

Escritura de gramáticas libres de contexto

Plan de la clase

1 Escritura de gramáticas

- Repaso
- Nociones de asociatividad y precedencia
- Otros ejemplos

2 Reducción

- Definiciones
- Algoritmos
 - Símbolos activos
 - Símbolos alcanzables
 - Símbolos anulables
- Reglas de renombramiento
- Recursividad a izquierda
- Factorización
- Orden

- “Practicamente todos los lenguajes que se usan en computación son independientes de contexto.”
- Como vimos, reconocer si una cadena pertenece a un lenguaje libre de contexto de manera determinista puede resultar no eficiente. (Los autómatas de pila no deterministas no siempre pueden pasarse a deterministas.)
- Vamos a ver distintas técnicas de parsing. Es decir, dada una gramática, tener un algoritmo que nos diga si una determinada cadena pertenece al lenguaje de esa gramática, y obtener un árbol de derivación de la misma, de manera que además comprendamos su estructura sintáctica.
- Y despues vamos a ver como agregar acciones a los parsers de manera de aprovechar que tenemos la estructura para procesar la cadena de entrada en forma sistemática.

Repaso

Ambigüedad: por ejemplo, la gramática

$$S \rightarrow aS | Sa | \lambda$$

es ambigua (ver dos árboles de derivación para a) pero es equivalente a la gramática no ambigua:

$$S \rightarrow aS | \lambda$$

Recordar: árboles de derivación sólo se aplican a gramáticas libres de contexto (y no a las sensibles al contexto o sin restricciones).

Asociatividad y precedencia

Supongamos que queremos definir un lenguaje de la lógica de primer orden (sólo predicados, sin funciones ni constantes). Una primera idea podría ser:

$$\begin{aligned} E &\rightarrow A \mid E \wedge E \mid E \vee E \mid E \Rightarrow E \mid \neg E \mid (E) \mid \forall \text{id}_V E \mid \exists \text{id}_V E \\ A &\rightarrow \text{id}_P (L) \\ L &\rightarrow \text{id}_V \mid L , L \end{aligned}$$

Donde id_V representa algún símbolo de variable (como x_1, x_2, \dots , y id_P es algún símbolo de relación (predicado), como $=, \leq, R$, etc.)

Asociatividad y precedencia

Una gramática no sólo nos sirve para definir un lenguaje, sino también para entender su estructura. En el proceso de reconocimiento de una cadena, obtenemos el *árbol de derivación*, que si la gramática es apropiada nos puede ayudar a entender el significado de la expresión.

Asociatividad y precedencia

La gramática definida es bastante ambigua...

Ejemplo:

$$A \vee A \wedge A$$

Hay dos árboles de derivación que producen esta cadena. Uno puede interpretarse como $(A \vee A) \wedge A$, y el otro como $A \vee (A \wedge A)$. Queremos una gramática equivalente en la cual únicamente sea posible esta última derivación-interpretación. Es decir, queremos representar la precedencia mayor del \wedge por sobre el \vee de alguna manera en la gramática. Se puede ver una intuición de que uno quiere que el \vee aparezca primero más arriba en el árbol de derivación. Surge la idea de jerarquizar niveles en la gramática dependiendo de la precedencia.

Asociatividad y precedencia

Por otro lado, si tenemos:

$$A \wedge A \wedge A$$

elegimos según asociatividad (digamos que a izquierda), de modo que se interprete:

$$(A \wedge A) \wedge A$$

Idea de cómo hacer esto con la gramática.

Asociatividad y precedencia

También con:

$$A \Rightarrow A \Rightarrow A$$

En este caso elegimos asociatividad a derecha

$$A \Rightarrow (A \Rightarrow A)$$

(Otro ejemplo de asociatividad a derecha sería la exponenciación:
 2^{3^4} se interpreta como $2^{(3^4)}$)

Asociatividad y precedencia

Recordar que asociatividad a izquierda y a derecha no tiene vínculo con derivación más a la izquierda o más a la derecha.

