

مروری بر مبانی نظریه محاسبه

امیرحسین رجبی کوروش مسلمی

۱ گرامرهای مستقل از متن

گرامرها^۱ یک مدل محاسباتی هستند و برخلاف اتوماتاهای متناهی که کلمات یک زبان را با قرار گرفتن در استیت اکسپت پذیرش می‌کنند، کلمات زبانی را به کمک تعدادی قاعده^۲ تولید می‌کنند.

تعریف ۱. گرامر مستقل از متن یک چهارتایی $G = (V, \Sigma, S, P)$ است که به ترتیب V مجموعه متناهی متغیرهای گرامر، Σ مجموعه متناهی حروف الفبا، S متغیر شروع گرامر ($S \in V$) و P مجموعه متناهی از قواعد که به فرم $A \rightarrow \alpha$ هستند که $A \in V$ و $\alpha \in (V \cup \Sigma)^*$.

مثال: گرامر زیر زبان $L = \{a^n b^n \mid n \in \mathbb{N}\}$ را تولید می‌کند. (زبان $AnBn$)

$$G = (\{S\}, \{a, b\}, S, \{S \rightarrow aSb, S \rightarrow \Lambda\})$$

$$S \rightarrow aSb$$

$$S \rightarrow \Lambda$$

یا به طور خلاصه $S \rightarrow aSb \mid \Lambda$. واضح است که $\Lambda \in (V \cup \Sigma)^*$ و به معنای قابلیت جایگزینی S با رشته تهی است؛ در نتیجه تولید یا پروداکشن می‌تواند خاتمه یابد. همچنین به حروف Σ ترمینال نیز گفته می‌شود.

برای نشان دادن نحوه تولید یک کلمه به کمک زبان از نمادهای \Rightarrow ، \Rightarrow^* و \Rightarrow^n استفاده می‌شود.

مثال: گرامر G که در بالا توصیف شد را در نظر بگیرید، رشته $aabb$ به صورت زیر از قواعد گرامر مشتق^۳ می‌شود:

$$S \Rightarrow aSb \Rightarrow aaSbb \Rightarrow aa\Lambda bb \Rightarrow aabb$$

یا به طور خلاصه $S \Rightarrow_G^* aabb$ نوشته می‌شود که یعنی با تعدادی مرحله به رشته مذکور می‌رسیم. همچنین با زیرنویس G مشخص می‌کنیم که از قواعد گرامر G استفاده شده است.

تعریف ۲. اگر $G = (V, \Sigma, S, P)$ یک گرامر مستقل از متن باشد، منظور از زبان تولید شده توسط G ،

$$L(G) = \{x \in \Sigma^* \mid S \Rightarrow_G^* x\}$$

است. زبانی مستقل از متن است اگر گرامر مستقل از متنی وجود داشته باشد که آن را تولید کند.

مثال: زبان حاوی کلمات پالیندروم، زبانی مستقل از متن است. (زبان Pal)

¹Context Free Grammar

²Production Rule

³Derivation

گرامر مستقل از متن زیر را داریم: (با فرض $\Sigma = \{a, b\}$)

$$S \rightarrow aSa \mid bSb \mid \Lambda$$

مثال: زبان حاوی کلماتی که تعداد a در آنها با تعداد b برابر است، زبانی مستقل از متن است. (زبان $AEqB$)

گرامر مستقل از متن زیر را داریم: (با فرض $\Sigma = \{a, b\}$)

