	6	7	8	9	A	B	C	D	E	F
8	—		Г	г	Л	л	Т	т	Т	т	+	■	■	■	■	■
128	129	130	131	132	133	134	135	136	137	138	139	140	141	142	143	
9	▤	▥	▦	г	■	●	√	≈	≤	≥	nbsp	Ј	•	²	•	÷
144	145	146	147	148	149	150	151	152	153	154	155	156	157	158	159	
A	=		F	ё	П	Г	г	П	П	Е	Л	Ј	Ј	Ј	Ј	
160	161	162	163	164	165	166	167	168	169	170	171	172	173	174	175	
B			г	ё	П	Г	г	П	П	Е	Л	Ј	Ј	Ј	Ј	©
176	177	178	179	180	181	182	183	184	185	186	187	188	189	190	191	
C	ю	а	б	ц	д	е	ф	г	х	и	й	к	л	м	н	о
192	193	194	195	196	197	198	199	200	201	202	203	204	205	206	207	
D	п	я	р	с	т	у	ж	в	ь	ы	з	ш	э	щ	ч	б
208	209	210	211	212	213	214	215	216	217	218	219	220	221	222	223	
E	Ю	А	Б	Ц	Д	Е	Ф	Г	Х	И	Й	К	Л	М	Н	О
224	225	226	227	228	229	230	231	232	233	234	235	236	237	238	239	
F	П	Я	Р	С	Т	У	Ж	В	Ь	Ы	З	Ш	Э	Щ	Ч	Б
240	241	242	243	244	245	246	247	248	249	250	251	252	253	254	255	

Кодировки

cp866

	0	1	2	3	4	5	6	7	8	9	A	B	C	D	E	F
8	А	Б	В	Г	Д	Е	Ж	З	И	Й	К	Л	М	Н	О	П
9	Р	С	Т	У	Ф	Х	Ц	Ч	Ш	Щ	Ъ	Ы	Ь	Э	Ю	Я
A	а	б	в	г	д	е	ж	з	и	й	к	л	м	н	о	п
B	␣	␣	␣													
C																
D																
E	р	с	т	у	ф	х	ц	ч	ш	щ	ъ	ы	ь	э	ю	я
F	Ё	ё	Є	є	İ	ı	Ÿ	ÿ	°	•	•	√	№	¤	■	nbsp

Кодировки

cp1251

	0	1	2	3	4	5	6	7	8	9	A	B	C	D	E	F
8	Ъ	ґ	,	ґ	„	…	†	‡	І	%	Љ	<	Њ	Ќ	Ћ	Ќ
128	129	130	131	132	133	134	135	136	137	138	139	140	141	142	143	
9	ђ	‘	’	“	”	•	—	І	™	љ	>	њ	ќ	ћ	џ	
144	145	146	147	148	149	150	151	152	153	154	155	156	157	158	159	
A	nbsp	Ў	ў	Ј	Ѡ	Г	І	Ѕ	Ё	Є	«	¬	shy	®	Ї	
160	161	162	163	164	165	166	167	168	169	170	171	172	173	174	175	
B	°	±	І	і	г	μ	¶	•	ё	№	є	»	ј	Ѕ	ѕ	ї
176	177	178	179	180	181	182	183	184	185	186	187	188	189	190	191	
C	А	Б	В	Г	Д	Е	Ж	З	И	Й	К	Л	М	Н	О	П
192	193	194	195	196	197	198	199	200	201	202	203	204	205	206	207	
D	Р	С	Т	У	Ф	Х	Ц	Ч	Ш	Щ	Ъ	Ы	Ь	Э	Ю	Я
208	209	210	211	212	213	214	215	216	217	218	219	220	221	222	223	
E	а	б	в	г	д	е	ж	з	и	й	к	л	м	н	о	п
224	225	226	227	228	229	230	231	232	233	234	235	236	237	238	239	
F	р	с	т	у	ф	х	ц	ч	ш	щ	ъ	ы	ь	э	ю	я
240	241	242	243	244	245	246	247	248	249	250	251	252	253	254	255	

