

Author of Solutions: Manzoor Ali

Email: manzooralis29@gmail.com

Natural Language Processing

Chapter 04:

NAIVE BAYES AND SENTIMENT CLASSIFICATION

Book: Speech and language processing Book

(Authors of book: Daniel Jurafsky and James H.Martin)

1. Assume the following likelihoods for each word being part of a positive or negative movie review, and equal prior probabilities for each class.

	pos	neg
I	0.09	0.16
always	0.07	0.06
like	0.29	0.06
foreign	0.04	0.15
films	0.08	0.11

What class will Naive bayes assign to the sentence "I always like foreign films."?

Solution:

We ignore the equal prior probabilities in our computation because it's having equal effect on both results.

$$P(\text{"I always like foreign films"}|\text{pos}) = 0.09 \times 0.07 \times 0.29 \times 0.04 \times 0.08 = 0.000005846$$

$$P(\text{"I always like foreign films"}|\text{neg}) = 0.16 \times 0.06 \times 0.06 \times 0.15 \times 0.11 = 0.000009504$$

$P(s|\text{neg}) > P(s|\text{pos})$ so the Naive bayes assigns a **neg** class.

2. Given the following short movie reviews, each labeled with a genre, either comedy or action:

- a. 1. fun, couple, love, love **comedy**
- b. 2. fast, furious, shoot **action**
- c. 3. couple, fly, fast, fun, fun **comedy**
- d. 4. furious, shoot, shoot, fun **action**
- e. 5. fly, fast, shoot, love **action**

and a new document D:

- f. fast, couple, shoot, fly

compute the most likely class for D. Assume a naive Bayes classifier and use add-1 smoothing for the likelihoods.

Solution:

Vocabulary (V) = [fun, couple, love, fast, furious, shoot, fly]

len(V) = 7

Sentence (S) = fast, couple, shoot, fly

Classes (C) = [comedy(c), action(a)]

bigdoc[c] = [fun, couple, love, love, couple, fly, fast, fun, fun] = 9

bigdoc[a] = [fast, furious, shoot, furious, shoot, shoot, fun, fly, fast, shoot, love] = 11

We will only find the probabilities of the test sentence (S)

$P(c) = \frac{2}{9}$

$P(a) = \frac{3}{11}$

$$P(\text{fast}|c) = \frac{1+1}{9+7} = \frac{2}{16}$$

$$P(\text{fast}|a) = \frac{2+1}{11+7} = \frac{3}{18}$$

$$P(\text{couple}|c) = \frac{2+1}{9+7} = \frac{3}{16}$$

$$P(\text{couple}|a) = \frac{0+1}{11+7} = \frac{1}{18}$$

$$P(\text{shoot}|c) = \frac{0+1}{9+7} = \frac{1}{16}$$

$$P(\text{shoot}|a) = \frac{4+1}{11+7} = \frac{5}{18}$$

$$P(\text{fly}|c) = \frac{1+1}{9+7} = \frac{2}{16}$$

$$P(\text{fly}|a) = \frac{1+1}{11+7} = \frac{2}{18}$$

$$c = \frac{2}{9} \times \frac{2}{16} \times \frac{3}{16} \times \frac{1}{16} \times \frac{2}{16} = 0.000073242$$

$$a = \frac{1}{18} \times \frac{3}{18} \times \frac{1}{18} \times \frac{5}{18} \times \frac{2}{18} = 0.000171468$$

The selected class will be = $\text{argmax}(c, a) = a = \text{Action}$

3. Train two models, multinomial naive Bayes and binarized naive Bayes, both with add-1 smoothing, on the following document counts for key sentiment words, with positive or negative class assigned as noted.

doc	“good”	“poor”	“great”	(class)
d1.	3	0	3	pos
d2.	0	1	2	pos
d3.	1	3	0	neg
d4.	1	5	2	neg
d5.	0	2	0	neg

Use both naive Bayes models to assign a class (pos or neg) to this sentence:

A good, good plot and great characters, but poor acting.

Do the two models agree or disagree?

Solution:

Vocabulary (V) = [good, poor, great]

$\text{len}(V) = 3$

Sentence (S) = A good, good plot and great characters, but poor acting.

Classes (C) = [pos, neg]

Prior probabilities:

$P(\text{pos}) = \frac{1}{5}$

$P(\text{neg}) = \frac{4}{5}$

a) Multinomial naive Bayes

$\text{bigdoc}[\text{pos}] = [\text{good}, \text{good}, \text{good}, \text{great}, \text{great}, \text{great}, \text{poor}, \text{great}, \text{great}] = 9$

$\text{bigdoc}[\text{neg}] = [\text{good}, \text{poor}, \text{poor}, \text{poor}, \text{good}, \text{poor}, \text{poor}, \text{poor}, \text{poor}, \text{poor}, \text{great}, \text{great}, \text{poor}, \text{poor}] = 14$

$$P(\text{good}|\text{pos}) = \frac{3+1}{9+3} = \frac{4}{12} \quad P(\text{good}|\text{neg}) = \frac{2+1}{14+3} = \frac{3}{17}$$

$$P(\text{great}|\text{pos}) = \frac{5+1}{9+3} = \frac{6}{12} \quad P(\text{great}|\text{neg}) = \frac{2+1}{14+3} = \frac{3}{17}$$

$$P(\text{poor}|\text{pos}) = \frac{1+1}{9+3} = \frac{2}{12} \quad P(\text{poor}|\text{neg}) = \frac{10+1}{14+3} = \frac{11}{17}$$

$$[\text{pos}] = 0.4 \times (4 \div 12) \times (6 \div 12) \times (2 \div 12) = 0.011111111$$

$$[\text{neg}] = 0.6 \times (3 \div 17) \times (3 \div 17) \times (11 \div 17) = 0.012090372$$

So $P(\text{neg}) > P(\text{pos})$ we assert that multinomial naive bayes will assign class negative **(neg)**

b) Binarized naive Bayes

$$\text{bigdoc}[\text{pos}] = [\text{good}, \text{great}, \text{great}, \text{poor}] = 4$$

$$\text{bigdoc}[\text{neg}] = [\text{good}, \text{poor}, \text{good}, \text{poor}, \text{great}, \text{poor}] = 6$$

$$P(\text{good}|\text{pos}) = \frac{1+1}{4+3} = \frac{2}{7} \quad P(\text{good}|\text{neg}) = \frac{2+1}{6+3} = \frac{3}{9}$$

$$P(\text{great}|\text{pos}) = \frac{2+1}{4+3} = \frac{3}{7} \quad P(\text{great}|\text{neg}) = \frac{1+1}{6+3} = \frac{2}{9}$$

$$P(\text{poor}|\text{pos}) = \frac{1+1}{4+3} = \frac{2}{7} \quad P(\text{poor}|\text{neg}) = \frac{3+1}{6+3} = \frac{4}{9}$$

$$[\text{pos}] = 0.4 \times (2 \div 7) \times (3 \div 7) \times (2 \div 7) = 0.013994169$$

$$[\text{neg}] = 0.6 \times (3 \div 9) \times (2 \div 9) \times (4 \div 9) = 0.019753086$$

Hence, $P(\text{neg}) > P(\text{pos})$ we assert that binarized naive Bayes will assign class negative **(neg)**

In this case, both models are agreed.