$$S \rightarrow aSbS \mid bSaS \mid \Lambda$$

اثبات. باید نشان دهیم هر کلمه مثل x که $n_a(x) = n_b(x)$ توسط این قواعد تولید می‌شود و همچنین باید ثابت کنیم هر کلمه‌ای که توسط این قواعد تولید می‌شود چنین ویژگی را دارد. گزاره دوم واضح است. گزاره اول را ثابت می‌کنیم. به استقرای قوی حکم را روی $|x|$ ثابت می‌کنیم. پایه اسقرا $x = \Lambda$ است که توسط گرامر تولید می‌شود. اکنون فرض کنید حکم برای همه رشته‌های به طول زوج کوچکتر از $|x|$ برقرار باشد، حکم را برای x نشان می‌دهیم. تابع $f: \Sigma^* \rightarrow \mathbb{Z}$ را با ضابطه $f(x) = n_a(x) - n_b(x)$ تعریف می‌کنیم. همچنین فرض می‌کنیم $x = x_1 x_2 \dots x_n$ ($x_i \in \Sigma$). بدون کاسته شدن از کلیت مسئله فرض کنید $x_1 = a$. اولین i که $f(x_1 x_2 \dots x_i) = 0$ را در نظر می‌گیریم. در این صورت نتیجه می‌شود $x_i = b$. پس $x_2 \dots x_{i-1} \in AEqB$ و همچنین $x_{i+1} \dots x_n \in AEqB$. طبق فرض استقرا این دو کلمه به کمک قواعد G از متغیر S قابل تولیدند. با توجه به قاعده $S \rightarrow aSbS$ کلمه x نیز توسط گرامر قابل تولید است. \square

قضیه ۱. اگر L_1 و L_2 زبان‌هایی مستقل از متن باشند، زبان‌های $L_1 \cup L_2$ ، $L_1 L_2$ و L_1^* نیز مستقل از متن هستند.

اثبات. فرض کنید $G_1 = (V_1, \Sigma, S_1, P_1)$ و $G_2 = (V_2, \Sigma, S_2, P_2)$ به ترتیب گرامرهای مستقل از متن برای زبان‌های L_1 و L_2 باشند.

اگر $L(H_1) = L_1 \cup L_2$ ، آنگاه $H_1 = (V_1 \cup V_2, \Sigma, S, P_1 \cup P_2 \cup \{S \rightarrow S_1 \mid S_2\})$

اگر $L(H_2) = L_1 L_2$ ، آنگاه $H_2 = (V_1 \cup V_2, \Sigma, S, P_1 \cup P_2 \cup \{S \rightarrow S_1 S_2\})$

اگر $L(H_3) = L_1^*$ ، آنگاه $H_3 = (V_1, \Sigma, S, P_1 \cup P_2 \cup \{S \rightarrow SS_1 \mid \Lambda\})$

\square

۱.۱ گرامرهای منظم و خطی راست و چپ

تعریف ۳. گرامر مستقل از متن $G = (V, \Sigma, S, P)$ را گرامر خطی راست ^۴ می‌گوییم اگر همه قواعد آن به فرم $A \rightarrow \sigma B$ یا $A \rightarrow \Lambda$ باشد که در آن $\sigma \in \Sigma$ و $A, B \in V$.

تعریف ۴. گرامر مستقل از متن $G = (V, \Sigma, S, P)$ را گرامر خطی چپ ^۵ می‌گوییم اگر همه قواعد آن به فرم $A \rightarrow B\sigma$ یا $A \rightarrow \Lambda$ باشد که در آن $\sigma \in \Sigma$ و $A, B \in V$.

تعریف ۵. گرامر مستقل از متن $G = (V, \Sigma, S, P)$ را گرامر منظم ^۶ می‌گوییم اگر خطی راست یا خطی چپ باشد.

قضیه ۲. زبان L منظم است اگر و تنها اگر گرامر منظم G وجود داشته باشد که $L = L(G)$.

