

زبان‌های مستقل از متن

CONTEXT-FREE LANGUAGES



۱-۵ گرامرها و زبان‌های مستقل از متن

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

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

باشد.

تعریف

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

زبان مستقل از متن زبان L مستقل از متن است، اگر و فقط اگر یک گرامر مستقل از متن G وجود داشته باشد که $L = L(G)$ باشد.

تعریف

- ◀ هر گرامر منظم، یک گرامر مستقل از متن است \Leftarrow هر زبان منظم، یک زبان مستقل از متن است.
- ◀ خانواده‌ی زبان‌های منظم، زیرمجموعه‌ی خانواده‌ی زبان‌های مستقل از متن است.
- ◀ مستقل از متن یعنی این که جایگزینی متغیرهای سمت چپ قواعد را می‌توان در هر زمانی که آن متغیر در یک شکل جمله‌ای ظاهر می‌شود انجام داد و این به بقیه‌ی فرم جمله‌ای وابستگی ندارد (نتیجه‌ی وجود تنها یک متغیر در سمت چپ هر قاعده).

گرامر خطی گرامر $G = (V, T, S, P)$ خطی است که در سمت راست تمامی قواعد تولید آن حداکثر یک ناپایانه وجود داشته باشد. یعنی تمامی قواعد P به صورت

$$A \rightarrow xBy \quad \text{یا} \quad A \rightarrow x \quad , x, y \in T^*, A, B \in V$$

باشد.

تعریف

مثال

گرامر $G = (S, a, b, S, P)$ با قواعد $S \rightarrow aSa, S \rightarrow bSb, S \rightarrow \lambda$ مستقل از متن است. یک نمونه اشتقاق از این گرامر به صورت زیر است.

$$S \Rightarrow aSa \Rightarrow aaSaa \Rightarrow aabSbaa \Rightarrow aabbbaa$$

می‌توان با استقرا ثابت کرد که گرامر فوق زبان زیر را تولید می‌کند:

$$L(G) = \{ww^R : w \in \{a, b\}^*\}$$

مشاهده می‌شود که مطابق تعریف، گرامر فوق یک گرامر خطی است.

مثال

زبان $L = \{a^n b^m : n \neq m\}$ مستقل از متن است. ملاحظه می‌کنیم که زبان فوق می‌تواند به صورت اجتماع دو زبان نوشته شود:

$$L = \{a^n b^m : n > m\} \cup \{a^n b^m : n < m\}$$

و به این ترتیب گرامر مستقل از متن تولید کننده‌ی آن می‌شود:

$$\begin{aligned} S &\rightarrow S_1 \mid S_2 \\ S_1 &\rightarrow AS' & S_2 &\rightarrow S''b \\ S' &\rightarrow aS'b \mid \lambda & S'' &\rightarrow aS''b \mid \lambda \\ A &\rightarrow aA \mid a & B &\rightarrow bB \mid b \end{aligned}$$

مشاهده می‌شود که مطابق تعریف، گرامر فوق یک گرامر خطی نیست.

مثال

گرامر با قواعد $S \rightarrow aSb \mid SS \mid \lambda$ مستقل از متن است، ولی خطی نیست. زبان این گرامر عبارت است از:

$$L = \{w \in \{a, b\}^* : n_a(w) = n_b(w) \wedge n_a(v) \geq n_b(v), w \text{ هر زیررشته‌ای از } v\}$$

که با در نظر گرفتن a به عنوان پرانتز باز و b به عنوان پرانتز بسته ساختار پرانتزهای تودرتوی صحیح را مشخص می‌کند.

۱-۱-۵ اشتقاق در گرامرهای مستقل از متن

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

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

◀ **نکته** در صورت وجود اشتقاق، هم اشتقاق راست و هم اشتقاق چپ وجود خواهد داشت.

مثال

گرامر با قواعد $S \rightarrow aAB$, $A \rightarrow bBb$, $B \rightarrow A|\lambda$ را در نظر بگیرید. برای رشته‌ی $abbbb$:
یک اشتقاق چپ‌ترین:

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

یک اشتقاق راست‌ترین:

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

تعریف

درخت اشتقاق اگر $G = (V, T, S, P)$ 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 $T \cup \{\lambda\}$ باشد.
- ۳) هر یک از گره‌های داخلی دارای برچسبی از V باشد.
- ۴) اگر گره‌ای دارای برچسب $A \in V$ باشد و فرزندان آن از چپ به راست به صورت

$$a_1, a_2, \dots, a_n$$

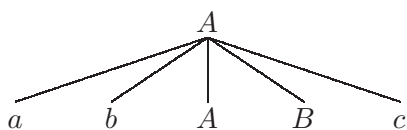
برچسب‌گذاری شده باشد، آن‌گاه P حاوی قاعده‌ای به شکل $A \rightarrow a_1 a_2 \dots a_n$ باشد.

