

Positive tweets

I am happy because I am Learning NLP

I am happy, not sad.

Negative tweets

I am sad, I am not Learning NLP

I am sad, not happy.

Frequency table

word	Pos	Neg
I	3	3
am	3	3
happy	2	1
because	1	0
Learning	1	1
NLP	1	1
sad	1	2
not	1	2
N class	13	12

word	Pos	Neg
I	$\frac{3}{13} = 0.23$	0.25
am	0.23	0.25
happy	0.15	0.08
because	0.08	0
Learning	0.08	0.08
NLP	0.08	0.08
sad	0.08	0.17
not	0.08	0.17

I am happy today; I am learning

$$\prod_{i=1}^m \frac{P(W_i | \text{Pos})}{P(W_i | \text{neg})} = \frac{0.23}{0.25} \frac{0.23}{0.25} \frac{0.15}{0.08} \cdot 1 \cdot \frac{0.23}{0.25} \frac{0.23}{0.25} \frac{0.08}{0.08}$$

$$= 1.34 > 1 \quad \text{positive!}$$

Laplacian Smoothing.

usually we compute

$$P(w_i | \text{class}) = \frac{\text{freq}(w_i, \text{class})}{N_{\text{class}}}, \quad \text{class} \in \{\text{Positive}, \text{Negative}\}$$

However, if a word does not appear in the training, it gets a probability of 0. to fix this, we add smoothing

$$P(w_i | \text{class}) = \frac{\text{freq}(w_i, \text{class}) + 1}{N_{\text{class}} + V}$$

V : number of unique words in vocabulary

N : number of total frequency in each class

Likelihood ratio

$$\begin{aligned} \text{For a word } w_i, \quad \text{ratio} = (w_i) &= \frac{P(w_i | \text{Pos})}{P(w_i | \text{Neg})} \\ &\approx \frac{\text{freq}(w_i, 1) + 1}{\text{freq}(w_i, 0) + 1} \end{aligned}$$

For N_{positive}
 $\approx N_{\text{negative}}$

To infer sentiment of a tweet, we can compute

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

Log likelihood

$$\log \frac{P(\text{Pos})}{P(\text{neg})} + \sum_{i=1}^m \log \frac{P(w_i | \text{Pos})}{P(w_i | \text{neg})} \quad \left\{ \begin{array}{ll} > 0 & \text{positive} \\ < 0 & \text{negative} \end{array} \right.$$