

Văn phạm phi ngữ cảnh

Trần Vĩnh Đức



Trường Đại Học Bách Khoa Hà Nội

Ngày 12 tháng 3 năm 2019

Thuật ngữ

- ▶ Context-Free Language (CFL) : Ngôn ngữ phi ngữ cảnh
- ▶ Context-Free Grammar (CFG): Văn phạm phi ngữ cảnh

Nội dung

Giới thiệu

Định nghĩa

Ví dụ

Thiết kế văn phạm

Tính nhập nhằng

Dạng chuẩn Chomsky

Thuật toán Cocke-Younger-Kasami (CYK)

Nội dung

Giới thiệu

Định nghĩa

Ví dụ

Thiết kế văn phạm

Tính nhập nhằng

Dạng chuẩn Chomsky

Thuật toán Cocke-Younger-Kasami (CYK)

Giới thiệu

- ▶ Ta đã xem xét hai phương pháp khác nhau, nhưng tương đương, để mô tả ngôn ngữ: Các otomat hữu hạn và các biểu thức chính quy.
- ▶ Ta đã biết rằng có những ngôn ngữ đơn giản nhưng không thể mô tả theo cách trên. Ví dụ ngôn ngữ

$$B = \{0^n 1^n \mid n \geq 0\}.$$

- ▶ Trong chương này ta xem xét một phương pháp mạnh hơn để mô tả ngôn ngữ: Văn phạm phi ngữ cảnh (viết tắt là CFG).

Văn phạm phi ngữ cảnh

- ▶ Được sử dụng đầu tiên trong nghiên cứu ngôn ngữ tự nhiên. Nó giúp hiểu quan hệ giữa các khái niệm *câu*, *động từ*, và *mệnh đề*.
- ▶ Có thể mô tả một số cấu trúc đệ quy. Điều này rất có ích trong nhiều ứng dụng, đặc biệt trong đặc tả và biên dịch các ngôn ngữ lập trình.
- ▶ Lớp ngôn ngữ gắn với CFG gọi là ngôn ngữ phi ngữ cảnh. Chúng chứa thực sự lớp ngôn ngữ chính quy.

Ví dụ

Dưới đây là một CFG

$$A \rightarrow 0A1$$

$$A \rightarrow B$$

$$B \rightarrow \#$$

- ▶ bao gồm tập *quy tắc thay thế*, hay còn gọi là *sản xuất*.
- ▶ A, B gọi là *biến*
- ▶ $0, 1, \#$ gọi là ký hiệu *kết thúc*
- ▶ Có một biến *bắt đầu*, ví dụ A .

Mô tả ngôn ngữ bằng CFG

CFG mô tả ngôn ngữ bằng cách sinh ra mỗi xâu theo cách sau đây:

1. Viết ra biến bắt đầu. Thường là biến bên trái của quy tắc đầu tiên.
2. Tìm một biến trong xâu vừa viết và chọn một quy tắc có vế trái là biến đó. Thay thế biến đó bởi vế phải của quy tắc đã chọn.
3. Lặp lại bước 2. cho đến khi không còn biến nào trong xâu đã viết.

Dãy phép thay thế để đạt được xâu gọi là *dẫn xuất*.

Ví dụ

Xét CFG G_1 dưới đây

$$A \rightarrow 0A1$$

$$A \rightarrow B$$

$$B \rightarrow \#$$

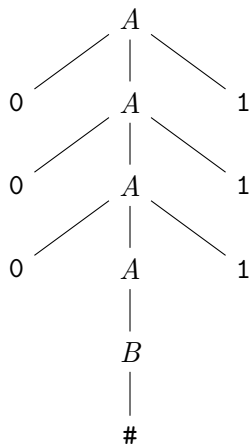
- ▶ Một dẫn xuất của xâu $000\#111$ trong G_1 là

$$A \Rightarrow 0A1 \Rightarrow 00A11 \Rightarrow 000A111 \Rightarrow 000B111 \Rightarrow 000\#111$$

