

Дрождин В.В., Баканов А.Б. Грамматика описания домена фамилий. – Вопросы радиоэлектроники. Серия «Электронная вычислительная техника». – 2007, Вып.1. – С. 77-82.

В.В. Дрождин
А.Б. Баканов

УДК 681.3

ГРАММАТИКА ОПИСАНИЯ ДОМЕНА ФАМИЛИЙ

Рассматривается формализация представления данных в домене фамилий на основе формальной грамматики. Предложен критерий и способ его вычисления, позволяющий сравнивать языки и грамматики.

Использование семантических свойств в представлении и организации данных в базе данных (БД) позволяет существенно повысить эффективность функционирования систем баз данных и автоматизированных информационных систем (АИС), построенных на их основе. Легко установить, что основные семантические свойства классов объектов отражаются в структуре синтаксического представления информации об этих объектах. Поэтому серьезным препятствием для обеспечения высокой эффективности функционирования систем БД являются:

- использование в них универсальных языков представления различных данных на внутреннем (физическом) уровне;
- использование универсальных (единых для всей системы и для различных данных) алгоритмов обработки данных.

Например, в реляционной модели данных, получившей в настоящее время наибольшее распространение, данные рассматриваются как атомарные (целостные, неделимые) объекты, вследствие чего объем данных в БД возрастает прямо пропорционально объему накапливаемой информации о предметной области. В сетевой и иерархической моделях данных объем одних и тех же данных в БД несколько меньше, чем в реляционной, но громоздкость обрабатываемых структур и отсутствие мощного языка запросов теоретико-множественного типа сводят на нет это преимущество.

Доменно-ориентированный подход позволяет снять эти проблемы, рассматривая домены как отдельные самостоятельные системы БД, взаимодействующие между собой на основе единого стандартизованного интерфейса. Обособленность и специализация доменов позволяют в полном объеме реализовать принцип физической независимости данных в БД.

Домен содержит информацию о множестве объектов одного типа и, следовательно, все данные в домене имеют одинаковую или незначительно различающуюся структуру. Это позволяет разрабатывать внутреннюю организацию домена специально ориентированную на хранение и обработку данных именно с требуемой структурой, что повышает эффективность функционирования домена на несколько порядков.

Для выявления структуры данных элементов домена может использоваться аналитический подход, разработанный в математической

лингвистике для анализа естественных и формальных языков. При этом любые системы данных считаются открытыми как вверх (создание суперсистем, в которых данные системы являются элементами), так и вниз (любой элемент является сложным, состоящим из элементов более низкого уровня, вплоть до отдельных символов). Поэтому на основе анализа элементов домена можно выявлять их более или менее устойчивую структуру и формировать существенно более эффективную внутреннюю организацию домена.

Рассмотрим внутреннюю организацию домена фамилий, в основе которой будет находиться формальная грамматика.

1. Неформальное описание домена фамилий

Домен фамилий представляет собой некоторое множество слов-фамилий. Для конкретности будем рассматривать русские фамилии, представляемые в виде:

$$Z_1 Z_2 \dots Z_n \text{ или } Z_1 Z_2 \dots Z_n - Z_{n+1} Z_{n+2} \dots Z_m - Z_{m+1} Z_{m+2} \dots \quad (1)$$

где $Z_1, Z_{n+1}, Z_{m+1}, \dots$ – прописная буква кириллицы;

$Z_2, \dots, Z_n, Z_{n+2}, \dots, Z_m, Z_{m+2}, \dots$ – строчная буква кириллицы.

Домен может содержать также иностранные фамилии, записанные на русском языке в виде (1) или в виде:

$$Z_1 Z_2 \dots Z_n _\underline{Z}_{n+1} Z_{n+2} \dots Z_m _\underline{Z}_{m+1} Z_{m+2} \dots \quad (2)$$

где ‘ ’ – символ пробела.

Каждой фамилии сопоставим бинарный признак-пол со значениями ‘ж’ – женский и ‘м’ – мужской.

Так как между фамилией и полом задано простое отношение ассоциации (соответствия), то домен фамилий относится к простым доменам.

Таким образом, полное значение элемента домена фамилий представляется в виде:

$$Z_1 Z_2 \dots Z_n \delta \# \text{ или } Z_1 Z_2 \dots Z_n \gamma Z_{n+1} Z_{n+2} \dots Z_m \gamma Z_{m+1} Z_{m+2} \dots \delta \#, \quad (3)$$

где γ – символ ‘–’ или ‘ ’;

δ – признак-пол со значениями ‘ж’ или ‘м’;

– финальный (закрывающий) символ элемента домена.

