

Лекции по теории формальных языков

Лекция 9.

Синтаксический анализ для LL(1)-грамматик (окончание).
LL(k)-грамматики.

Синтаксический анализ методом рекурсивного спуска.
Общая схема восходящего синтаксического анализа

Александр Сергеевич Герасимов

<http://gas-teach.narod.ru>

Кафедра математических и информационных технологий
Санкт-Петербургского академического университета
Российской академии наук.
Весенний семестр 2010/11 учебного года

8 апреля 2011 г.

План

- 1 Вычисление множеств выбора правил: повторение
- 2 Вычисление множеств выбора правил: новый материал
- 3 Обработка синтаксических ошибок при LL(1)-анализе
- 4 LL(k)-грамматики
- 5 Синтаксический анализ методом рекурсивного спуска
- 6 Общая схема восходящего синтаксического анализа

План

- 1 Вычисление множеств выбора правил: повторение
- 2 Вычисление множеств выбора правил: новый материал
- 3 Обработка синтаксических ошибок при LL(1)-анализе
- 4 LL(k)-грамматики
- 5 Синтаксический анализ методом рекурсивного спуска
- 6 Общая схема восходящего синтаксического анализа

Вычисление $\text{SELECT}((A \rightarrow \alpha)$ и $\text{FIRST}(X)$

Пусть дана LL(1)-грамматика $G = (\Sigma, \Gamma, P, S)$.

$$\text{SELECT}(A \rightarrow \alpha) = \begin{cases} \text{FIRST}(\alpha), & \text{если } \varepsilon \notin \text{FIRST}(\alpha), \\ (\text{FIRST}(\alpha) \setminus \{\varepsilon\}) \cup \text{FOLLOW}(A), & \text{если } \varepsilon \in \text{FIRST}(\alpha). \end{cases}$$

Для вычисления $\text{SELECT}(A \rightarrow \alpha)$ достаточно уметь вычислять $\text{FIRST}(\alpha)$ и $\text{FOLLOW}(A)$.

Как вычислить $\text{FIRST}(X)$?

- $\text{FIRST}(a) = \{a\}$.
- $\text{FIRST}(A) = \bigcup_{(A \rightarrow \alpha) \in P} \text{FIRST}(\alpha)$.
- Пусть $L \subseteq_t M$ означает, что $L \cap \Sigma \subseteq M$.
- Рассмотрим $\alpha = X_1 \dots X_n$, $n > 0$.
 - ▶ $\text{FIRST}(X_1) \subseteq_t \text{FIRST}(\alpha)$;
 - ▶ $\text{FIRST}(X_2) \subseteq_t \text{FIRST}(\alpha)$, если $\varepsilon \in \text{FIRST}(X_1)$;
 - ▶ $\text{FIRST}(X_3) \subseteq_t \text{FIRST}(\alpha)$, если $\varepsilon \in \text{FIRST}(X_1) \cap \text{FIRST}(X_2)$; ...;
 - ▶ $\text{FIRST}(X_n) \subseteq_t \text{FIRST}(\alpha)$, если $\varepsilon \in \text{FIRST}(X_1) \cap \dots \cap \text{FIRST}(X_{n-1})$;
 - ▶ $\varepsilon \in \text{FIRST}(\alpha)$, если $\varepsilon \in \text{FIRST}(X_1) \cap \dots \cap \text{FIRST}(X_n)$.

Алгоритмы вычисления $\text{FIRST}(A)$ и $\text{FIRST}(\alpha)$

Вход. LL(1)-грамматика $G = (\Sigma, \Gamma, P, S)$.

Выход. Массив множеств $\text{FIRST}(A)$ для каждого $A \in \Gamma$.

