

04.09.19

LR-анализ

Сегодня посмотрим только на готовый анализатор. Как его создавать посмотрим позднее.

Самая важная тема курса!

Готовый анализатор — автомат с магазинной памятью. Главное отличие автомата от обычного — возможность вынуть более одного символа из стека за раз. Но заглянуть в стек больше, чем на один символ, по-прежнему нельзя! Но стековый символ может нести в себе информацию о символах ниже.

Заучивать не надо — КСГ распознаются МПА.

С какими грамматиками будем работать:

- однозначными, так как нам нужен детерминированный алгоритм
- приведённые

Можем работать с аннулирующими правилами!

Будем рассматривать LR(k) грамматики.

Вспомним нашу арифметическую грамматику.

$$E \rightarrow E^{(1)} + T^{(2)} | T^{(2)}$$

$$T \rightarrow T^{(1)} * F^{(4)} | F^{(4)}$$

$$F \rightarrow (E)^{(5)} | x^{(6)}$$

в GOTO — нетерминалы грамматики. Она говорит о том, что нужно делать во время свёртки

Без цифр — перенос. С цифрой — свёртка по правилу из грамматики.

	A	C	T	I	O	N	GO		TO
	+	*	x	()	\neg	E	T	F
E^1	+					✓			
T^1	(2)	*			(2)	(2)			
F^1	(4)	(4)			(4)	(4)			
(x	(E^2	T^1	F^1
x	(6)	(6)			(6)	(6)			
+			x	(T^2	F^1
*			x	(F^2
E^2	+)				
T^2	(1)	*			(1)	(1)			
F^2	(3)	(3)			(3)	(3)			
)	(5)	(5)			(5)	(5)			
∇			x	(E^1	T^1	F^1

Если происходит перенос, то вершина на стеке остаётся!

Анализатор смотрит на вершину на стеке, на символ, делает то, что написано

Пример

Вершину стека пишем справа!

$(x + x) * x$

Стек	Остаток входной строки	Комментарий
∇	$(x + x) * x \dashv$	
$\nabla($	$x + x) * x \dashv$	
$\nabla(x$	$+x) * x \dashv$	Свёртка по (6)
$\nabla(F^1$	$+x) * x \dashv$	F^1 — потому что на стеке была открытая скобка, а свернулись к F. См. GOTO
$\nabla(T^1$	$+x) * x \dashv$	Свернулись по (4)
$\nabla(E^2$	$+x) * x \dashv$	Свернулись по (2). Не было бы скобки, было бы E^1
$\nabla(E^2 +$	$x) * x \dashv$	
$\nabla(E^2 + x$	$) * x \dashv$	
$\nabla(E^2 + F^1$	$) * x \dashv$	Свёртка по (6)
$\nabla(E^2 + T^2$	$) * x \dashv$	Свёртка по (4). Взяли F, сворачиваемся к T, на стеке +
$\nabla(E^2$	$) * x \dashv$	Свёртка по (1). Его длина — 3, поэтому и из стека выкидываем 3 символа.
$\nabla(E^2)$	$*x \dashv$	
∇F^1	$*x \dashv$	
∇T^1	$*x \dashv$	
$\nabla T^1 *$	$x \dashv$	
$\nabla T^1 * x$	\dashv	
$\nabla T^1 * F^2$	\dashv	
∇T^1	\dashv	
∇E^1	\dashv	
✓	✓	

Внешняя аксиома S'

$$G = (\Sigma, \Gamma, P, S)$$

$$G' = (\Sigma, \Gamma \cup \{S'\}, P \cup \{S' \rightarrow S\}, S')$$

Свёртка по добавленному правилу говорит, что свёртка произошла по первой аксиоме. S' в правых частях отсутствует, поэтому она — начало. Свёртка к ней — это команда допуска.

Распознаватель — таблица, которую мы нарисовали выше. *Автомат* — другая штука, который будем рассматривать далее. Состояния этого автомата — стековый алфавит анализатора.

Автомат LR(0)-пунктов

Опр. LR(k)-**пункт** в $G = (\Sigma, \Gamma, P, S)$ это четвёрка $[A \rightarrow \beta_1 \cdot \beta_2, v]$, где

- $A \rightarrow \beta_1 \beta_2 \in P$
- $|v| = k$ или $|v| < k$ и $v[[v]] = \perp$

Пункт — разделение правила на две части. Он указывает, какую часть продукции мы уже посмотрели в данной точке в процессе синтаксического анализа.

