

АДАПТИВНЫЙ СИНТАКСИЧЕСКИЙ АНАЛИЗАТОР

И. А. Волкова, М. Г. Мальковский, Н. В. Одинцев
Московский государственный университет им. М. В. Ломоносова
volkova@cs.msu.su, malk@cs.msu.su, nickone@yandex.ru

Ключевые слова: корпус текстов, модель управления, синтаксический анализатор.

В рамках концепции адаптивных лингвистических процессоров исследуется задача создания адаптивного синтаксического анализатора. Предлагаемый синтаксический анализатор использует описание языка в виде моделей управления, настроенное на работу в требуемой предметной области и полученное в результате анализа корпуса текстов этой предметной области. Для первоначального анализа корпуса текстов и осуществления такой настройки используется отдельный синтаксический анализатор, работающий с описанием языка в виде сетевой грамматики. В ходе этого анализа каждой выделяемой модели управления приписывается частотность, характеризующая вероятность использования этой модели управления для новых текстов данной предметной области.

Преимуществами адаптивного синтаксического анализатора являются высокая точность описания языка предметной области, высокая скорость работы (точнее, скорость получения первого, наиболее вероятного варианта синтаксической структуры), а также возможность работы с неграмматичными конструкциями. Кроме того, знание общих стилистических особенностей различных текстов позволяет решать ряд практических проблем, например, проблемы определения авторства текста.

1. Введение

Понятие адаптивности [1] системы подразумевает возможность настройки на новые области применения (адаптируемость) и способность к самостоятельному пополнению базы лингвистических и проблемно-ориентированных знаний в рабочем режиме (адаптивность и обучаемость).

Как правило ([1] и др.), адаптация анализатора подразумевает пополнение его словаря, пополнение набора используемых моделей управления, смягчение условий проверки грамматических признаков и т. д.; алгоритмы работы при этом принципиально не изменяются. Важным отличием предлагаемого анализатора является его разделение на два компонента, один из которых использует только базовые лингвистические знания, и результатом его работы является адаптированное для данной предметной области описание языка, а другой компонент использует это адаптированное описание, представляющее собой открытые лингвистические знания.

2. Базовые лингвистические знания и их использование

Источником данных для адаптации системы являются результаты анализа достаточно больших корпусов текстов. Базовые знания о синтаксисе языка, используемые для такого рода анализа, описаны в виде сетевой грамматики - расширенной сети переходов (РСП) [2, 3]; для работы с ними предназначен отдельный компонент системы (РСП-анализатор).

РСП-анализатор состоит из сетевой грамматики, описывающей входной язык, и программы, выполняющей анализ текста. Таким образом, сама программа не зависит от входного языка.

Ниже приводится пример одной подсети (описывающей именную группу с главным словом – существительным).

Именная группа с главным словом - существительным

```

$NPNoun
(
0: $AP 0 @ToLocGender @ToLocNumber @ToLocCase
   $Pt1 0 @ToGlobVp1SubjectInfo @ToLocGender @ToLocNumber @ToLocCase
   $NPPersonalPronoun 2 @CheckNonAccusCase
   $Adv 0 @CheckNoAdvO
   <> 1
;
1: <Сущ> 2 @ToLocNpFromNoun
;
2: $AP 2 @ToLocGender @ToLocNumber @ToLocCase
   $Pt1 2 @ToLocGender @ToLocNumber @ToLocCase @CheckSimplePt1
   $NP 2 @CheckNonAccusCase
   $PP 2
   $Adv 2 @CheckNoAdvO
   $Inf 2
   <Запятая> 3
   <Кавычки> 5
   <> *
;
3: $Pt1 4 @ToLocGender @ToLocNumber @ToLocCase @CheckComplexPt1
;
4: <Запятая> 2
   <Запятая> *
;
5: $NP 6 @CheckNomCase
;
6: <Кавычки> 2
) @ToGlobGender @ToGlobNumber @ToGlobCase @ToGlobAnimate @ToGlobPerson

```

В каждом состоянии (в приведенном примере – 0, 1, 2, 3, 4, 5 и 6) ожидается появление одной из указанных подсетей или терминальных символов (в состоянии 0 - \$AP, \$Pt1, \$NPPersonalPronoun, \$Adv, т. е. адъективной, причастной, именной с главным словом – личным местоимением и адвербиальной группы). Анализ начинается с состояния 0. При успешном разборе подсети осуществляется проверка выполнения контекстных условий (операторы с префиксом «@»), и переход в указанное состояние. Переход в состояние «*» соответствует успешному анализу описываемой группы при условии выполнения контекстных операторов.

