



اصول طراحی کامپایلر

۲

درس نامه‌ی

کاظم فولادی

<http://kazim.fouladi.ir>

ویراست اول: ۱۳۸۸

ویراست دوم: ۱۳۹۳

فهرست مطالب

۱	۴	تحلیل نحوی
۱	۱-۴	مقدمه
۲	۲-۴	بازنمایی نحو زبان با گرامرهای مستقل از متن
۲	۱-۲-۴	گرامر
۲		قواعد تولید
۳		اشتقاق
۴		زبان تولید شده توسط یک گرامر
۵		فرم‌های جمله‌ای یک گرامر
۵	۲-۲-۴	هم‌ارزی گرامرها
۵	۳-۲-۴	گرامر مستقل از متن
۶		قراردادهای نمادگذاری
۶		نمادگذاری بی. ان. اف (BNF)
۶		اشتقاق‌های قانونمند
۷		درخت اشتقاق
۸		درخت اشتقاق جزیی
۹	۳-۴	ابهام در گرامر
۹	۱-۳-۴	زبان ذاتاً مبهم

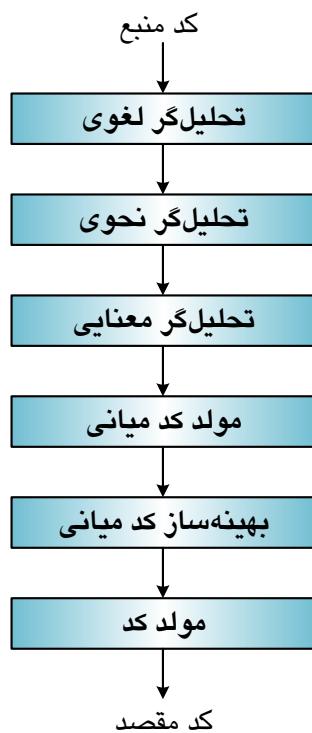
۱۰	تقدیم عملگرها برای رفع ابهام از گرامر عبارات ریاضی	۲-۳-۴
۱۱	شرکت‌پذیری عملگرها: شرکت‌پذیری از راست	۳-۳-۴
۱۲	شرکت‌پذیری عملگرها: شرکت‌پذیری از چپ	۴-۳-۴
۱۳	گرامر دستورهای شرطی	۵-۳-۴
۱۴		۴-۴ تبدیل‌های گرامرها
۱۴	گرامرها کاہش‌یافته، بدون دور و خالی از تهی	۱-۴-۴
۱۴	ایجاد گرامرها کاہش‌یافته	۱-۴-۴
۱۵	نایانه (متغیر)، قاعده و گرامر بازگشتی	۲-۴-۴
۱۵	حذف قواعد بازگشتی از چپ مستقیم	۱-۴-۴
۱۶	حذف بازگشتی از چپ غیر مستقیم	۱-۴-۴
۱۷	پیشوند مشترک	۳-۴-۴
۱۷	فاکتورگیری از چپ برای حذف پیشوند مشترک	۱-۴-۴
۱۸	الگوریتم رفع پیشوندهای مشترک با فاکتورگیری از چپ	۱-۴-۴
۱۹		۵-۴ تجزیه‌گر (Parser)
۱۹	روش‌های تجزیه	۱-۵-۴
۲۰	طبقه‌بندی روش‌های تجزیه	۱-۵-۴
۲۰		۶-۴ طبقه‌بندی چامسکی
۲۱		* تمرین

تحلیل نحوی

SYNTAX ANALYSIS

۱-۴ مقدمه

تحلیل نحوی دومین مرحله‌ی فرایند کامپایل است.



شکل ۱-۴: تحلیل نحوی به عنوان دومین مرحله‌ی فرایند کامپایل

- در تحلیل نحوی، توکن‌ها در ساختار گرامری زبان گروه‌بندی می‌شوند.
- در این مرحله ساختار سلسله‌مراتبی کد منبع به دست می‌آید.



شکل ۲-۴: مرحله‌ی تحلیل نحوی

۲-۴ بازنمایی نحو زبان با گرامرهای مستقل از متن

هر زبان برنامه‌سازی، قواعدی دارد که ساختار نحوی برنامه‌های درست را مشخص می‌کند.

- عبارت‌های منظم برای بازنمایی نحو محدودیت دارند (مانند بازنمایی پرانتزگذاری صحیح).
- گرامرهای مستقل از متن برای بازنمایی نحو زبان‌های برنامه‌سازی مناسب هستند.

۱-۲-۴ گرامر

گرامر، ابزاری برای تولید رشته‌های یک زبان است

گرامر G یک چهارتایی مرتب به صورت

$$G = (N, T, S, P)$$

است که در آن

N مجموعه‌ای متناهی از ناپایانه‌ها (non-terminal) یا متغیرها (variable)

T مجموعه‌ای از پایانه‌ها (terminal) (الفبای گرامر، مجموعه‌ی توکن‌ها)

S یک عنصر خاص از N با نام نماد شروع (start symbol)

P مجموعه‌ای متناهی از قواعد تولید (production rule)

با فرض اینکه N و T ناتهی و مجزا هستند ($N \cap T = \emptyset$)

قواعد تولید

قواعد تولید اساس تعريف یک گرامر است.

هر قاعده‌ی تولید (قاعده) یک زوج مرتب به صورت $(x, y) \in P$ است که به شکل

$$x \rightarrow y$$

نمایش داده می‌شود که در آن

$x \in (N \cup T)^+$ (هر ترکیبی از پایانه‌ها و ناپایانه‌های گرامر بجز رشته‌ی تهی)

$y \in (N \cup T)^*$ (هر ترکیبی از پایانه‌ها و ناپایانه‌های گرامر)

از قواعد تولید در عمل اشتقاق (Derivation) برای تولید رشته استفاده می‌شود.

مثال

$$G = (N, T, S, P)$$

$$T = \{\circ, \text{۱}\}, \quad N = \{S\}, \quad P = \{S \rightarrow \circ S \text{۱}, S \rightarrow \circ \text{۱}\}$$

$$\begin{aligned} S &\Rightarrow \circ \text{۱} \\ S &\Rightarrow \circ S \text{۱} \Rightarrow \circ \circ \text{۱} \\ S &\Rightarrow \circ S \text{۱} \Rightarrow \circ \circ S \text{۱} \text{۱} \Rightarrow \circ \circ \circ \text{۱} \text{۱} \text{۱} \\ &\dots \\ S &\Rightarrow^* \circ^n \text{۱}^n \end{aligned}$$

(*) اشتقاق طی صفر یا چند مرحله را نشان می‌دهد.)

$$L(G) = \{\circ^n \text{۱}^n : n \geq 1\}$$



اشتقاق

اشتقاق (Derivation): اعمال یک قاعده‌ی تولید بر روی یک فرم جمله‌ای اگر رشته‌ی $w = uxv$ به صورت w باشد و قاعده‌ی تولید $y \rightarrow x$ را داشته باشیم، می‌توانیم آن را بر روی این رشته اعمال کنیم و رشته‌ی $z = uyz$ را به دست آوریم. به این عمل اشتقاق می‌گوییم و آن را به صورت

$$w \Rightarrow z$$

نشان می‌دهیم و می‌گوییم w, z را مشتق می‌کند یا z از w مشتق می‌شود.
می‌توان قواعد تولید را به طور پی در پی و به تعداد دلخواه استفاده کرد تا رشته‌ی مورد نظر به دست آید:

$$w_1 \Rightarrow w_2 \Rightarrow \dots \Rightarrow w_n$$

و گفته می‌شود که w_n از w_1 مشتق می‌شود. می‌نویسیم

$$w_1 \Rightarrow^* w_n$$

که به معنی آن است که w_n طی صفر یا چند مرحله از w_1 تولید می‌شود.