Зачем это нужно? Чтобы отслеживать, насколько мы готовы к свёртке.

Например, пункт $A \rightarrow \cdot XYZ$ указывает, что во входном потоке мы ожидаем встретить строку, порождаемую XYZ . Пункт $A \rightarrow X \cdot YZ$ указывает, что нами уже просмотрела строка, порождаемая X , и мы ожидаем получить из входного потока строку, порождаемую YZ . Пункт $A \rightarrow XYZ \cdot$ говорит о том, что уже обнаружено тело XYZ и что, возможно, пришло время свернуть XYZ в A .

LR(0)-пункт — просто правило с точкой.

$E \rightarrow E + T$:

$E \rightarrow \cdot E + T \quad E \rightarrow E \cdot + T \quad E \rightarrow E + \cdot T \quad E \rightarrow E + T \cdot$

Теперь построим автомат пунктов, который будем использовать для принятия решений в процессе синтаксического анализа — **LR(0)-автомат**. Как этот автомат используется? Если в нём есть переход по текущему символу входной строки, то осуществляется перенос. Если нет, то произведём свёртку по правилу, о котором сообщает пункт.

Опр. **Автоматом LR(0)-пунктов** расширенной грамматики $G = (\Sigma, \Gamma, P, S')$ называется λ -НКА

$I_G = (\Sigma \cup \Gamma, I, \delta, i_0, I)$, где :

- I — множество LR(0)-пунктов грамматики

I — начальные состояния, совпадают с конечными

- $i_0 = [S' \rightarrow \cdot S]$
- δ — множество переходов вида: $[A \rightarrow \beta_1 \cdot x \beta_2] \xrightarrow{x} [A \rightarrow \beta_1 x \cdot \beta_2]$ — базисные
 $[A \rightarrow \beta_1 \cdot B \beta_2] \xrightarrow{\lambda} [B \rightarrow \cdot \beta]$

Если после точки стоит терминал, то из этого пункта по этому терминалу можно перейти в следующее состояние, в котором точка передвинется на шаг вправо.

Если после точки стоит нетерминал, то можно перейти по лямбде в следующее множество пунктов: все правила для этого нетерминала, в самом начале правых частей которых стоит точка.

// Переход по лямбде == замыкание

Базисные пункты — $[A \rightarrow \beta_1 \cdot \beta_2], \beta_1 \neq \lambda$ и i_0

Начальный пункт и все пункты, у которых точки расположены не у левого края

Опр. **Активный префикс** — префикс r -формы, не выходящий за правый конец основы.

== префикс, который может находиться в стеке нашего анализатор.

В стеке не может лежать основа и ещё что-то над ней, потому что как только мы видим основу, мы её сворачиваем

Опр. Пункт $[A \rightarrow \beta_1 \cdot \beta_2, v]$ **допустим** для активного префикса γ , если \exists правый вывод

$S \Rightarrow^+ \gamma' A w \Rightarrow \gamma' \beta_1 \beta_2 w \Rightarrow^* u w$, где $\gamma = \gamma' \beta_1$, v — префикс $w \dashv$

LR(0)-автомат **распознаёт** активные префиксы

Основная теорема LR-анализа

LR(0)-пункт i допустим для активного префикса $\gamma \iff$ в I_G существует пусть из i_0 в i , помеченный γ

Лемма 1

Для активного префикса γ существует допустимый для него базисный пункт.

Д-во

Давайте рассмотрим вывод, в котором появляется данный активный префикс. Нас интересует первое появление. Как только встретили форму, у которой появился такой префикс

$S \Rightarrow^* \gamma \alpha \Rightarrow^* w$

Перед появлением префикса было применено какое то правило, значит, префикс откуда-то взялся, и его можно свернуть до этого нетерминала. Если основа лежит в альфе, то предыдущая форма тоже начиналась с гаммы, а значит это не первое появление. Теперь покажем, что основа лежит на границе гаммы и альфы — $\beta_1 \beta_2$.

