

ĐẠI HỌC QUỐC GIA HỒ CHÍ MINH  
TRƯỜNG ĐẠI HỌC BÁCH KHOA  
KHOA KHOA HỌC VÀ KỸ THUẬT MÁY TÍNH



## XỬ LÝ NGÔN NGỮ TỰ NHIÊN (BT) (CO3085)

---

**KỲ 251: LAB 3**

**BÀI TẬP TOÁN**

**Giáo viên:** Bùi Khánh Vĩnh  
**Sinh viên:** Nguyễn Mạnh Dũng 2210583 (L01)

THÀNH PHỐ HỒ CHÍ MINH, 2025



## Contents

1 Bài 1	3
2 Bài 2	3
3 Bài 3	3
4 Bài 4	4
5 Bài 5	4
6 Bài 6	4

## 1 Bài 1

Công thức trigram:

$$P(w_i | w_{i-2}, w_{i-1}) = \frac{\text{count}(w_{i-2}, w_{i-1}, w_i)}{\text{count}(w_{i-2}, w_{i-1})}$$

$$\begin{aligned} P(I | \langle s \rangle, \langle s \rangle) &= \frac{2}{3} = 0.67 \\ P(am | \langle s \rangle, I) &= \frac{1}{2} = 0.5 \\ P(Sam | I, am) &= \frac{1}{2} = 0.5 \\ P(\langle /s \rangle | am, Sam) &= \frac{1}{1} = 1 \\ P(Sam | \langle s \rangle, \langle s \rangle) &= \frac{1}{3} = 0.33 \\ P(I | \langle s \rangle, Sam) &= \frac{1}{1} = 1 \\ P(am | Sam, I) &= \frac{1}{1} = 1 \\ P(\langle /s \rangle | I, am) &= \frac{1}{2} = 0.5 \\ P(do | \langle s \rangle, I) &= \frac{1}{2} = 0.5 \\ P(not | do, I) &= \frac{1}{1} = 1 \\ P(like | do, not) &= \frac{1}{1} = 1 \\ P(green | like, not) &= \frac{1}{1} = 1 \\ P(eggs | green, like) &= \frac{1}{1} = 1 \\ P(and | green, eggs) &= \frac{1}{1} = 1 \\ P(Sam | eggs, and) &= \frac{1}{1} = 1 \\ P(\langle /s \rangle | and, Sam) &= \frac{1}{1} = 1 \end{aligned}$$

## 2 Bài 2

Công thức bigram:

$$P(w_1 w_2 \dots w_n) \approx \prod_i P(w_i | w_{i-1})$$

Xác suất	Không biến đổi Laplace	Có biến đổi Laplace
$P(i   \langle s \rangle)$	0.19	0.19
$P(want   i)$	0.33	0.21
$P(chinese   want)$	0.0065	0.0029
$P(food   chinese)$	0.52	0.52
$P(\langle /s \rangle, food)$	0.40	0.40
$P(câu)$	$8.477 \cdot 10^{-5}$	$2.4068 \cdot 10^{-5}$

Bảng 1: Xác suất sử dụng bigram

## 3 Bài 3

Xác suất unsmoothed cao hơn xác suất smoothed, vì smoothing làm phẳng phân phối, giảm xác suất ở những bigram quan sát được để phân bổ cho các bigram chưa quan sát. Nếu tất cả bigram đều tồn tại ( $p > 0$ ), xác suất unsmoothed sẽ cao hơn xác suất smoothed

## 4 Bài 4

$V = 11$

Số lần xuất hiện từ "am": 3

Số lần xuất hiện cặp "am Sam": 2

Với bigram có dùng add-one smoothing, xác suất là:

$$\frac{2+1}{3+11} \approx 0.2143$$

## 5 Bài 5

Độ dài của corpus: 25

Số lần xuất hiện từ "Sam": 4

Số lần xuất hiện từ "am": 3

Số lần xuất hiện cặp "am Sam": 2

Công thức tính xác suất:

$$P(\text{Sam} \mid \text{am}) = \lambda_1 P_{\text{bi}} + \lambda_2 P_{\text{uni}} = \frac{1}{2} \cdot \frac{2}{3} + \frac{1}{2} \cdot \frac{4}{25} = \frac{31}{75} \approx 0.4133.$$

## 6 Bài 6

Xác suất trong mô hình unigram được tính theo công thức:

$$P(w_i) = \frac{\text{số lần xuất hiện của } w_i}{\text{độ dài của corpus}}$$

$$P(0) = \frac{91}{100}$$

$$P(3) = \frac{1}{100}$$

Xác suất:  $P = (0.91)^8 \times (0.01)^1$ .

Vậy  $PP = P^{-\frac{1}{9}} \approx 1.81$