Кодовое пространство Unicode

- ▶ обозначения: U+xxxx, U+xxxxx, U+xxxxxx
- ▶ пространство разделено на 17 плоскостей по 2^{16} символов
- ▶ первые 128 символов совпадают с ASCII
- ▶ плоскость 0 (base multilingual plane) содержит основные символы
- ▶ остальные плоскости содержат символы редких письменностей
- ▶ 2048 кодов U+DC00 - U+DFFF заняты под “суррогатные пары”

Всего символов в Unicode

$$17 \cdot 2^{16} - 2048 = 1112064$$

Кодировки

UTF-8

- ▶ Нужен для передачи Unicode по однобайтовым каналам связи
- ▶ Начало аналогично первой половине таблицы ASCII
- ▶ Обладает свойством самосинхронизации
- ▶ Мультибайтная кодировка

0xxxxxxx

110xxxxx 10xxxxxx

1110xxxx 10xxxxxx 10xxxxxx

11110xxx 10xxxxxx 10xxxxxx 10xxxxxx

Кодировки

UTF-16

- ▶ Символы из base multilingual plane кодируются как есть
- ▶ Символы из других плоскостей кодируются парой слов с помощью “суррогатных пар” U+D800 - U+DFFF
- ▶ Важен порядок байт в слове!
- ▶ BOM (byte order mark)
- ▶ big-endian (от старших разрядов к младшим): U+FEFF (UTF-16BE)
- ▶ little-endian (от старших разрядов к младшим): U+FFFE (UTF-16LE)

110110xx xxxxxxxx 110111xx xxxxxxxx

Кодировки

Объединение и дублирование символов

Для многих символов в Unicode есть отдельные коды. Это имеет место по историческим причинам, так как они присутствуют в национальных кодировках.

Новые составные символы в Unicode не добавляются, их нужно “конструировать” из нескольких кодов:

À 00C5	:	Ä 0041	◌̊ 030A
Ô 00F4	:	Ö 006F	◌̈ 0302

Кодировки

Классы эквивалентности в Unicode

Cyrillic extensions

0400	Ё	CYRILLIC CAPITAL LETTER IE WITH GRAVE ≡ 0415 Ё 0300 ѐ
0401	ѐ	CYRILLIC CAPITAL LETTER IO ≡ 0415 Ё 0308 ѐ
0402	Ђ	CYRILLIC CAPITAL LETTER DJE
0403	Ѓ	CYRILLIC CAPITAL LETTER GJE ≡ 0413 Г 0301 ѓ
0404	Є	CYRILLIC CAPITAL LETTER UKRAINIAN IE
0405	Ѕ	CYRILLIC CAPITAL LETTER DZE
0406	І	CYRILLIC CAPITAL LETTER BYELORUSSIAN- UKRAINIAN I → 0049 І latin capital letter i → 0456 і cyrillic small letter byelorussian- ukrainian i → 04C0 І cyrillic letter pabochka
0407	Ї	CYRILLIC CAPITAL LETTER YI ≡ 0405 І 0308 й
0408	Ј	CYRILLIC CAPITAL LETTER JE
0409	Љ	CYRILLIC CAPITAL LETTER LJE
040A	Њ	CYRILLIC CAPITAL LETTER NJE
040B	Ѡ	CYRILLIC CAPITAL LETTER TSHE
040C	Ѣ	CYRILLIC CAPITAL LETTER KJE ≡ 041A К 0301 ѣ
040D	Ѥ	CYRILLIC CAPITAL LETTER I WITH GRAVE ≡ 0418 И 0300 ѥ
040E	Ѧ	CYRILLIC CAPITAL LETTER SHORT U ≡ 0423 У 0306 Ѧ
040F	Ѣ	CYRILLIC CAPITAL LETTER DZHE

