

УДК 655.25+601.4

АВТОМАТИЧЕСКОЕ ВЫЯВЛЕНИЕ СТРУКТУРЫ ПАРАДИГМАТИЧЕСКИХ
ОТНОШЕНИЙ ТЕКСТА

М.К.Тимофеева

I. Известно, что каждая единица языковой системы противопоставляется некоторой другой единице, вступая с ней в парадигматические отношения. В связи с этим представляется интересным автоматическое построение формальной модели парадигматики произвольного текста, представляющего некоторый язык. Предлагаемый нами метод использует наиболее универсальные свойства известных языковых систем и состоит в автоматическом выявлении конкретных средств реализации этих свойств в предложенных текстах. Возможность полной автоматизации процесса обусловлена выбором таких универсальных свойств, которые опираются не на семантику текста, а на его статистико-комбинаторные характеристики. Это, в свою очередь, приводит к тому, что выявляются не все, а только наиболее типичные свойства парадигматики текста. Степень адекватности модели действительной системе парадигматических отношений языка зависит от величины и информативности текста, его представляющего.

Итеративное применение метода позволяет строить иерархию отношений единиц текста. Исходный материал для первой итерации - 1-граммный спектр текста [I] (частотный словарь встречаемости в нем подцепочек различной длины), для каждой последующей - 1-граммный спектр текста, полученного из предыдущего путем замены каждой выделенной в нем единицы кодом соответствующей ей парадигмы.

Предлагаемый метод анализа текста может быть полезен при разработке формы представления текста для работы со стенорайтерами, при автоматизации редакционно-издательских работ, при сжатии информации, хранимой в памяти ЭВМ, а также при типологическом изучении языков и выявлении тенденций их развития.

2. Известен ряд работ [2,3], предлагающих формальные методы выделения значимых подцепочек текста на основе использования универсальных свойств языков.

В [2] выделение значимых подцепочек текста производится с помощью анализа спектра 1-грамм, но задача установления парадигматических отношений между этими подцепочками не ставится. Выделенное множество значительно шире множества грамматических единиц текста. При больших объемах текстов производить отбор таких единиц вручную было бы затруднительно. Для автоматизации этого процесса необходимо ввести дополнительные критерии выделения грамматических единиц, учитывающие отношения между ними.

В [3] ставится задача выявления парадигматических отношений, но накладываются ограничения на допустимые позиции единиц и их длину. Основным критерием выделения единиц служит предположение о том, что в каждой парадигме найдется хотя бы одна единица, для которой условная вероятность встречаемости в тексте (в заданной позиции или в заданном окружении) существенно отличается от независимой. Введение такого критерия значительно уменьшает перебор подцепочек текста в ходе выделения единиц, но при снятии ограничений на позицию и длину этих единиц его использование становится трудоемким. Кроме того, добавление еще одного критерия, не связанного с выявлением отношений, несколько сужает область применения метода, так как указанное предположение может выполняться не для всех языков.

Благодаря использованию закономерности, обнаруженной на спектрах 1-грамм, появилась возможность достаточно быстрой организации перебора подцепочек текста. Основным критерием выделения единиц - наличие парадигматических отношений между ними. Остальные критерии служат для отбрасывания тех единиц, которые в традиционной лингвистике обычно не рассматриваются как самостоятельные. После окончания работы алгоритма нахождение отброшенных единиц, если это необходимо, не представляет затруднений.

3. Рассмотрим свойства языковых систем, положенные в основу метода. Пусть задан текст T в алфавите A , через A^* обозначим множество всех цепочек, состоящих из символов алфавита A ; Λ - пустая цепочка, W - множество всех подцепочек T , обладающих заданным набором формальных признаков. Причем W должно быть таким, чтобы для любых двух $\omega_1, \omega_2 \in W$ не существовало непустых $g, \omega'_1, \omega'_2 \in A^*$ таких, что $\omega_1 = \omega'_1 g$, $\omega_2 = g \omega'_2$ и $\omega'_1 g \omega'_2$ - подцепочка T .