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

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

مثال

درخت اشتقاق برای قاعده‌ی $A \rightarrow abABc$:



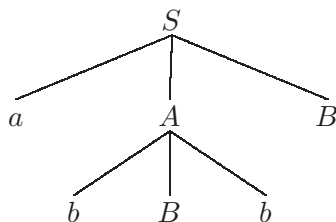
تعریف

درخت اشتقاق جزئی درخت اشتقاق جزئی همانند درخت اشتقاق است با این تفاوت که به جای شرط (۲)، شرط «هر برگ دارای برچسبی از $V \cup T \cup \{\lambda\}$ است» را دارد.

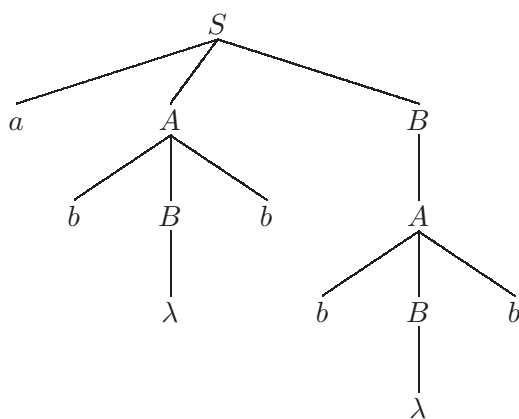
حاصل درخت اشتقاق جزئی از پیمایش عمق - اول برگ‌ها به دست می‌آید و یک فرم جمله‌ای از گرامر است.

مثال

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



یک درخت اشتقاق برای جمله‌ی $abbbb$:



قضیه

اگر $G = (V, T, S, P)$ یک گرامر مستقل 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 جزئی برای G با ریشه‌ی S باشد، آن‌گاه حاصل t_G یک شکل جمله‌ای از G است.

◀ اثبات.

برای اثبات اینکه برای هر فرم جمله‌ای از $L(G)$ یک درخت اشتقاق جزئی وجود دارد، از استقرا بر روی تعداد گام‌های اشتقاق استفاده می‌کنیم:

- پایه: این گزاره برای هر فرم جمله‌ای که در یک گام مشتق می‌شود، صحیح است (زیرا اگر $S \rightarrow u \in P$ ، آن‌گاه $S \Rightarrow^1 u$).
- فرض: برای هر فرم جمله‌ای که در n گام به دست می‌آید یک درخت اشتقاق جزئی وجود دارد ($S \Rightarrow^n w$).
- حکم: هر رشته‌ی w که در $n + 1$ قدم مشتق می‌شود باید به گونه‌ای باشد که در n گام داشته باشیم:

$$S \Rightarrow^* xAy, \quad x, y \in (V \cup T)^*, \quad A \in V$$

$$xAy \Rightarrow xa_1a_2 \dots a_my = w, \quad a_i \in V \cup T$$

بر اساس فرض استقرا یک درخت اشتقاق جزئی با حاصل xAy وجود دارد و چون گرامر باید قاعده‌ای مانند $A \rightarrow a_1a_2 \dots a_m$ داشته باشد، با توسعه‌ی برگ A یک درخت اشتقاق جزئی به دست می‌آوریم که حاصل آن $xa_1a_2 \dots a_my$ است.

بنابراین با استفاده از استقرا ادعا می‌کنیم که نتیجه برای همه‌ی فرم‌های جمله‌ای صحیح است. به همین ترتیب می‌توانیم نشان دهیم که هر درخت اشتقاق یک فرم جمله‌ای را نشان می‌دهد.

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 تأثیری ندارد.

۲-۵ تجزیه و ابهام

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

تعریف

۱-۲-۵ جستجوی جامع به عنوان یک الگوریتم عضویت

برای تشخیص $w \in L(G)$ می‌توانیم

- ابتدا اشتقاق‌های یک مرحله‌ای یعنی $S \Rightarrow x$ را بررسی کنیم،
- سپس اشتقاق‌های دو مرحله‌ای،
- ...
- اشتقاق‌های $n = |w|$ مرحله‌ای و ...

