

Интеллектуальные информационные системы

Тема 2

Морфологический анализ

Павел Исаакович Браславский
pb@imach.uran.ru
весенний семестр 2006

План

- Зачем нужна морфология?
- Типы морфологической обработки в задачах ИП
- Подходы
- Морфологические модули для русского языка

Зачем нужна морфология?

- Классы эквивалентности ключевых слов при поиске:

кошка, кошки, кошку, кошкой, кошке...

- Последующая обработка (синтаксический анализ, семантический анализ...)

Типы анализа

- Стемминг – выделение основы
лесной, лес, лесистый, леса → лес
система, системный, систематизировать → систем
- Приведение к словарной форме
лесного, лесному → лесной
леса → лес
танцующая → танцевать

Типы анализа – 2

- POS-tagging (part-of-speech)

Танцующая <V> в <PREP> темноте <N>

- Полная морфологическая информация

Танцующая <V, прич, несоверш, наст., ед., жен., им.> в <PREP> темноте <N, жен., неод., ед., предл.>

Части речи	Грамматические категории (в скобках приведены сокращенные названия их значений)
Существительное (сущ.)	Род (м., ж., ср., р), число (ед., мн.), падеж (им., род., дат., вин., тв., пр.), одушевленность (од., неод., о)
Полное прилагательное (прил. полн.)	Пассивность (пасс., акт., п), время (прош., наст., буд.), род (м., ж., ср., р), число (ед., мн.), падеж (им., род., дат., вин., тв., пр.), одушевленность (од., неод., о), вид (сов, псв, вид)
Краткое прилагательное (прил. кр.)	Пассивность (пасс., акт.), время (прош., наст., буд.), род (м., ж., ср., р), число (ед., мн.), вид (сов., псв., вид)
Глагол (глагол.)	Пассивность (акт., пасс.), время (прош., наст., буд.), род (м., ж., ср., р), число (ед., мн.), вид (сов, псв, вид)
Инфинитив (инф.)	Пассивность (акт., пасс.), род (м., ж., ср., р), число (ед., мн., ч)
Деепричастие (деепр.)	Пассивность (акт., пасс.), время (прош., наст.), вид (сов, псв, вид)
Наречие (нареч.)	Тип: обстоятельное (обст.), определительное (опр.)
Количественное числительное (числ.)	Тип: «1», «2», «5», дробное (дробн.), неопределенное (неопр.), именованное (именов.)
Местоимение (мест.)	Класс: притяжательное (прит.), указательное (указ.), возвратное (возвр.), возвратно-атрибутивное (возвр.-атр.), третьего лица (3 л.); падеж (им., род., дат., вин., тв., пр., п), число (ед., мн., ч), род (м., ж., ср., р)
Союз	Тип: сочинительный (соч.), подчинительный (подч.)
Предлог (предл.)	Падеж (род., дат., вин., тв., пр.)
Частичка	Тип: вопросительная (вопр.), отрицательная (отр.)
Синтаксический знак (синт. зн.)	—

Попов, 1982

Грамматическая омонимия

Объект анализа – отдельное слово →
неоднозначность

падали → падаль (сущ.), падать (гл.)

печь → печь (сущ.), печь (гл.)

черепах → череп (сущ., муж. род), черепаха (сущ., жен. род.)

стекла → стекло (сущ.), стекать (гл.)

ученый → учить (гл.), ученый (сущ.)

Английский: 1,2 – 1,5 тэга на словоформу

Разрешение – частота, учет контекста,
синтаксис

Пример: mystem

Он сделал это так неловко, что задел образок моего ангела, висевший на дубовой спинке кровати, и что убитая муха упала мне прямо на голову.

Л.Н.Толстой, «Детство»

Он{он=S, сред, неод=(им, ед|им, мн|род, ед|род, мн|дат, ед|дат, мн|вин, ед|вин, мн|твор, ед|твор, мн|пр, ед|пр, мн)|он=S, ед, муж, од=им}

сделал{сделать=V, сов=прош, ед, изъяв, муж}

это{это=S, ед, сред, неод=(им|вин)|этот=A=(им, ед, сред|вин, ед, сред)|это=PART=}

так{так=ADV=|так=PART=|так=CONJ=}

неловко{неловкий=A=ед, кр, сред|неловко=ADV=}

что{что=CONJ=|что=S, ед, сред, неод=(им|вин)}

задел{задевать=V=прош, ед, изъяв, муж, сов|задел=S, муж, неод=(им, ед|вин, ед)}

образок{образок=S, муж, неод=(им, ед|вин, ед)}

моего{мой=A=(род, ед, муж|род, ед, сред|вин, ед, муж, од)}

ангела{ангел=S, муж, од=(род, ед|вин, ед)}

висевший{висеть=V, несов=(прош, им, ед, прич, муж|прош, вин, ед, прич, муж, неод)}

на{на=PR=|на=PART=}

дубовой{дубовый=A=(род, ед, жен|дат, ед, жен|твор, ед, жен|пр, ед, жен)|дубова=S, жен, од=(род, ед|дат, ед|твор, ед|пр, ед)}

спинке{спинка=S, жен, неод=(дат, ед|пр, ед)}

2006

П.И. Браславский - Интеллектуальные ИС

9

Методы

- Процедурный
- Табличный
- Статистический
- Различные комбинации

Алгоритм Портера

- Самый распространенный стеммер для английского языка
- 5 циклов усечения
- Каждый цикл – набор команд
- В первую очередь выполняется операция над самым длинным суффиксом