1. для каждого $a \in \Sigma$ $\text{FIRST}(a) := \{a\}$;
2. для каждого $A \in \Gamma$
3. **если** $((A \rightarrow \varepsilon) \in P)$ $\text{FIRST}(A) := \{\varepsilon\}$;
4. **иначе** $\text{FIRST}(A) := \emptyset$;
5. **пока** все множества $\text{FIRST}(A)$ не стабилизировались, **повторять**
6. **для** каждого $(A \rightarrow X_1 \dots X_n) \in P$, где $n > 0$, **выполнить**
7. $i := 1$;
8. $\text{FIRST}(A) := \text{FIRST}(A) \cup (\text{FIRST}(X_i) \cap \Sigma)$;
9. **если** $(\varepsilon \in \text{FIRST}(X_i))$
10. **если** $(i < n)$
11. $i := i + 1$; **перейти на 8**;
12. **иначе** $\text{FIRST}(A) := \text{FIRST}(A) \cup \{\varepsilon\}$

Вычисление $\text{FIRST}(\alpha)$ для $\alpha = X_1 \dots X_n$ ($n > 0$): инициализировать $\text{FIRST}(\alpha) := \emptyset$ и выполнить строки 7–12, где A заменён на α .

План

- 1 Вычисление множеств выбора правил: повторение
- 2 Вычисление множеств выбора правил: новый материал**
- 3 Обработка синтаксических ошибок при LL(1)-анализе
- 4 LL(k)-грамматики
- 5 Синтаксический анализ методом рекурсивного спуска
- 6 Общая схема восходящего синтаксического анализа

Вычисление FOLLOW(A)

- $\dagger \in \text{FOLLOW}(S)$.
- Какие символы из $\Sigma \cup \{\dagger\}$ могут стоять непосредственно за нетерминалом B в (каких угодно, необязательно левых) выводах из аксиомы?
- Рассмотрим каждое вхождение B в правую часть каждого правила: $A \rightarrow \alpha B \beta$. Символ a может стоять непосредственно за B в указанных выводах в случаях:
 - ▶ $\beta \Rightarrow^* a\beta'$, тогда $a \in \text{FIRST}(\beta)$;
 - ▶ $\beta \Rightarrow^* \varepsilon$, тогда $a \in \text{FOLLOW}(A)$.

Пример вычисления FIRST(A) и FOLLOW(A)

Грамматика GA_3 :

- (1) $E \rightarrow TE'$, (2) $E' \rightarrow +TE'$, (3) $T \rightarrow FT'$, (4) $T' \rightarrow *FT'$,
 (5) $F \rightarrow (E)$, (6) $F \rightarrow x$, (7) $E' \rightarrow \varepsilon$, (8) $T' \rightarrow \varepsilon$.

Построение FIRST	Построение FOLLOW
(8) $\varepsilon \mapsto T'$	$\dagger \mapsto E$
(7) $\varepsilon \mapsto E'$	(1) $\dagger \mapsto E'$
(6) $x \mapsto F$	(1) $+ \mapsto T$
(5) $(\mapsto F$	(1) $\dagger \mapsto T$
(4) $* \mapsto T'$	(3) $+, \dagger \mapsto T'$
(3) $x, (\mapsto T$	(3) $* \mapsto F$
(2) $+ \mapsto E'$	(3) $+, \dagger \mapsto F$
(1) $x, (\mapsto E$	(5) $) \mapsto E$
	(1) $) \mapsto E'$
	(1) $) \mapsto T$
	(3) $) \mapsto T'$
	(3) $) \mapsto F$

Итрм.	FIRST	FOLLOW
E	$(, x$	$), \dagger$
E'	$+, \varepsilon$	$), \dagger$
T	$(, x$	$+,), \dagger$
T'	$*, \varepsilon$	$+,), \dagger$
F	$(, x$	$+, *,), \dagger$

План

- 1 Вычисление множеств выбора правил: повторение
- 2 Вычисление множеств выбора правил: новый материал
- 3 Обработка синтаксических ошибок при LL(1)-анализе**
- 4 LL(k)-грамматики
- 5 Синтаксический анализ методом рекурсивного спуска
- 6 Общая схема восходящего синтаксического анализа

Ошибки, обнаруживаемые при синтаксическом анализе

- Хороший синтаксический анализатор, кроме сообщения о найденной во входной цепочке ошибке, должен восстановить своё состояние («исправить ошибку») и продолжить анализ.
- МПА обнаруживает ошибку, когда у него нет команды с левой частью (входной символ, верхний символ стека).
- Возможные синтаксические ошибки в цепочке, которая должна быть арифметическим выражением, задаваемым грамматикой GA_3 :
 - ▶ отсутствует операнд;
 - ▶ отсутствует оператор;
 - ▶ правая скобка без парной ей левой скобки;
 - ▶ левая скобка без парной ей правой скобки.