Если $\beta_1 \neq \lambda$, тогда $\gamma \neq \lambda$ и $S \Rightarrow^* \gamma' A \alpha; \Rightarrow \gamma' \beta_1 \beta_2 \alpha' \Rightarrow^*$, где $\gamma' \beta_1 = \gamma$, $\beta_2 \alpha' = \alpha$ и пункт $[A \rightarrow \beta_1 \cdot \beta_2]$ — допустим

Если $\gamma = \lambda$, то i_0 допустим для γ

■

Лемма 2

Пункт $[B \rightarrow \cdot \beta]$ допустим для активного префикса $\gamma \iff$ он достижим по λ -переходу из некоторого базисного пункта, допустимого для γ .

Д-во

⇐ Рассмотрим допустимый для γ пункт $[A \rightarrow \beta_1 \cdot B\beta_2]$. Надо перейти к пункту из условия. Существует правый вывод $S' \Rightarrow^* \gamma' Aw \Rightarrow \gamma' \beta_1 \cdot B\beta_2 w \Rightarrow uw$.

Где то тут использовали первую лемму, чтобы обосновать возможность таких переходов в выводе.

Мы не заботились о том, как преобразовывали формы. Но, так как грамматика приведённая, из нетерминала B можно вывести разные терминальные цепочки:

$$S' \Rightarrow^* \gamma' Aw \Rightarrow \gamma' \beta_1 B\beta_2 w \Rightarrow^* \gamma' \beta_1 Bvw \Rightarrow \gamma' \beta_1 \beta w \Rightarrow^* vw$$

⇒

$[B \rightarrow \cdot \beta]$, значит, существует правый вывод: $S' \Rightarrow^* \gamma Bw \Rightarrow \gamma \beta w \Rightarrow^* uw$. По первой лемме для γ существует допустимый базисный пункт $[A \rightarrow \beta_1 \cdot \beta_2]$, который мы хотим найти.

Распишем тот же самый вывод, уточняя переходы через беты.

$$S' \Rightarrow^* \gamma' Au \Rightarrow \gamma' \beta_1 \beta_2 u = \gamma \beta_2 u = [\gamma Bw \Rightarrow \gamma \beta w] \Rightarrow^* \gamma Bw \Rightarrow \gamma \beta w \dots$$

$$1) \gamma \beta_2 u = \gamma Bv$$

$$[A \rightarrow \beta_1 \beta_2] \xrightarrow{\lambda} [B \rightarrow \cdot \beta]$$

$$2) \beta_2 \text{ разворачивается в } B\alpha$$

β_2 начинается с нетерминала C , из которого выводится цепочка, начиная с B

Существенно, что вывод правосторонний! Если бы это было не так, то могла бы случиться такая ситуация: $\beta \rightarrow DC$

$$\gamma \beta_2 u = \gamma C\alpha_1 u \Rightarrow^* \gamma Cu_1 u \Rightarrow \gamma C_2 \alpha_2 u_1 u \Rightarrow \gamma C_2 u_2 u_1 u \Rightarrow^* \dots \gamma Bu_k \dots u_1 u$$

$$A \rightarrow \beta_a C\alpha_1$$

$$C \rightarrow C_2 \alpha_2$$

$$C_2 \rightarrow C_3 \alpha_3$$

$$C_{k-1} \rightarrow B\alpha_k$$

Все лежат в P

$$[A \rightarrow \beta_1 C\alpha_1] \xrightarrow{\lambda} [C \rightarrow \cdot C_2 \alpha_2] \xrightarrow{\lambda} [C_2 \rightarrow \cdot C_3 \alpha_3] \xrightarrow{\lambda} \dots \xrightarrow{\lambda} [C_{k-1} \rightarrow \cdot B\alpha_k] \xrightarrow{\lambda} [B \rightarrow \cdot \beta]$$

■

Доказательство теоремы

⇐

Индукция по длине γ

БИ. $|\gamma| = 0$

$i_0 = [S' \rightarrow \cdot S]$ допустима для $\gamma = \lambda$ по определению

Всё, что достижимо из i_0 , по λ -переходу достижимо для $\gamma = \lambda$ по второй лемме

ШИ. Ненулевая длина, значит на конце есть какой-то символ

$$\gamma = \bar{\gamma}X$$

Последний базисный переход в пути, помеченном γ

$$[A \rightarrow \beta_1 \cdot X\beta_2] \xrightarrow{X} [A \rightarrow \beta_1 X \cdot \beta_2]$$

$[A \rightarrow \beta_1 \cdot X\beta_2]$ допустим для $\bar{\gamma}$ по ПИ $\Rightarrow \exists$ правосторонний вывод

$$S' \Rightarrow^* \gamma' Aw \Rightarrow \gamma' \beta_1 X\beta_2 w \Rightarrow^* ww \Rightarrow [A \rightarrow \beta_1 \cdot X\beta_2] \text{ допустим для } \gamma$$

$$\bar{\gamma} = \gamma' \beta_1 \quad \gamma = \gamma' \beta_1 X$$

Все пункты, достижимые из $[A \rightarrow \beta_1 \cdot X\beta_2]$ по λ -переходам, допустим для γ по второй лемме.

