

CSE431  
Final Exam

Tasnim Sakib Apon

ID : 20241068

Old ID : 18301297

1. Naive bayes algorithm in unknown words:

A simple way to deal with unknown words is simply to add an extra word to the vocabulary.

Let  $w_u$  = "unknown word"

$$\hat{p}(w_u|e) = \frac{\text{Count}(w_u, e) + 1}{\left( \sum_{w \in V} \text{Count}(w, e) \right) + |V| + 1}$$

$$= \frac{1}{\left( \sum_{w \in V} \text{Count}(w, e) \right) + |V| + 1}$$

Unknown words will be modeled with the equation above.

And normally stopwords are removed before applying naive bayes. However if stopwords are found, naive bayes simply ignore the words. Stop words are removed both from test set and train set.

## 2 5 usages of text categorization

- ⇒ language identification
- ⇒ Authorship identification
- ⇒ Age/gender identification
- ⇒ Spam detection.
- ⇒ Assigning subject categories, topic or genres.

## 3

	category	Document	Words	Words in same category
Train	Positive	I like walking	3	1
	Negative	This is not good	4	1
Test	?	I like drinking coffee	4	

4 prior probability: Total Note = 2; Positive = 1; Negative = 1.

$$p(-) = \frac{1}{2} \text{ (Negative)}$$

$$p(+) = \frac{1}{2} \text{ (Positive)}$$

5 Vocabulary size  $|V| = \{ \text{I, like, walking, This, is, not, good} \}$

Theme:

$$6 \quad p(\text{like} | \text{negative}) = \frac{0+1}{4+2} = \frac{1}{11}$$

$$p(\text{like} | \text{positive}) = \frac{1+1}{3+2} = \frac{2}{10} = \frac{1}{5}$$

$$p(\text{drinking} | \text{negative}) = \frac{0+1}{4+2} = \frac{1}{11}$$

$$p(\text{drinking} | \text{positive}) = \frac{0+1}{3+2} = \frac{1}{10}$$

$$p(\text{coffee} | \text{negative}) = \frac{0+1}{4+2} = \frac{1}{11}$$

$$p(\text{coffee} | \text{positive}) = \frac{0+1}{3+2} = \frac{1}{10}$$

$$\therefore p(\text{negative} | \text{test} | \text{negative})$$

$$\Rightarrow \frac{1}{11} \times \frac{1}{11} \times \frac{1}{11}$$

$$\Rightarrow 7.51 \times 10^{-4}$$

$$p(\text{positive} | \text{test} | \text{positive})$$

$$\Rightarrow \frac{1}{5} \times \frac{1}{10} \times \frac{1}{10}$$

$$\Rightarrow 2 \times 10^{-3}$$

$\therefore$  So our test document was positive as  
 $p(\text{positive}) > p(\text{negative})$



8.

Document 1  $\Rightarrow$  Dialogue system is computer system intended to converse with a human.

Document 2  $\Rightarrow$  SIRI is a very common chatbot now a days

Document 3  $\Rightarrow$  Chatbot is also known as dialogue system.

8.

	Document 1	Document 2	Document 3
Chatbot	0	1	1
Dialogue System	1	0	1

9.