Armamos entonces la tabla de precedencia (las primeras líneas tienen la menor precedencia y las últimas la mayor):

\Rightarrow	(Asoc. a derecha)
\vee	(Asoc. a izq.)
\wedge	(Asoc. a izq.)
unarios:	($\forall \exists \neg$)

Sí, la precedencia de los operadores unarios también debe ser especificada; tienen la mayor precedencia. Ejemplo:

$$\neg A \wedge A$$

Ambigüedad

¿Cómo eliminar la ambigüedad y tener el árbol que queremos? No siempre se puede: hay lenguajes inherentemente ambiguos, y en general no se puede determinar si una gramática es ambigua o no. En casos como este, si empezamos por los operadores de menor precedencia podemos pensar que nos subdividen la entrada en subárboles.

Ejemplo:

$$A \vee A \Rightarrow A \wedge A \Rightarrow A \wedge A \vee A$$

Ambigüedad

Ahora, separamos de menor a mayor precedencia, poniendo símbolos no terminales y reglas de producción para cada nivel de precedencia. Nos queda la siguiente gramática:

$$\begin{aligned} E &\rightarrow D \Rightarrow E \mid D \\ D &\rightarrow D \vee C \mid C \\ C &\rightarrow C \wedge U \mid U \\ U &\rightarrow \forall \text{id}_V U \mid \exists \text{id}_V U \mid \neg U \mid A \\ A &\rightarrow \text{id}_P (L) \mid (E) \\ L &\rightarrow \text{id}_V \mid L , L \end{aligned}$$

Ambigüedad

Todavía podemos mejorar esto: en realidad, no hace falta separar U y A (de alguna forma, son todas operaciones unarias):

$$U \rightarrow \forall \text{id}_V U \mid \exists \text{id}_V U \mid \neg U \mid \text{id}_P (L) \mid (E)$$

(Esto motiva la escritura de gramáticas, que veremos a continuación.)

Además, la parte de la lista de variables sigue siendo ambigua, así que hacemos el siguiente cambio:

$$L \rightarrow \text{id}_V \mid L , \text{id}_V$$

Otros ejemplos

Con aritmética:

$$E \rightarrow E + T \mid T$$

$$T \rightarrow T * F \mid F$$

$$F \rightarrow \text{id} \mid (E)$$

es una gramática no ambigua

Tampoco hay problema en que haya dos operadores con la misma precedencia:

$$E \rightarrow E + T \mid E - T \mid T$$

$$T \rightarrow T * F \mid T / F \mid F$$

$$F \rightarrow \text{id} \mid (E)$$

Sigue siendo una gramática no ambigua. Observar que la resta es asociativa a izquierda: $3 - 2 - 1$ interpretamos como $(3 - 2) - 1$ y no como $3 - (2 - 1)$.

Definiciones

$X \in V_N$ es *útil* si participa en alguna derivación que genere una cadena de terminales. Es decir, si hay una derivación:

$$S \Rightarrow^* \alpha X \beta \Rightarrow^* \alpha \gamma \beta$$

con $\alpha, \beta, \gamma \in V_T^*$.

Notar que X es útil si suceden dos cosas:

- $X \in V_N$ es *alcanzable* si aparece en alguna forma sentencial, es decir, si hay $\alpha, \beta \in V^*$ tales que:

$$S \Rightarrow^* \alpha X \beta$$

- $X \in V_N$ es *activo* si puede generar una cadena de terminales, es decir, si existe alguna derivación:

$$X \Rightarrow^* \gamma$$

con $\gamma \in V_T^*$.

Definiciones

Como uno esperaría, se define que X es *inalcanzable* si no es alcanzable, X es *inactivo* si no es activo, y es *inútil* si no es útil.

Ejemplo

Ejemplos de gramáticas con símbolos...

- alcanzables pero inactivos
- activos pero inalcanzables
- inalcanzables inactivos
- alcanzables y activos (útiles)

Una gramática puede reducirse a otra equivalente sin símbolos inútiles.

Algoritmos

Es fácil determinar si hay símbolos inútiles, calculando los conjuntos de símbolos activos y de símbolos alcanzables, y viendo que ambos conjuntos sean iguales a V_N .