- ▶ Mọi xâu sinh ra theo cách này gọi là *ngôn ngữ của văn phạm* G_1 , ký hiệu $L(G_1)$.
- ▶ Ta có thể kiểm tra

$$L(G_1) = \{0^n\#1^n \mid n \geq 0\}$$

Cây dẫn xuất



Hình: Dẫn xuất

$A \Rightarrow 0A1 \Rightarrow 00A11 \Rightarrow 000A111 \Rightarrow 000B111 \Rightarrow 000\#111$

Mô tả ngôn ngữ Tiếng Việt

- ▶ Một câu Tiếng Việt viết hoặc nói ra chỉ gồm các ký tự thuộc bộ chữ cái tiếng Việt, nhưng quá trình viết hoặc nói được hướng dẫn bởi các quy tắc ngữ pháp
- ▶ Hướng dẫn về cấu trúc đơn giản của một câu thường có dạng

$$\langle \text{câu} \rangle \rightarrow \langle \text{chủ ngữ} \rangle \langle \text{vị ngữ} \rangle$$

nghĩa là một câu đúng có thể bắt đầu bằng một chủ ngữ, tiếp theo là vị ngữ.

- ▶ về chủ ngữ

$$\langle \text{chủ ngữ} \rangle \rightarrow \langle \text{đại từ} \rangle \mid \langle \text{danh từ} \rangle \mid \langle \text{câu} \rangle$$

nghĩa là chủ ngữ có thể là đại từ hoặc danh từ hoặc câu.

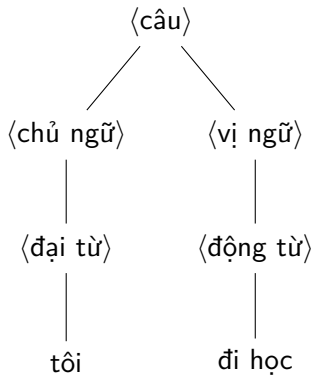
Văn phạm Tiếng Việt giản lược

Ví dụ

Văn phạm sau đây hướng dẫn cách viết đúng nhiều câu Tiếng Việt đơn giản

$\langle \text{câu} \rangle \rightarrow \langle \text{chủ ngữ} \rangle \langle \text{vị ngữ} \rangle$
 $\langle \text{chủ ngữ} \rangle \rightarrow \langle \text{đại từ} \rangle \mid \langle \text{danh từ} \rangle \mid \langle \text{câu} \rangle$
 $\langle \text{vị ngữ} \rangle \rightarrow \langle \text{động từ} \rangle \mid \langle \text{tính từ} \rangle$
 $\langle \text{đại từ} \rangle \rightarrow \text{tôi} \mid \text{anh} \mid \text{nó} \mid \dots$
 $\langle \text{động từ} \rangle \rightarrow \text{đi học} \mid \text{đi chơi} \mid \dots$

Hãy đưa ra một cây dẫn xuất cho câu “tôi đi học”.



Hình: Cây dẫn xuất cho câu “tôi đi học”.

Nội dung

Giới thiệu

Định nghĩa

Ví dụ

Thiết kế văn phạm

Tính nhập nhằng

Dạng chuẩn Chomsky

Thuật toán Cocke-Younger-Kasami (CYK)

Định nghĩa

Một *Văn phạm phi ngữ cảnh* là bộ bốn (V, Σ, R, S) trong đó

1. V là một tập hữu hạn, mỗi phần tử thuộc nó được gọi là *biến*,
2. Σ là một tập hữu hạn, phân biệt với V , mỗi phần tử thuộc nó được gọi là ký hiệu *kết thúc*,
3. R là một tập hữu hạn, mỗi phần tử của R gọi là một *quy tắc* có dạng $A \rightarrow u$ trong đó $A \in V$ và u là một xâu trên $V \cup \Sigma$, và
4. $S \in V$ là biến bắt đầu.