اثبات. (شهودی) اثبات معادل بودن گرامرهای خطی راست با اتوماتاهای متناهی ساده است. فرض کنید $M = (Q, \Sigma, q_0, A, \delta)$ یک اتوماتای متناهی قطعی ^۷ باشد. گرامر خطی راست $G = (V, \Sigma, S, P)$ را این‌گونه می‌سازیم. قرار می‌دهیم $V = Q$ و $S = q_0$. همچنین برای q $\delta(p, \sigma) = q$ که

⁴Right Linear Grammars (RLG)

⁵Left Linear Grammars (LLG)

⁶Regular Grammars

⁷Deterministic Finite Automata (DFA)

$p, q \in Q$ و $\sigma \in \Sigma$ قاعده σq نیز در P قرار خواهد داشت. ضمناً اگر $r \in A$ آنگاه قاعده $r \rightarrow \Lambda$ نیز در P را دارد. به سادگی می‌توان دید که $L(G) = L(M)$. از طرفی اگر $G = (V, \Sigma, S, P)$ گرامر خطی راست باشد، اتوماتای متناهی غیر قطعی^۹ $M = (Q, \Sigma, q_0, A, \delta)$ را این‌گونه می‌سازیم. قرار می‌دهیم $Q = V$ و $q_0 = S$. همچنین تابع $\delta : Q \times \Sigma \rightarrow Q$ برای هر قاعده $X \rightarrow \sigma Y$ در P ، این‌گونه تعریف می‌کنیم $\delta(X, \sigma) = Y$ و برای هر قاعده به فرم $X, X \rightarrow \Lambda$ در A خواهد بود. دوباره به سادگی می‌توان دید که $L(M) = L(G)$.

برای معادل بودن گرامرهای خطی چپ کار کمی دشوارتر است. اگر هر قاعده در یک گرامر خطی چپ مثل $X \rightarrow Y\sigma$ با قاعده $X \rightarrow \sigma Y$ جایگزین شود، زبان تولید شده توسط این گرامر شامل عکس تمام رشته‌های زبان گرامر اولیه خواهد بود. (نه فقط شامل آنها بلکه دقیقاً همان‌ها می‌شود.) از طرفی در صورت داشتن یک اتوماتای متناهی قطعی می‌توان یک اتوماتای متناهی دیگر ساخت که دقیقاً عکس کلمات پذیرفته شده توسط اتوماتای نخست را پذیرش کند. همچنین هر اتوماتای غیرقطعی معادل یک اتوماتای قطعی است. پس به کمک این دو گزاره، فرآیند زیر برای اثبات قضیه طی می‌شود:

اتوماتای پذیرنده معکوس کلمات \longleftrightarrow اتوماتای قطعی \longleftrightarrow اتوماتای غیرقطعی \longleftrightarrow گرامر خطی راست \longleftrightarrow برعکس کردن قواعد \longleftrightarrow گرامر خطی چپ

□

پس بنابر قضیه بالا گرامرهای منظم قدرت محاسباتی بیشتری از اتوماتاهای متناهی ندارند.

۲.۱ ابهام و درخت اشتقاق

مثال: گرامر $G = (\{S\}, \{a, +, *\}, S, P)$ با قواعد زیر مبهم است:

$$S \rightarrow S + S$$

$$S \rightarrow S * S$$

$$S \rightarrow a$$

دو درخت اشتقاق^{۱۰} برای کلمه $a + a * a$ در زیر کشیده شده‌اند:



پس درخت اشتقاق به صورت تصویری نشان می‌دهد برای تولید یک کلمه چه قواعدی استفاده شده‌اند.

تعریف ۶. گرامر مستقل از متن G مبهم است اگر کلمه $x \in L(G)$ موجود باشد به طویکه x دارای حداقل دو درخت اشتقاق باشد.

در نتیجه گرامر فوق مبهم است. رفع ابهام نیازمند کمی خلاقیت است. در اینجا می‌توان گرامر زیر را ارائه کرد:

$$S \rightarrow S + T \mid T$$

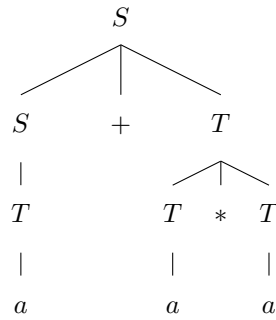
$$T \rightarrow T * T \mid a$$

^۹چرا اتوماتای غیر قطعی تعریف می‌کنیم؟ چرا مستقیم سراغ قطعی نمی‌رویم؟

^۹Nondeterministic Finite Automata

^{۱۰}Derivation tree

با درخت اشتقاق زیر برای کلمه $a + a * a$:



۳.۱ فرم‌های نرمال

به مسئله تصمیم در گرامرهای مستقل از متن فکر کنید. یعنی با داشتن پروداکشن‌های گرامر G و رشته x مشخص کنیم x توسط قواعد G تولید می‌شود یا خیر. دانستن نکاتی درباره فرم پروداکشن‌ها به ما در پاسخ به این سوال کمک می‌کند. در مراحل اشتقاق^{۱۱} یک رشته به کمک قواعد گرامر تعداد ترمینال‌های رشته (t) کاهش نمی‌یابند. اگر تصور کنیم گرامر هیچ لاندا-قاعده‌ای^{۱۲} نداشته باشد، (یعنی هیچ قاعده‌ای به فرم $X \rightarrow \Lambda$) در این صورت طول رشته (l) نیز در مراحل اشتقاق کاهش نمی‌یابد. (در تعبیر درخت اشتقاق، تعداد رئوس در هر سطح درخت کاهش نمی‌یابد.) همچنین اگر بدانیم گرامر هیچ یونیت-پروداکشنی^{۱۳} (یعنی هیچ قاعده‌ای به فرم $X \rightarrow Y$ که Y متغیر باشد) نیز نداشته باشد، می‌توانیم بگوییم تعداد رئوس در هر سطح درخت اشتقاق اکیدا صعودی است. یعنی $l + t$ اکیدا صعودی است. اما وقتی رشته x تولید شده باشد داریم $2|x| = |x| + |x| = l + t$. از طرفی هنگام شروع تولید از متغیر S داریم $1 = 0 + 1 = l + t$. در نتیجه تولید رشته x بیش از $2|x| - 1$ مرحله نیاز ندارد. پس با بررسی همه درخت‌های اشتقاق با حداکثر عمق $2|x| - 1$ می‌توانیم درباره وجود x در $L(G)$ تصمیم بگیریم.

در ادامه نشان می‌دهیم هر گرامر را می‌توان به گرامری تبدیل کرد که هیچ لاندا-قاعده و هیچ یونیت پروداکشنی نداشته باشد. سپس پا را فرا تر گذاشته و نشان می‌دهیم می‌توان قاعده‌ها را محدود به قاعده‌هایی به فرم $X \rightarrow YZ$ کرد که به فرم نرمال چامسکی^{۱۴} معروف است. در نهایت به معرفی الگوریتم CYK^{۱۵} می‌پردازیم که در مرتبه چندجمله‌ای عضویت رشته دلخواه x در $L(G)$ را بررسی می‌کند.

قضیه ۳. برای هر گرامر G ، گرامر G' وجود دارد که هیچ لاندا قاعده‌ای نداشته باشد و همچنین $L(G') = L(G) - \{\Lambda\}$.

اثبات. مراجعه شود به صفحه ۱۵۰ کتاب مارتین. □

قضیه ۴. برای هر گرامر بدون لاندا-قاعده G ، گرامر G' وجود دارد که هیچ یونیت پروداکشنی نداشته باشد و همچنین $L(G') = L(G)$.

اثبات. مراجعه شود به صفحه ۱۵۲ کتاب مارتین. □

تعریف ۷. گرامر G به فرم نرمال چامسکی است اگر هر قاعده گرامر به فرم $A \rightarrow BC$ یا $A \rightarrow \sigma$ باشد که $\sigma \in \Sigma$.

قضیه ۵. برای هر گرامر مستقل از متن G ، گرامر G' به فرم نرمال چامسکی وجود دارد که $L(G') = L(G) - \{\Lambda\}$.