قاعده‌ی اشتقاق هرگاه سمت چپ یک قاعده با یک زیررشته از یک فرم جمله‌ای تطابق پیدا کند، می‌توانیم سمت راست آن قاعده را با آن زیررشته جایگزین کنیم.

$$w \Rightarrow z \quad \text{iff} \quad w = w_1 uw_2, z = w_1 vw_2, w_1, w_2 \in (N \cup T)^*, u \rightarrow v \in P$$

نمادهای عمل اشتقاق برای بیان جزئیات عمل اشتقاق از نمادگذاری‌های زیر استفاده می‌شود:

	اشتقاق توسط گرامر G	\Rightarrow_G
	اشتقاق در یک مرحله G	\Rightarrow
(دارای خاصیت بازتابی و تراگذری)	اشتقاق در صفر مرحله یا بیشتر G	\Rightarrow^*
(دارای خاصیت تراگذری)	اشتقاق در یک مرحله یا بیشتر G	\Rightarrow^+
	اشتقاق در n مرحله	\Rightarrow^n
	اشتقاق با قاعده‌ی r	$\xrightarrow{(r)}$

زبان تولید شده توسط یک گرامر

اگر $G = (N, T, S, P)$ یک گرامر باشد، آنگاه مجموعه‌ی

$$L(G) = \{w \in T^* : S \Rightarrow^* w\}$$

زبان تولید شده توسط G نام دارد.

مثال

$$G = (N, T, S, P)$$

$$T = \{a, b\}, \quad N = \{S, A, B\},$$

$$\begin{array}{ll} S & \rightarrow AB \\ A & \rightarrow aA \\ A & \rightarrow \epsilon \\ B & \rightarrow bB \\ B & \rightarrow \epsilon \end{array} \quad \begin{array}{l} (1) \\ (2) \\ (3) \\ (4) \\ (5) \end{array}$$

$$S \xrightarrow{(1)} AB \xrightarrow{(2)} aAB \xrightarrow{(1)} aaAB \xrightarrow{(2)} aaaAB \xrightarrow{(3)} aaaB \xrightarrow{(4)} aaabB \xrightarrow{(5)} aaabbB \xrightarrow{(6)} aaabb$$

$$L(G) = \{a^m b^n : m, n \geq 0\}$$

مثال

$$G = (N, T, S, P)$$

$$T = \{a\}, \quad N = \{S, N, Q, R\}$$

$$\begin{array}{ll} S & \rightarrow QNQ \\ QN & \rightarrow QR \\ RN & \rightarrow NNR \\ RQ & \rightarrow NNQ \\ N & \rightarrow a \\ Q & \rightarrow \epsilon \end{array} \quad \begin{array}{l} (1) \\ (2) \\ (3) \\ (4) \\ (5) \\ (6) \end{array}$$

$$S \xrightarrow{(1)} QNQ \xrightarrow{(2)} QRQ \xrightarrow{(3)} QNNQ \xrightarrow{(4)} QRNQ \xrightarrow{(5)} QNNRQ \xrightarrow{(6)} QNNNNQ \Rightarrow^* aaaa$$

$$L(G) = \{a^{(\forall^n)} : n \geq 0\}$$



فرم‌های جمله‌ای یک گرامر

اگر $w \in L$ باشد، آنگاه

$$S \Rightarrow w_1 \Rightarrow w_2 \Rightarrow \dots \Rightarrow w_n \Rightarrow w$$

را یک اشتقاق از جمله‌ی w می‌گویند.

(sentential
رشته‌های S, w_1, w_2, \dots, w_n و w که حاوی پایانه‌ها و ناپایانه‌ها هستند به فرم جمله‌ای form) موسوم هستند.

مجموعه‌ی همه‌ی فرم‌های جمله‌ای قابل اشتقاق از یک گرامر را با $SF(G)$ نمایش می‌دهیم و به صورت

$$SF(G) = \{W \in (T \cup N)^*: S \Rightarrow^* W\}$$

تعریف می‌کنیم. بدیهی است که

$$L(G) \subset SF(G)$$

۲-۲-۴ همارزی گرامرها

دو گرامر G_1 و G_2 معادل هستند، اگر و فقط اگر هر دو یک زبان واحد را تولید کنند.

$$G_1 \equiv G_2 \quad \text{iff} \quad L(G_1) = L(G_2)$$

۳-۲-۴ گرامر مستقل از متن

گرامر $G = (N, T, S, P)$ را مستقل از متن می‌گوییم، اگر و فقط اگر همه‌ی قواعد P به صورت

$$A \rightarrow x \quad , x \in (N \cup T)^*, \quad A \in N$$

باشد.

در گرامر مستقل از متن، سمت چپ هر قاعده فقط یک ناپایانه وجود دارد.

قراردادهای نمادگذاری

از قراردادهای زیر برای خلاصه نویسی گرامرها استفاده می‌کنیم:

a, b, c, \dots	حروف کوچک ابتدای الفبای انگلیسی	پایانه‌ها
\mathbf{id}, \dots	رشته‌های سیاه (bold)	
$1, 2, 3, \dots$	ارقام	
$+,-,(,),\dots$	عملگرها و نمادهای خاص	
A, B, C, \dots	حروف بزرگ ابتدای الفبای انگلیسی	ناپایانه‌ها
	حروف بزرگ S	
$stmt, expr, \dots$	رشته‌های ایتالیک	
\dots, X, Y, Z	حروف بزرگ انتهای الفبای انگلیسی	نمادهای گرامر (پایانه یا ناپایانه)
$\alpha, \beta, \gamma, \dots$	حروف کوچک یونانی	رشته‌های از نمادهای گرامر
\dots, w, x, y, z	حروف کوچک انتهای الفبای انگلیسی	رشته‌های از پایانه‌ها
	نایانه‌ی سمت چپ اولین قاعده	نماد شروع

نمادگذاری بی. ان. اف (BNF)

Backus-Naur Form

این فرم نمایشی از گرامرها مستقل از متن، در توصیف گرامرها شامل تعداد زیادی پایانه و نایانه به کار می‌رود.