Исправление ошибок: пример

	x	+	*	()	⊖
E	TE'			TE'		
E'		+TE'			ε	ε
T	FT'			FT'		
T'		ε	*FT'		ε	ε
F	x			(E)		

Позиция указателя	Содержимое стека	Исправление
◇)(x + x)(*x + x ⊖	E∇	пропустить)
)(x + x)◇ (*x + x ⊖	T'E'∇	вставить *
)(x + x)(◇ *x + x ⊖	E)T'E'∇	пропустить *
)(x + x)(*x + x ◇ ⊖)T'E'∇	вставить)

Позиция указателя	Содержимое стека	Исправление
(x + x)◇)(*x + x ⊖	T'E'∇	
(x + x)◇)(*x + x ⊖	E'∇	
(x + x)◇)(*x + x ⊖	∇	пропустить) ?

Общая стратегия восстановления после ошибок: режим паники

- Пропуск входных символов до тех пор, пока не встретится символ, (предположительно) позволяющий продолжить анализ.
- Если в момент обнаружения ошибки на вершине стека находится нетерминал A , то можно пропустить все входные символы до первого из $\text{FOLLOW}(A)$ и снять A с вершины стека.
- Также можно пропустить все входные символы до первого из $\text{FIRST}(A)$ и продолжить анализ, оставив A на вершине стека.
- Если при обнаружения ошибки на вершине стека находится терминал, то можно вставить этот терминал во входную цепочку.
- Замечание. Ошибочный входной символ может повредить стек до обнаружения ошибки.

План

- 1 Вычисление множеств выбора правил: повторение
- 2 Вычисление множеств выбора правил: новый материал
- 3 Обработка синтаксических ошибок при LL(1)-анализе
- 4 LL(k)-грамматики**
- 5 Синтаксический анализ методом рекурсивного спуска
- 6 Общая схема восходящего синтаксического анализа

LL(k)-грамматики: определения

- Пусть $G = (\Sigma, \Gamma, P, S)$ — КС-грамматика, $\alpha \in (\Sigma \cup \Gamma)^*$. Тогда
$$\text{FIRST}_k(\alpha) = \{u \in \Sigma^* \mid (\exists \beta (\alpha \Rightarrow^* u\beta) \wedge |u| = k) \vee (\alpha \Rightarrow^* u \wedge |u| < k)\}.$$
- Пусть k — целое неотрицательное число. КС-грамматика G с аксиомой S называется *LL(k)-грамматикой*, если из существования двух (левых) выводов
 - ▶ $S \Rightarrow^* wA\alpha \Rightarrow w\beta\alpha \Rightarrow^* wx$ и
 - ▶ $S \Rightarrow^* wA\alpha \Rightarrow w\gamma\alpha \Rightarrow^* wy$,для которых $\text{FIRST}_k(x) = \text{FIRST}_k(y)$, следует, что $\beta = \gamma$.
- Грамматика называется *LL-грамматикой*, если она является LL(k)-грамматикой для некоторого целого неотрицательного числа k .
- Язык называется *LL(k)-языком*, если существует LL(k)-грамматика, порождающая этот язык.
- Язык называется *LL-языком*, если он является LL(k)-языком для некоторого k .

LL(k)-грамматики: факты и примеры

- Класс всех LL(k)-грамматик \subsetneq класс всех LL($k + 1$)-грамматик.
- Грамматика $G_1 = \{S \rightarrow a^k b | a^k c\}$ не является LL(k)-грамматикой, но является LL($k + 1$)-грамматикой.
- Грамматика

$$G_2 = \{S \rightarrow A|B, A \rightarrow aAb|0, B \rightarrow aBbb|1\}$$

не является LL-грамматикой:

- ▶ $S \Rightarrow^* S \Rightarrow \beta = A \Rightarrow^* x = a^k 0 b^k,$
- ▶ $S \Rightarrow^* S \Rightarrow \gamma = B \Rightarrow^* y = a^k 1 b^{2k},$

$\text{FIRST}_k(x) = \text{FIRST}_k(y)$, но $\beta \neq \gamma$.