Basic Russian alphabet

0410	А	CYRILLIC CAPITAL LETTER A
0411	Б	CYRILLIC CAPITAL LETTER BE → 0183 б latin small letter b with topbar
0412	В	CYRILLIC CAPITAL LETTER VE
0413	Г	CYRILLIC CAPITAL LETTER GHE
0414	Д	CYRILLIC CAPITAL LETTER DE
0415	Е	CYRILLIC CAPITAL LETTER IE
0416	Ж	CYRILLIC CAPITAL LETTER ZHE

0428	Ы	CYRILLIC CAPITAL LETTER YERU
042C	Ь	CYRILLIC CAPITAL LETTERS OFT SIGN
042D	Ѣ	CYRILLIC CAPITAL LETTER E
042E	Ю	CYRILLIC CAPITAL LETTER YU
042F	Я	CYRILLIC CAPITAL LETTER YA
0430	а	CYRILLIC SMALL LETTER A
0431	б	CYRILLIC SMALL LETTER BE
0432	в	CYRILLIC SMALL LETTER VE
0433	г	CYRILLIC SMALL LETTER GHE
0434	д	CYRILLIC SMALL LETTER DE
0435	е	CYRILLIC SMALL LETTER IE
0436	ж	CYRILLIC SMALL LETTER ZHE
0437	з	CYRILLIC SMALL LETTER ZE
0438	и	CYRILLIC SMALL LETTER I
0439	й	CYRILLIC SMALL LETTER SHORT I ≡ 0438 и 0306 й
043A	к	CYRILLIC SMALL LETTER KA
043B	л	CYRILLIC SMALL LETTER EL
043C	м	CYRILLIC SMALL LETTER EM
043D	н	CYRILLIC SMALL LETTER EN
043E	о	CYRILLIC SMALL LETTER O
043F	п	CYRILLIC SMALL LETTER PE
0440	р	CYRILLIC SMALL LETTER ER
0441	с	CYRILLIC SMALL LETTER ES
0442	т	CYRILLIC SMALL LETTER TE
0443	у	CYRILLIC SMALL LETTER U
0444	ф	CYRILLIC SMALL LETTER EF
0445	х	CYRILLIC SMALL LETTER HA
0446	ц	CYRILLIC SMALL LETTER TSE
0447	ч	CYRILLIC SMALL LETTER CHE
0448	ш	CYRILLIC SMALL LETTER SHA
0449	щ	CYRILLIC SMALL LETTER SHCHA
044A	ъ	CYRILLIC SMALL LETTER HARD SIGN
044B	ы	CYRILLIC SMALL LETTER YERU → A651 ы cyrillic small letter yeru with back yer
044C	ь	CYRILLIC SMALL LETTER SOFT SIGN → 0185 ь latin small letter tone six → A64F ь cyrillic small letter neutral yer

Кодировки

Полезные инструменты

- CLD (Compact Language Detector) – C++, Python
 - <http://code.google.com/p/chromium-compact-language-detector/>
- LanguageDetection – Java
 - <http://code.google.com/p/language-detection/>
- Apache Tika
 - <http://tika.apache.org/>
- Видеолекция (Яндекс, RuSSIR 2012)
http://videlectures.net/russir2012_grigoriev_language/

Определение языка

Подходы

- ▶ Графематический
- ▶ N-граммный
- ▶ Лексический

Графематический подход

Система письменности

- ▶ Кириллица
- ▶ Латиница
- ▶ ...

Алфавит

- ▶ Русский А..Я
- ▶ Украинский - не используются Ё, Ъ, Ы, Э, но есть Г', І с точками и т.д.
- ▶ Казахский