$<\text{expr}>$	درون برآکت گوشیدار	ناپایانه‌ها
number	رشته‌های کاراکتری	پایانه‌ها
	$::=$	نماد جایگزینی

1	$<\text{goal}> ::= <\text{expr}>$
2	$<\text{expr}> ::= <\text{expr}> \text{ } \langle\text{op}\rangle \text{ } <\text{expr}>$
3	number
4	id
5	$\langle\text{op}\rangle ::= +$
6	-
7	*
8	/

در فرم گسترش یافته (EBNF-Extended BNF)، سمت راست قواعد می‌تواند با یک عبارت منظم مشکل از نمادهای گرامر (پایانه یا نایانه) نمایش داده شود.

اشتقاق‌های قانونمند

• اشتقاق چپ‌ترین (Left-most Derivation)

در هر قدم از اشتقاق، سمت چپ‌ترین متغیر در فرم جمله‌ای جایگزین می‌شود (\Rightarrow_{LM}).

• اشتقاق راستترین: (Right-most Derivation)

در هر قدم از اشتقاق، سمت راستترین متغیر در فرم جمله‌ای جایگزین می‌شود (\Rightarrow_{RM}).

در صورت وجود اشتقاق، هم اشتقاق راستترین و هم اشتقاق چپترین وجود خواهد داشت.

▼ مثال

گرامر با قواعد

$$\begin{array}{lcl} S & \rightarrow & aAB \\ A & \rightarrow & bBb \\ B & \rightarrow & A \mid \epsilon \end{array}$$

را در نظر بگیرید. برای رشته‌ی $:abbbb$ یک اشتقاق چپترین:

$$S \Rightarrow a\underline{A}B \Rightarrow ab\underline{B}bB \Rightarrow ab\underline{A}bB \Rightarrow abb\underline{B}bbB \Rightarrow abbbb\underline{B} \Rightarrow abbbb$$

یک اشتقاق راستترین:

$$S \Rightarrow a\underline{A}B \Rightarrow a\underline{A} \Rightarrow ab\underline{B}b \Rightarrow ab\underline{A}b \Rightarrow abb\underline{B}bb \Rightarrow abbbb$$



درخت اشتقاق

اگر $G = (N, T, S, P)$ یک گرامر مستقل از متن باشد،

درخت ریشه‌دار مرتب t_G یک درخت اشتقاق (derivation tree) برای G است، اگر و فقط اگر

۱) ریشه دارای برچسب S باشد.

۲) هر یک از برگ‌ها دارای برچسبی از $\{\epsilon\} \cup T$ باشد.

۳) هر یک از گره‌های داخلی دارای برچسبی از N باشد.

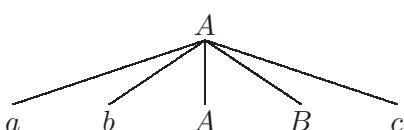
۴) اگرگره‌ای دارای برچسب $A \in N$ باشد و فرزندان آن از چپ به راست به صورت a_1, a_2, \dots, a_n برچسب‌گذاری شده باشد، آنگاه P حاوی قاعده‌ای به شکل $A \rightarrow a_1 a_2 \dots a_n$ باشد.

۵) گره‌ای که دارای فرزندی با برچسب ϵ باشد، هیچ فرزند دیگری نداشته باشد.

حاصل درخت اشتقاق از پیماش عمق - اول برگ‌ها به دست می‌آید و یک رشته در زبان گرامر است.



مثال



درخت اشتقاق برای قاعده‌ی $:A \rightarrow abABC$

درخت اشتقاق جزیی

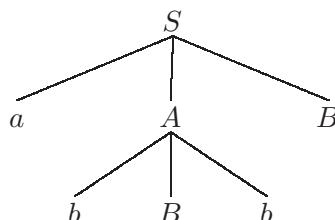
درخت اشتقاق جزیی همانند درخت اشتقاق است با این تفاوت که «هر برگ دارای برجستگی از $\{\epsilon\}$ است»

حاصل درخت اشتقاق جزیی از پیمایش عمق - اول برگ‌ها به دست می‌آید و یک فرم جمله‌ای از گرامر است.

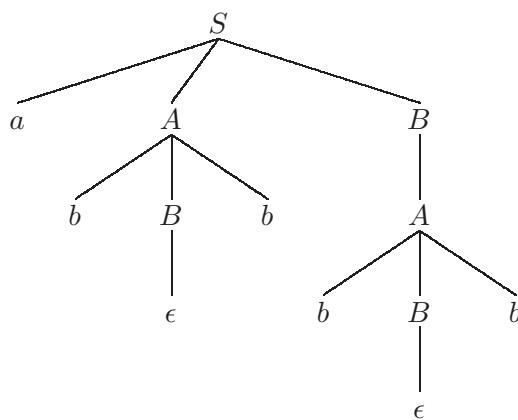
- اگر $G = (N, T, S, P)$ یک گرامر مسقبل از متن باشد، در این صورت به ازای هر $w \in L(G)$ یک درخت اشتقاق برای G داریم که حاصل آن w است و برعکس حاصل درخت اشتقاق در $L(G)$ است.
- اگر t_G یک درخت اشتقاق جزیی برای G باشد، آنگاه حاصل t_G یک فرم جمله‌ای از G است.
- درخت اشتقاق نشان می‌دهد که چه قواعدی برای به دست آوردن یک جمله استفاده شده است، اما ترتیب استفاده از آنها را نشان نمی‌دهد.
- به عبارت دیگر ترتیب به کارگیری قواعد در هر مرحله در نتیجه‌ی نهایی تاثیری ندارد.

مثال

گرامر G با قواعد $S \rightarrow aAB, A \rightarrow bBb, B \rightarrow A|\epsilon$ را در نظر بگیرید:
یک درخت اشتقاق جزیی برای فرم جمله‌ای $abBbB$ از $:G$



یک درخت اشتقاق برای جمله‌ی $abbbb$:



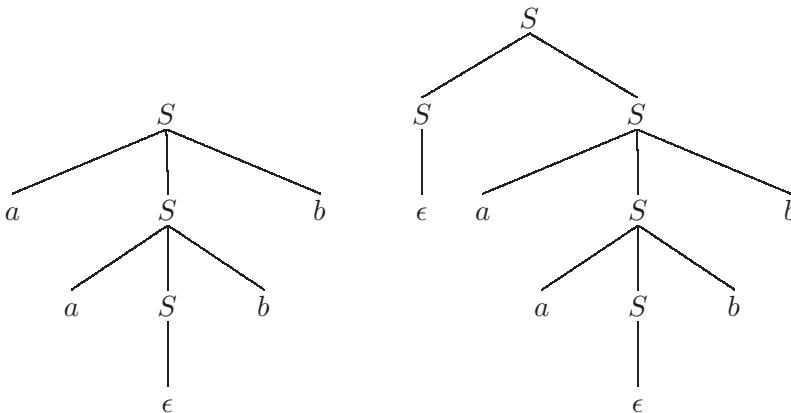
۳-۴ ابهام در گرامر

گرامر مستقل از متن G را **مبهم** (ambiguous) می‌نامیم اگر و فقط اگر حداقل یک $w \in L(G)$ وجود داشته باشد که حداقل دو درخت اشتقاق متفاوت داشته باشد.

ابهام به معنی وجود دو یا چند اشتقاق چپ‌ترین متفاوت یا راست‌ترین متفاوت برای G است.