► Ta đã xét CFG $G_1 = (\{A, B\}, \{0, 1, \#\}, R, A)$ với tập quy tắc

$$R = \{A \rightarrow 0A1, A \rightarrow B, B \rightarrow \#\}.$$

Ngôn ngữ của văn phạm

- ▶ Xét u, v , và w là xâu trên $V \cup \Sigma$, và $A \rightarrow w$ là một quy tắc của văn phạm. Khi đó ta nói uAv *dẫn trực tiếp* $uwxv$ và viết $uAv \Rightarrow uwxv$.
- ▶ Ta nói rằng u *dẫn* ra v và viết

$$u \Rightarrow^* v$$

nếu $u = v$ hoặc nếu tồn tại một dãy u_1, u_2, \dots, u_k , với $k \geq 0$ và

$$u \Rightarrow u_1 \Rightarrow u_2 \Rightarrow \dots \Rightarrow u_k = v.$$

- ▶ *Ngôn ngữ của văn phạm* G là

$$L(G) = \{w \in \Sigma^* \mid S \Rightarrow^* w\}.$$

Nội dung

Giới thiệu

Định nghĩa

Ví dụ

Thiết kế văn phạm

Tính nhập nhằng

Dạng chuẩn Chomsky

Thuật toán Cocke-Younger-Kasami (CYK)

Ví dụ

Xét văn phạm $G_3 = (\{S\}, \{a, b\}, R, S)$. Tập quy tắc R là

$$S \rightarrow aSb \mid SS \mid \varepsilon.$$

- ▶ Văn phạm này sinh ra xâu như abab, aaabbb, và aababb.
- ▶ Xem a như dấu ngoặc trái "(" và b như dấu ngoặc phải ")"
- ▶ $L(G_3)$ là ngôn ngữ gồm mọi xâu "cân bằng ngoặc".

Ví dụ

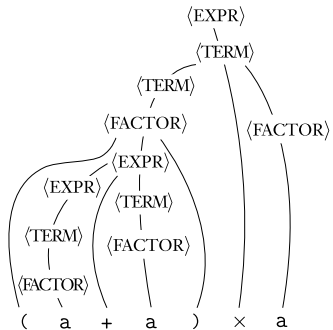
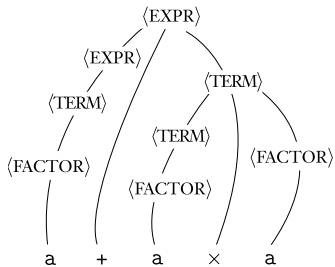
Xét văn phạm $G_4 = (V, \Sigma, R, \langle \text{EXPR} \rangle)$ trong đó

- ▶ $V = \{ \langle \text{EXPR} \rangle, \langle \text{TERM} \rangle, \langle \text{FACTOR} \rangle \}$
- ▶ $\Sigma = \{ a, +, \times, (,) \}$
- ▶ Các quy tắc là

$$\langle \text{EXPR} \rangle \rightarrow \langle \text{EXPR} \rangle + \langle \text{TERM} \rangle \mid \langle \text{TERM} \rangle$$

$$\langle \text{TERM} \rangle \rightarrow \langle \text{TERM} \rangle \times \langle \text{FACTOR} \rangle \mid \langle \text{FACTOR} \rangle$$

$$\langle \text{FACTOR} \rangle \rightarrow (\langle \text{EXPR} \rangle) \mid a$$



Nội dung

Giới thiệu

Định nghĩa

Ví dụ

Thiết kế văn phạm

Tính nhập nhằng

Dạng chuẩn Chomsky

Thuật toán Cocke-Younger-Kasami (CYK)

Hợp của nhiều văn phạm

Ví dụ

Xây dựng CFG cho ngôn ngữ $\{0^n 1^n \mid n \geq 0\} \cup \{1^n 0^n \mid n \geq 0\}$.

- Xây dựng CFG

$$S_1 \rightarrow 0S_11 \mid \varepsilon$$

cho ngôn ngữ $\{0^n 1^n \mid n \geq 0\}$ và CFG

$$S_2 \rightarrow 1S_20 \mid \varepsilon$$

cho ngôn ngữ $\{1^n 0^n \mid n \geq 0\}$.

- Kết hợp lại ta được CFG

$$S \rightarrow S_1 \mid S_2$$

$$S_1 \rightarrow 0S_11 \mid \varepsilon$$

$$S_2 \rightarrow 1S_20 \mid \varepsilon$$

Với ngôn ngữ chính quy

Ta xây dựng DFA đoán nhận ngôn ngữ đó. Sau đó xây dựng CFG tương đương như sau:

- ▶ Các biến R_i tương ứng với các trạng thái q_i của DFA
- ▶ Thêm quy tắc

$$R_i \rightarrow aR_j$$

nếu $\delta(q_i, a) = q_j$.