این روش یک اشکال دارد و آن این است که اگر $w \notin L(G)$ این روال خاتمه نمی‌یابد. اگر گرامر مورد نظر قاعده‌ی λ و قاعده‌ی یک $A \rightarrow B$ را نداشته باشد، این روال پس از $|w|$ مرحله متوقف می‌شود.

اگر $G = (V, T, S, P)$ یک گرامر مستقل از متن باشد که قواعدی به شکل $A \rightarrow B$ یا $A \rightarrow \lambda$ را نداشته باشد، آن‌گاه روش جستجوی کامل می‌تواند به صورت الگوریتمی در آید که برای هر $w \in \Sigma^*$ یا یک تجزیه برای w تعیین می‌کند و یا به ما می‌گوید که تجزیه ممکن نیست.

قضیه

◀ اثبات.

برای هر فرم جمله‌ای، طول و تعداد پایانه‌های آن را در نظر می‌گیریم. هرگام اشتقاق حداقل یکی از این دو مورد را افزایش می‌دهد (به دلیل وجود قواعد یک‌ه و تهی). از آنجا که نه طول فرم جمله‌ای و نه تعداد پایانه‌ها نمی‌تواند بیش از $|w|$ باشد، یک اشتقاق نمی‌تواند بیش از $۲|w|$ مرحله داشته باشد و تا آنجا یا تجزیه به طور موفق به پایان رسیده است و یا w نمی‌تواند توسط آن گرامر تولید شود.

$$\text{تعداد دورها} \leq |w| + |w| = ۲|w|$$

کاربرد عملی جستجوی جامع برای تجزیه بسیار محدود است، زیرا تعداد فرم‌های جمله‌ای بسیار زیاد است. برای محاسبه‌ی حد بالای تعداد فرم‌های جمله‌ای با فرض استفاده از اشتقاق چپ‌ترین، در دور اول حداکثر $|P|$ فرم جمله‌ای، در دور دوم حداکثر $|P|^2$ فرم جمله‌ای،
 \dots
 در حداکثر دور $|w|$ حداکثر $|P|^{|w|}$ فرم جمله وجود دارد.
 در مجموع حداکثر تعداد کل فرم‌های جمله‌ای می‌شود:

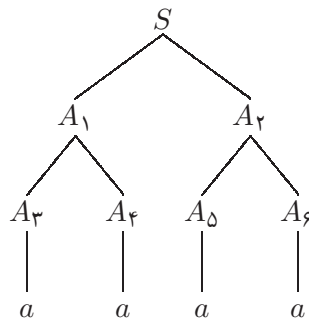
$$M = |P| + |P|^2 + \dots + |P|^{|w|} = |P| \frac{1 - |P|^{|w|}}{1 - |P|}$$

مثال

مثالی از بدترین حالت تعداد فرم‌های جمله‌ای را می‌توان در گرامر زیر مشاهده کرد:

$$\begin{aligned} S &\rightarrow A_1 A_2 \\ A_1 &\rightarrow A_3 A_4 \\ A_2 &\rightarrow A_5 A_6 \\ A_3 &\rightarrow a \\ A_4 &\rightarrow a \\ A_5 &\rightarrow a \\ A_6 &\rightarrow a \end{aligned}$$

که درخت اشتقاق آن برای رشته‌ی $aaaa$ در شکل زیر رسم شده است:



قضیه

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

اثبات این قضیه بر پایه‌ی روش CYK برای تجزیه است که در فصل بعدی تشریح می‌شود.

گرامر ساده گرامر مستقل از متن $G = (V, T, S, P)$ یک گرامر ساده (simple grammar, *s-grammar*) نام دارد اگر و فقط اگر تمامی قواعد آن به فرم

$$A \rightarrow ax \quad , A \in V, a \in T, x \in V^*$$

باشد و هر زوج (A, a) حداکثر یک مرتبه در P ظاهر شود.

تعریف

مثال

گرامر $S \rightarrow aSbSSc \mid b$ یک گرامر ساده است ولی $C \mid aSS \mid bSS \mid aS$ این گونه نیست.