N-граммный подход

Russian	Ukrainian	English	French
^п 1.91 ^по 0.84	^п 1.97 ^на 0.85	^t 3.17 ^th 2.00	es 2.31 es\$ 1.77
^с 1.71 ^пр 0.68	^в 1.75 на\$ 0.73	th 2.48 the 1.62	le 1.97 ^de 0.98
^в 1.68 ^на 0.66	^н 1.68 ^по 0.72	^a 2.41 he\$ 1.44	^d 1.84 le\$ 0.82
^н 1.55 ^и\$ 0.61	на 1.45 ^пр 0.63	he 2.24 ed\$ 0.78	^l 1.74 de\$ 0.76
ст 1.43 ^в\$ 0.60	^з 1.40 ^за 0.59	in 1.94 nd\$ 0.73	on 1.70 ^le 0.72
то 1.29 ^не 0.56	^с 1.25 ^не 0.56	er 1.60 ing 0.73	re 1.48 re\$ 0.68
но 1.23 ть\$ 0.48	ро 1.13 ого 0.54	an 1.54 ^an 0.72	^c 1.46 nt\$ 0.58

- ▶ Ранговый
- ▶ Марковский

Пословный подход

- ???
 - án került vagy től majd új ami ő kategória ben szerint amikor hogy amerikai két ezt mint alatt magyar itt második már
- ???
 - cel cod său cu cea l după ro va județul această în către sunt pe toate astfel ani prin ca departamentul din timpul într
- ???
 - ayrıca iklimi gibi tarafından olu kültür birlikte ula yol tarihinde veya iyi sonra türk bulunan kar çalı göre oldu

Пословный подход

- Hungarian
 - án került vagy től majd új ami ő kategória ben szerint amikor hogy amerikai két ezt mint alatt magyar itt második már
- Romanian
 - cel cod său cu cea l după ro va județul această în către sunt pe toate astfel ani prin ca departamentul din timpul într
- Turkish
 - ayrıca iklimi gibi tarafından olu kültür birlikte ula yol tarihinde veya iyi sonra türk bulunan kar çalı göre oldu

Кореференция: синонимы

Различные способы названия одного и того же объекта

- ▶ Синонимы: [“ШАВЕРМА”, “ШАУРМА”]
- ▶ Аббревиатуры: [“БМП”, “БОЕВАЯ МАШИНА ПЕХОТЫ”]
- ▶ Транслитерация: [“PLAZMA”, “ПЛАЗМА”]
- ▶ Грамматические замены: [“ПОЗДРАВЛЕНИЕ”, “ПОЗДРАВИТЬ”]
- ▶ Переводы: [“ВОЗДУШНАЯ ТЮРЬМА”, “CON AIR”]
- ▶ Джойны: [“АУДИО КОДЕКИ”, “АУДИОКОДЕКИ”]

Примеры

поиск@mail.ru айфон x q

Интернет Картинки Видео Приложения Новости Ответы

iPhone – Apple (RU)

apple.com/ru/iphone

Откройте для себя мир iPhone. Взгляните на iPhone 6s, iPhone 6 и iPhone SE. Зайдите на сайт Apple, чтобы изучить информацию, купить и получить поддержку.

поиск@mail.ru мгу x q

Интернет Картинки Видео Приложения Новости Ответы

Московский государственный университет имени М.В. Ломоносова

msu.ru

Информация о факультетах, институтах, центрах, руководстве МГУ. Учеба: высшее, дополнительное, дистанционное образование, практическое обучение...

Москва, м-р. Ленинские горы, д. 1 ☎ +7 (495) 939-10-00

Об университете

Поступающим

Образование

Университетская жизнь

поиск@mail.ru гиппопотам x q

Интернет Картинки Видео Приложения Новости Ответы



Обыкновенный бегемот — Википедия

ru.wikipedia.org/wiki/...

Обыкновенный бегемот, или гиппопотам — млекопитающее из отряда парнокопытных, подотряда свинообразных (неквачных), семейства бегемотовых, единственный современный вид рода Hippopotamus. Характерной особенностью бегемота является его полуводный образ жизни — большую часть времени он проводит...

