

# COMS W4705: Natural Language Processing

## Written Homework 1

Qianrui Zhao (qz2338)

September 29, 2018

### Problem 1

(a)

$$P(\textit{Spam}) = \frac{3}{5}$$

$$P(\textit{Ham}) = \frac{2}{5}$$

(b)  $P(\textit{Word}|\textit{Class})$  in table

| $P(\textit{Word} \textit{Class})$ | Spam           | Ham           |
|-----------------------------------|----------------|---------------|
| "buy"                             | $\frac{1}{12}$ | 0             |
| "car"                             | $\frac{1}{12}$ | $\frac{1}{7}$ |
| "Nigeria"                         | $\frac{1}{6}$  | $\frac{1}{7}$ |
| "profit"                          | $\frac{1}{6}$  | 0             |
| "money"                           | $\frac{1}{12}$ | $\frac{1}{7}$ |
| "home"                            | $\frac{1}{12}$ | $\frac{1}{7}$ |
| "bank"                            | $\frac{1}{6}$  | $\frac{2}{7}$ |
| "check"                           | $\frac{1}{12}$ | 0             |
| "wire"                            | $\frac{1}{12}$ | 0             |
| "fly"                             | 0              | $\frac{1}{7}$ |

(c) Predict class label

(a) Predict class of Nigeria

$$P(\text{Spam}|\text{Nigeria}) = \frac{P(\text{Spam}) \cdot P(\text{Nigeria}|\text{Spam})}{P(\text{Nigeria})} = \frac{\frac{3}{5} \cdot \frac{1}{6}}{P(\text{Nigeria})} = \frac{\frac{1}{10}}{P(\text{Nigeria})}$$

$$P(\text{Ham}|\text{Nigeria}) = \frac{P(\text{Ham}) \cdot P(\text{Nigeria}|\text{Ham})}{P(\text{Nigeria})} = \frac{\frac{2}{5} \cdot \frac{1}{7}}{P(\text{Nigeria})} = \frac{\frac{2}{35}}{P(\text{Nigeria})}$$

Because  $P(\text{Nigeria})$  is a positive number and  $P(\text{Spam}|\text{Nigeria}) > P(\text{Ham}|\text{Nigeria})$

Therefore, Predicted class for Nigeria is Spam.

(b) Predict class of Nigeria, Home

$$P(\text{Spam}|\text{Nigeria}, \text{home}) = \frac{P(\text{Spam}) \cdot P(\text{Nigeria}|\text{Spam}) \cdot P(\text{Home}|\text{Spam})}{P(\text{Nigeria}) \cdot P(\text{Home})}$$

$$= \frac{\frac{3}{5} \cdot \frac{1}{6} \cdot \frac{1}{12}}{P(\text{Nigeria}) \cdot P(\text{Home})} = \frac{\frac{1}{120}}{P(\text{Nigeria}) \cdot P(\text{Home})}$$

$$P(\text{Ham}|\text{Nigeria}, \text{Home}) = \frac{P(\text{Ham}) \cdot P(\text{Nigeria}|\text{Ham}) \cdot P(\text{Home}|\text{Ham})}{P(\text{Nigeria}) \cdot P(\text{Home})}$$

$$= \frac{\frac{2}{5} \cdot \frac{1}{7} \cdot \frac{1}{7}}{P(\text{Nigeria}) \cdot P(\text{Home})} = \frac{\frac{2}{245}}{P(\text{Nigeria}) \cdot P(\text{Home})}$$

Because  $P(\text{Nigeria}) \cdot P(\text{Home})$  is a positive number and  $P(\text{spam}|\text{Nigeria}, \text{Home}) > P(\text{Ham}|\text{Nigeria}, \text{Home})$

Therefore, predicted class for Nigeria, Home is Spam.

