

## Analytical Component

### Problem 1

1. For finding prior probabilities , we have 3 spam emails and 2 ham emails.

So  $P(\text{spam}) = 3/5$

$P(\text{ham}) = 2/5$

2. For spam emails

buy – 1

car -1

Nigeria – 2

profit -2

money-1

home-1

bank-2

check-1

wire-1

Total words = 12

$P(\text{buy} \mid \text{spam}) = P(\text{car} \mid \text{spam}) = P(\text{money} \mid \text{spam}) = P(\text{home} \mid \text{spam}) = P(\text{check} \mid \text{spam}) = P(\text{wire} \mid \text{spam}) = 1/12$

$P(\text{Nigeria} \mid \text{spam}) = P(\text{profit} \mid \text{spam}) = P(\text{bank} \mid \text{spam}) = 2/12 = 1/6$

For ham emails

money -1

bank-1

home-2

car-1

Nigeria-1

fly-1

Total words = 7

$P(\text{money} \mid \text{ham}) = P(\text{bank} \mid \text{ham}) = P(\text{car} \mid \text{ham}) = P(\text{Nigeria} \mid \text{ham}) = P(\text{fly} \mid \text{ham}) = 1/7$

$P(\text{home} \mid \text{ham}) = 2/7$

3. For Nigeria

$P(\text{spam} \mid \text{Nigeria}) = (P(\text{Nigeria} \mid \text{Spam}) \cdot P(\text{spam})) / P(\text{Nigeria})$

Since the denominator will remain constant in both cases of ham and spam, we are ignoring the denominator during our following calculations.

$$\begin{aligned}
 P(\text{spam} \mid \text{Nigeria}) &= P(\text{Nigeria} \mid \text{Spam}) \cdot P(\text{spam}) \\
 &= \frac{2}{12} \cdot \frac{3}{5} \\
 &= 0.1
 \end{aligned}$$

$$\begin{aligned}
 P(\text{ham} \mid \text{Nigeria}) &= P(\text{Nigeria} \mid \text{ham}) \cdot P(\text{ham}) \\
 &= \frac{1}{7} \cdot \frac{2}{5} \\
 &= 0.0571
 \end{aligned}$$

Since  $P(\text{spam} \mid \text{Nigeria})$  is greater, predicted class = spam

For Nigeria home

$$\begin{aligned}
 P(\text{spam} \mid \text{Nigeria home}) &= P(\text{Nigeria home} \mid \text{spam}) P(\text{spam}) \\
 &= P(\text{Nigeria} \mid \text{spam}) \cdot P(\text{home} \mid \text{spam}) P(\text{spam}) \quad [\text{assuming conditional independence of words}] \\
 &= \frac{2}{12} \cdot \frac{1}{7} \cdot \frac{3}{5} \\
 &= 0.01428
 \end{aligned}$$

$$\begin{aligned}
 P(\text{ham} \mid \text{Nigeria home}) &= P(\text{Nigeria home} \mid \text{ham}) P(\text{ham}) \\
 &= P(\text{Nigeria} \mid \text{ham}) \cdot P(\text{home} \mid \text{ham}) P(\text{ham}) \\
 &= \frac{1}{7} \cdot \frac{2}{7} \cdot \frac{2}{5} \\
 &= 0.01632
 \end{aligned}$$

So, predicted class = ham

For home bank money

$$\begin{aligned}
 P(\text{spam} \mid \text{home bank money}) &= P(\text{home} \mid \text{spam}) \cdot P(\text{bank} \mid \text{spam}) \cdot P(\text{money} \mid \text{spam}) \cdot P(\text{spam}) \\
 &= \frac{1}{12} \cdot \frac{2}{12} \cdot \frac{1}{12} \cdot \frac{3}{5} \\
 &= 0.000694
 \end{aligned}$$

$$\begin{aligned}
 P(\text{ham} \mid \text{home bank money}) &= P(\text{home} \mid \text{ham}) \cdot P(\text{bank} \mid \text{ham}) \cdot P(\text{money} \mid \text{ham}) \cdot P(\text{ham}) \\
 &= \frac{1}{7} \cdot \frac{1}{7} \cdot \frac{2}{7} \cdot \frac{2}{5} \\
 &= 0.00233
 \end{aligned}$$

So, predicted class = ham

## Problem 2

Let  $P(n)$  denote sum of probabilities of all sentences of length  $n$ .

For  $P(1)$

Let vocab = {a,b,c}, then sentences are -

START a (probability is  $\frac{1}{3}$ )

START b (probability is  $\frac{1}{3}$ )

START c (probability is  $\frac{1}{3}$ )

So total sum is 1. So  $P(1) = 1$

Lets take generalized case of all sentences of length  $n-1$ .

Our vocab is  $\{a,b,c\}$

Lets say we have  $k$  possible sentences of length  $n-1$ .

Probability of each sentence (using bigram model) =  $1/k$  (as we can see from  $n=1$ , and  $n=2$ )

Now if we move to all sentence of length  $n$ .

Our total no. of sentences =  $3k$  ( since our vocab size is 3)

In each of sentences , the probability product upto  $n-1$ th word =  $1/k$

And the probability of last word  $P(w_n | w_{n-1}) = 1/3$

So, probability of the entire sentence =  $1/k * 1/3 = 1/3k$

So sum of probability of all sentences of length  $n = P(n) = 1/3k * 3k = 1$

Hence proved,  $\sum P(w_1, w_2, \dots, w_n) = 1$