\Rightarrow

Индукция по $|\gamma|$

БИ. $\gamma = \lambda$

i_0 допустим для γ

Воспользуемся леммой два. Все переходы по пустому слову из i_0 достижимы и допустимы.

ШИ.

$$\gamma = \bar{\gamma}X$$

$[A \rightarrow \beta_1 \cdot \beta_2]$ допустим для $\gamma \Rightarrow \exists$ правосторонний вывод:

$$S' \Rightarrow^* \gamma' Aw \Rightarrow \gamma' \beta_1 X\beta_2 w \Rightarrow^* ww$$

$$\gamma = \gamma' \beta_1 = \bar{\gamma}X$$

$$\beta_1 = \beta'_1 X$$

$$[A \rightarrow \beta_1 \cdot \beta_2] = [A \rightarrow \beta'_1 X \cdot \beta_2]$$

$[A \rightarrow \beta'_1 \cdot X\beta_2]$ достижим в I_G по пути, помеченному $\bar{\gamma}$ (по ПИ)

$$[A \rightarrow \beta_1 \cdot X\beta_2] \xrightarrow{X} [A \rightarrow \beta_1 X \cdot \beta_2]$$

■

19.09.19

LR(0)-автомат

Следствие 1. Автомат I_G распознаёт язык активных префиксов грамматики G

Автомат, который мы построили — недетерминированный из-за лямбда-переходов. Поэтому построим эквивалентный ДКА, **LR(0)-автомат** — A_G . Обычно он рассматривается как неполный, поэтому все состояния — терминальные.

Следствие 2. Состояние i автомата A_G , достижимое из i_0 , по пути, помеченному γ , совпадает с множеством пунктов, допустимых для активного префикса γ

Следствие 3. Если в состоянии автомата A_G есть пункт $A \rightarrow \beta_1 X \cdot \beta_2$ и $B \rightarrow \beta'_1 Y \cdot \beta'_2$, то $X = Y$.

Как построить состояния?

Опр. $\text{CLOSURE}(M)$ — минимальное по включению множество пунктов, содержащее M , такое, что если в нём содержится пункт вида $[A \rightarrow \beta_1 \cdot B\beta_2]$, то в $\text{CLOSURE}(M)$ содержатся все пункты вида $[B \rightarrow \cdot \beta]$

Добавляем все пункты, которые можно получить переходом по лямбде

$$\text{CLOSURE}(M) = \{i \mid i \in M \text{ и } \exists (i, j)\text{-путь, помеченный } \gamma\} \cup M$$