مثال

گرامر $S \rightarrow aSb|SS|\epsilon$ است، زیرا جمله‌ای $aabb$ دو درخت اشتقاق متفاوت دارد:



۱-۳-۴ زبان ذاتاً مبهم

زبان L را یک زبان ذاتاً مبهم (inherently ambiguous) می‌گوییم، اگر همه‌ی گرامرهای تولیدکننده‌ی L مبهم باشند.

مثال

زبان ذاتاً مبهم $L = \{a^n b^n c^m\} \cup \{a^n b^m c^m\}$ است.

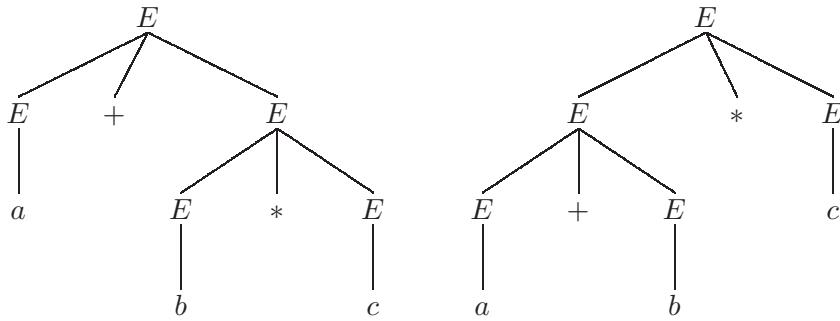
- زبان L را غیرمبهم می‌گوییم اگر و فقط اگر حداقل یک گرامر غیرمبهم برای آن موجود باشد.
- مساله‌ی تشخیص ابهام در یک گرامر مستقل از متن، تصمیم‌ناپذیر است، یعنی هیچ الگوریتمی برای تشخیص یا رفع ابهام وجود ندارد.

مثال

گرامر $G = (\{E\}, \{a, b, c, +, *, (,)\}, E, P)$

$E \rightarrow E + E \mid E * E \mid (E) \mid a \mid b \mid c$

برای عبارت‌های حسابی، مبهم است، زیرا: برای رشته‌ی $a + b * c$ دو درخت اشتقاق وجود دارد:



مثال

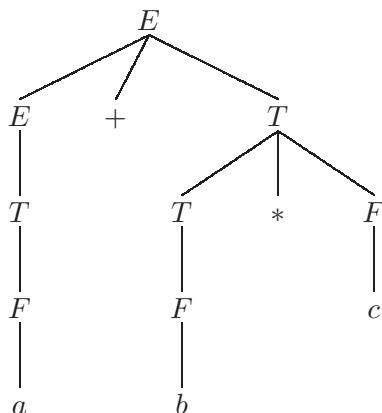
(رفع ابهام از گرامر) برای گرامر مبهم

$$E \rightarrow E + E \mid E * E \mid (E) \mid a \mid b \mid c$$

می‌توان با بازنویسی گرامر، یک گرامر معادل غیر مبهم پیدا کرد:

$$\begin{aligned} E &\rightarrow E + T \mid T \\ T &\rightarrow T * F \mid F \\ F &\rightarrow (E) \mid a \mid b \mid c \end{aligned}$$

در این گرامر رشته‌ی $a + b * c$ تنها یک درخت اشتقاق دارد:



۲-۳-۴ تقدم عملگرها برای رفع ابهام از گرامر عبارات ریاضی

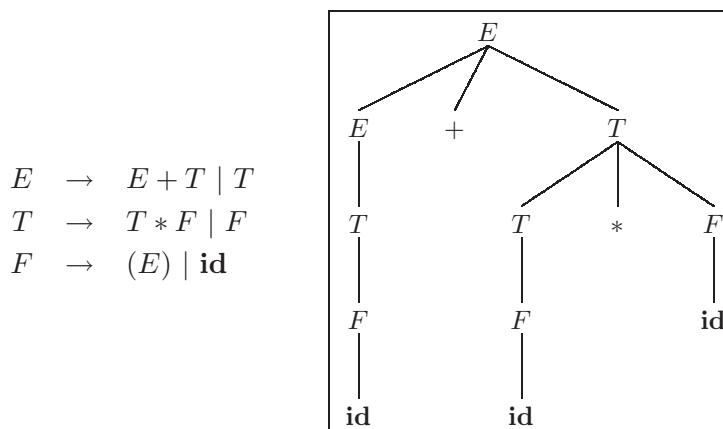
از دو عملگر که در یک سطح پرانتزگذاری ظاهر شده‌اند، آنکه تقدم بالاتری دارد، باید در عمق بیشتری نسبت به دیگری در درخت تجزیه ظاهر شود.

برای این منظور باید ناپایانه‌های جدیدی را وارد گرامر کنیم و به کمک آنها قواعد گرامر را طوری بازنویسی کنیم که قاعده‌ی فوق پیاده شود.

مثال

$$E \rightarrow E + E \mid E * E \mid (E) \mid \text{id}$$

⇓



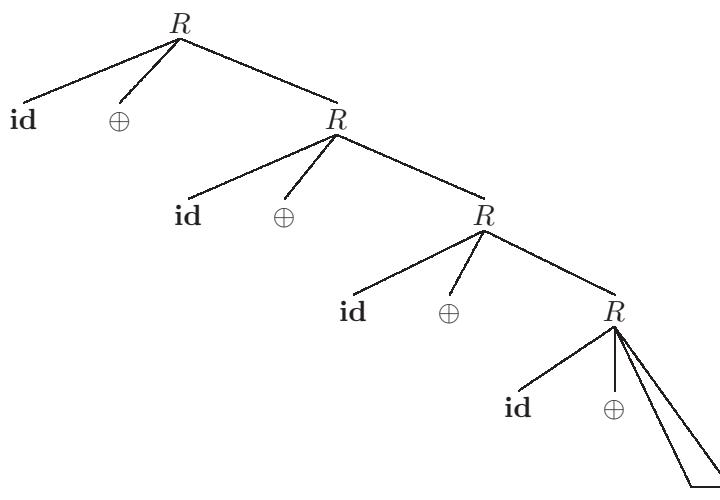
$\text{id} + \text{id} * \text{id}$

۳-۳-۴ شرکت‌پذیری عملگرها: شرکت‌پذیری از راست

عملگرها بی که شرکت‌پذیر از راست (right-associative) هستند، با یک قاعده‌ی بازگشتی از راست (right-recursive) تولید می‌شوند.

$$R \rightarrow \text{id} \oplus R$$

درخت تجزیه برای قاعده‌ی بازگشتی از راست، از راست رشد می‌کند.