Множество $q_{f_1, f_j} \in W$ образует квадрат^{*)} для $f_1, f_j \in A^*$, если

1) $f_1 \neq f_j, f_1 \cdot f_j \neq \Lambda$;

2) найдутся такие непустые и несовпадающие цепочки, $\Delta_1 \tilde{X}_1, \Delta_2 \tilde{X}_2 \in A^*$, что q_{f_1, f_j} представимо в виде:

$$q_{f_1, f_j} = \{ \Delta_1 f_1 \tilde{X}_1, \Delta_1 f_j \tilde{X}_1, \Delta_2 f_1 \tilde{X}_2, \Delta_2 f_j \tilde{X}_2 \}.$$

Упорядоченные пары $(\Delta_1, \tilde{X}_1), (\Delta_2, \tilde{X}_2)$ - контексты для f_1, f_j ; $K(f_1, f_j)$ - множество всех контекстов для f_1, f_j ; $K(f_1) = \cup K(f_1, f_j)$ - множество контекстов для f_1 ; $k_{f_1, f_j} = |K(f_1, f_j)|$; S - некоторая заданная константа больше 2.

Цепочки $f_1, f_j \in A^*$ взаимозаменяемы, если $k_{f_1, f_j} \geq S$.

Подцепочка f_1 текста T участвует в системе парадигматических отношений, если найдется цепочка f_j , с которой она взаимозаменяема. Обозначим через F множество всех таких подцепочек текста T и через P_1, \dots, P_n - максимальные множества попарно взаимозаменяемых элементов из F .

Для каждого $f_1 \in F$ упорядочим некоторым образом контексты $K(f_1)$. Вхождением f_1 - это пара (f_1, N) , где N - номер контекста из $K(f_1)$.

Рассмотрим такие вхождения $(f_1, N_1), (f_j, N_2)$ разных цепочек f_1, f_j , для которых $\Delta_1 f_1 \tilde{X}_1 = \Delta_2 f_j \tilde{X}_2$, где (Δ_1, \tilde{X}_1) - N_1 -й контекст из $K(f_1)$, (Δ_2, \tilde{X}_2) - N_2 -й контекст из $K(f_j)$, $n_1 = |\Delta_1|$, $n_2 = |\Delta_1 f_1|$, $n_3 = |\Delta_2|$, $n_4 = |\Delta_2 f_j|$. Вхождение (f_1, N_1) вложено во вхождение (f_j, N_2) , если $n_1 \geq n_3, n_2 \leq n_4$. Вхождения (f_1, N_1) и (f_j, N_2) пересекаются, если либо $n_3 \leq n_1 < n_4$, либо $n_3 < n_2 \leq n_4$.

Элемент $f_1 \in F$ минимален^{**)}, если для любого P_r , содержащего f_1 , не существует таких цепочек $g, h \in A^*$ ($g \cdot h \neq \Lambda$), что любое $f_j \in P_r$ представимо в виде $f_j = g f_1' h, f_1' \in A^*$.

Далее, обозначим через $K(P_r)$ множество контекстов, поставленное в соответствие множеству P_r и состоящее из таких контекстов (Δ_1, \tilde{X}_1) , для которых существует хотя бы одно $f_1 \in P_r$, такое, что $(\Delta_1, \tilde{X}_1) \in K(f_1)$ и f_1 не входит ни в одно из $P_j, j \neq r$.

*) Это понятие опирается на обобщенный "квадрат Гринберга" [4].

**) Сходные свойства минимальности используются в [3].

Контексты множества $K(f_i)$ минимальны, если для любого P_r , содержащего f_i , не существует таких цепочек $g, h \in A^*(g \cdot h \neq \Lambda)$, что каждый контекст $(\Delta_1, \chi_1) \in K(P_r)$ представим в виде $(\Delta'_1 g, h \chi'_1)$, где $\Delta'_1, \chi'_1 \in A$ и элементы множества $P'_r = \{gf_i h | f_i \in P_r\}$ минимальны.

Грамматические единицы - минимальные вступающие в парадигматические отношения подцепочки текста T , входящие которых не пересекаются и контексты минимальны.