Общая информация

Облик и строение

Название

Происхождение и систематика

Подвиды

Ареал и численность

Источники расширений

- ▶ Словари синонимов
- ▶ Энциклопедические сайты
 - ▶ Википедия
 - ▶ Тематические сайты (kinopoisk)
- ▶ Скобочные написания в документах
- ▶ Логи запросов с кликами
- ▶ Переформулировки запросов
- ▶ Грамматические преобразования
- ▶ Языковые модели и дистрибутивная семантика

Дистрибутивная семантика

Дистрибутивная гипотеза

- ▶ Значение лингвистической единицы складывается только из ее употребления, использования.
- ▶ В мозге хранится сумма всех тех контекстов, в рамках которых мы слышали или видели то или иное слово.
- ▶ Это и есть его смысл. Без знания типичных соседей никакой семантики нет.

Вывод:

- ▶ Слова с похожими типичными контекстами имеют схожее значение

Дистрибутивная семантика

Счетные модели

- ▶ Совмместная встречаемость
- ▶ Косинусная близость

Predictive models

- ▶ word2vec

Дистрибутивная семантика

Пример модели

► <http://ling.go.mail.ru/dsm/ru/>

Ближайшие по контекстам слова для *дефицит* (с косинусной близостью)

1. [нехватка](#) 0.65354729
2. [профицит](#) 0.62004328
3. [госдолг](#) 0.61431211
4. [несбалансированность](#) 0.59282905
5. [дефицитность](#) 0.58463109
6. [дисбаланс](#) 0.57497036
7. [недопоступление](#) 0.57044709

Ближайшие по контекстам слова для *сингл* (с косинусной близостью)

1. [альбом](#) 0.91364485
2. [донглей](#) 0.88309568
3. [пластинка](#) 0.85516018
4. [кавер](#) 0.84436756
5. [студийник](#) 0.83440500

Усечение окончаний

- ▶ Отсекаем самое длинное возможное окончание от слова, надеемся, что это не сильно ухудшит результат по сравнению с лемматизацией
- ▶ Для каждого языка свои таблицы окончаний
- ▶ Для некоторых языков другие аналогичные методы
- ▶ Результат выходит забавный :)
 - ▶ Сочи и сочиться приводятся к одной форме

Примеры

Sample text: Such an analysis can reveal features that are not easily visible from the variations in the individual genes and can lead to a picture of expression that is more biologically transparent and accessible to interpretation

Lovins stemmer: such an analysis can reveal features that are not easily visible from the variations in the individual genes and can lead to a picture of expression that is more biologically transparent and accessible to interpretation

Porter stemmer: such an analysis can reveal features that are not easily visible from the variations in the individual genes and can lead to a picture of expression that is more biologically transparent and accessible to interpretation

Paice stemmer: such an analysis can reveal features that are not easily visible from the variations in the individual genes and can lead to a picture of expression that is more biologically transparent and accessible to interpretation

Преимущества и недостатки

- ▶ Усечение окончаний улучшает одни запросы и ухудшает другие
- ▶ Дает в большинстве случаев увеличение полноты
- ▶ При этом часто ухудшает точность
- ▶ Алгоритм Портера определяет следующий класс эквивалентности
 - ▶ operate operating operates operation operative operatives operational
- ▶ Ясно что в запросах станет хуже
 - ▶ operational AND research
 - ▶ operating AND system
 - ▶ operative AND densistry

Лемматизация

Цель: привести все разные формы одного слова к начальной (каноничной).

Примеры:

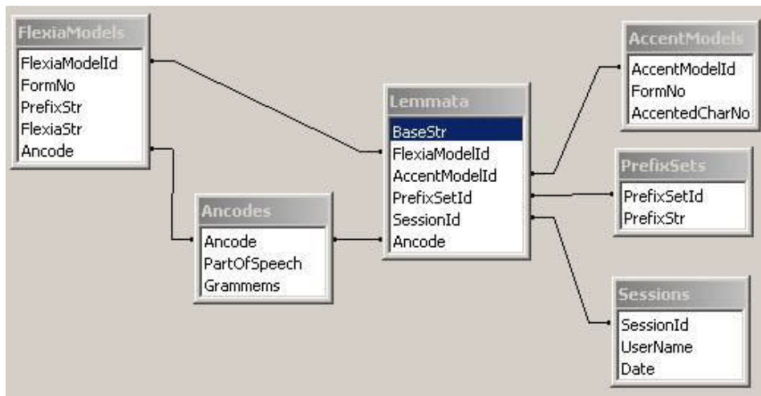
- ▶ am, are, is → be
- ▶ car, cars, car's, cars' → car
- ▶ the boy's cars are different colors → the boy car be different color