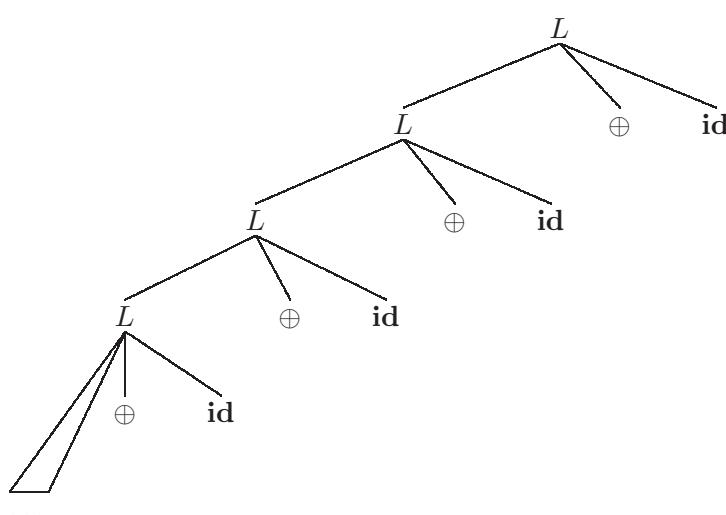
$\text{id} \oplus \text{id} \oplus \text{id} \oplus \text{id} \oplus \dots$

۴-۳-۴ شرکت‌پذیری عملگرها: شرکت‌پذیری از چپ

عملگرهایی که شرکت‌پذیر از چپ (left-associative) هستند، با یک قاعده‌ی بازگشتی از چپ (left-recursive) تولید می‌شوند.

$$L \rightarrow L \oplus \text{id}$$

درخت تجزیه برای قاعده‌ی بازگشتی از چپ، از چپ رشد می‌کند.



$\dots \oplus \text{id} \oplus \text{id} \oplus \text{id} \oplus \text{id}$

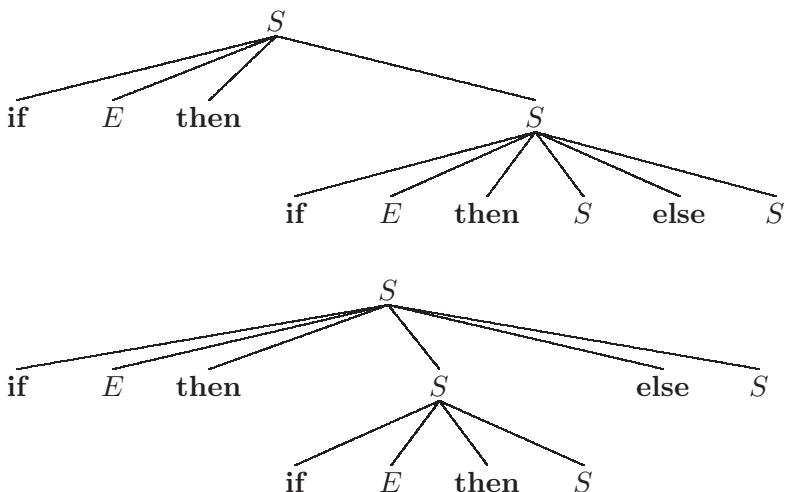
۵-۳-۴ گرامر دستورهای شرطی

$$\begin{array}{l} S \rightarrow \text{if } E \text{ then } S \\ | \\ \text{if } E \text{ then } S \text{ else } S \\ | \\ \text{other} \end{array}$$

این گرامر مبهم است، زیرا فرم جمله‌ای متنه‌ی به رشتی نهایی

if E then if E then S else S

دارای دو درخت اشتقاق متفاوت است:

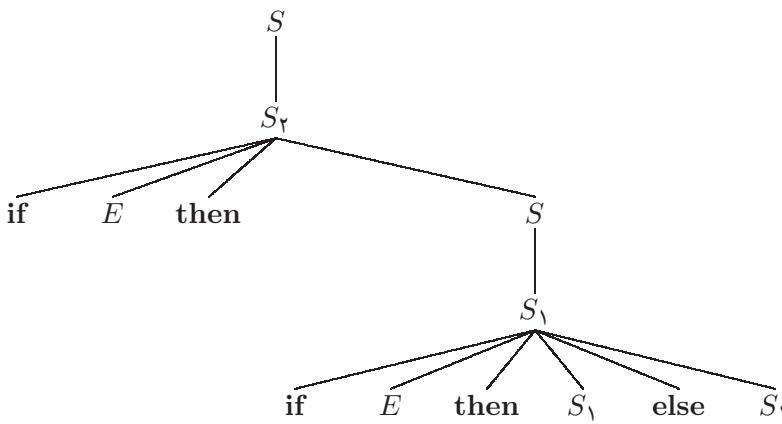


رفع ابهام از گرامر دستورهای شرطی برای این منظور، از قاعده‌ی معنایی زیر استفاده می‌کنیم که در بیشتر زبان‌های برنامه‌سازی در نظر گرفته می‌شود:

هر else باید با نزدیک‌ترین then تطابق‌نیافته، تطابق یابد.

$$\begin{array}{l} S \rightarrow S_1 \mid S_2 \\ S_1 \rightarrow \text{if } E \text{ then } S_1 \text{ else } S_1 \\ | \\ \text{other} \\ S_2 \rightarrow \text{if } E \text{ then } S \\ | \\ \text{if } E \text{ then } S_1 \text{ else } S_1 \end{array}$$

if E then if E then S else S



۴-۴ تبدیل‌های گرامرها

۱-۴-۴ گرامرها کاهش یافته، بدون دور و خالی از تهی

- گرامر کاهش یافته (Reduced) گرامری است که در آن:

- ۱) قواعدی به صورت $A \rightarrow A$ وجود ندارد.
 - ۲) تمام ناپایانه‌ها مولد هستند $x \in T^*$
 - ۳) همهی ناپایانه‌ها دسترس پذیر هستند $\forall A \in N \exists \alpha, \beta \in (N \cup T)^* : S \Rightarrow^* \alpha A \beta$
- گرامر بدون دور (Cycle-free) گرامری است که در آن اشتاقاچی مانند $A \Rightarrow^* A$ موجود نباشد.
 - گرامر خالی از تهی (ϵ -free) گرامری است که قاعده‌ای مانند $\epsilon \rightarrow A$ ندارد.

ایجاد گرامرها کاهش یافته

به ترتیب:

- حذف قواعد $A \rightarrow A$
- حذف قواعد شامل متغیرهای نامولد
- حذف قواعد شامل متغیرهای دسترس ناپذیر

مثال

متغیرهای نامولد

$$\begin{array}{lcl} S & \rightarrow & X \mid Y \\ X & \rightarrow & () \\ Y & \rightarrow & (YY) \end{array}$$

متغیر Y نامولد است و در نتیجه قواعد استفاده کننده از آن بی‌استفاده (useless) هستند.

هر متغیر بازگشته که آلترناتیوی مختوم به یک رشته‌ی پایانی نداشته باشد، نامولد است.

مثال

متغیرهای دسترس ناپذیر

$$\begin{array}{lcl} S & \rightarrow & AB \\ A & \rightarrow & + \mid - \mid \epsilon \\ B & \rightarrow & \text{digit} \mid B \text{ digit} \\ C & \rightarrow & .B \end{array}$$

متغیر C دسترس ناپذیر است و در نتیجه قاعده‌ی آخر بی‌استفاده (useless) است.



جز S هر ناپایانه‌ای که در سمت راست هیچ قاعده‌ی تولیدی قرار نداشته باشد، دسترس ناپذیر است.

