

COMS 4705 Natural Language Processing (2020 Spring)

Homework 1

Shusen Xu - `sx2261@columbia.edu`

February 18, 2020

Problem 1 Text Classification with Naive Bayes

1. Prior Probability of a random email

$$P_{mle}(Class = spam) = \frac{3}{5}$$

$$P_{mle}(Class = ham) = \frac{2}{5}$$

2. Estimate the conditional probability distribution

Bag of Words:

| Indexes | Word | Count |
|---------|---------|-------|
| 0 | buy | 1 |
| 1 | car | 2 |
| 2 | Nigeria | 3 |
| 3 | profit | 2 |
| 4 | money | 2 |
| 5 | home | 3 |
| 6 | bank | 3 |
| 7 | check | 1 |
| 8 | wire | 1 |
| 9 | fly | 1 |

Conditional Probability distribution:

| Word | 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 | $\frac{1}{12}$ | $\frac{2}{7}$ |
| bank | $\frac{1}{6}$ | $\frac{1}{7}$ |
| check | $\frac{1}{12}$ | 0 |
| wire | $\frac{1}{12}$ | 0 |
| fly | 0 | $\frac{1}{7}$ |

3. Naive Bayes Classifier

Since $y^* = \arg \max_y P(y) \prod_i P(x_i|y)$ (Naive Bayes classifier)

(a) Nigeria

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

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

Since $\frac{1}{10} > \frac{2}{35}$, the class of "Nigeria" is spam.

(b) Nigeria home

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

$$P(\text{ham}) \cdot P(\text{Nigeria}|\text{ham}) \cdot P(\text{home}|\text{ham}) = \frac{2}{5} \cdot \frac{1}{7} \cdot \frac{2}{7} = \frac{4}{245}$$

Since $\frac{1}{120} < \frac{4}{245}$, the class of "Nigeria home" is ham.

(c) home bank money

$$P(\text{spam}) \cdot P(\text{home}|\text{spam}) \cdot P(\text{bank}|\text{spam}) \cdot P(\text{money}|\text{spam}) = \frac{3}{5} \cdot \frac{1}{12} \cdot \frac{1}{6} \cdot \frac{1}{12} = \frac{1}{1440}$$

$$P(\text{ham}) \cdot P(\text{home}|\text{ham}) \cdot P(\text{bank}|\text{ham}) \cdot P(\text{money}|\text{ham}) = \frac{2}{5} \cdot \frac{2}{7} \cdot \frac{1}{7} \cdot \frac{1}{7} = \frac{4}{1715}$$

Since $\frac{1}{1440} < \frac{4}{1715}$, the class of "home bank money" is ham.

Problem 2 Bigram Models

Goal: Prove that, if you sum up the probabilities of all sentence of length n under a bigram language model, this sum is exactly 1 by induction.

$$\sum_{w_1, w_2, \dots, w_n} P(w_1, w_2, \dots, w_n) = \sum_{w_1, w_2, \dots, w_n} P(w_1 | \text{start}) \cdot P(w_2 | w_1) \cdots P(w_n | w_{n-1}) = 1$$

1. Initialization

when the length of sentence is 1:

since $\sum_{w_1} p(w_1)$ sum over all possible w_1 , the value of it obviously is 1. So, combining with the definition of bigram model, we can easily get the following equation.

$$\sum_{w_1} p(w_1) = \sum_{w_1} p(w_1 | \text{start}) = 1$$

2. Maintenance:

- (a) First, assume when the length of sentence is equal to n , the following equation holds true.

$$\sum_{w_1, w_2, \dots, w_n} P(w_1, w_2, \dots, w_n) = \sum_{w_1, w_2, \dots, w_n} P(w_1 | \text{start}) \cdot P(w_2 | w_1) \cdots P(w_n | w_{n-1}) = 1$$

- (b) we need to prove when the length of sentence is equal to $n+1$, the following equation holds true.

$$\sum_{w_1, w_2, \dots, w_{n+1}} P(w_1, w_2, \dots, w_n, w_{n+1}) = \sum_{w_1, w_2, \dots, w_n} P(w_1 | \text{start}) \cdot P(w_2 | w_1) \cdots P(w_{n+1} | w_n) = 1$$

Firstly, according to the bigram model definition which based on Markov assumption, we can get the following equation:

$$\sum_{w_1, w_2, \dots, w_{n+1}} P(w_1, w_2, \dots, w_n, w_{n+1}) = \sum_{w_1, w_2, \dots, w_n} P(w_1 | \text{start}) \cdot P(w_2 | w_1) \cdots P(w_{n+1} | w_n)$$

Secondly, let's denote sentence (w_1, w_2, \dots, w_n) as event A . Notes: since $\sum_{w_{n+1}} P(w_{n+1} | A)$ summer over all possible w_{n+1} under the condition that the previous words sequence is A , so obviously the equation is equal to 1. so we can get:

$$\begin{aligned} \sum_{w_1, w_2, \dots, w_{n+1}} P(w_1, w_2, \dots, w_n, w_{n+1}) &= \sum_{A, w_{n+1}} P(A, w_{n+1}) \\ &= \sum_{A, w_{n+1}} P(A) \cdot P(w_{n+1} | A) \\ &= \sum_A P(A) \cdot \left(\sum_{w_{n+1}} P(w_{n+1} | A) \right) \end{aligned}$$

$$= \sum_A P(A) \cdot 1 = \sum_{w_1, w_2, \dots, w_n} P(w_1, w_2, \dots, w_n) = 1$$

3. we proved that: the equation is true when length of sentence is equal to $n+1$. By induction, we can conclude that for any length of sentence, this equation always holds true.