اثبات. مراجعه شود به صفحه ۱۵۲ و ۱۵۳ کتاب مارتین. □

الگوریتم CYK برای بررسی وجود یک کلمه در زبان تولیدشده توسط یک گرامر مستقل از متن استفاده می‌شود. کاربرد آن در بخش تجزیه^{۱۶} کامپایلرها است. این الگوریتم همچنین نحوه اشتقاق کلمه را نیز نشان می‌دهد یعنی اینکه در هر مرحله از کدام قاعده استفاده شود تا کلمه تولید شود. برای ورودی، این الگوریتم نیاز دارد تا گرامر ورودی به فرم نرمال چامسکی باشد.

¹¹Derivation steps

¹² Λ -production

¹³Unit production

¹⁴Chomsky normal form

¹⁵Cocke-Younger-Kasami algorithm

¹⁶Parser

مثال: گرامر زبان $AEqB$ را در نظر بگیرید:

$$S \rightarrow aSbS \mid bSaS \mid \Lambda$$

ابتدا لاند-قاعده‌های آن را حذف می‌کنیم:

$$S \rightarrow abS \mid aSbS \mid aSb \mid ab$$

$$S \rightarrow baS \mid bSaS \mid bSa \mid ba$$

این گرامر به طور طبیعی شامل هیچ یونیت-پروداکشنی نیست، بنابراین می‌توانیم مراحل تبدیل به فرم چامسکی را آغاز کنیم. برای این کار ابتدا حروف ترمینال a و b را در قاعده‌هایی که سمت راست آن‌ها ترمینال نیست، به ترتیب با A و B تعویض می‌کنیم. بنابراین خواهیم داشت:

$$S \rightarrow ABS \mid ASBS \mid ASB \mid AB$$

$$S \rightarrow BAS \mid BSAS \mid BSA \mid BA$$

$$A \rightarrow a$$

$$B \rightarrow b$$

گرامر فوق طبق تعریف ۷ در فرم نرمال چامسکی قرار ندارد بنابراین آن را به صورت زیر اصلاح می‌کنیم:

$$S \rightarrow AY \mid XY \mid YX \mid BX$$

$$S \rightarrow AB \mid BA \mid XB \mid YA$$

$$X \rightarrow AS$$

$$Y \rightarrow BS$$

اجرای الگوریتم CYK برای رشته $aababb$: ^{۱۷} ^{۱۸}

$\{A\}$	\emptyset	$\{X\}$	\emptyset	$\{X\}$	$\{S\}$
	$\{A\}$	$\{S\}$	$\{X\}$	$\{S\}$	\emptyset
		$\{B\}$	$\{S\}$	$\{Y\}$	\emptyset
			$\{A\}$	$\{S\}$	\emptyset
				$\{B\}$	\emptyset
					$\{B\}$

فرض کنید رشته مورد بررسی این الگوریتم $x = x_1x_2 \dots x_n$ باشد. در این صورت جدولی $n \times n$ تشکیل می‌شود و در درایه سطر i ام و ستون j ام آن، مجموعه متغیرهایی قرار می‌گیرند که می‌توانند زیررشته $x_i \dots x_j$ را تولید کنند. این مجموعه را X_{ij} با می‌نامیم. به وضوح مقادیر زیر قطر اصلی جدول معنایی ندارند. مقادیر قطر اصلی را به کمک قاعده‌های $C \rightarrow \sigma$ پر می‌کنیم ($\sigma \in \Sigma$)؛ یعنی در صورتی که $x_i = \sigma$ و قاعده $C \rightarrow \sigma$ در گرامر باشد، آنگاه

^{۱۷} با تشکر فراوان از کوروش مسلمی که در تکمیل بخش الگوریتم CYK کمک کرد.

^{۱۸} این الگوریتم مشابه الگوریتم ضرب زنجیر ماتریس‌ها (Chain Matrix Multiplication) است که در آن نحوه پرانتز گذاری برای ضرب n ماتریس $A_1A_2 \dots A_n$ با تعداد سطر و

ستون‌های زیر

$$A_1(r_1, c_1), A_2(r_2 = c_1, c_2), \dots, A_n(r_n = c_{n-1}, c_n)$$

مطلوب است که کل تعداد عملیات‌های ضرب کمینه شود.