Лемматизация заключается в поиске начальной формы (леммы) в словаре

Лемматизация

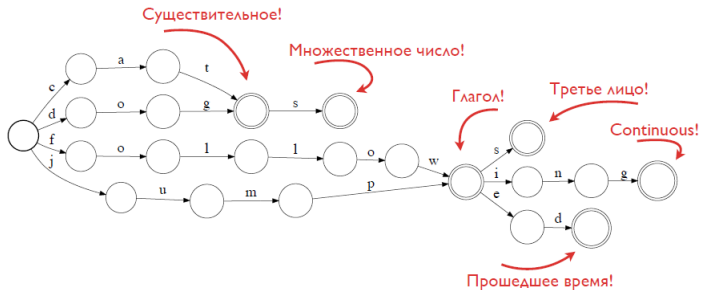
- Слово = машинная основа + парадигма
- Парадигма
 - Парадигм а
 - Парадигм ы
 - Парадигм е
 - Парадигм у
 - Парадигм ой
 - Парадигм е

Лемматизация



- (с) Сокирко А.В.,
<http://aot.ru/docs/sokirko/Dialog2004.htm>

Лемматизация



Лемматизация

- ▶ Что делать со словами, которых нет в словаре?
- ▶ Ищем похожие!

Предиктор

Input Your text:

микроблог

☐ English ☒ Russian ☐ German

☐ With paradigms

Submit Request

Found	Dict ID	Lemma	Grammems	
-	но,	МИКРОБЛОГ	С мр,вн,им,ед,	АНАЛОГ
-		МИКРОБЛГИЙ	КР_ПРИЛ но,од,мр,ед,	НЕДОЛГИЙ

Откуда взять морфологию?

Английская

- ▶ Стеммер Портера (Porter)
- ▶ Стеммер Ловинса (Lovins)

Русская

- ▶ aot.ru
- ▶ keva.ru (СтемКа)
- ▶ MyStem
(<http://company.yandex.ru/technologies/mystem/>)
- ▶ PyMorphy

Стеммер можно сделать самостоятельно

- ▶ Snowball - фреймворк для алгоритмов стемминга

Какие проблемы?

Языки сильно отличаются

- ▶ Изолирующие/Аналитические
- ▶ Синтетические
 - ▶ Флективные
 - ▶ Агглютинативные
- ▶ Полисинтетические

Флективные языки

Индо-европейские

- ▶ Русский
- ▶ Испанский
- ▶ Немецкий

Агглютинативные языки

Тюркские языки

Turkish	English
<i>ev</i>	(the) house
<i>evler</i>	(the) houses
<i>evin</i>	your (sing.) house
<i>eviniz</i>	your (pl./formal) house
<i>evim</i>	my house
<i>evimde</i>	at my house
<i>evlerinizin</i>	of your houses
<i>evlerinizden</i>	from your houses

Полисинтетические языки

Полисинтетические языки — языки, в которых все члены предложения или некоторые компоненты словосочетания соединяются в единое целое без формальных показателей у каждого из них.

- ▶ Чукотско-камчатские
- ▶ Эскимосско-алеутские

Тымэйнылевтпыгтыркын

(t-ə-mejŋ-ə-levt-pəyt-ə-rkən)

У меня сильно болит голова.