- Класс всех LL(k)-языков \subsetneq класс всех LL($k + 1$)-языков.
- $L(G_2)$ не является LL-языком.

План

- 1 Вычисление множеств выбора правил: повторение
- 2 Вычисление множеств выбора правил: новый материал
- 3 Обработка синтаксических ошибок при LL(1)-анализе
- 4 LL(k)-грамматики
- 5 Синтаксический анализ методом рекурсивного спуска**
- 6 Общая схема восходящего синтаксического анализа

Метод рекурсивного спуска

- Применяется для синтаксического анализа цепочек языка, порождаемого нелеворекурсивной грамматикой.
- Программа синтаксического анализатора состоит из набора процедур, по одной для каждого нетерминала грамматики.
- Процедура $A()$ для нетерминала A распознаёт во входной цепочке подцепочку, выводимую из A .

```
процедура  $A()$  {  
    выбрать  $A$ -правило  $A \rightarrow X_1 \dots X_k$ ;  
    for ( $i = 1, \dots, k$ ) {  
        if ( $X_i$  — нетерминал)  
            вызвать процедуру  $X_i()$ ;  
        else if ( $X_i =$  текущий входной символ)  
            перейти к следующему входному символу;  
        else возврат в точку выбора правила или обработка ошибки;  
    }  
}
```

Выбор A -правила в процедуре $A()$

- Если имеется более одного A -правила, то в общем случае требуется возврат и выбор другого A -правила (и повторное чтение части входной цепочки!).
- Для LL(1)-грамматики возвраты не требуются: выбирается такое A -правило $A \rightarrow \alpha$, что текущий входной символ принадлежит $\text{SELECT}(A \rightarrow \alpha)$.

План

- 1 Вычисление множеств выбора правил: повторение
- 2 Вычисление множеств выбора правил: новый материал
- 3 Обработка синтаксических ошибок при LL(1)-анализе
- 4 LL(k)-грамматики
- 5 Синтаксический анализ методом рекурсивного спуска
- 6 Общая схема восходящего синтаксического анализа**

Восходящий анализ. Соглашения

- Синтаксический анализ называется *восходящим анализом* (*анализом снизу вверх*), если дерево вывода анализируемой цепочки строится снизу (от листьев) вверх. При этом цепочка терминалов «сворачивается» в аксиому грамматики и строится (от конца к началу) правый вывод этой цепочки.
- Далее считаем все рассматриваемые КС-грамматики приведёнными, неукорачивающими и однозначными, выводы — правыми, а формы — r -формами, если не будет оговорено иное.

Основа r -формы

- Пусть дан (правый) вывод в грамматике с аксиомой S :

$$S = \alpha_0 \Rightarrow \alpha_1 \Rightarrow \dots \Rightarrow \alpha_n = w .$$

Тогда для каждого $k = 1, \dots, n$ *основой* r -формы α_k называют такое вхождение цепочки β в цепочку α_k , что $\alpha_{k-1} = \gamma_1 A \gamma_2$ и $\alpha_k = \gamma_1 \beta \gamma_2$ для некоторых γ_1, γ_2 и A .

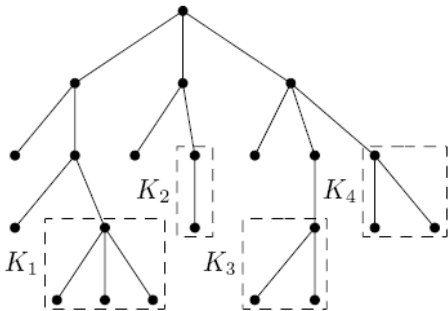
- Часто об основе формы α_k говорят как о цепочке β , подразумевая вхождение цепочки β в цепочку α_k .
- Основа r -формы определяется единственным образом в силу однозначности грамматики.
- Пример. Грамматика

$$G_1 = \{S \rightarrow aFSd|c, F \rightarrow Fb|b\}.$$

Если цепочка $aFbcd$ является r -формой, то её основа — Fb .
 $S \Rightarrow aFSd \Rightarrow aFcd \Rightarrow aFbcd \Rightarrow abbcd$

Куст дерева

Поддерево K дерева T назовём *кустом*, если K состоит из некоторого узла дерева T и всех сыновей этого узла, причём все эти сыновья являются листьями дерева T .