۲-۴-۴ ناپایانه (متغیر)، قاعده و گرامر بازگشته

- قاعده‌ی $A \rightarrow \alpha A \beta \rightarrow \alpha A_1 \beta$ که در آن متغیر سمت چپ در سمت راست قاعده نیز ظاهر شده است، یک قاعده‌ی بازگشته (recursive) نام دارد.
- قاعده‌ی $A \rightarrow A\alpha \rightarrow A\alpha_1 \beta$ که در آن متغیر سمت چپ در منتهی الیه چپ سمت راست قاعده نیز ظاهر شده است، یک قاعده‌ی بازگشته از چپ (left recursive) نام دارد.
- قاعده‌ی $A \rightarrow \alpha A \rightarrow \alpha A_1 \beta$ که در آن متغیر سمت چپ در منتهی الیه راست سمت راست قاعده نیز ظاهر شده است، یک قاعده‌ی بازگشته از راست (right recursive) نام دارد.
- متغیری که برای آن حداقل یک قاعده‌ی بازگشته وجود داشته باشد، متغیر بازگشته نام دارد.
- گرامری که در آن حداقل یک قاعده‌ی بازگشته وجود داشته باشد، گرامر بازگشته نام دارد.

حذف قواعد بازگشته از چپ مستقیم

گاهی اوقات وجود قواعد بازگشته از چپ در گرامر برای ما نامطلوب است. برای مثال در تجزیه‌گرهای بالا به پایین وجود این قواعد باعث بروز حلقه‌ی بینهایت می‌شود. در چنین مواردی لازم است این قواعد با قواعد مناسبی که بازگشته از چپ نیستند، جایگزین شوند. به این فرآیند، حذف بازگشته از چپ گفته می‌شود.

برای حذف قواعد بازگشته از چپ متغیر بازگشته A ، قواعد A را به دو بخش بازگشته از چپ و غیر بازگشته از چپ تقسیم می‌کنیم:

$$\begin{array}{lcl} A & \rightarrow & A\alpha_1 \mid A\alpha_2 \mid \dots \mid A\alpha_n \\ A & \rightarrow & \beta_1 \mid \beta_2 \mid \dots \beta_m \end{array}$$

حال این مجموعه قواعد را با مجموعه قواعد زیر جایگزین می‌کنیم:

$$\begin{array}{lcl} A & \rightarrow & \beta_1 A' \mid \beta_2 A' \mid \dots \mid \beta_m A' \\ A' & \rightarrow & \alpha_1 A' \mid \alpha_2 A' \mid \dots \mid \alpha_n A' \mid \epsilon \end{array}$$

گرامر حاصل معادل با گرامر اولیه است.

در حالت ساده، قواعد $A \rightarrow A\alpha \mid \beta A' \mid \dots \mid \alpha A' \mid \epsilon$ با قواعد $A \rightarrow \alpha A' \mid \beta A' \mid \dots \mid \alpha A' \mid \epsilon$ جایگزین می‌شود.

مثال

حذف قواعد بازگشتی از چپ مستقیم:

$$\frac{\begin{array}{l} E \rightarrow E + T \mid T \\ T \rightarrow T * F \mid F \end{array}}{F \rightarrow (E) \mid \text{id}}$$

جایگزینی با:

$$\frac{\begin{array}{l} E \rightarrow TE' \\ E' \rightarrow +TE' \mid \epsilon \\ T \rightarrow FT' \\ T' \rightarrow *FT' \mid \epsilon \end{array}}{F \rightarrow (E) \mid \text{id}}$$

- متغیر A یک متغیر بازگشتی غیر مستقیم است اگر و فقط اگر اشتقاق $\alpha A \beta \Rightarrow^n A$ در آن گرامر ممکن باشد ($n \geq 2$).

- متغیر A یک متغیر بازگشتی از چپ غیر مستقیم است اگر و فقط اگر اشتقاق $A \alpha \Rightarrow^n A$ در آن گرامر ممکن باشد ($n \geq 2$).

- متغیر A یک متغیر بازگشتی از راست غیر مستقیم است اگر و فقط اگر اشتقاق $\alpha A \Rightarrow^n A$ در آن گرامر ممکن باشد ($n \geq 2$).

حذف بازگشتی از چپ غیر مستقیم

برای حذف بازگشتی از چپ غیر مستقیم، به صورت زیر عمل می‌کنیم:

- یک ترتیب بر روی متغیرهای گرامر در نظر می‌گیریم.
- بازگشتی از چپ مستقیم را برای متغیر اول حذف می‌کنیم.
- برای متغیرهای بعدی:

قواعد متغیر قبلی را جایگزین می‌کنیم و بازگشتی از چپ مستقیم را برای قواعد متغیر فعلی حذف می‌کنیم.

- **Input:** grammar G without cycles and ϵ -productions.
- **Output:** An equivalent grammar without left recursion.
- **Number the nonterminals in some order** A_1, A_2, \dots, A_n
- **for** $i = 1$ **to** n **do**
 - **for** $j = 1$ **to** $i - 1$ **do**
 - ▷ replace $A_i \rightarrow A_j \gamma$
 - ▷ with $A_i \rightarrow \delta_1 \gamma | \dots | \delta_k \gamma$
 - ▷ where $A_j \rightarrow \delta_1 | \dots | \delta_k$ are all the current A_j -productions.
 - **Eliminate immediate left-recursion for** A_i
 - ▷ New nonterminals generated above are numbered A_{i+n}

مثال

حذف بازگشتی از چپ غیر مستقیم:

- **Original Grammar:**
 - (1) $S \rightarrow Aa | b$
 - (2) $A \rightarrow Ac | Sd | e$
- **Ordering of nonterminals:** $S \equiv A_1$ **and** $A \equiv A_2$.
- $i = 1$
 - do nothing as there is no immediate left-recursion for S
- $i = 2$
 - replace $A \rightarrow Sd$ by $A \rightarrow Aad | bd$
 - hence (2) becomes $A \rightarrow Ac | Aad | bd | e$
 - after removing immediate left-recursion:
 - ▷ $A \rightarrow bdA' | eA'$
 - ▷ $A' \rightarrow cA' | adA' | \epsilon$

۳-۴-۴ پیشوند مشترک

در مجموعه قواعد

$$A \rightarrow \alpha\beta_1 | \alpha\beta_2 | \dots | \alpha\beta_m$$

α پیشوند مشترک نام دارد.

فاکتورگیری از چپ برای حذف پیشوند مشترک

برای حذف پیشوند مشترک،

$$A \rightarrow \alpha\beta_1 | \alpha\beta_2 | \dots | \alpha\beta_m$$

این مجموعه قواعد را با فاکتورگیری از چپ (left-factoring)، با مجموعه قواعد زیر جایگزین می‌کنیم:

$$\begin{array}{lcl} A & \rightarrow & \alpha B \\ B & \rightarrow & \beta_1 | \beta_2 | \dots | \beta_m \end{array}$$

گرامر حاصل معادل با گرامر اولیه است.