2. Теоретико-множественная модель домена фамилий

Для записи домена фамилий используется конечный **алфавит** символов:

$$A_\phi = A_1 \cup A_2 \cup A_3 \cup \{\#\}, \quad (4)$$

где $A_1 = \{‘А’…‘Я’\}$ – алфавит прописных букв кириллицы;

$A_2 = \{‘а’…‘я’\}$ – алфавит строчных букв кириллицы;

$A_3 = \{‘-’, ‘_\’\}$ – алфавит специальных символов;

‘#’ – символ, являющийся признаком завершения фамилии.

На A_ϕ задается **свободная полугруппа** \mathcal{F}^* [1] в виде множества всех конечных последовательностей символов (слов), на котором определена ассоциативная и некоммутативная бинарная операция конкатенации

(сцепления). Из \mathcal{F}^* выделим такое подмножество \mathcal{F} , в котором каждое слово $x_\phi \in \mathcal{F}$ представляется в формате:

$$x_\phi = z_1 z_2 \dots z_n \delta \# \text{ или } x_\phi = z_1 z_2 \dots z_n \gamma z_{n+1} z_{n+2} \dots z_m \gamma z_{m+1} z_{m+2} \dots \delta \#,$$

где $z_1, z_{n+1}, z_{m+1}, \dots \in A_1$;

$z_2, \dots, z_n, z_{n+2}, \dots, z_m, z_{m+2}, \dots \in A_2$;

$\gamma \in A_3$;

$\delta \in \{\text{'ж'}, \text{'м'}\} \subset A_2$.

Тогда **язык** L_ϕ , представляющий конкретный домен фамилий, является подмножеством:

$$L_\phi \subseteq \mathcal{F}. \quad (5)$$

Множество \mathcal{F} есть универсальный язык представления всех возможных фамилий, а полугруппа \mathcal{F}^* – универсальный язык над A_ϕ .

Учитывая форматы слов $x_\phi \in L_\phi$, в \mathcal{F}^* можно выделить свободную подполугруппу $\mathcal{F}_2^* \subset \mathcal{F}^*$ со свободным порождающим множеством A_2 .

Для количественной оценки полугруппы \mathcal{F}^* или любого ее подмножества введем оценочную функции φ . Например, $\varphi(\mathcal{F}^*)$ – количество всех символов в полугруппе \mathcal{F}^* , $\varphi(L_\phi)$ – количество всех символов в языке L_ϕ , $\varphi(x_\phi)$ – количество символов в слове x_ϕ .

Таким образом, домен фамилий представляется формальным языком $L_\phi = \{x_{\phi i} \mid 0 \leq i \leq |L_\phi|\}$, находящемся в отношении $L_\phi \subseteq \mathcal{F} \subset \mathcal{F}^*$. Здесь $|L_\phi|$ – количество слов-фамилий в языке L_ϕ .

3. Контекстная грамматика, порождающая домен фамилий

При представлении домена фамилий с помощью некоторой порождающей грамматики будем использовать **алфавит символов**:

$$\tilde{A}_\phi = T_\phi \cup N_\phi,$$

где $T_\phi = A_\phi$ – конечное фиксированное множество терминальных символов;
 N_ϕ – множество нетерминальных символов.

Правила подстановки

$$P_\phi = \{p_\phi\}$$

являются множеством продукции вида $p_\phi : x \rightarrow y$, каждая из которых на основе слова x порождает слово y .

Начальный символ грамматики S_ϕ является последовательностью слов вида $S_\phi = \{\lambda_i \mid 0 \leq i \leq n\}$,

где λ_i – слово, позволяющее восстановить одно или несколько слов $x_\phi \in L_\phi$;
 n – количество слов в последовательности S_ϕ .

Грамматика G_ϕ , порождающая язык $L(G_\phi)$, задается в виде [2]:

$$G_\phi = (S_\phi, T_\phi, N_\phi, P_\phi), \quad (6)$$

Если в объявлении S_ϕ параметр $n = 0$, то $S_\phi = \emptyset$ и грамматика G_ϕ порождает пустой язык $L(G_\phi) = \emptyset$.