Куст дерева вывода и основа формы

- Пусть (правый) вывод

$$S = \alpha_0 \Rightarrow \alpha_1 \Rightarrow \dots \Rightarrow \alpha_n = w$$

в грамматике (с аксиомой S) представлен деревом вывода T ; T_i, T_{i+1} — стандартные поддеревья дерева T , представляющие выводы форм α_i и α_{i+1} соответственно.

- Тогда T_i получается из T_{i+1} удалением всех листьев самого левого куста.
- На листьях самого левого куста стандартного поддерева дерева вывода написана основа r -формы, вывод которой представлен этим поддеревом.

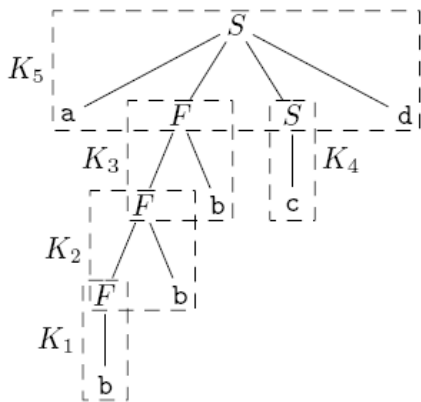
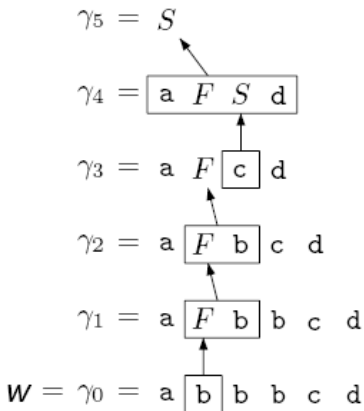
Схема восходящего анализа

- Предполагаем, что заданная цепочка терминалов $w = \gamma_0$ выводима из аксиомы S .
- Находим предполагаемую основу β цепочки γ_0 .
- Если в грамматике нет правила вида $A \rightarrow \beta$, то заключаем, что γ_0 невыводима.
- Иначе выберем правило вида $A \rightarrow \beta$ и произведём *свёртку по этому правилу*, т. е. заменим в цепочке γ_0 её основу на A . Обозначим полученную в результате этой замены цепочку через γ_1 .
- Описанный процесс повторяем для цепочки γ_1 и т. д. (Конкретный алгоритм анализа может предписывать возврат в точки выбора правил для испытания других свёрток.)
- Если получили цепочку $\gamma_n = S$, то имеем (правый) вывод цепочки w :

$$S = \gamma_n \Rightarrow \gamma_{n-1} \Rightarrow \dots \Rightarrow \gamma_1 \Rightarrow \gamma_0 = w .$$

Пример восходящего анализа

Грамматика $G_1 = \{S \rightarrow aFSd|c, F \rightarrow Fb|b\}$.



$S \Rightarrow aFSd \Rightarrow aFcd \Rightarrow aFbcd \Rightarrow aFbbcd \Rightarrow abbbcd = w$

Реализация восходящего анализа при помощи стека (алгоритм типа «перенос-свёртка»)

- Вначале стек пуст.
- Входная цепочка посимвольно переносится в стек до тех пор, пока наверху стека не окажется некоторая основа (правая часть некоторого правила $A \rightarrow \beta$).
- Затем в стеке производится свёртка по правилу $A \rightarrow \beta$.
- Описанный процесс повторяется, пока вся входная цепочка не перенесена в стек и не сделаны все возможные свёртки после завершения этого переноса.
- Если в стеке находится только аксиома грамматики, то входная цепочка допускается.
- Почему каждая основа окажется наверху стека?
- Соглашение: при записи содержимого стека пишем символ дна ∇ слева, а символ на вершине стека — справа и считаем символ на вершине стека последним символом цепочки, представляющей содержимое стека.

