

Log Likelihoods

words can have many shades of emotional meaning, but for the purpose of sentiment classification, they are simplified into three categories: neutral, positive & negative.

All can be identified by using their conditional probabilities →

$$\text{ratio}(w_i) = \frac{P(w_i | \text{Pos})}{P(w_i | \text{Neg})}$$

word	Pos	Neg	ratio		
I	0.20	0.20	1	∞	Positive
am	0.20	0.20	1		
Happy	0.14	0.10	1.4	1	Neutral
Because	0.10	0.10	1		
learning	0.10	0.10	1		
NLP	0.10	0.10	1		
sad	0.10	0.15	0.6	0	negative
not	0.10	0.15	0.6		

The larger the ratio the more positive the word is going to be.

Now, on the other hands, negative words have a ratio smaller than 1.

$$\approx \frac{\text{freq}(w_i, 1) + 1}{\text{freq}(w_i, 0) + 1}$$

Naïve Bayes Inference

Class $\in \{\text{pos}, \text{neg}\}$

$\omega \rightarrow$ set of m words in a tweet

$$\prod_{i=1}^m \frac{P(\omega_i | \text{Pos})}{P(\omega_i | \text{Neg})} > 1 \quad \frac{P(\text{Pos})}{P(\text{Neg})}$$

↓
likelihood

- A simple, fast and powerful Baseline.
- A probabilistic method used for classification

Log Likelihood → • sentiments probability can requires multiplication of numbers with values between 0 and 1.

- Carrying out such multiplications on a computer the risk of numerical underflow when the number returned is so small if it can't be stored device

- mathematical trick to solve this -
- $$\log(a * b) = \log(a) + \log(b)$$

$$\log \left(\frac{P(\text{Pos})}{P(\text{Neg})} \prod_{i=1}^m \frac{P(\omega_i | \text{Pos})}{P(\omega_i | \text{Neg})} \right)$$

log (prior multiplied by likelihood)

$\log C_{prior} + \log C_{likelihood}$

or

$$\log \frac{P(POS)}{P(neg)} + \sum_{i=1}^n \log \frac{P(w_i | Pos)}{P(w_i | neg)}$$

Calculating Lambda —

$$\lambda(I) = \log \frac{0.05}{0.05} \Rightarrow 0$$

	pos	neg	λ
I	0.05	0.05	0
am	0.04	0.04	0
happy	0.09	0.01	2.2
use	0.01	0.01	0
sing	0.03	0.01	1.1
is	0.02	0.02	0
not	0.01	0.09	-2.2
it	0.02	0.03	-0.4

$$\lambda(w) = \log \frac{P(w|Pos)}{P(w|neg)}$$

$$g_{\text{lambda}}(w) = \frac{P(w|Pos)}{P(w|neg)}$$

word
sentiment

$$\lambda(w) = \log \frac{P(w|Pos)}{P(w|neg)}$$

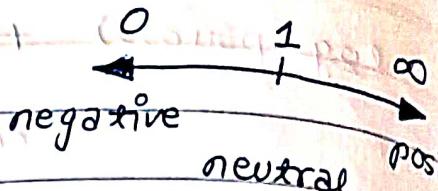
I am happy Because I am learning

$\downarrow \downarrow \downarrow \downarrow \downarrow \downarrow \downarrow$

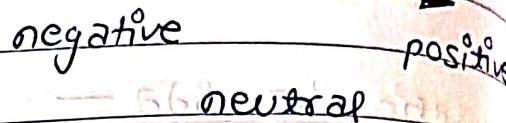
$$\rightarrow 0 + 0 + 2.2 + 0 + 0 + 0 + 1.1$$

$$\Rightarrow 3.3$$

$$\sum_{i=1}^m \frac{P(w_i^o | pos)}{P(w_i^o | neg)} > 1$$

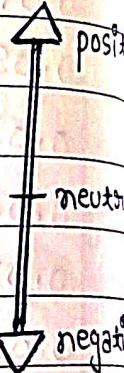


$$\sum_{i=1}^m \log \frac{P(w_i^o | pos)}{P(w_i^o | neg)} > 0$$



Tweet sentiment

$$\log \prod_{i=1}^m \text{ratio}(w_i^o) = \sum_{i=1}^m \lambda(w_i^o) > 0$$



Training Naive Bayes

p0: Collect and Annotate Corpus

positive tweets

: am Happy because I am Learning NLP
 : am happy, not sad. @ NLP

negative tweets

I am sad, I am not learning NLP
 I am sad, not happy !!

p1: Data-Preprocessing

- Lowercase
- Remove punctuation, urls, names
- Remove stopwords
- Stemming
- Tokenize sentences

positive tweets

[happy, Because, learn, NLP]
 [happy, not, sad]

negative tweets

[sad, not, learn, NLP]
 [sad, not, happy].

p2: word-count

ice, you have a clean corpus of processed tweets, we will start by computing the vocabulary for each word in class,

$\text{freq}(w, \text{class})$

word	Pos	Neg
happi	2	1
because	1	0
Learn	1	1
NLP	1	1
sad	1	2
not	1	2
N_{class}	7	7