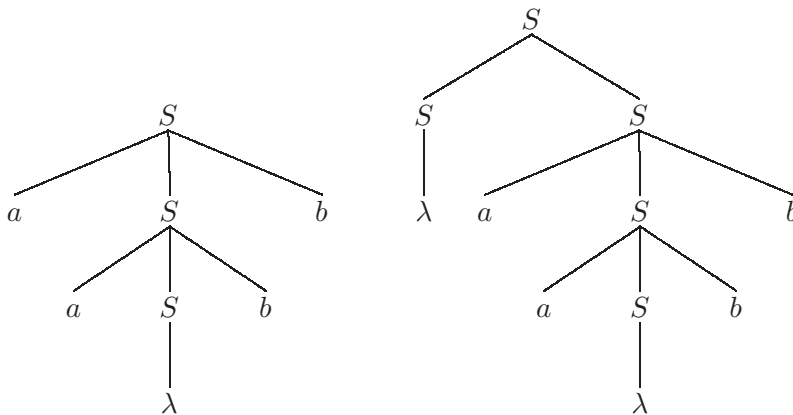
۲-۲-۵ ابهام در گرامر و زبان

گرامر مبهم 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. به عبارت دیگر، ابهام به معنی وجود دو یا چند اشتقاق چپ‌ترین یا راست‌ترین برای G است.

تعریف

مثال

گرامر $S \rightarrow aSb \mid SS \mid \lambda$ مبهم است: جمله‌ی $aabb$ دو درخت اشتقاق متفاوت دارد:



بنابراین حداقل دو اشتقاق چپ‌ترین (راست‌ترین) برای این رشته وجود خواهد داشت.

تعریف

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

مثال

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

تعریف

زبان غیرمبهم زبان L را غیرمبهم می‌گوییم اگر و فقط اگر حداقل یک گرامر غیرمبهم برای آن موجود باشد.

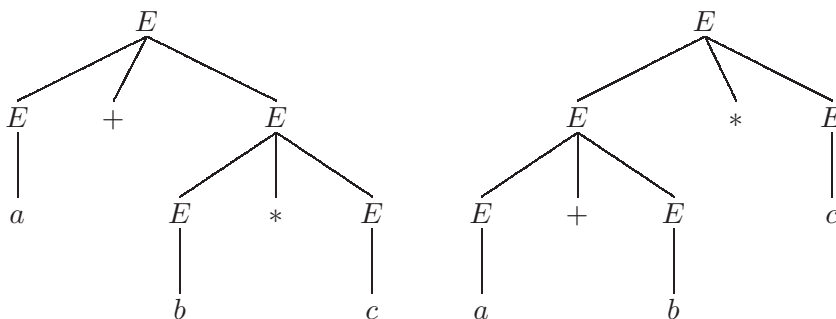
- ◀ زبان مبهم نداریم (زبان ذاتاً مبهم داریم).
- ◀ مسأله‌ی تشخیص ابهام در یک گرامر مستقل از متن، تصمیم‌ناپذیر است، یعنی هیچ الگوریتمی برای تشخیص یا رفع ابهام وجود ندارد.
- ◀ هیچ زبان منظمی ذاتاً مبهم نیست.

مثال

گرامر $G = (E, a, b, c, +, *, (,), E, P)$ با قواعد

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

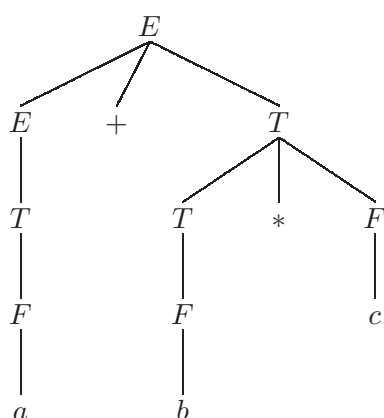
برای عبارت‌های حسابی، مبهم است: برای رشته‌ی $a + b * c$ دو درخت اشتقاق وجود دارد:



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

$$E \rightarrow E + T | T, \quad T \rightarrow T * F | F, \quad F \rightarrow (E) | a | b | c$$

در این گرامر رشته‌ی یاد شده تنها یک درخت اشتقاق دارد:



تشخیص ابهام یک گرامر مستقل از متن به عنوان یک مسأله‌ی تصمیم‌ناپذیر

گراف گرامر اگر G یک گرامر مستقل از متن باشد، گراف جهت دار $g_G = (N, E)$ که در آن N مجموعه‌ی راس‌ها، مجموعه‌ی فرم‌های جمله‌ای گرامر G و E مجموعه‌ی یال‌های برجسب‌دار گراف به صورت

