



УДК 004.82

ЛОГИКО-ЛИНГВИСТИЧЕСКАЯ МОДЕЛЬ ИДЕНТИФИКАЦИИ СЕМАНТИЧЕСКИХ ОТНОШЕНИЙ СУЩНОСТЕЙ СРЕДСТВАМИ АЛГЕБРЫ КОНЕЧНЫХ ПРЕДИКАТОВ

Хайрова Н.Ф., Узлов Д.Ю., Шаронова Н.В.

*Национальный технический университет “Харьковский политехнический институт”,
г. Харьков, Украина*

nina_khajrova@yahoo.com

popucik@mail.ru

nvsharonova@mail.ru

В работе предлагается логико-лингвистическая модель извлечения слабоструктурированных фактов из естественно языковых текстов. Для идентификации факта в тексте определяются некоторые сущности, выраженные лексическими единицами, и семантические связи между ними. Семантические связи определяются семантическими функциями партиципантов предложения, которые описаны предикатами алгебры конечных предикатов. Модель применяется на семантическом этапе лингвистического процессора информационной подсистемы идентификации криминалистически значимых фактов в слабоструктурированных текстах.

Ключевые слова: слабоструктурированные факты; семантические функции; алгебра конечных предикатов; лингвистический процессор.

Введение

Знания о некоторой предметной области (ПО) представляют собой совокупность сведений об объектах/субъектах данной ПО, их существенных свойствах и связывающих отношениях, а также о фактах, семантически объединяющих партиципанты ПО и их отношения в триаду субъект – атрибут – значение (или субъект – отношение – объект).

Для извлечения фактов, представленных в хорошо структурированных текстовых документах, существуют достаточно надежные алгоритмы [Baeza-Yates, 1999].

К плохо структурированной фактографической информации относятся сведения, представленные различными нерегламентированными словесными конструкциями на естественном языке. Задача извлечения плохо структурированных фактов из произвольных текстов до сих пор не имеет сколь-нибудь общего решения [Ландэ, 2009].

Для идентификации некоторого знания, представленного в форме слабоструктурированного факта, необходимо извлечь из текстовой информации некоторые объекты/субъекты,

выраженные лексическими единицами, и определить семантические отношения между ними.

Так как факт выражается на естественном языке в форме законченного высказывания, то необходимо построить некий шаблон, отображающий семантические (или понятийные) связи партиципантов предложения (участников действия, выраженных существительными).

Для задания таких смысловых связей предлагается использовать семантические функции, выражаемые через отношения морфологических и семантических категорий партиципантов предложения средствами алгебры конечных предикатов (АКП).

1. Описание используемой модели.

Введем конечное множество грамматических и семантических характеристик партиципантов предложения $M = \{m_1, \dots, m_n\}$, где n количество указанных характеристик. Отношения между характеристиками можно представить в виде $m_i * m_j * \dots * m_k$, где $m_i, m_j, \dots, m_k \in M$, а знак $*$ – обозначает, что данные характеристики соответствуют существительному, выполняющему некоторую семантическую функцию.

На множестве M введем систему предикатов S так, чтобы любой предикат $P(q_m) \in S$, обращался в 1 на множестве существительных с грамматико-семантической информацией, соответствующей определенной семантической функции, и был равен 0 в противном случае. Таким образом, множество предикатов S можно сопоставить с множеством семантико-грамматических характеристик приписанных партиципantu предложения.

Для формализации семантических функций партиципantов предложения русского языка и их явного представления средствами поверхностной структуры были выделены и описаны предметными переменными морфологические (грамматический падеж) и семантические категории существительных [Бондаренко, 2007]. Рассматривались семантические категории: живой / неживой, инструмент, часть тела, объемное пространство, пункт назначения, место отправления, плоскость/точка, механизм, определенный час, период, месяц/сезон.