$$i_0 = [S' \rightarrow \cdot S]$$

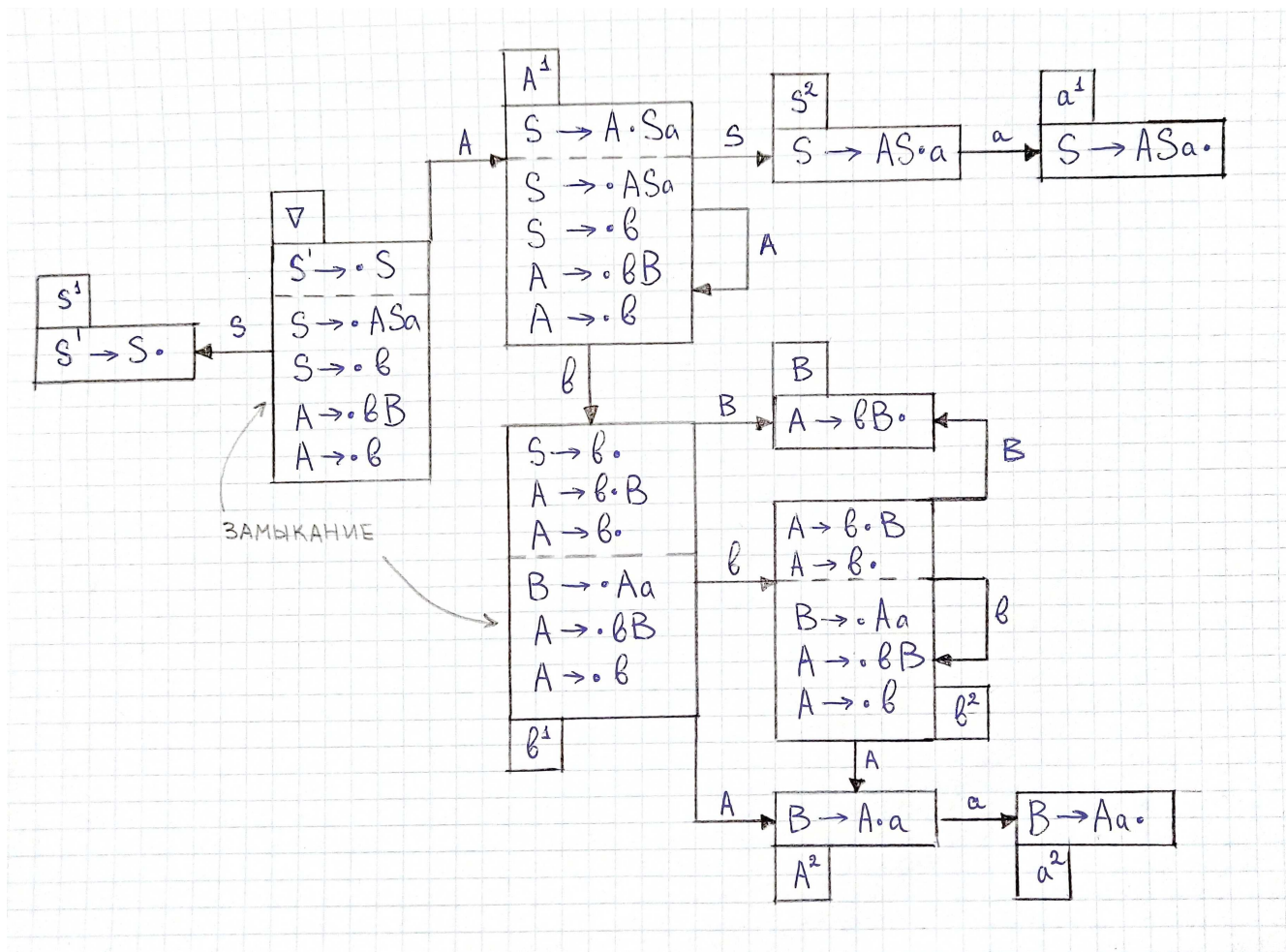
В A_G начальное состояние — $\text{CLOSURE}(\{i_0\})$

$\text{GOTO}(M, X)$ — новые состояния, функция перехода в ДКА

$$\text{GOTO}(M, X) = \text{CLOSURE}(\{[A \rightarrow \beta_1 X \cdot \beta_2] \mid [A \rightarrow \beta_1 \cdot X\beta_2] \in M\})$$

Пример

$$S' \rightarrow S \quad S \rightarrow ASa \mid b \quad A \rightarrow bB \mid b \quad B \rightarrow Aa$$



Опр. LR(0)-грамматика — грамматика, LR(0)-автомат которой не содержит конфликтов. То есть любое состояние, содержащее пункт вида $[A \rightarrow \beta \cdot]$, содержит ровно 1 пункт ???

Конфликты

β и β_1 кончаются на один символ (см. следствие 3)

1. Перенос-свёртка:

$$A \rightarrow \beta \cdot$$

$$B \rightarrow \beta_1 \cdot a \beta_2$$

Если точка стоит перед нетерминалом — будем замыкать, и всё равно появится терминал. И либо первый конфликт, либо второй