Грамматика $G_1 = \{S \rightarrow aFSd|c, F \rightarrow Fb|b\}$, цепочка $w = abbbcd$.

Такт	Содержимое стека	Позиция указателя
1	∇	$\diamond abbbcd \uparrow$
2	∇a	$a \diamond bbbcd \uparrow$
3	∇ab	$ab \diamond bbcd \uparrow$
4	∇aF	$ab \diamond bbcd \uparrow$
5	∇aFb	$abb \diamond bcd \uparrow$
6	∇aF	$abb \diamond bcd \uparrow$
7	∇aFb	$abbb \diamond cd \uparrow$
8	∇aF	$abbb \diamond cd \uparrow$
9	∇aFc	$abbbc \diamond d \uparrow$
10	∇aFS	$abbbc \diamond d \uparrow$
11	$\nabla aFSd$	$abbbcd \diamond \uparrow$
12	∇S	$abbbcd \diamond \uparrow$

Конкатенация содержимого стека и необработанного суффикса цепочки w является r -формой в выводе цепочки w :

$S \Rightarrow aFSd \Rightarrow aFcd \Rightarrow aFbcd \Rightarrow aFbbcd \Rightarrow abbbcd = w$.

Предложение о восходящем анализе при помощи стека

Предложение

Пусть восходящий анализ производится при помощи стека. Тогда верхний символ стека либо входит в основу текущей r -формы, либо входит в текущую r -форму левее этой основы.

Доказательство.

- Индукция по порядковому номеру свёртки.
- База индукции верна: перед первой свёрткой основа находится в стеке и её последний символ находится на вершине стека.
- Индукционный переход. Рассмотрим два идущих подряд применения правил в правом выводе:

$$(1) S \Rightarrow^* \alpha A w \Rightarrow \alpha \underline{\underline{\beta V}} v w \Rightarrow \alpha \beta \underline{\gamma} v w$$

или

$$(2) S \Rightarrow^* \alpha B u A w \Rightarrow \alpha B u \underline{\underline{v}} w \Rightarrow \alpha \underline{\gamma} u v w.$$

Предложение о восходящем анализе при помощи стека: окончание доказательства

- В случае (1)

$$S \Rightarrow^* \alpha Aw \Rightarrow \alpha \underline{\underline{\beta Bvw}} \Rightarrow \alpha \beta \underline{\gamma} vw$$

сразу после свёртки по правилу $B \rightarrow \gamma$ в стеке будет $\alpha\beta B$;
верхний символ стека B принадлежит очередной основе βBv .

- В случае (2)

$$S \Rightarrow^* \alpha BuAw \Rightarrow \alpha Bu \underline{\underline{v}} w \Rightarrow \alpha \underline{\gamma} uvw$$

сразу после свёртки по правилу $B \rightarrow \gamma$ в стеке будет αB ; верхний
символ стека B находится левее очередной основы v .



Два типа алгоритмов восходящего анализа

Мы изучим два типа алгоритмов восходящего анализа.

1. Алгоритмы, основанные на *отношениях предшествования*.
Такой алгоритм находит основу, анализируя пары соседних символов в текущей цепочке.

2. Алгоритмы *LR-анализа*.

Для выделения основы используется информация о символах, расположенных до предполагаемой правой границы основы, и информация о k символах, идущих после правой границы (число k заранее фиксировано).

Литература

Основная литература

- Замятин А. П., Шур А. М. Языки, грамматики, распознаватели: Учебное пособие. Екатеринбург : Изд-во Урал. ун-та, 2007 (электронный вариант книги — на <http://elar.usu.ru>, поиск).

Дополнительная литература

- Ахо А., Лам М., Сети Р., Ульман Дж. Компиляторы: принципы, технологии и инструментарий. М.: ООО "И.Д. Вильямс", 2008.
- Ахо А., Ульман Дж. Теория синтаксического анализа, перевода и компиляции. М.: Мир, 1978.
- Мартыненко Б. К. Языки и трансляции: Учеб. пособие. СПб.: Издательство С.-Петербургского университета, 2004 (электронный вариант книги — на <http://www.math.spbu.ru/user/mbk>).