Область изменения введенных переменных формально задается следующим образом:

$$\begin{aligned} x^o \vee x^h &= 1, \\ z^h \vee z^p \vee z^l \vee z^b \vee z^t \vee z^n &= 1, \\ y^m \vee y^c \vee y^n \vee y^i \vee y^t \vee y^o \vee y^b \vee y^n \vee y^u &= 1, \end{aligned}$$

где x^i , z^j , y^k – предметные переменные, характеризующие: x^o – категорию живого, x^h – категорию неживого; y^m – наличие признака «механизм», y^n – наличие семантического признака «инструмент», y^i – наличие семантического признака «часть тела», y^t – наличие семантического признака «плоскость/точка», y^o – семантического признака «объемное пространство», y^b – семантического признака «определенное время», y^n – наличие семантического признака «период», y^u – наличие семантического признака «пункт назначения»; z^h , z^p , z^l , z^b , z^t , z^n — грамматическая категория падежа.

Семантическая функция существительного — партиципantа предложения описывается предикатом $P(x, y, z) = 1$, связывающим элементы семантического значения существительного x и y с его грамматическими значениями z [Хайрова, 2012]. Тогда, используя конъюнкцию предикатов, можно записать:

$$P(x, y, z) \rightarrow P(x) \bullet P(y) \bullet P(z), \quad (1)$$

где \bullet — операция конъюнкции.

Таким образом, отношения между морфологическими и семантическими признаками существительных предложения, выражающие семантические функции, можно записать логическим произведением:

$$\begin{aligned} P(x_n) * P(y_n) * P(z_n) &= \\ = \gamma_k(x_n, y_n, z_n) \bullet P(x_n) \bullet P(y_n) \bullet P(z_n), \end{aligned} \quad (2)$$

Логическое произведение предикатов $P(x_n)$, $P(y_n)$ и $P(z_n)$, описывает всевозможные отношения

морфологических и семантических характеристик, а предикат $\gamma_k(x_n, y_n, z_n)$ исключает часть связей, которые не реализуются в данной семантической функции. Предикат γ_k принимает значение 1, если морфосемантическая информация словоформы n формирует некоторую семантическую функцию лексемы, и значение 0 в противном случае.

Множество рассмотренных в системе семантических категорий значения существительного задается предикатом:

$$P(y_n) = y_n^m \vee y_n^c \vee y_n^n \vee y_n^i \vee y_n^t \vee y_n^o \vee y_n^b \vee y_n^n \vee y_n^n. \quad (3)$$

Множество значений грамматических категорий, определяющих грамматические падежи существительного, выражается предикатом:

$$P(z_n) = z_n^h \vee z_n^p \vee z_n^l \vee z_n^b \vee z_n^t \vee z_n^n. \quad (4)$$

Признак одушевленности выражается предикатом:

$$P(x_n) = x_n^o \vee x_n^h. \quad (5)$$

Семантическая функция агенса, представляющая субъект действия, обычно выступающего инициатором действия, лицо или предмет, имеющее потенцию на осуществление действий, выражается предикатом:

$$\gamma_A(x_n, y_n, z_n) = x_n^o z_n^h \vee z_n^h x_n^h y_n^m \vee z_n^h x_n^o y_n^c. \quad (6)$$

Семантическая функция инструменталиса, определяющая непосредственную причину действия, играющую определенную роль в совершении процесса, выражается предикатом:

$$\gamma_{II}(x_n, y_n, z_n) = z_n^t x_n^h y_n^n \vee z_n^t x_n^h y_n^i. \quad (7)$$

Семантическая функция локатива, выражающая характеристики месторасположения, пространственной ориентации действия или состояния, выражается предикатом:

$$\gamma_L(x_n, y_n, z_n) = z_n^l x_n^h y_n^t \vee z_n^l x_n^h y_n^m \vee z_n^l x_n^h y_n^i \vee z_n^l x_n^h y_n^o. \quad (8)$$

Семантическая функция обьектив, определяющая объект, над которым непосредственно осуществляется действие, выражается предикатом:

$$\gamma_O(x_n, y_n, z_n) = z_n^b x_n^h \vee z_n^b x_n^o. \quad (9)$$

Семантическая функция темпоралиса, выражающая временную характеристику действия, выражается предикатом:

$$\gamma_T(x_n, y_n, z_n) = z_n^b x_n^h y_n^b \vee z_n^l x_n^h y_n^n. \quad (10)$$

Семантические функции в различных естественных языках имеют разные формы формального выражения и соответственно определяются предикатами АПК различного вида. Например, в русском, украинском и белорусском языках они выражаются грамматической формой и семантическими категориями существительного,

стоящего после или перед определенным глаголом, как представлено в нашей модели.