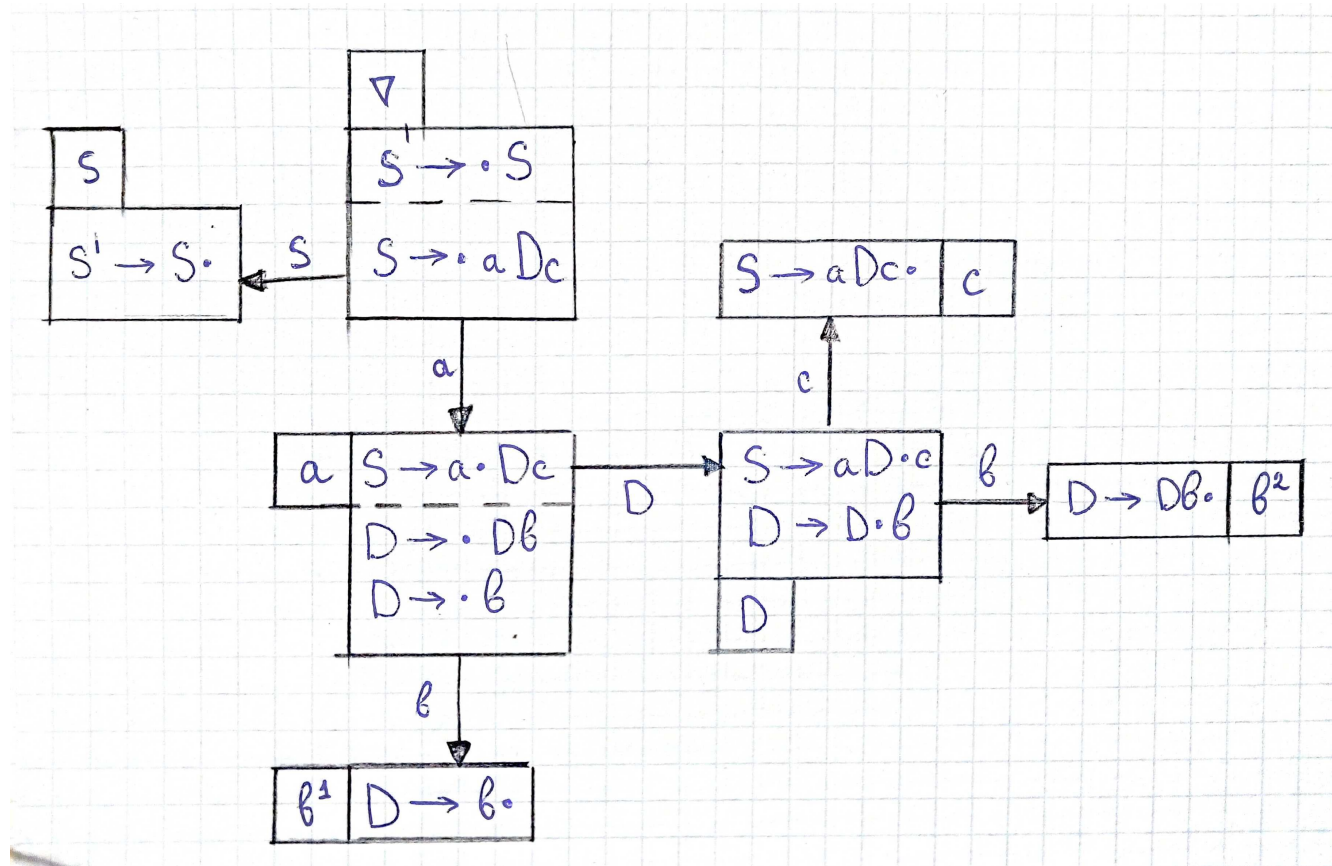
2. Свёртка-свёртка:

$$A \rightarrow \beta_1 \cdot$$

$$B \rightarrow \beta_2 \cdot$$

Другой пример

$S' \rightarrow S \quad S \rightarrow aDc \quad D \rightarrow Db|b$



Стек	Входная строка
∇	$abbbc \dashv$
∇a	$bbbc \dashv$
∇ab^1	$bbc \dashv$
откатываемся на длину основы (правой части)	
∇a	$bbc \dashv$
∇aD	$bbc \dashv$
∇aDb^2	$bc \dashv$
откатываемся на 2 элемента к состоянию D	
∇a	$bc \dashv$
∇aD	$bc \dashv$
∇aDb^2	$c \dashv$
∇a	$c \dashv$
∇aD	$c \dashv$
∇aDc	\dashv
∇	\dashv
∇S	\dashv
\checkmark	\dashv