تعریف

$$E = \{(x, r, y) : x \xrightarrow[r]{G} y, \quad x, y \in N\}$$

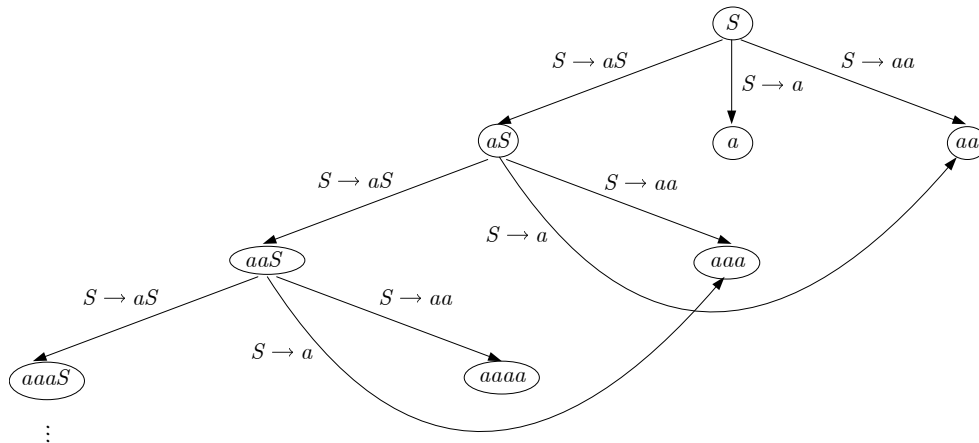
را گراف (چپ) گرامر G می‌گوییم.

- ◀ گراف گرامر با درخت تجزیه متفاوت است: درخت تجزیه چگونگی ایجاد یک رشته توسط گرامر را مشخص می‌کند،
- ◀ گراف گرامر کلیه‌ی اشتقاق‌های چپ (یا راست) ممکن درون یک گرامر را مشخص می‌کند.
- ◀ در درخت تجزیه هر گره معرف یک نماد گرامر است، در صورتی که در گراف گرامر هر گره معرف یک فرم جمله‌ای است.
- ◀ تجزیه‌گر (parser) برای یافتن چگونگی اشتقاق، فرم‌های جمله‌ای در گراف گرامر را جستجو می‌کند. بسته به نحوه‌ی جستجو، روش‌های مختلفی برای تجزیه خواهیم داشت:
 - بالا به پایین
 - پایین به بالا
 - عمق - اول
 - عرض - اول
- ◀ گرامر مبهم، در گراف خود حداقل دارای یک حلقه است.
- ◀ اگر گرامری مبهم نباشد، گراف آن یک درخت است.

گرامر G مبهم است اگر و فقط اگر گراف آن حداقل حاوی یک حلقه باشد.

مثال

گراف گرامر مبهم $S \rightarrow aS|a|aa$ حاوی حلقه است.



۳-۵ زبان‌های برنامه‌سازی و گرامرهای مستقل از متن

- ◀ تعریف زبان‌های برنامه‌سازی با گرامر متداول است.
- ◀ همه‌ی خواص زبان‌های برنامه‌سازی را نمی‌توان با گرامرهای مستقل از متن ساده توصیف کرد (\Leftrightarrow) نیاز به گرامرهای پیچیده‌تر مانند LL یا LR)
- ◀ همه‌ی خواص زبان‌های برنامه‌سازی را نمی‌توان با گرامرهای مستقل از متن توصیف کرد. برخی از این خواص با گرامرهای سطح بالاتر قابل توصیف هستند (\Leftrightarrow تحلیل معنایی)
- ◀ گرامر زبان‌های برنامه‌سازی باید غیر مبهم باشد.

۱-۳-۵ نمادگذاری BNF برای گرامرهای مستقل از متن

نمادگذاری BNF (Backus-Naur Form) از نمادهای زیر استفاده می‌کند:

- $\langle \text{variable} \rangle$ یک متغیر (ناپایانه) گرامر
- $::=$ به جای نماد \rightarrow
- پایانه‌ها بدون علامت اضافی به کار می‌روند.
- اسامی ناپایانه‌ها به شکل معنی‌دارتری انتخاب می‌شود.

مثال

نمادگذاری BNF برای چند قاعده‌ی گرامری زبان:

```
<if-statement> ::= if <expression> <then-clause> <else-clause>  
<expression>  ::= <term> | <expression> + <term>  
<term>        ::= <factor> | <term> * <factor>
```