- ▶ Thêm quy tắc

$$R_i \rightarrow \varepsilon$$

nếu q_i là trạng thái chấp nhận của DFA.

- ▶ Đặt R_0 là trạng thái bắt đầu, với q_0 là trạng thái bắt đầu của DFA.

Nội dung

Giới thiệu

Định nghĩa

Ví dụ

Thiết kế văn phạm

Tính nhập nhằng

Dạng chuẩn Chomsky

Thuật toán Cocke-Younger-Kasami (CYK)

Sự nhập nhằng

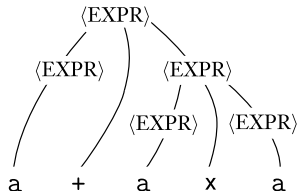
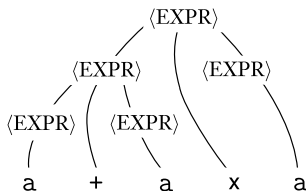
- ▶ Một CFG có thể sinh cùng một xâu theo nhiều cách khác nhau.
- ▶ Xâu này có nhiều cây dẫn xuất, nên nó có nhiều nghĩa.
- ▶ Điều này không mong muốn trong nhiều ứng dụng, như trong các ngôn ngữ lập trình: Các chương trình phải được diễn dịch một cách duy nhất.

Ví dụ

- ▶ Xét văn phạm G_5

$$\langle \text{EXPR} \rangle \rightarrow \langle \text{EXPR} \rangle + \langle \text{EXPR} \rangle \mid \langle \text{EXPR} \rangle \times \langle \text{EXPR} \rangle \mid (\langle \text{EXPR} \rangle) \mid a$$

- ▶ Văn phạm này sinh ra xâu $a + a \times a$ một cách nhập nhằng.



- ▶ Văn phạm không nhập nhằng G_4 có $L(G_4) = L(G_5)$.

Định nghĩa văn phạm nhập nhằng

- ▶ Hai dẫn xuất có thể khác nhau chỉ do thứ tự thay thế các biến trong xâu. Ta sẽ chỉ quan tâm đến dẫn xuất trái nhất.
- ▶ Một dẫn xuất là *trái nhất* nếu tại mọi bước ta luôn thay thế biến còn lại bên trái nhất.

Định nghĩa

- ▶ Một xâu w được dẫn xuất *một cách nhập nhằng* trong CFG G nếu nó có ít nhất hai dẫn xuất trái nhất.
- ▶ Văn phạm G là *nhập nhằng* nếu nó sinh một cách nhập nhằng một xâu nào đó.

Nhập nhằng cố hữu

- ▶ Đôi khi, khi cho một văn phạm nhập nhằng ta có thể tìm một văn phạm không nhập nhằng sinh ra cùng ngôn ngữ.
- ▶ Tuy vậy, có những CFL chỉ sinh bởi những văn phạm nhập nhằng. Những ngôn ngữ như vậy gọi là *nhập nhằng cố hữu*.
- ▶ Ví dụ, ngôn ngữ

$$\{a^i b^j c^k \mid i = j \text{ hoặc } j = k\}$$

là nhập nhằng cố hữu.

Nội dung

Giới thiệu

Định nghĩa

Ví dụ

Thiết kế văn phạm

Tính nhập nhằng

Dạng chuẩn Chomsky

Thuật toán Cocke-Younger-Kasami (CYK)

Dạng chuẩn Chomsky

Định nghĩa

Một CFG là ở *dạng chuẩn Chomsky* nếu mọi quy tắc của nó có dạng

$$A \rightarrow BC$$

$$A \rightarrow a$$

trong đó a là một ký hiệu kết thúc và A, B , và C là các biến—nhưng B và C không là biến bắt đầu. Thêm nữa, ta cho phép quy tắc

$$S \rightarrow \varepsilon$$

nếu S là biến bắt đầu.

Đưa một CFG về dạng chuẩn Chomsky

Định lý

Mỗi ngôn ngữ phi ngữ cảnh đều sinh bởi một văn phạm phi ngữ cảnh ở dạng chuẩn Chomsky.