Парадигмы - максимальные множества попарно взаимозаменяемых грамматических единиц текста T .

Омонимы - совпадающие по написанию грамматические единицы, входящие в разные парадигмы.

Качество выделения грамматических единиц зависит от выбора W и S . В первом варианте алгоритма предполагается, что в алфавите A выделено некоторое подмножество разделителей R , содержащее пробел и знаки препинания. Параметр S рассматривается как равный 2. Множество W состоит из всех подцепочек T , расположенных в нем между двумя символами из R и не содержащих их внутри себя. Если R неизвестно, то при выборе W используются спектры 1-грамм текста T . Для выработки критериев построения W в общем случае эксперименты проводятся на спектрах 1-грамм, хотя при известном R можно было бы использовать более простые данные. Алгоритмы анализа элементов W не зависят от способа его построения.

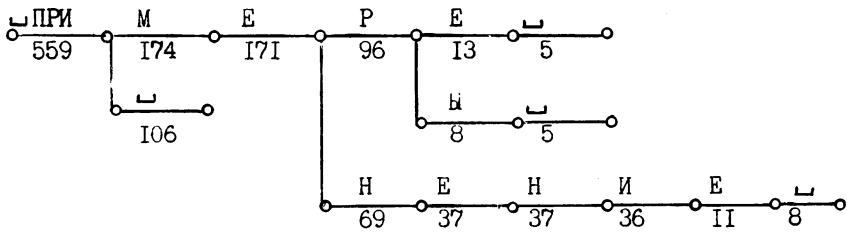
4. Анализ текста осуществляется посредством итеративной процедуры, на каждом шаге которой выделяются только те грамматические единицы, для которых не существует вложенных вхождений других грамматических единиц. Кодируются некоторым образом парадигмы, обнаруженные на каждом шаге. Производится обращение к алгоритму выявления минимальных общих признаков элементов каждого множества $K(P_r)$, отличающих их от элементов любого другого множества $K(P_i)$ [6]. Множества $K(P_r)$ дополняются контекстами из $K(f_i)$, которые по выделенным признакам могут быть однозначно отнесены к какой-либо парадигме. Все $\Delta_1, \chi_1, (\Delta_1, \chi_1) \in K(P_r)$ в тексте заменяются на цепочки $\Delta_1 P_r \chi_1$, где P_r - код парадигмы P_r . Указанное свойство элементов множества W обеспечивает однозначность такой перекодировки. На следующей итерации анализируется перекодированный текст.

Спектр 1-грамм текста и элементы W упаковываются в виде деревьев, построение которых можно иллюстрировать следующим примером.

Задан спектр 1-грамм:

- 1 = 4 ПРИ - 559
- 1 = 5 ПРИМ - 174; ПРИ - 106
- 1 = 6 ПРИМЕ - 171
- 1 = 7 ПРИМЕР - 96; ПРИМЕН - 69
- 1 = 8 ПРИМЕРЕ - 13; ПРИМЕР. - 5; ПРИМЕРЫ - 8; ПРИМЕНЕ - 37
- 1 = 9 ПРИМЕРЕ - 5; ПРИМЕРЫ, - 4; ПРИМЕНЕН - 37
- 1 = 10 ПРИМЕНЕНИ - 36
- 1 = 11 ПРИМЕНЕНИЕ - 11
- 1 = 12 ПРИМЕНЕНИЕ - 8

Ему ставим в соответствие дерево:



Если в результате анализа такого дерева грамматические единицы не выявлены, то рассматривается обратное дерево, в котором элементы W упакованы не с начала, как в данном, а с конца.

Звено дерева - путь ρ , образованный вершинами v_1, \dots, v_m , в котором $m \geq 2$ и

- 1) из v_1 выходит хотя бы одна дуга, ведущая не в v_2 ;
- 2) если из v_r ($1 < r < m$) выходит только одна дуга, то путь, образованный вершинами v_r, \dots, v_m , - линейный;
- 3) путь ρ не может быть продолжен без потери двух первых свойств.

Используются следующие способы ускорения просмотра деревьев в ходе выявления грамматических единиц.