در حالت ساده، قواعد $A \rightarrow \alpha B \rightarrow \alpha\beta_1 \mid \alpha\beta_2$ و $B \rightarrow \beta_1 \mid \beta_2$ با قواعد $A \rightarrow \alpha\beta_1 \mid \alpha\beta_2$ جایگزین می‌شود.

الگوریتم رفع پیشوندهای مشترک با فاکتورگیری از چپ

برای رفع پیشوندهای مشترک، آن قدر روال فوق را تکرار می‌کنیم تا دیگر هیچ دو آلتنتاتیو یک متغیر دارای پیشوند مشترک نباشند.

- **Input:** context free grammar G
- **Output:** equivalent **left-factored** context-free grammar G'
- **for each nonterminal A do**
 - find the longest non- ϵ prefix α that is common to right-hand sides of two or more productions
 - replace
 - $A \rightarrow \alpha\beta_1 \mid \dots \mid \alpha\beta_n \mid \gamma_1 \mid \dots \mid \gamma_m$
with
 - ▷ $A \rightarrow \alpha A' \mid \gamma_1 \mid \dots \mid \gamma_m$
 - ▷ $A' \rightarrow \beta_1 \mid \dots \mid \beta_n$
 - repeat the above process until A has no two productions with a common prefix.

مثال

با فاکتورگیری از چپ از گرامر:

$$S \rightarrow (S) \mid ()$$

و جایگزینی با:

$$\begin{array}{l} S \rightarrow (B \\ B \rightarrow S) \mid () \end{array}$$

در یک مرحله هیچ دو آلتنتاتیو با پیشوند مشترک باقی نمی‌ماند.

مثال

فاکتورگیری از چپ و حذف بازگشتی از چپ

- **Original grammar:**
 $S \rightarrow (S) \mid SS \mid ()$
- **To remove immediate left-recursion, we have**
 - $S \rightarrow (S)S' \mid ()S'$
 - $S' \rightarrow SS' \mid \epsilon$
- **To do left-factorizing, we have**
 - $S \rightarrow (S'' \mid ()S'$
 - $S'' \rightarrow S)S' \mid)S'$
 - $S' \rightarrow SS' \mid \epsilon$

۵-۴ تجزیه‌گر (Parser)

جزیه‌گر (parser)، الگوریتمی است که برای رشته‌ی w یک اشتقاق می‌یابد و یا می‌گوید اشتقاق ممکن نیست.

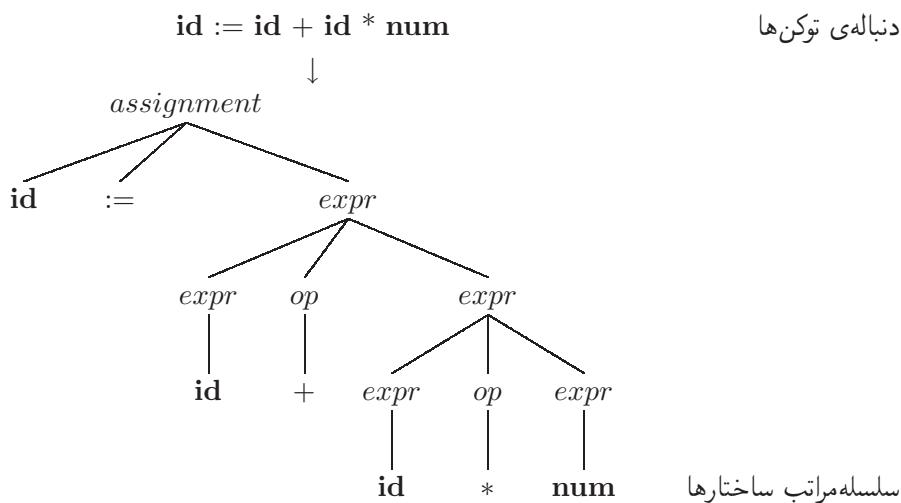
خروجی تجزیه‌گر، نوعی بازنمایی از درخت تجزیه (درخت اشتقاق) است.

به ازای هر گرامر مستقل از متن، الگوریتمی وجود دارد که هر رشته‌ی $w \in L(G)$ را در تعداد مراحلی که متناسب با $|w|^3$ است، تجزیه می‌کند.

مثال

نمونه‌ای ساده از تحلیل نحوی:

$$\begin{array}{lcl} \textit{assignment} & \rightarrow & \text{id} := \textit{expr} \\ \textit{expr} & \rightarrow & \text{id} \mid \text{num} \mid \textit{expr} \text{ op } \textit{expr} \mid (\textit{expr}) \\ \textit{op} & \rightarrow & + \mid - \mid * \mid / \end{array}$$



۱-۵-۴ روش‌های تجزیه

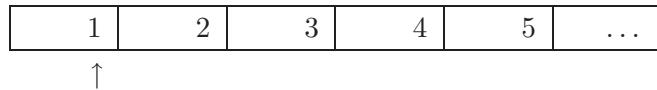
- روش‌های بالا به پایین (Top-Down Parsing)
 - درخت تجزیه از بالا به پایین (از ریشه به سمت برگ‌ها) ساخته می‌شود.

- روش‌های پایین به بالا (Bottom-Up Parsing)
 - درخت تجزیه از پایین به بالا (از برگ‌ها به سمت ریشه) ساخته می‌شود.

طبقه‌بندی روش‌های تجزیه

روش‌های بالا به پایین	
Recursive Descent with Backtracking	- روش نزولی بازگشتی با عقبگرد
Recursive Descent	- روش نزولی بازگشتی
	- روش‌های LL(k)
روش‌های پایین به بالا	
Operator Precedence	- روش تقدم عملگر
Simple Precedence	- روش تقدم ساده
Simple LR	- روش SLR ساده
Lookahead LR	- روش LALR
Canonical LR	- روش CLR
Earley, CYK, ...	- روش‌های عمومی

نماد پیش‌بینی (Lookahead) نماد پیش‌بینی، نمادی است که باید برای تشخیص قاعده‌ی مناسب، از پیش آن را ببینیم. پرسش این است که به چند توکن بعد با شروع از توکن جاری باید نگاه کنیم؟



۶-۴ طبقه‌بندی چامسکی

زبان	گرامر	قالب قواعد گرامر	پذیرنده
زبان‌های شمارش‌پذیر بازگشتی	گرامرهای بدون محدودیت	$\alpha \rightarrow \beta$	ماشین‌های تورینگ
نوع صفر	Unrestricted	$\alpha \neq \epsilon$	Turing Machine (TM)
زبان‌های حساس به متن	گرامرهای حساس به متن	$\alpha \rightarrow \beta$ $ \alpha \leq \beta $	آutomاتی کراندار خطی
نوع یک	Context-Sensitive		Linear Bounded Automata (LBA)
زبان‌های مستقل از متن	گرامرهای مستقل از متن	$\alpha \rightarrow \beta$ $\alpha \in N, \alpha = 1$	آtomاتی پسته‌ای
نوع دو	Context-Free		Push-down Automata (PDA)
زبان‌های منظم	گرامرهای منظم	$\alpha \rightarrow \beta$ $\beta = aB$ با $\beta = a$	آtomاتی متنه‌ای
نوع سه	Regular	$\beta = Ba$ یا $\beta = a$	Finite Automata (FA)

$$\mathcal{L}_3 \subset \mathcal{L}_2 \subset \mathcal{L}_1 \subset \mathcal{L}_0.$$

◀ تمرین

۱. گرامر مستقل از متن زیر را در نظر بگیرید:

$$S \rightarrow SS + \mid SS * \mid a$$

آ) نشان دهید که رشته‌ی $aa + a*$ را می‌توان توسط این گرامر تولید نمود.