Step 3% $P(CN/\text{class})$

$$\sqrt{\text{class}} = 6$$

$\text{freq}(w, \text{class}) + 1$

$$N_{\text{class}} + \sqrt{\text{class}}$$

Step 4%
Get Lambda(λ)

$$\lambda(w) = \log \frac{p(w/\text{POS})}{p(w/\text{neg})}$$

word	Pos	Neg
happi	0.23	0.15
because	0.15	0.07
Learn	0.08	0.08
NLP	0.08	0.08
sad	0.08	0.17
not	0.08	0.17

word	pos	neg	λ
happi	0.23	0.15	0.43
because	0.15	0.07	0.6
Learn	0.08	0.08	0
NLP	0.08	0.08	0
sad	0.08	0.17	-0.75
not	0.08	0.17	-0.75

Step 5: Get the log prior

D_{pos} = Nos of positive tweets

D_{neg} = Nos of negative tweets

$$\log \text{prior} = \log \frac{D_{pos}}{D_{neg}}$$

If dataset is balanced, $D_{pos} = D_{neg}$ and $\log \text{prior} = 0$

Testing Naïve Bayes →

- Predict using a Naïve Bayes Model

- Using your validation set to compute model dev.

word	λ	log-likelihood dictionary
I	-0.01	$\lambda(w) = \log \frac{P(w pos)}{P(w neg)}$
the	-0.01	
happy	0.63	
because	0.01	$\log \text{prior} = \log \frac{D_{\text{pos}}}{D_{\text{neg}}} = 0$
pass	0.5	
NLP	0	
sad	-0.75	• Tweet :
not	-0.75	(check)

[I, pass, the, NLP, interview]
 ↓ ↓ ↓ ↓ ↓

$$\text{score} = -0.01 + 0.5 - 0.01 + 0 + 0 + \log \text{prior}$$

$$\Rightarrow 0.48$$

$\text{pred} = \text{score} > 0$ (positive)

• Xval Yval λ logprior

$\text{score} = \text{predict}(X\text{val}, \lambda, \log \text{prior})$

$\text{pred} = \text{score} > 0$

$$\begin{bmatrix} 0.5 \\ -1 \\ 1.3 \\ \vdots \\ \vdots \\ \text{Score}_m \end{bmatrix} > 0 \Rightarrow \begin{bmatrix} 0.5 > 0 \\ -1 > 0 \\ 1.3 > 0 \\ \vdots \\ \vdots \\ \text{Score}_m > 0 \end{bmatrix}$$

produces a vector populated with 0's & 1's indicating if the predicted sentiment is negative or positive.

$$\begin{bmatrix} 0 \\ \frac{1}{2} \\ 1 \\ \frac{1}{2} \\ 0 \\ \vdots \\ \text{predm} \end{bmatrix} = \begin{bmatrix} 0 \\ \frac{1}{2} \\ 1 \\ \frac{1}{2} \\ 0 \\ \vdots \\ \varphi_{\text{valm}} \end{bmatrix}$$

$$\begin{array}{l} 0-0 \rightarrow 1 \\ 1-0 \rightarrow 0 \\ 1-1 \rightarrow 1 \end{array}$$

$$\frac{1}{m} \sum_{i=1}^m (\text{pred}_i = \varphi_{\text{val}_i}) \Rightarrow \begin{bmatrix} 1 \\ 0 \\ 1 \\ \vdots \\ \text{predm} = \varphi_{\text{valm}} \end{bmatrix}$$

summary • X_{val} , φ_{val} → performance on unseen data

predict using λ and logprior for each new tweet

$$\text{Accuracy} \rightarrow \frac{1}{m} \sum_{i=1}^m (\text{pred}_i = \varphi_{\text{val}_i})$$

what about words that don't appear in $\lambda(\text{cw})$?

applications of Naïve Bayes:

$$P(\text{pos} | \text{tweet}) = P(\text{pos}) P(\text{tweet} | \text{pos})$$

(1) Twitter
Sentiment Analysis

$$P(\text{neg} | \text{tweet}) = P(\text{neg}) P(\text{tweet} | \text{neg})$$

$$\frac{P(\text{pos} | \text{tweet})}{P(\text{neg} | \text{tweet})} \approx \frac{P(\text{pos})}{P(\text{neg})} \prod_{i=1}^m \frac{P(w_i | \text{pos})}{P(w_i | \text{neg})}$$

priors likelihood

estimating the probability for each class by using
the joint probability of the words in classes

② Identification system

Author identification: $P(C \text{ "shakespeare" } | \text{Book})$
 $P(C \text{ "meadley" } | \text{Book})$

If you have two large corpora, each written by two different authors, you could train a model to recognize whether a new document was written by one or another.

③ Word disambiguation

Consider that you have only two possible interpretations of a given word with text.

Let, say you don't know if the word Bank in reading is referring to the Bank of a river or to financial institution.

To disambiguate your word calculate the score of documents, given that it refers to each one of possible meanings.

In this case, if the text refers to the concept of money instead of the concept of river, then the score will be bigger than one.

$$P(\text{river} | \text{text})$$

$$P(\text{money} | \text{text})$$