

ТЕОРЕТИКО-МНОЖЕСТВЕННОЕ МОДЕЛИРОВАНИЕ СЕМАНТИЧЕСКИХ ПОЛЕЙ СЛОВ

Черепанова Ю. Ю.

Харьковский государственный технический университет
радиоэлектроники

The article deals with set-theoretical modeling of semantic fields. The thesaurus of semantic fields is one of tools for the description of semantics of natural language that can be used in the systems of natural language understanding for solving some problems. The received set-theoretical model of the thesaurus of semantic fields can be used as a basis for the further modeling of semantic fields.

Тезаурус семантических полей является одним из инструментов для описания семантики и анализа естественного языка, позволяющим решить некоторые проблемы, возникающие в процессе обработки естественно-языковых текстов. Семантические поля представляют собой структуры, включающие в себя лексические единицы различных уровней и связывающие их семантические отношения. Такая структура семантических полей дает возможность моделирования их с помощью математического аппарата теории множеств, теории графов и теории категорий.

Формализуем лингвистические модели семантических полей с помощью теоретико-множественного аппарата.

Определим следующие множества:

- $U = \{v_1, v_2, \dots, v_M\}$ – множество слов языка;
- $A = \{a_1, a_2, \dots, a_K\}$ – множество семантических признаков;
- $T = \{t_1, t_2, \dots, t_N\}$ – множество семантических типов.

Множества U , A , T представляют собой лексические единицы различных уровней. Для отражения существующих между ними семантических связей введем следующие отношения.

Определим отображение $f_{\text{код}} : U \rightarrow A$, где $U = \{v_1, v_2, \dots, v_M\}$ – множество слов языка, $A = \{a_1, a_2, \dots, a_K\}$ – множество семантических признаков. Данное отображение представляет собой кодирование слов в квазиосновы, семантические признаки. Это отображение является всюду определенным, сюръективным. Так как один признак может соответствовать нескольким словам (отображение не является инъективным), то $K \leq M$.

Отношение $R_{\text{см}} \subseteq A \times A$ содержательно означает «иметь семантический множитель». Это отношение всюду определенное, рефлексивное, транзитивное. Тогда сечения $R_{\text{см}}(a_i)$ будут представлять собой семантические объемы слов $v_i = f_{\text{код}}^{-1}(a_i)$.

При построении семантических полей все множество слов U разбивается на множество семантических полей U_1, U_2, \dots, U_N . Семантическое поле U_i – некоторое множество слов $\{v_{i1}, v_{i2}, \dots, v_{Mi}\}$, связанных между собой парадигматическими и синтагматическими семантическими отношениями, соответственно разделяются парадигматические и синтагматические семантические поля. К парадигматическим полям относятся различные классы лексических единиц, имеющие общие признаки в своих значениях. То есть если каждому слову v_i поставить в соответствие множество $W_i = R_{cm}(a_i) = \{a_{i1}, a_{i2}, \dots, a_{iK}\}$ (семантический объем), где a_{iL} – семантический признак слова v_i , то слова v_i, v_j объединяются в одно парадигматическое поле U_z при условии, что $W_i \cap W_j \neq \emptyset$.

Таким образом, можно задать отношение $R_{\Pi} \subseteq U \times U$, определяющее наличие парадигматических связей между словами. Отношение R_{Π} всюду определенное, сюръективное, рефлексивное, симметричное. Сечения $R_{\Pi}(v_i)$ будут представлять собой парадигматические поля слов v_i .

Синтагматическими называются отношения между сочетаемыми единицами языка. Таким образом, синтагматические семантические поля слова представляют собой связи слов либо их абстрактных категорий в текстах. Для отражения этих связей необходимо ввести несколько отношений.

Введем всюду определенное, сюръективное отношение $R_{ct} \subseteq U \times T$ означающее «сочетаться с семантическим типом».

Отношение $R_{tt} \subseteq T \times T$ означающее «сочетаться», всюду определено, сюръективно, симметрично, антирефлексивно.

Введем отображение $f_T: U \rightarrow T$, содержательно означающее «относиться к семантическому типу».

Тогда отношение $R_{cn} \subset R_{ct} \cdot f_T^{-1}$ будет определять синтагматические связи между словами. Сечения $R_{tt}(t_j)$ будут представлять собой синтагматические связи категории t_j , сечения $R_{ct}(v_i)$ – синтагматические связи слов v_i с абстрактными категориями в текстах, а сечения $R_{cn}(v_i)$ – связи слов v_i с другими словами.

Разумеется, данная модель является очень упрощенной. Отношения R_{Π} и R_{cn} отражают лишь наличие соответствующих связей между словами. На самом деле парадигматические и синтагматические связи должны отражаться множеством отношений. Кроме того, данная модель отражает многозначность лексических единиц лишь в том случае, если разные значения слов будут отображаться как разные слова v_i, v_j .

Полученную теоретико-множественную модель тезауруса семантических полей будем использовать как основу для дальнейшего моделирования семантических полей.