Результатом анализа является синтаксическая структура анализируемого фрагмента, представленная в виде дерева. Например, при анализе именной группы *конечного числа возможностей* мы получим следующую структуру:

```

$NP (
  $NPSimple (
    $NPNoun (
      $AP (

```

\$APSimple (
 \$APQualAdj (
 конечного <Прил, Средн, Ед, Род>)))
 числа <Сущ, Неод, Средн, Ед, Род>
 \$NP (
 \$NPSimple (
 \$NPNoun (
 возможностей <Сущ, Неод, Жен, Множ, Род>))))

3. Открытые лингвистические знания

Формализм, используемый для хранения открытых лингвистических знаний, должен, во-первых, предоставлять возможность фиксировать непосредственные зависимости между словами и, во-вторых, задавать частотные характеристики, определяющие использование таких зависимостей. Выполнение этих условий позволило бы осуществлять адаптацию описания языка для различных категорий текстов, однако РСП этим условиям не удовлетворяет.

Предлагается описывать синтаксис языка с помощью моделей управления (МУ). Введем следующую иерархию МУ [4]:

- модели управления 3-го уровня – атомарные модели управления, связывающие конкретное слово и словоформу (например, {*обогнать* {сущ., вин. п., ж. р., ед. ч., неод., *машину*} {частотная характеристика}});
- модели управления 2-го уровня, связывающие слово с определенной зависимой группой, на главное слово которой налагаются ограничения (например, {*обогнать* {сущ., вин. п.} {частотная характеристика}});
- модели управления 1-го уровня, связывающие абстрактное (принадлежащее некоторому классу) главное слово группы с зависимой группой (например, {глагол {сущ., вин. п.} {частотная характеристика}});
- модели управления 0-го уровня, связывающие описываемую группу с составляющими ее группами (например, {предложение {сущ., им. п.} {частотная характеристика}}).

Каждой обобщенной модели управления (вне зависимости от ее уровня) соответствует частотная характеристика, значение которой вычисляется при анализе корпуса текстов. В дальнейшем эта характеристика используется синтаксическим анализатором, основанным на МУ.

Для каждой категории текстов частотные характеристики (да и сами модели управления) будут отличаться.

4. Адаптация описания языка

Процесс адаптации системы для работы с определенной категорией текстов состоит из следующих стадий:

Выбор корпусов текстов. На этой стадии выбираются тексты, суммарный объем которых достаточен для поставленных целей. В некоторых случаях вопрос полноты корпуса тестов является тривиальным, например, для анализа произведений какого-либо писателя можно взять все его произведения, в других случаях следует использовать корпус возможно большего объема.

Разбиение текстов на предложения или другие более мелкие фрагменты. РСП-анализатор способен работать с простыми предложениями и некоторыми типами сложных. Таким образом, исходный текст может быть разбит на фрагменты, являющиеся предложениями или их составными частями. Размер фрагментов выбирается из следующих соображений: чем меньше фрагмент, тем быстрее и точнее его можно разобрать, но при этом повышается вероятность отсутствия в нем интересующей нас пары управляющего и управляемого слова. К сожалению, время работы анализатора, основанного на РСП, фактически является экспонентой от количества лексем во фрагменте. На практике это приводит к тому, что на анализ предложения из 20-25 слов может быть затрачено 20-25 минут (если оно вообще будет разобрано), что вынуждает в качестве объекта анализа выбирать более мелкие, чем предложение, фрагменты. При этом возникает другая проблема – как выделить из предложения синтаксически замкнутые фрагменты, не зная априори его синтаксической структуры? Для решения этой проблемы приходится использовать эвристики, позволяющие правдоподобно расставить границы фрагментов.

Синтаксический анализ фрагментов с помощью анализатора, основанного на РСП. На этой стадии формируется список синтаксических структур фрагментов, при этом каждому фрагменту может соответствовать несколько структур.

Анализ синтаксических структур фрагментов, выделение из них моделей управления и загрузка их в базу данных. Получаемые здесь модели управления – это модели управления 3-го уровня.

Автоматизированная обработка моделей управления. Целью этой обработки является, во-первых, удаление «неверных» моделей управления, и, во-вторых, получение более общих МУ. Удаление «неверных» моделей управления осуществляется в автоматизированном режиме: можно просмотреть МУ в поисках наиболее вероятных кандидатов на удаление с учетом частотных характеристик. Здесь нужно, правда, заметить, что если "неверная" модель редко встречается, то ее, в принципе, можно и не удалять, поскольку вероятность ее последующего применения, как функция от частоты использования, будет минимальной. Также для каждой МУ можно посмотреть примеры ее использования.

Заметим, что совершенно необязательно добиваться успешного и правильного анализа каждого фрагмента предложения (не говоря уже о том, что если бы можно было бы требовать абсолютно правильной работы от РСП-анализатора, то незачем было бы создавать анализатор, основанный на МУ). Некоторые формально правильные фрагменты могут быть отброшены, но, подчеркнем, задачей является не сам анализ как таковой, а выявление общих закономерностей использования сочетаний слов, представляемых моделями управления. Важно лишь не исключить слишком строгими правилами в РСП-описании языка целые классы допустимых словосочетаний.