Biến đổi CFG theo bốn bước sau:

1. Thêm biến bắt đầu mới để biến bắt đầu không xuất hiện ở vế phải của bất cứ quy tắc nào.
2. Loại bỏ các ϵ -quy tắc.
3. Loại bỏ các quy tắc đơn có dạng $A \rightarrow B$.
4. Chuyển các quy tắc còn lại về dạng đúng.

Bước 1: Xử lý biến bắt đầu

- ▶ Để biến bắt đầu S không xuất hiện vế phải của quy tắc nào.
- ▶ Ta thêm biến mới S_0 và quy tắc

$$S_0 \rightarrow S,$$

với S là biến khởi đầu.

Bước 2: Loại bỏ ε -quy tắc

Với mỗi quy tắc dạng $A \rightarrow \varepsilon$, ta thực hiện

- ▶ Loại bỏ quy tắc $A \rightarrow \varepsilon$
- ▶ Nếu có quy tắc dạng $R \rightarrow uAv$, thì ta thêm quy tắc $R \rightarrow uv$. Ta làm bước này với mọi xuất hiện của A . Có nghĩa rằng nếu có quy tắc dạng $R \rightarrow uAvAw$ thì ta phải thêm các quy tắc

$$R \rightarrow uAvw \mid uvAw \mid uvw$$

- ▶ Nếu có quy tắc $R \rightarrow A$, ta thêm quy tắc $R \rightarrow \varepsilon$ **trừ khi** ở các bước trước ta đã loại bỏ quy tắc $R \rightarrow \varepsilon$.

Bước 3: Loại bỏ quy tắc đơn

Với mỗi quy tắc đơn dạng $A \rightarrow B$, ta thực hiện

- ▶ Loại bỏ quy tắc $A \rightarrow B$.
- ▶ Với mỗi quy tắc dạng $B \rightarrow u$, ta thêm quy tắc

$$A \rightarrow u$$

trừ khi đây là quy tắc đơn đã bị loại bỏ ở các bước trước đó.

Bước 4: Chuyển các quy tắc còn lại về dạng đúng

- Thay thế các quy tắc dạng

$$A \rightarrow u_1 u_2 \dots u_k, \quad \text{với } k \geq 3 \text{ và } u_i \in V \cup \Sigma$$

bởi các quy tắc

$$A \rightarrow u_1 A_1, \quad A_1 \rightarrow u_2 A_2, \quad \dots, \quad A_{k-2} \rightarrow u_{k-1} u_k.$$

- Thay thế mọi ký hiệu kết thúc u_i trong các quy tắc vừa thay thế ở trên bởi biến mới U_i và thêm quy tắc $U_i \rightarrow u_i$.

Ví dụ

1. Đảm bảo biến bắt đầu của văn phạm bên trái dưới đây không xuất hiện ở vế phải của quy tắc nào.

$$S \rightarrow ASA \mid aB$$

$$A \rightarrow B \mid S$$

$$B \rightarrow b \mid \varepsilon$$

$$S_0 \rightarrow S$$

$$S \rightarrow ASA \mid aB$$

$$A \rightarrow B \mid S$$

$$B \rightarrow b \mid \varepsilon$$

2. Loại bỏ các ε -quy tắc: Loại bỏ $B \rightarrow \varepsilon$ sau đó loại bỏ $A \rightarrow \varepsilon$

$$S_0 \rightarrow S$$

$$S \rightarrow ASA \mid aB \mid a$$

$$A \rightarrow B \mid S \mid \varepsilon$$

$$B \rightarrow b$$

$$S_0 \rightarrow S$$

$$S \rightarrow ASA \mid aB \mid a \mid SA \mid AS \mid S$$

$$A \rightarrow B \mid S$$

$$B \rightarrow b$$

3a. Loại bỏ các quy tắc đơn: Loại bỏ $S \rightarrow S$ và sau đó là $S_0 \rightarrow S$

$$S_0 \rightarrow S$$

$$S \rightarrow ASA \mid aB \mid a \mid SA \mid AS$$

$$A \rightarrow B \mid S$$

$$B \rightarrow b$$

$$S_0 \rightarrow \textcolor{blue}{ASA} \mid \textcolor{blue}{aB} \mid \textcolor{blue}{a} \mid \textcolor{blue}{SA} \mid \textcolor{blue}{AS}$$