Для построения грамматики G_ϕ будем использовать аналитический подход [3], заключающийся в анализе языка L_ϕ с целью выявления структуры слов, состава их элементов и отношений между этими элементами.

В начальном состоянии будем считать, что для языка L_ϕ определено **единичное разбиение** E , в каждую клетку e_i которого помещается одно и только одно слово $x_{\phi i}$, т.е. $x_{\phi i} \in e_i$. При этом $E = \{e_i\}$, $e_i = \{x_{\phi i}\}$, $|E| = |L_\phi|$,

$$L_\phi = \bigcup_{i=0}^{|L_\phi|} e_i, \quad 0 \leq i \leq |L_\phi|. \quad \text{Язык } L_\phi \text{ становится языком с парадигматической}$$

структурой (L_ϕ , E), но без морфологии, так как слова языка являются целостными объектами, не имеющими внутренней структуры.

Представлению языка (L_ϕ , E) соответствует грамматика $G_{0\phi} = (S_{0\phi}, T_\phi, N_{0\phi}, P_{0\phi})$, в которой $S_{0\phi} = E$, $\lambda_i = e_i = x_{\phi i}$. $P_{0\phi}$ содержит только правила описания формата представления слов $x_{\phi i}$:

$$\begin{aligned} p_1 : x_{\phi i} &\rightarrow z\delta\#, \quad 1 \leq i \leq |L_\phi| \\ p_2 : z &\rightarrow z_1z_2z_3\dots z_n \mid z\gamma z_1z_2z_3\dots z_n \\ p_3 : z_1 &\rightarrow A \mid B \mid \dots \mid Y \\ p_4 : z_2z_3\dots z_n &\rightarrow a \mid b \mid \dots \mid y \\ p_5 : \gamma &\rightarrow - \mid \sqcup \\ p_6 : \delta &\rightarrow ж \mid м \end{aligned}$$

$N_{0\phi}$ включает нетерминальные символы из правил подстановки $p_1 \dots p_6$.

Грамматика $G_{0\phi}$ будет порождать язык $L(G_{0\phi}) = L_\phi$ и иметь размерность $|L(G_{0\phi})| = |L_\phi|$ и $\varphi(L(G_{0\phi})) = \varphi(L_\phi)$.

Учитывая, что функция φ дает количественную оценку последовательности символов, то ее можно применять как к отдельным правилам p_ϕ , так и ко всему множеству правил P_ϕ . Поэтому становится возможным количественное сравнение языков с порождающими их грамматиками, а также сравнение грамматик между собой. Например, количественное сравнение языка L_ϕ и порождающей грамматики $G_{0\phi}$ дает соотношение $\varphi(G_{0\phi}) = \varphi(L_\phi)$. Это означает, что грамматика $G_{0\phi}$ представляет собой запись слов языка L_ϕ в определенной последовательности.

Совершенствование грамматики $G_{0\phi}$ может осуществляться двумя путями:

- построением более крупных слов на основе обобщения слов $x_\phi \in L_\phi$;
- уменьшением размера слов $x_\phi \in L_\phi$ на основе декомпозиции на составные части и классификации этих составных частей.

Оценку степени совершенства грамматики будем производить на основе количества символов, необходимых для записи грамматики. Будем считать грамматику G'' более совершенной относительно грамматики G' , если $\varphi(G'') < \varphi(G')$ при $L(G'') = L(G')$. Совершенство грамматики G'' относительно грамматики G' тем выше, чем больше разность $\Delta n = \varphi(G') - \varphi(G'')$.

Оценивая возможные пути совершенствования грамматики $G_{0\phi}$ по критерию уменьшения размера Δp можно ожидать получение наиболее совершенной грамматики $G_{1\phi}$ при использовании обобщения слов, т.к. в этом случае строится более крупное разбиение E_1 , в каждую клетку $e_{1,j}$, которого помещаются слова из одной или нескольких клеток $e_i \in E$, $e_{1,j} = \{x_{\phi i} \mid x_{\phi i} \in e_i, p_k : \alpha_j \rightarrow x_{\phi i}\}$.

Библиография

1. Скорняков Л.А. Элементы алгебры: Уч. пос. для вузов. – М.: Наука, 1986. – 240 с.
2. Рейуорд-Смит В.Дж. Теория формальных языков. Вводный курс. – М.: Радио и связь, 1988. – 128 с.
3. Маркус С. Теоретико-множественные модели языков. - М.: Наука, 1970. – 332 с.