

# NLP assignment2 Detect Sentiment Polarity with Naive Bayes

姓名: 王涵 学号: 2015211984

---

Keywords:  
Naive Bayes, Sentiment Analysis

## Problem 1

---

Generally speaking, sentiment analysis aims to determine the attitude of a document. In this task, the goal is to be able to tell whether the review is positive or negative.

In this experiment, to solve this two-class classification problem, the thought is to deploy a **Naive Bayes Model** to compute the posterior probabilities of words shows in the dataset.

## Usage

---

### Software Requirement

```
python 2.7.6  
numpy=1.12  
nltk=3.2
```

### Source Code

[naive-bayes.py](#)

## Algorithms and Models

---

Naive Bayes is a probability-based model to classify the dataset. The model is based on Bayes' Theorem.

$$P(A | B) = \frac{P(B | A) P(A)}{P(B)},$$

In training step, the goal is to compute the probability that word  $w_i$ 's conditional probability of show in  $c_j$  type of text.

In this model,  $C = \{c_i | i \in N\} = \{\text{Pos}, \text{Neg}\}$ , the goal is to compute  $P(w_i | \text{Pos})$  and  $P(w_i | \text{Neg})$

The probabilities is generated by **Maximum Likelihood Estimation**.

$$\begin{aligned} \hat{c} &= \arg \max_k \{P(c_k | x_i)\} \\ &= \arg \max_k \left\{ \frac{P(c_k) P(x_i | c_k)}{P(x_i)} \right\} \\ &= \arg \max_k \{P(c_k) P(x_i | c_k)\} \\ &= \arg \max_k \{ \log P(c_k) + \log P(x_i | c_k) \} \\ &= \arg \max_k \left\{ \log P(c_k) + \log \prod_{j=1}^M P(t_j | c_k)^{n_{ij}} \right\} \end{aligned}$$

where  $n_{ij}$  is the count of term  $t_j$  in document  $x_i$ .

假设每个文档只属于一个类别  
Bayes法则  
假设terms服从多项式分布且相互独立

Why?

While in test step, probability for each class is computed respectivaly to find the class with largest probability. In doing so, the class of a test sample is labeled.

## Preproccessing

- Remove all the stopwords. (Proposed in nltk corpus)
- Skip the word that doesn't show up in training step.

## Perfomances

To test the accuracy on this model, I pull out a 5-fold cross-validation.

	<b>Acc.(with stopwords)</b>	<b>Acc. (without stopwords)</b>
1	77.5328 %	77.2983 %
2	77.0637 %	76.6885 %
3	78.0018 %	76.7354 %
4	76.5947 %	76.3133 %
5	76.1949 %	75.6794 %

## Analysis and Conclusion

---

From the charts showed above, the model that didn't remove the stopwords is doing slight better. The possible explanation may be that the distribution of stopwords on this dataset is uneven.

Overall, the accuracy on this dataset is slightly higher than what is expected to see with Naive Bayes Model by several paper. It may because the preprocessing method is working. And it helps with the classification.

The part-of-speech tags and the polarity of words are not used in this experiment.