$C \in X_{ii}$. سپس سراغ قطره‌های بالاتر می‌رویم. این بار برای پرکردن مجموعه X_{ij} ، باید زیررشته $x_i \dots x_j$ را به دو زیررشته $x_i \dots x_k$ و $x_{k+1} \dots x_j$ بشکنیم ($i \leq k \leq j-1$) و در صورت وجود قاعده‌ای در گرامر چون $U \rightarrow VW$ که $U \in X_{ij}$ باشد، آنگاه $VW \in \bigcup_{k=i}^{j-1} X_{ik}X_{k+1j}$. به عنوان مثال در اینجا قطر اصلی به کمک A و B پر شده‌اند چون تنها متغیرهایی هستند که تک کاراکتر تولید می‌کنند. سپس در قطر بالاتر فقط متغیر S قرار گرفته‌است چون $AB \mid BA \mid S \rightarrow AB$ و از طرفی AA و BB در سمت راست هیچ متغیری نیستند. اکنون مجموعه X_{25} را در نظر بگیرید. بنابر جدول $X_{25} = \{S\}$. برای تحقیق این موضوع $X_{22}X_{35} = \{A\} \cdot \{Y\}$ ، $X_{22}X_{35} = \{S\} \cdot \{S\}$ و $X_{24}X_{55} = \{X\} \cdot \{B\}$ باید بررسی شوند که AY و XB در سمت راست قواعد ما موجودند. در نهایت، در صورت وجود متغیر شروع گرامر (S) در مجموعه X_{1n} الگوریتم با خروجی YES پایان می‌یابد. در غیر این صورت خروجی NO خواهد بود. الگوریتم از مرتبه $\mathcal{O}(n^3|V|^2|P|)$ خواهد بود که P مجموعه پروداکشن‌های گرامر است و V مجموعه متغیرهای آن است. (چون P و V مستقل از n هستند الگوریتم از مرتبه $\mathcal{O}(n^3)$ است.) لازم به ذکر است که هنگام محاسبه X_{ij} شکستن رشته به دو زیررشته کافی است زیرا گرامر به فرم نرمال چامسکی است و قواعد آن دو متغیره هستند.

روش الگوریتمی تبدیل گرامر بدون لاندا-قاعده زبان AEB به فرم نرمال چامسکی به صورت زیر است:

$$\begin{aligned} S &\rightarrow X_a X_b S \mid X_a S X_b S \mid X_a S X_b \mid X_a X_b \\ S &\rightarrow X_b X_a S \mid X_b S X_a S \mid X_b S X_a \mid X_b X_a \\ X_a &\rightarrow a \\ X_b &\rightarrow b \end{aligned}$$

و در نهایت:

$$\begin{aligned} S &\rightarrow X_a Y_1 \mid X_a Y_2 \mid X_a Y_3 \mid X_a X_b \\ S &\rightarrow X_b Y_4 \mid X_b Y_5 \mid Y_2 Y_5 \mid X_b X_a \\ Y_1 &\rightarrow S X_b \\ Y_2 &\rightarrow X_b S \\ Y_3 &\rightarrow S Y_2 \\ Y_4 &\rightarrow S X_a \\ Y_5 &\rightarrow X_a S \end{aligned}$$

اکنون گرامر در فرم نرمال چامسکی قرار دارد و می‌توان الگوریتم CYK را اجرا کرد. در جدول زیر اجرای این الگوریتم برای رشته $aababb$ نشان داده شده است:

$\{X_a\}$	\emptyset	$\{Y_5\}$	\emptyset	\emptyset	$\{S\}$
	$\{X_a\}$	$\{S\}$	$\{Y_4, Y_5\}$	$\{S\}$	$\{Y_1\}$
		$\{X_b\}$	$\{S\}$	$\{Y_1, Y_2\}$	\emptyset
			$\{X_a\}$	$\{S\}$	$\{Y_1\}$
				$\{X_b\}$	\emptyset
					$\{X_b\}$