Тогда как в английском языке семантические функции во многом определяются синтаксическими категориями, в частности отношением с предлогом, стоящим после глагола. Такие категории могут быть как уникальными для определенных глаголов, так и общими, как, например, признак направления движения, определяемый предлогом (в случае его наличия) после глаголов go, run, drive, ride, transport, ship и т.д.

2. Идентификация криминалистически значимых фактов из текстов

Рассмотрим применение данной модели для идентификации криминалистически значимых фактов в источниках полнотекстовой информации. Из огромных информационных потоков текстовой информации, обрабатываемой в процессе оперативно-служебной деятельности (сводки, объяснительные/служебные записки, отчеты, газетные и интернет публикации, словесные портреты фигурантов и т. п.), и сопровождающего ее шума следователь (или иное процессуально-должностное лицо) должен извлечь факты конкретного уголовного или иного дела. В подавляющем числе случаев к таким фактам относятся: сведения о фигурантах, сведения об объекте посягательства, сведения о механизме и способе совершения преступления.

Для извлечения из неструктурированной текстовой информации фактов о дате и месте рождения некоторых персоналий (фигурантов дела), с целью исключения или привлечения потенциального кандидата, а так же о дате, месте, субъекте и объекте некоторого противоправного действия использовались семантические функции, выражающие информацию, соответствующую требованиям:

- темпоралис – временная характеристика события, позволяющая определить дату (в нашем случае: дату рождения человека или некоторого противоправного деяния);
- локатив – функция, характеризующая местонахождение, положение или состояние объекта, определяя место (в нашем примере, место рождения человека или некоторого преступного события);
- объектив – функция, определяющая сферы и продукты деятельности человека (в данном случае: сведения об объекте посягательства);
- агенс – семантическая функция, представляющая субъект действия, обычно выступающего инициатором действия (в нашем случае: лицо/субъект противоправного действия).

Были выбраны наиболее распространенные глаголы, соответствующие фактам преступного действия и идентифицирующие личность, а также определены семантические функции, являющиеся

центральной частью триплета субъект – атрибут – значение и предикаты (формулы 6-9), описывающие отношения морфологических и семантических категорий существительных, соответствующих внешним элементам данного триплета (табл. 1).

Таблица 1 – Формализм модели идентификации семантических отношений

Глагол	Падеж	Предикат семантической функции
родиться, похищать, похитить, убивать, убитый, красть, выкрасть, украсть, обмануть, обманывать, грабить, ограбить и др.	Темпоралис whenacted	10
	Локатив whereacted	8
	Объектив toactsmth	9
	Агенс tobeactedbysmth	6

3. Этапы работы лингвистического процессора рассмотренной задачи

Предложенная модель используется на семантическом этапе лингвистического процессора, включающего также графемный, морфологический, контекстный и семантический этапы анализа (рис. 1).



Рис. 1. Структурная схема лингвистического процессора подсистемы идентификации криминалистически значимых фактов в слабоструктурированных текстах.

Подавляющую часть субъектов криминалистически значимых фактов, составляют персоналии, компании и организации, кроме того, часто необходимо определять факт местоположения того или иного субъекта. Специальное графическое оформление таких сущностей позволяет для их выделения использовать формализмы графемного анализа. Графемный этап используется для выявления антропонимов, топонимов (включающих названия населенных мест, внутригородских объектов и др.), названий организаций, предприятий

и учреждений и базируется на базе имен собственных и специальных алгоритмах.

На этапе морфологического анализа осуществляется обработка словоизменяемых и словообразовательных форм, позволяющая, например, отнести глаголы «~~красть~~» и «~~украсть~~» к одному и тому же факту.