Была проведена автоматическая обработка различных корпусов текстов, в частности, романов Достоевского [5], а также архива газеты *Известия* за июнь 1997 года [6]. В качестве примера можно привести результаты автоматического выделения из текстов моделей управления для слова *находить* (таблица 1):

Таблица 1. Модели управления уровня 3 для глагола *находить*.

Уровень	Главное слово	Предлог	Падеж	Источник	Зависимое слово	Суммарная оценка
3	находить	в	Пред	Известия.97.6.2	детдоме	1.000008
3	находить	в	Пред	Известия.97.6.3	машинах	1
3	находить		Твор	Известия.97.6.1	порой	0.666668
3	находить	в	Пред	Известия.97.6.3	скалах	0.500004
3	находить		Вин	Известия.97.6.3	проходы	0.500004
3	находить		Вин	Известия.97.6.1	транспорт	0.5
3	находить		Вин	Известия.97.6.2	год	0.333336

После автоматической обработки моделей управления 3-го уровня можно получить модели управления 2-го уровня (таблица 2):

Таблица 2. Модели управления уровня 2 для глагола *находить*.

Уровень	Главное слово	Предлог	Падеж	Источник	Суммарная оценка
2	находить	в	Пред	Известия.97.6.a	3,000012
2	находить		Вин	Известия.97.6.a	1,33334

В любой момент времени доступна ручная коррекция данных. Например, из полученных моделей управления 2-го уровня была удалена модель управления {*находить* {сущ., твор. п.} {0.666668}} – результат интерпретации словоформы *порой* как существительного в творительном падеже, а не как наречия.

5. Анализатор, использующий модели управления

Главным преимуществом МУ-анализатора является используемый способ описания языка множеством моделей управления, от алгоритма анализа требовалось в первую очередь адекватное использование этого преимущества.

Итоговый алгоритм анализа фактически является алгоритмом поиска по дереву, использующим сильные эвристики. Таким образом, алгоритм должен допускать возможность существования нескольких альтернатив на каждом этапе анализа, кроме того, должен существовать механизм, обеспечивающий выбор наиболее приоритетного варианта.

Для описания принципов работы анализатора введем понятие точки разбора. Допустим, на некотором этапе анализа мы находимся в определенном состоянии - точке разбора - полученном в результате последовательного применения

моделей управления (в частности, связывания управляемого и управляющего слов). Допустим также, что в этой точке возможно связать несколько еще не разобранных слов с уже разобранными словами с помощью моделей управления. Таким образом, каждому возможному применению моделей управления уже разобранных слов соответствует новая точка разбора, и альтернативы на каждом этапе анализа – это применение различных обобщенных моделей управления.

Механизм, обеспечивающий порядок выбора допустимых моделей управления, т. е. определяющий, какую точку разбора следует строить в первую очередь, обеспечивается существованием частотной характеристикой каждой модели управления, полученной при автоматическом формировании множества МУ.

Строя точки разбора вышеописанным способом, мы получаем древовидную структуру, в корне которой находится неразобранное предложение, а в каждой вершине – точка разбора с разобранными и неразобранными словами. Ее дочерними вершинами являются точки разбора, полученные в результате однократного применения какой-либо модели управления. Лист дерева, в котором все слова разобраны, является искомым вариантом анализа.

Заметим, что базовые лингвистические знания непосредственно не используются. Тем не менее они присутствуют в виде моделей управления, полученных при анализе корпусов текстов.

6. Заключение

Преимуществами адаптивного синтаксического анализатора являются высокая точность описания языка предметной области, высокая скорость работы (точнее, скорость получения первого, наиболее вероятного варианта синтаксической структуры), а также облегчение работы с неграмматичными конструкциями. Кроме того, знание общих стилистических особенностей различных текстов позволяет решать ряд практических проблем, например, проблемы определения авторства текста.

Литература

1. Мальковский М.Г. Диалог с системой искусственного интеллекта. М.: МГУ, 1985.
2. Вудс В.А. Сетевые грамматики для анализа естественных языков. // Кибернетический сборник. Новая Серия. Вып. 13. М.: Мир, 1978. с. 120-158.
3. Волкова И.А., Головин И.Г. Синтаксический анализ фраз естественного языка на основе сетевой грамматики. // ДИАЛОГ'98, Труды межд. семинара. М., 1998. с. 438-447
4. Одинцев Н.В. Обобщенные модели управления. Синтаксический анализатор на основе обобщенных моделей управления. // ДИАЛОГ'2002, Труды межд. семинара. М., 2002. с. 401-406.
5. Достоевский Ф.М., <http://www.lib.ru/LITRA/DOSTOEWSKIJ>
6. Машинный фонд русского языка, <http://irlras-cfirl.rema.ru/>