ب) درخت تجزیه (parse tree) را برای این رشته بسازید.

پ) این گرامر چه زبانی را تولید می‌کند؟

۲. چه زبانی توسط هر یک از گرامرهای زیر تولید می‌شود؟

$$S \rightarrow ^\circ S^\circ \mid ^\circ \circ$$

$$S \rightarrow +SS \mid -SS \mid a$$

$$S \rightarrow S(S)S \mid \epsilon$$

$$S \rightarrow aSbS \mid bSaS \mid \epsilon$$

$$S \rightarrow a \mid S + S \mid SS \mid S * \mid (S)$$

۳. کدام یک از گرامرهای تمرین فوق مبهم است؟ چرا؟

۴. برای هر یک از زبان‌های زیر، یک گرامر مستقل از متن بسازید:

آ) عبارت‌های ریاضی با عملگرهای $+$, $*$ و توان به شکل پیشوندی (prefix)

ب) عبارت‌های ریاضی با عملگرهای $+$, $*$ و توان به شکل پسوندی (postfix)

پ) لیست‌های «شرکت‌پذیر از چپ» (left-associative) از شناسه‌ها که با کاما جدا شده‌اند.

ت) لیست‌های «شرکت‌پذیر از راست» (right-associative) از شناسه‌ها که با کاما جدا شده‌اند.

ث) عبارت‌های ریاضی از اعداد و شناسه‌ها با چهار عملگر دو عملوندی $/$, $*$, $-$, $+$

ج) به عملگرهای ریاضی (۴ث) جمع و تفریق تک عملوندی (unary) را اضافه نمایید.

۵. گرامرهای تمرین فوق را به گونه‌ای بازنویسی کنید که مبهم نباشند.

۶. گرامرهای تمرین فوق را به گونه‌ای بازنویسی کنید که پیشوند مشترک یا بازگشتی از چپ نداشته باشند.

۷. یک گرامر مستقل از متن برای اعداد رومی بنویسید. برای توضیحات بیشتر در مورد ساختار اعداد

رومی، به آدرس زیر مراجعه کنید:

http://fa.wikipedia.org/wiki/عدد_رومی-رومانی

۸. در ارتباط با گرامر زیر، پرسش‌های زیر را پاسخ دهید:

$$S \rightarrow Aa \mid aBb \mid aBbDc$$

$$A \rightarrow Sa \mid \epsilon$$

$$B \rightarrow CD$$

$$C \rightarrow Ba \mid \epsilon$$

$$D \rightarrow Db \mid E \mid \epsilon$$

$$E \rightarrow Ec \mid D \mid \epsilon$$

آ) صورت کاهش‌یافته‌ی گرامر فوق را بنویسید.

ب) بازگشته از چپ مستقیم و غیرمستقیم را از گرامر فوق رفع کنید و پس از فاکتورگیری از چپ، گرامر حاصل را بنویسید.

۹. بازگشته از چپ را از گرامر زیر برطرف کنید.

$$\begin{array}{lcl} S & \rightarrow & Sab \mid SaS \mid X \\ X & \rightarrow & Xc \mid a \mid b \end{array}$$

۱۰. گرامر زیر را در نظر بگیرید:

$$S \rightarrow S^{\circ\circ} \mid 11S \mid {}^{\circ}S1 \mid {}^{\circ}S \mid \epsilon$$

آ) نشان دهید که رشته‌ی $11{}^{\circ}0{}^{\circ}0$ به وسیله‌ی گرامر فوق قابل تولید است.

ب) چند درخت تجزیه (parse tree) متفاوت برای این رشته وجود دارد؟ همه‌ی آنها را بسازید.

پ) با اضافه کردن قواعد لازم، گرامر فوق را به گرامر مستقل از متنی تبدیل کنید که تمام رشته‌ها به طول زوج را تولید می‌کند.

۱۱. گرامر رو برو دو عملگر $\$$ و $\#$ را برای رشته‌های متشكل از حروف الفبای انگلیسی تعریف می‌کند. عملگر $\#$ دو عملوند خود را به هم می‌چسباند (concatenation) و عملگر $\$$ عملوند اول خود را معکوس (reverse) کرده و آن را به انتهای عملوند دوم خود می‌چسباند. مثال برای عملگر $\$$:

$$abc\$def = defcba$$

عملگر $\$$ نسبت به عملگر $\#$ اولویت بالاتری دارد. با توجه به گرامر زیر:

$$\begin{array}{lcl} S' & \rightarrow & S \\ S & \rightarrow & T\$S \mid T\#S \mid T \\ T & \rightarrow & string \end{array}$$

آ) حاصل عبارت $xyz\$hij\#cs$ چیست؟

ب) یک اشتقاء راست‌ترین، درخت تجزیه (parse tree) و درخت نحو (syntax tree) را برای رشته $a\#b\$c$ نشان دهید.

پ) با توجه به اینکه عملگر $\$$ نسبت به عملگر $\#$ اولویت بالاتری دارد، گرامر معادلی را بیابید که عاری از ابهام باشد.

۱۲. برای هر یک از زبان‌های زیر یک گرامر طراحی کنید:

آ) مجموعه‌ی تمامی رشته‌های متشكل از 0 و 1 که هر 0 بلافاصله با حداقل یک 1 دنبال شود.

ب) مجموعه‌ی تمامی رشته‌های متشكل از 0 و 1 به طوری که تعداد صفرها و یک‌ها مساوی باشد.

پ) مجموعه‌ی تمامی رشته‌های متشكل از 0 و 1 به طوری که زیر رشته 11 را شامل نشوند.

۱۳. برای هر یک از زبان‌های زیر یک گرامر غیر مبهم طراحی کنید:

آ) رشته‌های به شکل 1^* به طوری که تعداد صفرها از تعداد یک‌ها بیشتر باشد.

ب) رشته‌های آینه‌ای (palindrome) روی الفبای 0 و 1 (به رشته‌ای آینه‌ای گفته می‌شود که خواندن از چپ به راست با خواندن از راست به چپ یکسان باشد).

- پ) رشته‌های متشکل از پرانتزها و برآکت‌های متوازن، مانند: $(([])([]))$
۱۴. گرامر زیر را در نظر بگیرید که برای عملیات منطقی نوشته شده است. در رابطه با اولویت عملگرها
نسبت به یکدیگر چه نتیجه‌ای می‌توان گرفت؟

$$\begin{array}{lcl} X & \rightarrow & X \text{or } Y \mid Y \\ Y & \rightarrow & Y \text{and } Z \mid Z \\ Z & \rightarrow & \text{not } Z \mid (X) \mid \text{true} \mid \text{false} \end{array}$$