Качественное построение шаблона факта предполагает определение возможных имен сущностей, которые могут встречаться в текстах под разными знаками. Для повышения точности определения факта на этапе синтаксического анализа необходимо решить проблему кореферентности, в частности, ассоциировать местоимения со своими антецедентами, соотнося их с именуемой сущностью.

Задача построения связей шаблона факта представляет собой одну из центральных задач извлечения знаний из текстов, она направлена на распознавание в тексте набора возможных отношений между элементами шаблона, разработанными на предыдущем этапе. Для решения данной задачи на этапе семантического анализа лингвистического процессора используется разработанная логико-лингвистическая модель идентификации отношений между сущностями.

4. Практическая реализация разработанного метода

Задача извлечения криминалистически значимых фактов из полнотекстовой информации, обрабатываемой в процессе оперативно-служебной деятельности, была реализована для извлечения даты, места рождения персоналии, а также сведений о дате, месте, субъекте и объекте некоторого противоправного действия.

Программа представляет собой веб-приложение, анализирующее текст или список рассматриваемых текстовых файлов. Определенная фактографическая информация представляется пользователю в форме диалогового окна, в котором отображается извлеченный факт и предложение или несколько предложений, представляющие описание данного факта.

В разработанном приложении определяются семантические функции, выражающие информацию, соответствующую поставленной задаче: темпоралис - временная характеристика события, локатив - характеристика местонахождения, положения или состояния; агенс и объектив - участники заданного действия. В процессе реализации модели был определен набор глаголов, требующих выполнения данных семантических функций участников предложения (табл. 1).

Тестовая проверка разработанного приложения показала приемлемость использования предложенной модели.

Заключение

Таким образом, разработанная логико-лингвистическая модель идентификации семантических отношений сущностей позволила извлекать факты из слабоструктурированной текстовой информации в виде триплета субъект – отношение – объект. Отношения между сущностями в данном триплете описывается предикатом взаимосвязи морфологических и семантических категорий участника того или иного действия. Модель применяется на семантическом этапе лингвистического процессора информационной подсистемы идентификации криминалистически значимых фактов в слабоструктурированных текстах для чего к настоящему моменту описаны функции: темпоралис, локатив, объектив и агенс.

Библиографический список

- [Baeza-Yates, 1999] Baeza-Yates R., Ribeiro-Neto B. Modern Information Retrieval. — Addison-Wesley, 1999. — 340 p.
- [Ландэ, 2009] Ландэ Д. В. Интернетика: Навигация в сложных сетях: модели и алгоритмы: моногр. / Д. В. Ландэ, А. А. Снарский, И. В. Безсуднов — М.: Либроком (Editorial URSS), 2009. 264 с.
- [Бондаренко, 2007] Бондаренко М. Ф. Теория интеллекта: учебник/ Бондаренко М. Ф., Шабанов-Кушнарченко Ю. П. Харьков: Комп. СМІТ, 2007. — 576 с.
- [Хайрова, 2012] Хайрова Н. Ф. Використання логіко-алгебраїчної моделі семантичних відмінків для семантичного аналізу речення/ Н. Ф. Хайрова // Зб. наук. пр. Військового ін-ту Київ. нац. ун-ту. — К.: ВІКНУ, 2012. — Вип. № 38. — С. 239—245.

LOGICAL-LINGUISTIC MODEL FOR IDENTIFICATION OF SEMANTIC RELATIONSHIPS BETWEEN ENTITIES ON BASE OF ALGEBRA OF FINITE PREDICATES

Khairova N., Uzlov D., Sharonova N.

*National Technical University
“Kharkiv Polytechnic Institute”, Kharkiv, Ukraine*

nina_khajrova@yahoo.com
poputcik@mail.ru
nvsharonova@mail.ru

This paper proposes a logical-linguistic model extracting semi-structured facts of natural language texts. To identify the fact some entities expressed by lexical units as well as semantic relations between them are defined in the text. The semantic relations are expressed by semantic functions of sentence participants. The functions are described by predicates of algebra of finite predicates. The model is applied to the semantic stage of linguistic processor of information subsystem for facts identification, which are essential for criminalistics, in the framework a semi-structured texts.

Key words: semi-structured facts; semantic functions; algebra of finite predicates; linguistic processor.