(c) Predict class of Home, Bank, Money

$$P(\text{Spam}|\text{Home}, \text{Bank}, \text{Money}) = \frac{P(\text{Spam}) \cdot P(\text{Home}|\text{Spam}) \cdot P(\text{Bank}|\text{Spam}) \cdot P(\text{Money}|\text{Spam})}{P(\text{Home}) \cdot P(\text{Bank}) \cdot P(\text{Money})}$$

$$= \frac{\frac{3}{5} \cdot \frac{1}{12} \cdot \frac{1}{6} \cdot \frac{1}{12}}{P(\text{Home}) \cdot P(\text{Bank}) \cdot P(\text{Money})} = \frac{\frac{1}{1440}}{P(\text{Home}) \cdot P(\text{Bank}) \cdot P(\text{Money})}$$

$$P(\text{Ham}|\text{Home}, \text{Bank}, \text{Money}) = \frac{P(\text{Ham}) \cdot P(\text{Home}|\text{Ham}) \cdot P(\text{Bank}|\text{Ham}) \cdot P(\text{Money}|\text{Ham})}{P(\text{Home}) \cdot P(\text{Bank}) \cdot P(\text{Money})}$$

$$= \frac{\frac{2}{5} \cdot \frac{1}{7} \cdot \frac{2}{7} \cdot \frac{1}{7}}{P(\text{Home}) \cdot P(\text{Bank}) \cdot P(\text{Money})} = \frac{\frac{4}{1715}}{P(\text{Home}) \cdot P(\text{Bank}) \cdot P(\text{Money})}$$

Because  $P(\text{Home}) \cdot P(\text{Bank}) \cdot P(\text{Money})$  is a positive number

and  $P(\text{Spam}|\text{Home}, \text{Bank}, \text{Money}) < P(\text{Ham}|\text{Home}, \text{Bank}, \text{Money})$

Therefore, predicted class of Home, Bank, Money is Ham.

## Problem 2

Prove that sum of prob of all sentence with length == n is 1:

Assume vocabulary size is  $V$ , words  $W = \{w_1, w_2, \dots, w_V\}$ ,  $w^n$  means nth word in the sentence.

Lemma 01: Two words  $w^n$  and  $w^{n+1}$ ,

$$\sum_{w^{n+1} \in W} P(w^{n+1}, w^n) = (P(w_1, w^n) + P(w_2, w^n) + \dots + P(w_V, w^n)) \cdot P(w^n) = P(w^n)$$

Lemma 02: When sentence length == 1,

$$\sum_{w_i \in W} P(w^1 | \text{START}) = P(w_1 | \text{START}) + P(w_2 | \text{START}) \dots p(w_n | \text{START}) = 1$$

Induction formula can be built as:

$$\begin{aligned} & \sum_{w^1, w_2, \dots, w^{n+1} \in W} P(w^1 | \text{START}) P(w^2 | w^1) \dots P(w^{n+1} | P(w^n)) \\ &= \sum_{w^1, w_2, \dots, w^n \in W} P(w^1 | \text{START}) P(w^2 | w^1) \dots P(w^n | P(w^{n-1})) \left( \sum_{w^{n+1} \in W} P(w^{n+1} | w^n) \right) \\ &= \sum_{w^1, w_2, \dots, w^n \in W} \frac{P(w^1 | \text{START}) P(w^2 | w^1) \dots P(w^n | P(w^{n-1}))}{P(w^n)} \left( \sum_{w^{n+1} \in W} P(w^{n+1}, w^n) \right) \end{aligned}$$

According to Lemma 01, where  $\sum_{w^{n+1} \in W} P(w^{n+1}, w^n) = P(w^n)$

$$\begin{aligned} &= \sum_{w^1, w_2, \dots, w^n \in W} \frac{P(w^1 | \text{START}) P(w^2 | w^1) \dots P(w^n | P(w^{n-1}))}{P(w^n)} \cdot P(w^n) \\ &= \sum_{w^1, w_2, \dots, w^n \in W} P(w^1 | \text{START}) P(w^2 | w^1) \dots P(w^n | w^{n-1}) \end{aligned}$$

Therefore,  $\sum_{w^1, w_2, \dots, w^{n+1} \in W} P(w^1 | \text{START}) P(w^2 | w^1) \dots P(w^{n+1} | w^n) = \sum_{w^1 \in W} P(w^1 | \text{START})$

According to Lemma 2, where  $\sum_{w^1 \in W} P(w^1 | \text{START}) = 1$

$$\sum_{w^1, w_2, \dots, w^{n+1} \in W} P(w^1 | \text{START}) P(w^2 | w^1) \dots P(w^{n+1} | w^n) = 1$$