



TRƯỜNG ĐẠI HỌC CÔNG NGHỆ THÔNG TIN  
KHOA KHOA HỌC MÁY TÍNH

---

MÔN HỌC : XỬ LÝ NGÔN NGỮ TỰ NHIÊN

---

**Nguyễn Quốc Khánh**

22520646

22520646@gm.uit.edu.vn

---

BÀI TẬP QUÁ TRÌNH 1

---

## Bài tập 1.

**1.1 Chứng minh rằng: Không có từ kết thúc </s> :** Mô hình có thể sinh chuỗi dài vô hạn

$$\sum_{n=1}^{\infty} \sum_{x_{1:n}} P(x_{1:n}) = \infty$$

**Giải:**

$$\sum_{n=1}^{\infty} \sum_{x_{1:n}} P(x_{1:n}) = \sum_{n=1}^1 \sum_{x_{1:n}} P(x_{1:n}) + \sum_{n=2}^2 \sum_{x_{1:n}} P(x_{1:n}) + \sum_{n=3}^3 \sum_{x_{1:n}} P(x_{1:n}) + \dots + \sum_{n=\infty}^{\infty} \sum_{x_{1:n}} P(x_{1:n})$$

Vì không có từ kết thúc câu </s> nên mỗi 1 độ dài tự nó sẽ bao quát hết toàn bộ xác suất của nó khi tạo ra các câu có chiều dài n

$$\Rightarrow \sum_{n=1}^1 \sum_{x_{1:n}} P(x_{1:n}) = 1$$

$$\sum_{n=2}^2 \sum_{x_{1:n}} P(x_{1:n}) = \sum_{x' \in V} \sum_{x \in V} P(x | <s>, x') P(x' | <s>) = 1$$

⋮

$$\Rightarrow \sum_{n=1}^{\infty} \sum_{x_{1:n}} P(x_{1:n}) = 1 + 1 + 1 + 1 + \dots = \infty$$

**1.2 Chứng minh rằng: Có từ kết thúc </s> :** Mô hình ngôn ngữ phải dừng lại việc sinh chuỗi bằng từ kết thúc câu </s>.

$$\sum_{n=1}^{\infty} \sum_{x_{1:n}} P(x_{1:n}, </s>) = 1$$

**Giải:**

Ta có:  $P(x_{1:n}, </s>) = P(x_1)P(x_2|x_1)P(x_3|x_1, x_2)\dots P(</s> | x_{1:n})$

Mỗi chuỗi kết thúc </s>, để đảm bảo việc sinh chuỗi là hữu hạn

Mà mô hình bao gồm tất cả các khả năng xảy ra khi tạo 1 sequence thuộc độ dài n.

⇒ Tổng xác suất của tất cả các sequence có độ dài bất kỳ, kết thúc bằng </s> sẽ đúng bằng 1

→ Vì nó sẽ bao quát hết toàn bộ không gian xác suất.

$$\begin{aligned} \sum_{x_{1:n}} P(x_{1:n}, </s>) &= P(\text{các câu} \in \text{độ dài } n) \\ &\Rightarrow \sum_{n=1}^{\infty} P(\text{các câu} \in \text{độ dài } n) = 1 \end{aligned}$$

## Bài tập 2.

Cho tập dữ liệu gồm nhiều văn bản thuộc các lớp  $C = \{c_1, c_2, c_3, \dots, c_k\}$  và mỗi văn bản chứa các từ từ tập từ vựng  $V$ . Hãy sử dụng phương pháp MLE để tính:

### 2.1 Xác suất tiên nghiệm của lớp $c_j$

$$\hat{P}(c_j) = \frac{\text{count}(c_j)}{N_{doc}}$$

Trong đó:

- $\text{Count}(c_j)$  : số văn bản thuộc lớp  $c_j$
- $N_{doc}$  : tổng số văn bản

### Giải:

Gọi  $x_i$  là văn bản thuộc lớp  $c_j$ , với  $i \in [1, m]$ ;  $m$  là số lượng văn bản thuộc lớp  $j$

$$\Rightarrow \sum_{i=1}^m x_i^j = \text{Count}(c_j)$$

Xét các biến  $x_i$  là độc lập

Theo phân phối Categorical (trường hợp tổng quát của phân phối Bermoulli):

$$P(x_i | \hat{P}(c_j)) = \prod_{i=1}^m \hat{P}(c_j)^{x_i}$$

Đánh giá  $\hat{P}(c_j)$  dựa trên Maximum log-likelihood:

$$\begin{aligned} \mathcal{L} \hat{P}(c_j) &= \arg \max [P(x_1, x_2, x_3, \dots, x_m | \hat{P}(c_j))] = \arg \max \prod_{i=1}^m P(x_i | \hat{P}(c_j)) = \arg \max \prod_{i=1}^m \prod_{j=1}^k \hat{P}(c_j)^{x_i} \\ &= \arg \max \prod_{j=1}^k \hat{P}(c_j)^{\sum_{i=1}^m x_i} = \arg \max \prod_{j=1}^k \hat{P}(c_j)^{\text{Count}(c_j)} = \arg \max \sum_{j=1}^k \text{Count}(c_j) \log(\hat{P}(c_j)) \end{aligned}$$

