

Корректность алгоритма разбора LC(k)-грамматик

Шпилевой Денис Б05-153
Руководитель - Павел Ахтямов

Цель исследований

1. Формализация LC-разбора.
2. Как часть глобальной задачи - исследования возможности определения класса языков по некоторым данным, которые разбираются за линейное время.

План презентации

- Формальные языки, иерархия Хомского
- КС-грамматики
- Разбор
- LL(k)-грамматики
- LC(k)-грамматики
- План работы на семестр
- Источники

Формальные языки. Иерархия Хомского

Опр. Формальный язык - множество конечных слов над конечным алфавитом.

Иерархия языков Хомского:

- Регулярные (праволинейные) $A \rightarrow aB$ $A, B \in N, a \in \Sigma$
- Контекстно-свободные $A \rightarrow X$ $A \in N, X \in (N \cup \Sigma)^*$
- Контекстно-зависимые $XAY \rightarrow XZY$ $A \in N, X, Y, Z \in (N \cup \Sigma)^*$
- Порождающие $X \rightarrow Y$ $X, Y \in (N \cup \Sigma)^*$

КС-грамматики. Парсинг

Опр. $G = (N, \Sigma, P, S)$ - КС-грамматика, если все ее правила имеют вид:

$$A \rightarrow X, \text{ где } A \in N, X \in (N \cup \Sigma)^*$$

Парсинг (синтаксический анализ) - процесс сопоставления линейной слова формального языка с его формальной грамматикой. Результатом обычно является дерево разбора.

Алгоритм разбора - алгоритм, определенный на множестве слов некоторого языка, и биективно сопоставляющий каждому слову его разбор (последовательность примененных правил).

Алгоритм Эрли (1968)

Алгоритм Эрли - алгоритм синтаксического анализа цепочки по контекстно-свободной грамматике, основанный на методе динамического программирования.

Разбирает все КС-грамматики, но асимптотика - от квадрата до куба в зависимости от коэффициента неоднозначности грамматики.

Рассмотрим примеры алгоритмов, разбирающих не на все КС-грамматики, но выполняющих разбор за линейное время от длины слова.

LL(k)-грамматики (1968)

LL(k)-разбор - левый разбор по известным k символам после головки.

Опр. КС-грамматику $G = (N, \Sigma, P, S)$ называют LL(k)-грамматикой, если из существования двух левых выводов:

- $S \Rightarrow wAa \Rightarrow_1 wba \Rightarrow wx$
- $S \Rightarrow wAa \Rightarrow_1 wca \Rightarrow wy \quad w, x, y \in \Sigma^*, A \in N, a, b, c \in (N \cup \Sigma)^*$
- при этом $\text{FIRST}_k(x) = \text{FIRST}_k(y)$, тогда $b = c$.

LL(k)-грамматики (1968)

Критерий LL(k)-грамматики. КС-грамматика $G = (N, \Sigma, P, S)$ является LL(k)-грамматикой тогда и только тогда, когда для двух различных правил $A \Rightarrow x, A \Rightarrow y$ пересечение $\text{FIRST}_k(xa)$ и $\text{FIRST}_k(ya)$ пусто при всех таких wAa , что $S \Rightarrow wAa$, где $w \in \Sigma^*$, $A \in N$, $x, y \in (N \cup \Sigma)^*$.

Цель критерия - показать, что для определения того, какое правило необходимо применить, не нужно помнить всю цепочку w до рассматриваемых k символов после головки.

LL(k)-грамматики (1968)

Алгоритм разбора LL(k)-грамматик. Построение управляющей таблицы разбора LL(k)-грамматик.

Вход: LL(k)-грамматика $G = (N, \Sigma, P, S)$ и множество её LL(k)-таблиц Γ .

Выход: Корректная управляющая таблица разбора M .

Метод: Управляющая таблица M определяется $(\Gamma \cup \Sigma \cup \{\$\}) \times \Sigma^*k$ так:

1. Если $A \rightarrow x_0B_1x_1B_2x_2\dots B_mx_m$ - правило с номером i , и $T_{A,L} \in \Gamma$, то для всех u , для которых $T_{A,L}(u) = (A \rightarrow x_0B_1x_1B_2x_2\dots B_mx_m, \langle Y_1, \dots, Y_m \rangle)$ полагаем, что $M(T_{A,L}, u) = (x_0 T_{B_1,Y_1} x_1 T_{B_2,Y_2} x_2 \dots T_{B_m,Y_m} x_m, i)$.
2. $M(a, av) = \text{выброс}$, для всех $v \in \Sigma^{*(k-1)}$.
3. $M($, e) = \text{допуск}$.
4. В остальных случаях $M(X, u) = \text{ошибка}$.
5. $T_{\$, e}$ - начальная таблица.

LC(k)-грамматики (1970)

LC(k)-разбор - левый разбор по известным k символам после головки и известен вывод из левого участка.

Опр. КС-грамматику $G = (N, \Sigma, P, S)$ называют LC(k)-грамматикой, если она удовлетворяет таким условиям:

Пусть $S \Rightarrow^* wAx, w \in \Sigma^*, A \in N, x, y \in (N \cup \Sigma)^*$. Тогда для каждой цепочки $v \in \Sigma^*$ и вывода $A \Rightarrow^* v$ существует не более одного такого правила $B \Rightarrow a$, что:

1.
 - 1) Если $a = Cb$, где $a, b \in (N \cup \Sigma)^*$, $C \in N$, то первые k букв $u \in \text{FIRST}_k(byx)$
б) Если $C = A$, то $u \notin \text{FIRST}_k(x)$.
 - 2) Если a начинается терминалом, то $u \in \text{FIRST}_k(ayx)$

LC(k)-грамматики

Алгоритм разбора LC(1)-грамматик. Построение управляющей таблицы разбора LC(1)-грамматик.

Вход: LC(1)-грамматика $G = (N, \Sigma, P, S)$.

Выход: Корректная управляющая таблица разбора T .

1. Пусть $B \rightarrow x$ - правило с номером i .

а) Если $x = Cy$, где C - нетерминал, то $T([A, C], a) = (y[A, B], i)$ для всех $A \in N$, $x \in \text{FIRST}_1(yde)$, таких, что $S \Rightarrow wAe$ и $A \Rightarrow Bd$. $x, y, d, e \in (N \cup \Sigma)^*$, $a \in N$.

б) Если x не начинается нетерминалом, то $T(A, a) = (x[A, B], i)$ для всех $A \in N$, $x \in \text{FIRST}_k(xde)$, таких, что $S \Rightarrow wAe$ и $A \Rightarrow Bd$. $x, d, e \in (N \cup \Sigma)^*$, $a \in N$.

2. $T([A, A], a) = (e, e)$, для всех $A \in N$, $a \in \text{FIRST}_k(e)$, таких, что $S \Rightarrow wAe$.

3. $T(a, av) = \text{выброс}$, для всех $a \in \Sigma$.

4. $T($, e) = \text{допуск}$.

5. Иначе $T(X, a) = \text{ошибка}$.

План работы на семестр

1. Описать алгоритм разбора LC(k)-грамматик.
2. Ввести эквивалентный критерий.
3. Доказать корректность алгоритма.
4. Доказать линейную сложность алгоритма.

Источники и литература

- Ахо А. Ульман Дж., Теория синтаксического анализа, перевода и компиляции, том 1, Синтаксический анализ.
- 11th Annual Symposium on Switching and Automata Theory (swat 1970), deterministic left corner parsing, pages 139-152.

