



جلسه‌ی ۱۴: زبان‌های مستقل از متن

نگارنده: حمید پورربیع رودسری

مدیرس: دکتر شهرام خزائی

۱ مقدمه

گرامر مستقل از متن^۱ یک سیستم نمادگذاری برای کلاسی از زبان‌هاست که با نام زبان‌های مستقل از متن^۲ شناخته می‌شوند. این زبان‌ها کلاس بزرگتری را نسبت به کلاس زبان‌های منظم شامل می‌شوند و می‌توانند زبان‌هایی را توصیف کنند که منظم نیستند. این گرامرها به خاطر ساختار ذاتی خود برای توصیف رشته‌هایی که دارای ساختار تودرتو هستند مناسب هستند.

۲ چند مثال

بهتر است برای آشنایی بیشتر با زبان‌های مستقل از متن با چند مثال شروع کنیم:

مثال ۱ همانطور که در بحث زبان‌های منظم گفته شد زبان L به صورت

$$L = \{w \in \{0, 1\}^* : w = 0^n 1^n, n \geq 1\}$$

که شامل همه‌ی رشته‌هایی به صورت $0^n 1^n$ است را نمی‌توان با یک ماشین حالت محدود و یا یک عبارت منظم نمایش داد زیرا منظم نیست. ما اینجا نشان می‌دهیم که این زبان، یک زبان مستقل از متن است و می‌توان آن را با گرامر مستقل از متن تولید کرد. برای تعریف این زبان با استفاده از یک گرامر می‌توان از دو ویژگی زیر استفاده کرد:

۱. رشته‌ی 0^1 عضو این زبان است؛

۲. در صورتی که رشته‌ی w عضو این زبان باشد، $0w1$ نیز عضو این زبان است.

حال این دو شرط را در گرامر مستقل از متن می‌توان به صورت دو «قانون» زیر نوشت:

$$\begin{aligned} S &\rightarrow 0^1 \\ S &\rightarrow 0S1 \end{aligned}$$

با بیان بالا می‌توان تمام رشته‌های موجود در زبان را با استفاده از استقرا به دست آورد.

^۱Context-Free Grammers

^۲Context-Free Languages

مثال ۲ می خواهیم گرامری برای زبانی که شامل تمام رشته‌هایی که قلب مستوی^۳ است را به دست آوریم. برای راحتی فرض می کنیم تمام رشته‌ها از $\{0, 1\}$ تشکیل شده‌اند. در این صورت رشته‌هایی مثل 1101011 و 1010101 در زبان موجود هستند اما رشته‌هایی مثل 11001 یا 101010 در این زبان نیستند. این زبان را می توان به صورت $L = \{w \in \{0, 1\}^* : w = w^R\}$ نشان داد. با کمی تفکر در رشته‌های این زبان به دو نکته‌ی اساسی پی می بریم:

۱. رشته‌های 0 و 1 و ϵ در این زبان موجودند.
۲. در صورتی که رشته‌ی w در این زبان باشد رشته‌های $1w1$ و $w0$ نیز عضو این زبان هستند. با استفاده از نکات بالا می توان این زبان را با استفاده از قوانین زیر تولید کرد:

$$\begin{aligned} S &\rightarrow \epsilon \\ S &\rightarrow 0 \\ S &\rightarrow 1 \\ S &\rightarrow 1S1 \\ S &\rightarrow 0S0 \end{aligned}$$

۳ تعریف رسمی

تعریف ۱ یک گرامر مستقل از متن یک چهارتایی به صورت $G = (V, T, P, S)$ می باشد که در آن:

۱. V یک مجموعه‌ی متناهی از نمادها به نام متغیر^۴ است.
- هر کدام از این متغیرها بیانگر یک زبان (مجموعه‌ی ای از چند رشته) هستند در مثال‌های بالا V فقط شامل متغیر S است.
۲. T مجموعه‌ی متناهی از نمادها به نام پایانه (ترمینال)^۵ است. در مثال‌های ذکر شده T برابر است با $\{0, 1\}$.
۳. P یک مجموعه‌ی متناهی از قوانین تولید^۶ به صورت $A \rightarrow \alpha$ است که $A \in V$ و $\alpha \in (V \cup T)^*$.
۴. S متغیر شروع^۷ است که $S \in V$.

مثال ۳ گرامر قلب مستوی به صورت زیر تعریف می شود:

$$G = \{\{S\}, \{0, 1\}, P, \{S\}\}$$

که در آن P بیانگر قوانینی است که در مثال ۲ بیان شده است.