Омонимия

Разные по смыслу слова имеют одинаковое написание

Примеры:

белки бегали по лесу и ели орехи
(Лемма: белка, сущ. жен. род)

белки различаются по степени растворимости в воде
(Лемма: белок, сущ. муж. род)

- *Словоформа* - конкретная морфологическая разновидность слова

белка, белку, белкой, белке

Неоднозначность

Английский

- ▶ Leg
- ▶ Chair

Русский

- ▶ Лук
- ▶ Очки
- ▶ Лист

Снятие омонимии

Rule-based

```
//~ TN 26.1.a
//~ TC Если К - омоним с прилагательным или наречием
//~ ТА И справа - enough
//~ ТА И если омоним - переходный глагол, и справа от
enough не существительное
else if ((IsAdj(k) || IsAdverb(k))
    && CheckQuantitativeParticular(k + 1, QP_ENOUGH)
    && !CheckPrepParticular(k + 2, PP_OF)
    && !(IsTransitiveVerb(k) && IsNoun(k + 2)))
{
    //~ TD Тогда не глагол
    NOT_VERB(k);
}
```

Снятие омонимии

Rule-based

```
//~ TN 28.  
//~ TC Если k - омоним глагол/существительное и k не прилагательное и не глагол 'need',  
//~ TA и слева есть ИГ  
//~ TA и за ИГ следует глагол (не омоним)  
//~ TA который не управляет инфинитивом без to  
//~ TA не управляет прилагочным  
//~ TA и не является строго непереходным  
//~ TA и справа с пропуском наречий to + infinitive  
//~ TA Или K - соч. союз или предлог  
//~ TA Если K переходный И справа не начинается ИГ И справа не DO  
/* The Princess of Wales has opened a new Aids CENTRE in south-east London.  
else if (IsVerb(k) && IsNoun(k) && !IsAdj(k)  
    && !CheckVerbSemantic(k, VS_NEED) //~ не глагол 'need'  
    && (nTmp = SearchNGAtLeftBeg(LEFT_K, SINGLM_NORMAL))  
    && (nTmp2 = SkipAdvOmon(nTmp, 5, SAO_LEFT))  
    && (IsVerb(nTmp2)  
        && IsOnePartOfSpeech(nTmp2)  
        && !IsParticipleI(nTmp2)  
        && !IsVerbObjBareInfControl(nTmp2)  
        && CheckVGClaueGovernment(GetPrizm(nTmp2), VCG_NOT_SET)  
        && !IsVGStrictlyIntransitive(GetPrizm(nTmp2)))  
    && (nTmp = SkipAdvHomo(k + 1))  
    && (IsTo(nTmp) && IsInfinitive(nTmp + 1)  
        && !(IsVerbInfControl(k) && !IsNounInfControl(k))  
        || IsCoConj(k) || IsPrep(k)  
        || IsTransitiveVerb(k) && !NGCheck(NGM_BEGIN, NGO_ALL, k + 1) && !IsDo(nTmp)))  
{  
    //~ TD Тогда не глагол  
    NOT_VERB(k);  
}
```

Статистическое снятие омонимии

Скрытые марковские модели

Задача: Приписать наиболее вероятным образом каждому слову w_k в предложении тег t_k

По формула Байеса $P(A|B) = \frac{P(B|A)P(A)}{P(B)}$ получаем

$$P(t_1 \dots t_m | w_1 \dots w_m) = \frac{P(w_1 \dots w_m | t_1 \dots t_m) P(t_1 \dots t_m)}{P(w_1 \dots w_m)}$$

Нужно найти для данного предложения набор тегов, который делает эту вероятность максимальной.

Скрытые марковские модели

Предположение: Распределение тегов подчиняется марковскому свойству:

$$P(w_1 \dots w_m | t_1 \dots t_m) P(t_1 \dots t_m) \approx \prod_{k=1}^m P(w_k | t_k) P(t_k | t_{k-(n-1)} \dots t_{k-1})$$

Вероятности могут быть оценены как частоты по большому размеченному корпусу (<http://ruscorpora.ru>)

Для поиска оптимальных тегов используется динамический алгоритм Витерби.

Вопросы