$$S \rightarrow ASA \mid aB \mid a \mid SA \mid AS$$

$$A \rightarrow B \mid S$$

$$B \rightarrow b$$

3b. Loại bỏ quy tắc $A \rightarrow B$ và sau đó là $A \rightarrow S$

$$S_0 \rightarrow ASA \mid aB \mid a \mid SA \mid AS$$

$$S \rightarrow ASA \mid aB \mid a \mid SA \mid AS$$

$$A \rightarrow S \mid \textcolor{blue}{b}$$

$$B \rightarrow b$$

$$S_0 \rightarrow ASA \mid aB \mid a \mid SA \mid AS$$

$$S \rightarrow ASA \mid aB \mid a \mid SA \mid AS$$

$$A \rightarrow b \mid \textcolor{blue}{ASA} \mid \textcolor{blue}{aB} \mid \textcolor{blue}{a} \mid \textcolor{blue}{SA} \mid \textcolor{blue}{AS}$$

$$B \rightarrow b$$

3. Sau khi đã loại bỏ hết quy tắc đơn, ta được

$$S_0 \rightarrow ASA \mid aB \mid a \mid SA \mid AS$$

$$S \rightarrow ASA \mid aB \mid a \mid SA \mid AS$$

$$A \rightarrow b \mid ASA \mid aB \mid a \mid SA \mid AS$$

$$B \rightarrow b$$

4. Đưa các quy tắc còn lại về dạng đúng

$$S_0 \rightarrow AA_1 \mid UB \mid a \mid SA \mid AS$$

$$S \rightarrow AA_1 \mid UB \mid a \mid SA \mid AS$$

$$A \rightarrow b \mid AA_1 \mid UB \mid a \mid SA \mid AS$$

$$B \rightarrow b$$

$$A_1 \rightarrow SA$$

$$U \rightarrow a$$

Nội dung

Giới thiệu

Định nghĩa

Ví dụ

Thiết kế văn phạm

Tính nhập nhằng

Dạng chuẩn Chomsky

Thuật toán Cocke-Younger-Kasami (CYK)

Bài toán

Cho CFG G và một xâu w . Kiểm tra xem G có sinh ra xâu w .

Ý tưởng thuật toán

- ▶ Xét G là văn phạm ở dạng chuẩn Chomsky:

$$A \rightarrow a \quad \text{hoặc} \quad A \rightarrow BC.$$

- ▶ Thuật toán dựa trên kỹ thuật **quy hoạch động**.
- ▶ Xét chuỗi vào

$$w = w_1 w_2 \cdots w_n.$$

- ▶ Ta xây dựng bảng $n \times n$: Phần tử (i, j) (với $i \leq j$) của bảng là tập các biến sinh ra chuỗi

$$w_i w_{i+1} \cdots w_j.$$

Ví dụ

Xâu $w = \text{baba}$ có sinh bởi văn phạm sau không?

$$S \rightarrow RT$$

$$R \rightarrow TR \mid \mathbf{a}$$

$$T \rightarrow TR \mid \mathbf{b}$$

	1	2	3	4
1	T	R, T	S	S, R, T
2		R	S	S
3			T	R, T
3				R

Với xâu vào $w = w_1 \cdots w_n$:

1. Nếu $w = \varepsilon$ và có luật $S \rightarrow \varepsilon$, **chấp nhận**. [[xử lý ε]]
2. For $i = 1$ đến n : [[xử lý xâu con độ dài 1]]
3. For mỗi biến A :
4. Kiểm tra xem có luật $A \rightarrow b$ với $w_i = b$.
5. Nếu có, đặt A vào $table(i, i)$.
6. For $l = 2$ đến n : [[l là độ dài xâu con]]
7. For $i = 1$ đến $n - l + 1$: [[i là vị trí bắt đầu]]
8. Đặt $j = i + l - 1$. [[j là vị trí kết thúc]]
9. For $k = i$ đến $j - 1$: [[k là vị trí cắt xâu]]
10. For mỗi luật $A \rightarrow BC$:
11. Nếu $B \in table(i, k)$ và $C \in table(k + 1, j)$,
thì đặt A vào $table(i, j)$.
12. Nếu $S \in table(1, n)$, **chấp nhận**. Ngược lại, **bác bỏ**.