مثال ۴ می خواهیم زبان عبارات جبری ساده که از a و b و 1 و 0 تشکیل شده‌اند و فقط جمع و ضرب در بین

^۳ رشته‌های قلب مستوی (palindromes)، رشته‌هایی هستند که از هر دو طرف به یک صورت خوانده می شوند مانند گرگ در فارسی یا ("Madam, I'm adam") در انگلیسی.

^۴ Variable
^۵ Terminal
^۶ Productions
^۷ Start symbol

آن‌ها مجاز است را بنویسم (متغیرها باید با a یا b شروع شوند و بعد از آن‌ها می‌تواند عدد یا حرف باشد):

$$\begin{aligned} E &\rightarrow I \\ E &\rightarrow E + E \\ E &\rightarrow E * E \\ E &\rightarrow (E) \\ I &\rightarrow a \\ I &\rightarrow b \\ I &\rightarrow Ia \\ I &\rightarrow Ib \\ I &\rightarrow I \setminus \\ I &\rightarrow I \circ \end{aligned}$$

قانون اول بیانگر عبارات ساده‌ی جبری است که فقط شامل یک متغیر هستند. قانون دوم عباراتی که جمع در آن‌ها وجود دارد را توصیف می‌کند، قانون سوم ضرب و قانون چهارم پرانتز دور عبارات را. متغیر کمکی I هم زبانی را تولید می‌کند که شامل همه‌ی متغیرهای قابل قبول است که در واقع یک زبان منظم به صورت

$$(a + b)(a + b + 0 + 1)^*$$

است.

نکته: مرسوم است قوانین تولیدی که ابتدای آن‌ها متغیر A است به عنوان قوانین تولید A شناخت. می‌توان قوانین تولید را به صورت خلاصه به این صورت نوشت که همه‌ی متغیرها را یک بار در سمت چپ آورد و تمام قواعد مربوط به هر متغیر را در سمت راست نوشت و آن‌ها را با یک خط عمودی از هم جدا کرد. برای مثال می‌توان قوانین $A \rightarrow \alpha_1, A \rightarrow \alpha_2, A \rightarrow \alpha_3, \dots$ را به صورت خلاصه به شکل $A \rightarrow \alpha_1 \mid \alpha_2 \mid \alpha_3 \mid \dots$ نیز نوشت. شایان ذکر است که می‌توان قوانین نوشته شده در مثال ۴ را به صورت زیر خلاصه کرد:

$$\begin{aligned} E &\rightarrow I \mid E + E \mid E * E \mid (E) \\ I &\rightarrow a \mid b \mid Ia \mid Ib \mid I \setminus \mid I \circ \end{aligned}$$

۴ اشتقاق

تعریف ۲ فرض کنید $G = (V, T, P, S)$ یک گرامر مستقل از متن باشد. به ازای هر $A \in V$ و هر $\alpha, \beta \in (V \cup T)^*$ می‌گوییم $\alpha \gamma \beta$ یک اشتقاق از $\alpha A \beta$ است و می‌نویسیم:

$$\alpha A \beta \Rightarrow \alpha \gamma \beta$$

به عنوان مثال:

$$(E * E) * I \circ \Rightarrow (E * E) * a \circ \Rightarrow (I * E) * a \circ$$

تعریف ۳ (اشتقاق مکرر، اشتقاق چند مرحله‌ای) فرض کنید $G = (V, T, P, S)$ یک گرامر مستقل از متن باشد. می‌گوییم β اشتقاق مکرری از α است (که $\alpha, \beta \in (V \cup T)^*$) و با $\alpha \xrightarrow{*} \beta$ نشان می‌دهیم، اگر α و β در تعریف استقرایی زیر بگنجانند:

پایه:

$$\alpha \xrightarrow{*} \alpha$$

استقرا:

$$[\alpha \xrightarrow{*} \gamma \wedge \gamma \Rightarrow \beta] \implies [\alpha \xrightarrow{*} \beta]$$

۵ عرف نمادگذاری

معمولا برای نمایش متغیرهای تولید از حرف‌های ابتدایی و بزرگ الفبا (A, B, C, \dots) و همچنین حرف S ، برای ترمینال‌ها از حروف ابتدایی و کوچک الفبا (a, b, c, \dots)، برای رشته‌های روی ترمینال‌ها از حروف انتهایی و کوچک الفبا (\dots, w, x, y, z)، برای متغیر یا پایانه از حروف انتهایی و بزرگ الفبا (\dots, X, Y, Z) و برای رشته‌های روی ترمینال‌ها و متغیرها از (α, β, \dots) استفاده می‌کنیم.