Los algoritmos que damos calculan los símbolos que tienen cierta propiedad, para luego eliminar las producciones que los incluyen.

Los algoritmos que damos terminan cuando en una iteración el conjunto que obtenemos es el mismo que en el paso anterior. No son eficientes, sólo muestran que los conjuntos se pueden calcular.

Símbolos activos

$$N = \emptyset$$

repetir hasta que N no cambie

$$N = N \cup \{A \mid A \rightarrow \alpha \in P \wedge \alpha \in (N \cup V_T)^*\}$$

Gramática ejemplo:

$$S \rightarrow a \mid aA \mid aC$$

$$A \rightarrow AB$$

$$B \rightarrow b$$

$$C \rightarrow cC$$

Símbolos alcanzables

$$N = S$$

repetir hasta que N no cambie

$$N = N \cup \{X \mid A \rightarrow \alpha X \beta \in P \wedge A \in N\}$$

Gramática ejemplo:

$$S \rightarrow a \mid aA \mid aC$$

$$A \rightarrow AB$$

$$B \rightarrow b$$

$$C \rightarrow cC$$

Símbolos inútiles

Para eliminar inútiles:

- Eliminar inactivos.
- Eliminar inalcanzables.

Hay que hacerlo en este orden. ¿Por qué?

Porque tal vez un símbolo era alcanzable a través de uno inactivo.

En cambio, si un símbolo para ser activo usa otros inalcanzables es porque él mismo era inalcanzable...

Símbolos inútiles

Para eliminar inútiles:

- Eliminar inactivos.
- Eliminar inalcanzables.

Ejemplo:

$$\begin{aligned} S &\rightarrow b \mid I \\ I &\rightarrow A B \\ A &\rightarrow b \end{aligned}$$

Si eliminamos los inalcanzables primero, no se elimina ninguno.
Entonces, al eliminar los inactivos (I, B), quedaría

$$\begin{aligned} S &\rightarrow b \\ A &\rightarrow b \end{aligned}$$

Símbolos anulables

También es práctico saber cuáles son los símbolos *anulables* de una gramática, es decir, los símbolos X tales que $X \Rightarrow^* \lambda$. Los podemos calcular de forma similar a los anteriores:

$$N = \{A \mid A \rightarrow \lambda \in P\}$$

repetir

hasta que N no cambie

$$N = N \cup \{A \mid A \rightarrow \alpha \in P \wedge \alpha \in N^*\}$$

Para eliminar las producciones λ : reemplazar cada regla de la forma $B \rightarrow \alpha_1 A_1 \alpha_2 A_2 \alpha_3 \dots \alpha_k A_k \alpha_{k+1}$, con $k \geq 1$ y A_1, A_2, \dots, A_k anulables, por todas las reglas de la forma

$B \rightarrow \alpha_1 \beta_1 \alpha_2 \beta_2 \alpha_3 \dots \alpha_k \beta_k \alpha_{k+1}$ tales que $\beta_j = A_j$ o λ para todo $1 \leq j \leq k$, y eliminar cada regla $A_j \rightarrow \lambda$, excepto $S \rightarrow \lambda$ si aparece, en ese caso S no debe aparecer a derecha.

Símbolos anulables

Gramática ejemplo:

$$\begin{aligned} S &\rightarrow d \mid aA \\ A &\rightarrow bA \mid BBC \\ B &\rightarrow bB \mid \lambda \\ C &\rightarrow cC \mid \lambda \end{aligned}$$

$$N = \{A, B, C\}$$

Resulta:

$$\begin{aligned} S &\rightarrow d \mid aA \mid a \\ A &\rightarrow bA \mid b \mid BBC \mid B \mid BB \mid C \mid BC \\ B &\rightarrow bB \mid b \\ C &\rightarrow cC \mid c \end{aligned}$$

Reglas de renombramiento

Una regla de renombramiento es una producción de la forma $A \rightarrow B$, con $A, B \in V_N$.