1) Так как $s \geq 2$, то достаточно рассматривать только те контексты, частота которых не менее двух.

2) В силу требования минимальности грамматических единиц каждая парадигма P_r содержит хотя бы два элемента f_1 и f_2 , первые символы которых различны. Такие элементы выявляются в первую очередь с помощью анализа точек ветвления дерева. При этом будут обнаружены все остальные элементы парадигмы (так как они взаимозаменяемы с f_1, f_2), но будут выявлены не все контексты, в которых

они взаимозаменяемы друг с другом. Отношения между этими элементами парадигмы P_r устанавливаются с помощью анализа их вхождений, не рассмотренных при выявлении f_i, f_j .

3) Шаг итерации состоит из нескольких этапов, на каждом из которых анализируются пути, содержащие r звеньев ($r \in \{1, 2, \dots, \dots, r_{\max}\}$, где r_{\max} - максимальное число звеньев, встречающихся на одном пути в дереве). Первоначально $r=1$. Пути, на которых выявлены грамматические единицы, на данном шаге больше не рассматриваются.

Причиной такого разбиения дерева послужила закономерность, согласно которой границы между значимыми единицами языка часто совпадают с границами между звеньями [5]. Единицы, для которых эта закономерность не выполняется, выявляются во вторую очередь. При этом число путей дерева, подлежащих анализу, существенно меньше первоначального.

4) Множество W разбивается на части в зависимости от начальных символов содержащихся в нем подцепочек. Каждая часть упаковывается в отдельное дерево. В результате анализа одного из деревьев, содержащих максимальное число вершин, выделяются некоторые множества грамматических единиц и парадигм. Рассматривается следующее по числу вершин дерево, в результате чего полученные множества дополняются. Анализ деревьев прекращается, как только добавление нового дерева не приводит к изменению множеств грамматических единиц и парадигм. Таким образом, если для каждого P_r пересечение множеств $K(f_i, f_j)$, где $f_i, f_j \in P_r$, велико, то достаточно рассмотреть небольшое число деревьев 1-грамм.

5. Реализована первая часть алгоритма, состоящая из построения деревьев 1-грамм и выявления грамматических единиц. Использовались спектры 1-грамм для технического текста на русском языке (см. [1]). Проведен один шаг итерации, на котором рассматривались пути, содержащие одно звено, и такие контексты $(\Delta_1, \bar{\Delta}_1)$, в которых $\bar{\Delta}_1 = \Lambda$.

В результате анализа одного дерева выделены звенья:

ПОЛОЖИТЕЛЬН: ОГО, ОЕ, ОИ, ЫМИ, ЫХ	ПОЛИНОМ: ОВ, Ам
ПОЛУЧЕН: НОГО, НЫЕ, НИИ	ПОКАЗАТЕЛ: Е, Ъ
ПОТЕНЦИАЛ: А, л, ЪНОИ, ОБ	ПОПЕРЕЧН: ОГО, ОЕ
ПОВЕРХНОСТ: И, Ъ, НОИ, НАИ, НИИ, НИХ, ЯМИ, ЯМ	ПЕРЕПОЛНЕНИ: л, Е
ПОСТОЯН: НО, НЬМИ, НИХ, НИЕ, СТВА	ПЕРЕНОС: л, Е
ПРОИЗВОЛЬН: УЮ, ОГО, ОИ, ОЕ, ЫК, ЪЕ, ЫМ	ПЕРЕМЕНН: ОИ, ОГО

ПРОСПЕКТИВН:ЫЕ,ЫХ,АЯ,ОЙ
 ПРОСТРАНСТВ:А,О,Е,ЕННЫХ
 ПРОПОРЦИОНАЛЬН:А,О,ОСТИ
 ПРЕДЫДУЩ:ИХ,ЕГО,ЕМ
 ПРЕДСТАВ:ЛЯЕТ,ЛЯТЬ,ЛЕНО,ИТЬ,ИМЫХ
 ПРЕОБРАЗОВАНИ:И,ЕМ
 ПАРАЛЛЕЛЬН:ЫХ,ОГО
 ПРОВОДНИК:ОМ,А,А,
 ПОСЛЕДОВАТЕЛЬ:НОЕ,НОСТЬ,НОСТИ,НО,НОГО,НЫМ,НЫХ.