Построение LR(0)-распознавателя

Для расширенной грамматики

1. Построить A_G
2. Описать таблицы ACTION, GOTO
3. Проиндексировать строки таблицы (стековые алфавиты) состояниями A_G
4. Проиндексировать столбцы ACTION символами из $\Sigma \cup \dashv$
5. Проиндексировать столбцы GOTO символами из Γ
6. $\forall \delta(I, a) : ACTION(I, a) = \delta(I, a), a \in \Sigma$
7. $ACTION([S' \rightarrow S\cdot], \dashv) = \checkmark$

SLR(1)-грамматики

26.09.19

Если видим умножение — то переносим, потому что есть такой переход в автомате. Если видим символы, которые есть в FOLLOW нетерминала слева, то можно к нему свернуться

| ACTION |

	+	*	()	x	\vdash
E^1	+					✓
T^1	(2)	*		(2)		(2)
T^2	(1)	*		(1)		(1)

Ещё пример

$S \rightarrow ac|bDc|Da \quad D \rightarrow a$

Фото автомата от 26.09 ~13:32

Есть конфликт в состоянии a . Либо сворачиваемся к D , либо переносим. $FOLLOW(D) = \{a, c\}$. Если на входе окажется c , то у нас конфликт, непонятно, что делать. НО. Этот пункт оказался из-за переноса замыкания в предыдущем состоянии — $D \rightarrow \cdot a$. А это замыкание было замыкания другого пункта, $S \rightarrow \cdot Da$. Значит, после D должно быть a , и если после него окажется c , то это ерунда. Значит, нужно символ переносить, но по одному FOLLOW мы это понять не можем. Вот так мы и переходим к LR(1)-грамматикам

Добавим к пунктам символы.

LR(1) анализ

Опр. LR(1)-пункт — $[A \rightarrow \beta_1 \cdot \beta_2, a]. \uparrow$ — **ядро**

LR(1)-пункт **допустим**, если \exists правый вывод $S \Rightarrow^+ \gamma' Aw \Rightarrow \gamma' \beta_1 \beta_2 w \Rightarrow^* uw$, где:

- $\gamma = \gamma' \beta_1$ — активный префикс
- v — префикс $w \vdash$
- a — первый символ $w \vdash$

То, с чего должна начинаться оставшаяся строка, чтобы можно было свернуться до A

Опр. **Автоматом LR(1)-пунктов** расширенной грамматики $G = (\Sigma, \Gamma, P, S')$ называется λ -НКА $I_G^1 = (\Sigma \cup \Gamma, I', \delta, i'_0, I')$, где :

- $i'_0 = [S' \rightarrow \cdot S, \vdash]$
- δ — множество переходов вида:
 - $[A \rightarrow \beta_1 \cdot x \beta_2, a] \xrightarrow{x} [A \rightarrow \beta_1 x \cdot \beta_2, a]$ — базисные
 - $[A \rightarrow \beta_1 \cdot B \beta_2, a] \xrightarrow{\lambda} [B \rightarrow \cdot \beta, b], \quad b \in FIRST(\beta_2 a)$ — λ -переходы

Прежде чем сделать переход, надо что-то свернуть. За тем, что выводится из B , должно следовать то, с чего начинается $\beta_2 a$

Детерминированный автомат теперь будет огромный, так как пункты с одинаковыми ядрами но разными символами должны быть разнесены по разным состояниям.

Построим автомат для грамматики, на которой не сработал SLR(1) автомат

$$S' \rightarrow \cdot S, \vdash$$

Начинаем замыкать:

$$S \rightarrow \cdot ac, \vdash S \rightarrow \cdot bDc, \vdash S \rightarrow \cdot Da, \vdash$$

$$D \rightarrow \cdot a, a \text{ — взяли FIRST}(a)$$

Фото 26.09 ~13:58

Надо быть осторожнее с состояниями, в которых возможны и перенос, и свёртка. Когда переносим, нам пофиг на символ, смотрим только на ядро.

A_G^1 — LR(1)-ДКА

$$[A \rightarrow \beta \cdot, a] \in I \Rightarrow ACTION(I, a) = \# (A \Rightarrow \beta)$$

| ACTION | GOTO |

	a	b	c	\vdash	S	D
S				✓		
a^1	(4)		c^1			
a^2			(4)			
a^3				(3)		
b	a^2					D^2
c^1				(1)		
c^2				(2)		
D^1	a^3					
D^2			c^2			
∇	a^1	b			S	D^1