



ĐẠI HỌC BÁCH KHOA HÀ NỘI
VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

Phân tích cú pháp xác suất

Viện Công nghệ Thông tin và Truyền thông

Làm cách nào chọn cây đúng?

- Ví dụ:

I saw a man with a telescope.

- Khi số luật tăng, khả năng nhập nhằng tăng
- Tập luật NYU: bộ PTCP Apple pie : 20,000-30,000 luật cho tiếng Anh
- Lựa chọn luật AD: V DT NN PP

(1) VP → V NP PP

NP → DT NN

(2) VP → V NP

NP → DT NN PP

Kết hợp từ (bigrams pr)

Ví dụ:

Eat ice-cream (high freq)

Eat John (low, except on Survivor)

Nhược điểm:

- P(John decided to bake a) có xác suất cao

• Xét:

$$P(w_3) = P(w_3|w_2w_1) = P(w_3|w_2)P(w_2|w_1)P(w_1)$$

Giả thiết này quá mạnh: chủ ngữ có thể quyết định bổ ngữ trong câu

Clinton admires honesty

➤ sử dụng cấu trúc ngữ pháp để dừng việc lan truyền

- Xét Fred watered his mother's small garden. Từ garden có ảnh hưởng như thế nào?

- $Pr(\text{garden} | \text{mother's small})$ thấp \Rightarrow mô hình trigram không tốt
- $Pr(\text{garden} | X)$ là thành phần chính của bổ ngữ cho động từ *to water* cao hơn

➤ sử dụng bigram + quan hệ ngữ pháp

Kết hợp từ (bigrams pr)

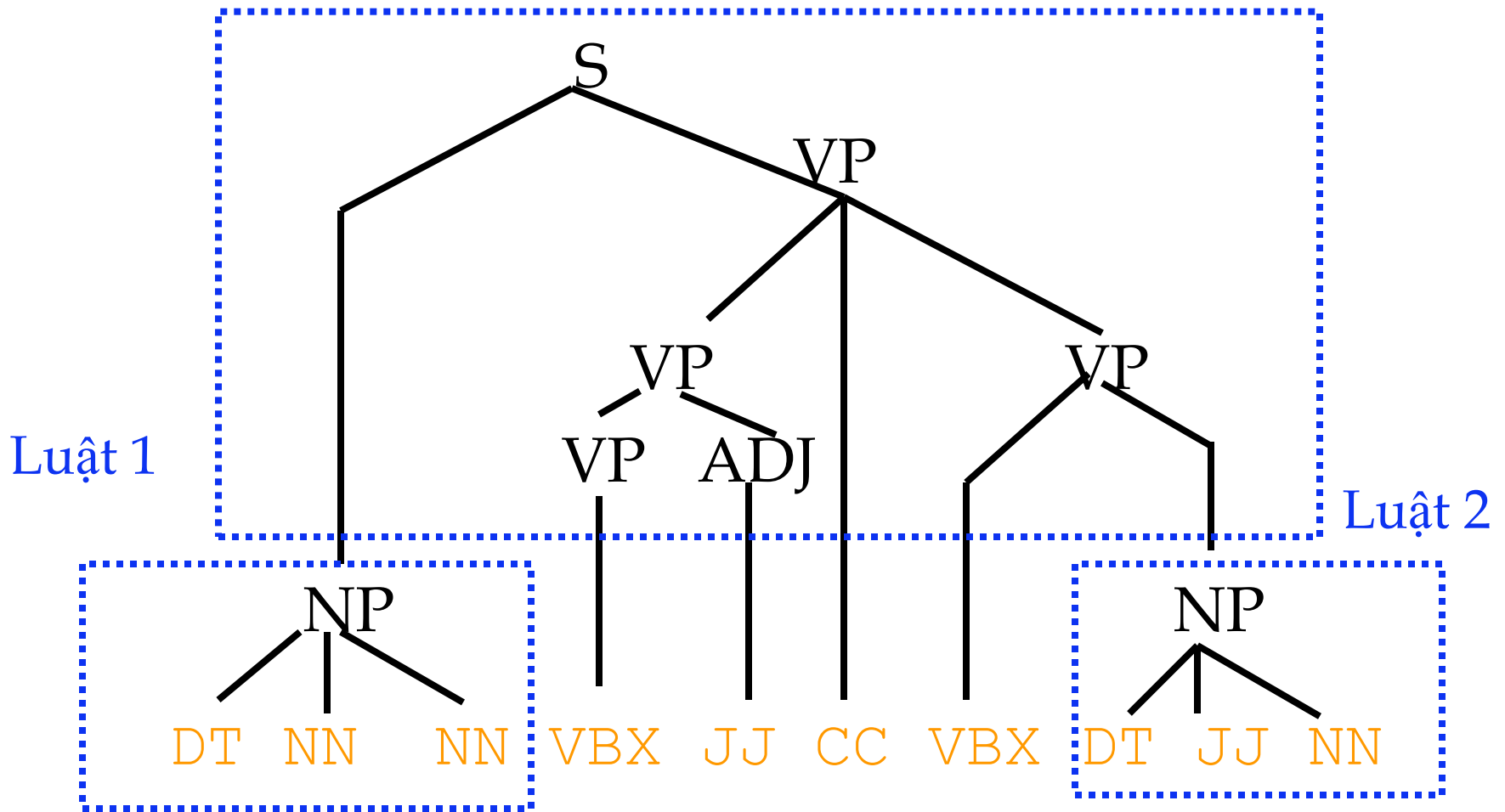
- V có một số loại bổ ngữ nhất định
⇒ Verb-with-obj, verb-without-obj
- Sự tương thích giữa chủ ngữ và bổ ngữ:
John admires honesty
Honesty admires John ???