Điều kiện chuẩn :  $\sum_{j=1}^k \hat{P}(c_j) = 1$

$\Rightarrow$  Áp dụng phương pháp nhân tử Lagrange:

$$\mathcal{L}(\hat{P}(c_j), \Lambda) = \sum_{j=1}^k \text{Count}(c_j) \log(\hat{P}(c_j)) + \lambda(1 - \sum_{j=1}^k \hat{P}(c_j))$$

Lấy đạo hàm riêng cho  $\mathcal{L}$

$$\frac{\partial \mathcal{L}(\hat{P}(c_j), \lambda)}{\partial \hat{P}(c_j)} = \frac{\text{Count}(c_j)}{\hat{P}(c_j)} = 0$$

$$\Rightarrow \hat{P}(c_j) = \frac{\text{Count}(c_j)}{\lambda} \quad (1)$$

$$\frac{\partial \mathcal{L}(\hat{P}(c_j), \lambda)}{\partial \lambda} = 1 - \sum_{j=1}^k \hat{P}(c_j) = 0$$

$$\Rightarrow \sum_{j=1}^k = 1 \quad (2)$$

Thế (1) vào (2) :

$$\Leftrightarrow \sum_{j=1}^k \frac{\text{Count}(c_j)}{\lambda} = 1$$

$$\Leftrightarrow \frac{1}{\lambda} \sum_{j=1}^k \text{Count}(c_j) = 1$$

$$\Leftrightarrow \frac{N_{doc}}{\lambda} = 1$$

$$\Rightarrow \lambda = N_{doc}$$

$$\Rightarrow \hat{P}(c_j) = \frac{\text{Count}(c_j)}{\lambda} = \frac{\text{Count}(c_j)}{N_{doc}} (dpcm)$$

## 2.2 Xác suất có điều kiện của từ $w_i$ trong lớp $c_j$

$$P(w_i | c_j) = \frac{\text{count}(w_i, c_j)}{\sum_{w \in V} \text{count}(w, c_j)}$$

Trong đó:

- $\text{count}(w_i, c_j)$ : số lần từ  $w_i$  xuất hiện trong lớp  $c_j$ .
- $\sum_{w \in V} \text{count}(w, c_j)$  : tổng tất cả các từ xuất hiện trong lớp  $c_j$ .

## Giải

Xét các từ  $w_i$  là độc lập.

Theo phân phối categorical:

$$P(count(w_i, c_j) | \hat{P}(w_i | c_j)) = \prod_{w_i \in V} \hat{P}(w_i | c_j)^{count(w_i, c_j)}$$

Đánh giá  $\hat{P}(w_i | c_j)$  dựa trên Maximum log-likelihood:

$$\mathcal{L} \hat{P}(w_i | c_j) = \arg \max [P(count(w_i, c_j) | \hat{P}(w_i | c_j))] = \arg \max [\prod_{w_i \in V} \hat{P}(w_i | c_j)^{count(w_i, c_j)}]$$

$$= \arg \max \sum_{w_i \in V} count(w_i, c_j) \log(\hat{P}(w_i | c_j))$$

Điều kiện chuẩn :  $\sum_{w_i \in V} \hat{P}(w_i | c_j) = 1$

$\Rightarrow$  Áp dụng phương pháp nhân tử Lagrange:

$$\mathcal{L}(\hat{P}(w_i | c_j), \lambda) = \sum_{w_i \in V} count(w_i, c_j) \log(\hat{P}(w_i | c_j)) + \lambda [1 - \sum_{w_i \in V} \hat{P}(w_i | c_j)]$$

Lấy đạo hàm riêng của L :

$$\frac{\partial \mathcal{L}(\hat{P}(w_i | c_j), \lambda)}{\partial \hat{P}(w_i | c_j)} = \frac{count(w_i, c_j)}{\hat{P}(w_i | c_j)} - \lambda = 0$$

$$\Rightarrow \hat{P}(w_i | c_j) = \frac{count(w_i, c_j)}{\lambda} \quad (1)$$

$$\frac{\partial \mathcal{L}(\hat{P}(w_i | c_j), \lambda)}{\partial \lambda} = 1 - \sum_{w_i \in V} \hat{P}(w_i | c_j) = 0$$

$$\Rightarrow \sum_{w_i \in V} \hat{P}(w_i | c_j) = 1 \quad (2)$$

Thế (1) vào (2):

$$\Rightarrow \sum_{w_i \in V} \frac{count(w_i, c_j)}{\lambda} = 1$$

$$\Leftrightarrow \frac{1}{\lambda} \sum_{w_i \in V} count(w_i, c_j) = 1$$

$$\Leftrightarrow \lambda = \sum_{w_i \in V} count(w_i, c_j)$$

$$\Rightarrow \hat{P}(w_i, c_j) = \frac{count(w_i, c_j)}{\lambda} = \frac{count(w_i, c_j)}{\sum_{w_i \in V} count(w_i, c_j)}$$

$$\text{hay } \hat{P}(w_i | c_j) = \frac{count(w_i, c_j)}{\sum_{w \in V} (count(w, c_j))} (dpcm)$$