ПРАВИЛЬН:ОЙ,ЫХ
 ПРАКТИЧЕСК:ОЙ,ОЕ
 ПРИНЦИП:А,ОМ
 ПРОИЗВЕДЕНИ:Я,Е,И
 ПРОЦЕСС:А,Е
 ПЛОТНОСТ:И,ЬЮ,Ь,ЕИ
 ПЛОСКОСТ:И,Ь,ЯМИ
 ПАРАГРАФ:Е,АХ

Первая цепочка символов, отделенная двоеточием, соответствует пути дерева, ведущему из его корня в первую вершину звена, цепочки после двоеточия – полным путям поддерева с корнем в этой вершине. При анализе дерева выделены следующие грамматические единицы: АЯ, ОЙ, ОЕ, ОГО, ОМ, ЫМ, ЫЕ, ЫХ, ЫМИ, А, И, Ъ, Я, О, Е, ЯМИ, А, ОСТИ. Из них неверно выделена только последняя единица, эта ошибка объясняется тем, что проверка вложенности единиц в данном эксперименте не проводилась. Часть грамматических единиц осталась невыявленной, так как использовался спектр 1-грамм, содержащий только 1-граммы неединичной частоты.

Рассмотрим основные типы ошибок в выделении окончаний и их парадигм:

1) Словам, относящимся к двум разным синтаксическим классам, но имеющим одну основу, сопоставляется одна парадигма вместо двух. Например, основе ПОСЛЕДОВАТЕЛЬН- как краткому прилагательному соответствует парадигма {О,А}, а как полному прилагательному – {ОЕ, ОГО, ОМУ, ...}. Алгоритмом будет выделена одна парадигма {О,А, ОЕ, ОГО, ОМУ, ...}.

2) Если имеется чередование символов в основе (в том числе и чередование с пустой цепочкой), то граница, отделяющая основу от окончания, будет смещена. Например, в словах ТОЧЕН, ТОЧНА, ТОЧНО выделяется окончания ЕН, НА, НО вместо А, А, О.

Большая часть ошибок указанных типов будет устранена в следующем варианте алгоритма при дополнении модели синтагматическими отношениями текста (т.е. отношениями между соседними его единицами), также выявляемыми автоматически, и использовании более тонких признаков парадигм.

Алгоритмы реализованы на ПК ОС ЕС.

Л и т е р а т у р а

1. ГУСЕВ В.Д., КОСАРЕВ Ю.Г., ТИТКОВА Т.Н. О задаче поиска повторяющихся отрезков текста. - В кн.: Вычислительные системы. Вып. 62. Ассоциативное кодирование. Новосибирск, 1975, с.49-71.
2. СУХОТИН Б.В. Оптимизационные методы исследования языка. - М.: Наука, 1976. - 120 с.
3. АНДРЕЕВ Н.Д. Статистико-комбинаторные методы в теоретическом и прикладном языковедении. - Л.: Наука, 1967. - 403 с.
4. ГРИНБЕРГ Дж. Квантитативный подход к морфологической типологии языков. - В кн.: Новое в лингвистике. Вып. 3. М., 1963, с. 60-95.
5. ГУСЕВ В.Д., КОСАРЕВ Ю.Г., ТИТКОВА Т.Н. Методы поиска и анализ статистических закономерностей в символьных последовательностях. - В кн.: Машинные методы обнаружения закономерностей. (Материалы Всесоюз. симпозиума, 5-7 апреля 1976). Новосибирск, 1976, с. 75-84.
6. КОСАРЕВ Ю.Г., ЧУЖАНОВА Н.А. Автоматический синтез алгоритмов классификации словоформ по типам словоизменительных парадигм. - В кн.: Структурная и математическая лингвистика. Киев, 1978, с.24-32.

Поступила в ред.-изд.отд.
29 октября 1980 года