Алгоритм Портера – фрагмент

- *sses* → *ss*
- *ies* → *i*
- *ational* → *ate*
- *tional* → *tion*

- Weight of word sensitive rules
- $(m > 1)$ *EMENT* →
 - *replacement* → *replac*
 - *cement* → *cement*

Manning&Raghavan, 2005

Морфологические типы существительных

	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16
Им. ед.	∅	∅	∅	а	я, ья	я	я	ь	ь, й	й	о	о	я	е	е, ье	е
Род. ед.	а	а	а	ы, п	и, ьи	и	и	и	и, я	я	а	а	и	а	я, ья	я
Дат. ед.	у	у	у	е	е, ье	е	е	е	ю	ю	у	у	и	у	ю, ью	ю
Вин. ед.				у	ю, ью	ю	ю	ь								
Тв. ед.	ом, ем	ом	ом	сй, ой	ей, бей	сй	ей	ью	ем	ем	ом	ем	ем	ем	ем, ьем	ем
Предл. ед.	е	е	е	е	е, ье	е	п	п	е	и	е	е	п	е	е, ье	и
Им. мн.	ы, п, ья	ы, п	а	ы, п	и, ьи	я	п	и	п	п	а	п, ья	а	а	я, ье	я
Род. мн.	ов, ев, ей, ьев	∅	ов	∅	ей	й, 1., ∅	й	ей	ей, ев	ев	∅, ов	∅, ов, ьев	∅	∅, ев	ий, ей	й
Дат. мн.	ам, ьям	ам	ам	ам	ям, ьям	ям	ям	ям	ям	ям	ам	ам, ьям	ам	ам	ям, ьям	ям
Вин. мн.																
Тв. мн.	амн, ьямн	амн	амн	амн	ьямн, ямн	ямн	ямн	ямн	ямн	ямн	амн	ьямн, амн	амн	амн	ямн, ьямн	ямн
Предл. мн.	ах, ьях	ах	ах	ах	ьях, ях	ях	ях	ях	ях	ях	ах	ах, ьях	ах	ах	ях, ьях	ях
Примеры	Профес, брат	Граам	Номер	Матрица	Ступня, статья	Идея	Липня	Милень	Забой, уголь	Салаторий	Тело	Дерево	Время	Пологале	Поле, устье	Валище

Процедурный подход

- Словарь основ
 - Словарь готовых форм (СГФ)
 - Поиск в СГФ
 - Выделение основы
 - Поиск основы в словаре
- (см. Попов, 1982, с. 234-235)

Табличный подход

- волка → волк (муж., од.; ед.ч., [р.п.|в.п.]
- не → не (частица)
- корми → кормить (несоверш.; повел. накл., ед. ч.)
- в → в (предлог)

Как сформировать таблицу?

Зализняк А.А. Грамматический словарь
русского языка

~100 тыс. входов

Модель русского словоизменения

Пример: лев мо 1*b (животное)

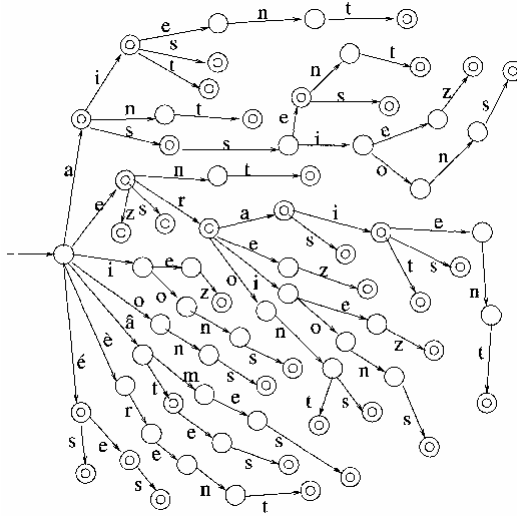
лев м 1а (денежная единица)

стричь нсв 8b (-г-)

прихожая ж (п 4а)

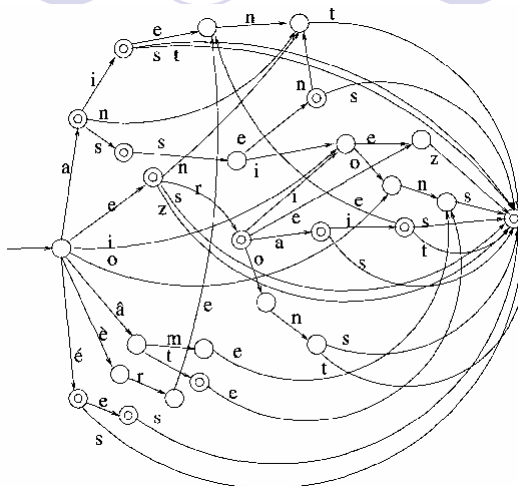
Основа большинства машинных морфологий РЯ
Автомобилестроения (мн.ч.), деревянное (ср.), при → пря

Как хранить? – Trie



Daciuk et al., 2000

Как хранить лучше? – FSA



Finite
State
Automata

Daciuk et al., 2000

Статистический стеммер

словарями → словарь → словар-ями → ар-ями
топорами → топор → топор-ами → ор-ами
летающего → лететь → лет-ящего → ет-ящего
летающего → летящий → летящ-его → ящ-его

+ правило: одна гласная в
основе

имя	ра	546
има	ро	154
огещя	те	12
оге	щя	12