- 1 Construir para cada A el conjunto $N_A = \{B \in V_N \mid A \Rightarrow^* B\}$

$$N_A = A$$

repetir hasta que N_A no cambie

$$N_A = N_A \cup \{C \mid B \rightarrow C \in P \wedge B \in N_A\}$$

- 2 Para cada $B \rightarrow \alpha$, agregar $A \rightarrow \alpha$ si $B \in N_A$.
- 3 Eliminar todas las producciones de la forma $A \rightarrow B$

Renombramientos

Gramática ejemplo:

$$\begin{aligned} S &\rightarrow A \\ A &\rightarrow a \mid B \\ B &\rightarrow A \end{aligned}$$

$$N_A = \{B \in V_N \mid A \Rightarrow^* B\}$$

Para cada $B \rightarrow \alpha$, agregar $A \rightarrow \alpha$ si $B \in N_A$.

Quedará

$$\begin{aligned} S &\rightarrow a \\ A &\rightarrow a \\ B &\rightarrow a \end{aligned}$$

Al eliminar inútiles, quedará

$$S \rightarrow a$$

Ciclos

Cómo eliminar ciclos, es decir la posibilidad de $A \Rightarrow^+ A$.
De no haber producciones λ , $A \Rightarrow^+ A$ sólo es posible si hay al menos una regla de renombramiento que involucre a A .
Por eso conviene previamente eliminar las producciones λ .

Recursividad a izquierda

Una gramática es recursiva a izquierda si existe un no terminal A tal que $A \Rightarrow^+ A\alpha$ para alguna cadena α . ¿Cómo eliminar esta condición?

Cambiar cada conjunto de producciones de la forma $A \rightarrow A\alpha \mid \beta$ por

$$\begin{aligned} A &\rightarrow \beta A' \\ A' &\rightarrow \alpha A' \mid \lambda \end{aligned}$$

agregando un nuevo $A' \notin V_N$.

Recursividad a izquierda

Gramática ejemplo:

$$\begin{aligned}E &\rightarrow E + T \mid T \\T &\rightarrow T * F \mid F \\F &\rightarrow (E) \mid \text{id}\end{aligned}$$

Quedará

$$\begin{aligned}E &\rightarrow TE' \\E' &\rightarrow +TE' \mid \lambda \\T &\rightarrow FT' \\T' &\rightarrow *FT' \mid \lambda \\F &\rightarrow (E) \mid \text{id}\end{aligned}$$

Después se puede hacer necesario eliminar reglas λ .

Factorización a izquierda

Cambiar cada conjunto de producciones de la forma

$A \rightarrow \alpha\beta_1 \mid \alpha\beta_2 \mid \dots \mid \alpha\beta_n \mid \gamma$, donde $\alpha \notin \text{Ini}(\gamma)$, por

$$\begin{aligned} A &\rightarrow \alpha A' \mid \gamma \\ A' &\rightarrow \beta_1 \mid \beta_2 \mid \dots \mid \beta_n \end{aligned}$$

agregando un nuevo $A' \notin V_N$.

Factorización a izquierda

Gramática ejemplo:

$$\begin{aligned}P &\rightarrow iEtP \mid iEtPeP \mid a \\E &\rightarrow b\end{aligned}$$

¿Qué representa en lenguajes de programación?

El **if ... then ... else ...**

Quedará

$$\begin{aligned}P &\rightarrow iEtPP' \mid a \\P' &\rightarrow eP \mid \lambda \\E &\rightarrow b\end{aligned}$$

Gramática reducida y propia

Una gramática es *reducida* si no tiene símbolos inútiles.

(Algunos piden también que símbolo inicial no aparezca en el lado derecho de ninguna producción. Para que ocurra esto, podría ser necesario introducir un renombramiento.)

Una gramática es *propia* si no tiene producciones λ , ni ciclos, ni símbolos inútiles.

Orden en la sala

¿En qué orden aplicar los algoritmos? Sin ser el único, un orden posible es:

- ① Eliminar recursividad a izquierda
- ② Eliminar producciones λ
- ③ Eliminar renombramientos
- ④ Eliminar inútiles
 - ① Eliminar inactivos
 - ② Eliminar inalcanzables