Nhược điểm:

- Kích thước tập ngữ pháp tăng
- Các bài báo của tạp chí Wall Street Journal trong 1 năm: 47,219 câu, độ dài trung bình 23 từ, gán nhãn bằng tay: chỉ có 4.7% hay 2,232 câu có cùng cấu trúc ngữ pháp
- Không thể dựa trên việc tìm các cấu trúc cú pháp đúng cho cả câu. Phải xây dựng tập các mẫu ngữ pháp nhỏ

Ví dụ

Luật 3



This apple pie looks good and is a real treat

Luật

1. NP → DT NN NN
2. NP → DT JJ NN
3. S → NP VBX JJ CC VBX NP
 - Nhóm (NNS, NN) thành NX; (NNP, NNPs) = NPX;
 - (VBP, VBZ, VBD) = VBX;
 - Chọn các luật theo tần suất của nó

Tính xác suất

$$\Pr(X \rightarrow Y) = \frac{\text{Số lượng } X \text{ chuyển thành } Y}{\text{Số lượng } X} = \frac{1470}{9711} = 0.1532$$

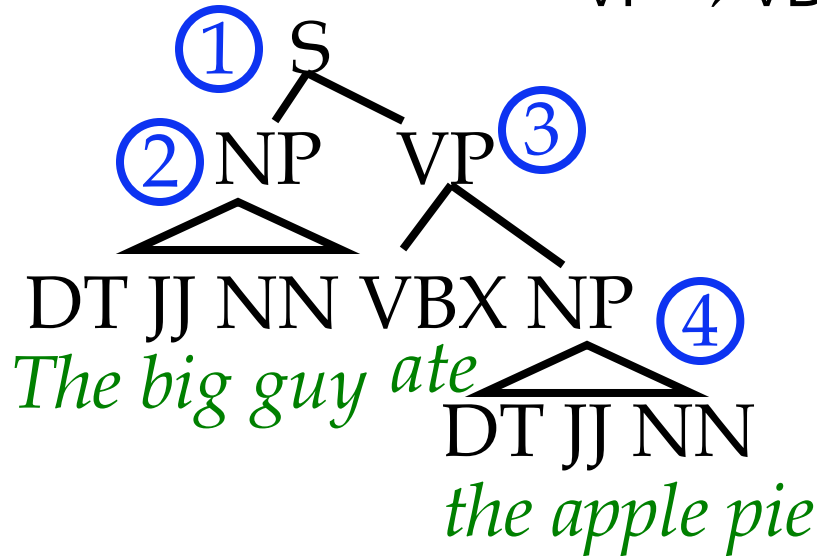
The diagram illustrates the calculation of the probability $\Pr(X \rightarrow Y)$. It shows a transition from state X to state Y . The number of transitions from X to Y is 1470. The total number of transitions from X to any state is 9711. The probability is calculated as $\frac{1470}{9711} = 0.1532$.

Tính Pr

$S \rightarrow NP VP; 0.35$

$NP \rightarrow DT JJ NN; 0.1532$

$VP \rightarrow VBX NP; 0.302$



Luật áp dụng

1 $S \rightarrow NP VP$

2 $NP \rightarrow DT JJ NN$

3 $VP \rightarrow VBX NP$

4 $NP \rightarrow DT JJ NN$

$Pr = 0.0025$

Chuỗi Pr

0.35

$0.1532 \times 0.35 = 0.0536$

$0.302 \times 0.0536 = 0.0162$

$0.1532 \times 0.0162 = 0.0025$

Văn phạm phi ngữ cảnh xác suất

- 1 văn phạm phi ngữ cảnh xác suất (Probabilistic Context Free Grammar) gồm các phần thông thường của CFG
- Tập ký hiệu kết thúc $\{w^k\}$, $k = 1, \dots, V$
- Tập ký hiệu không kết thúc $\{N^i\}$, $i = 1, \dots, n$
- Ký hiệu khởi đầu N^1
- Tập luật $\{N^i \rightarrow \zeta^j\}$, ζ^j là chuỗi các ký hiệu kết thúc và không kết thúc

- Tập các xác suất của 1 luật là:

$$\forall i \sum_j P(N^i \rightarrow \zeta^j) = 1$$

- Xác suất của 1 cây cú pháp:

$$P(T) = \prod_{i=1..n} p(r(i))$$

Các giả thiết

- **Độc lập vị trí:** Xác suất 1 cây con không phụ thuộc vào vị trí của các từ của cây con đó ở trong câu

$\forall k, P(N_{jk}(k+c) \rightarrow \zeta)$ là giống nhau

- **Độc lập ngữ cảnh:** Xác suất 1 cây con không phụ thuộc vào các từ ngoài cây con đó

$P(N_{jkl} \rightarrow \zeta | \text{các từ ngoài khoảng } k \text{ đến } l) = P(N_{jkl} \rightarrow \zeta)$

- **Độc lập tổ tiên:** Xác suất 1 cây con không phụ thuộc vào các nút ngoài cây con đó

$P(N_{jkl} \rightarrow \zeta | \text{các nút ngoài cây con } N_{jkl}) = P(N_{jkl} \rightarrow \zeta)$