

II Naive Bayes (assumption: words in sentence are INDEPENDENCE)

$$\textcircled{1} \quad P(X|Y) = P(Y|X) \frac{P(X)}{P(Y)} \quad (\text{Bayes Rule})$$

\textcircled{2} Overview Steps

Raw data
 $y=1: \text{pos tweet}$
 $=0: \text{neg } \uparrow$

Tweets	train-x	train-y
T ₁	1	
T ₂	0	
T ₃	1	
T _m	?	

+ remove stop words, punctuations, stemming

Tweets	y
[list of words] ⁽ⁿ⁾	1
↑ ↑ ↑ (12)	0
↑ ↑ (N)	1
⋮	⋮

Frequency (Count)

Words	Pos	Neg
Word ₁	9	80
Word ₂	100	0
⋮	⋮	⋮
Word _N	⋮	⋮

Sum V Words N pos N neg

unique words in vocab

- likelihood (of the same word w.r.t different classes)

$$\frac{P(\text{word}_{\text{pos}})}{P(\text{word}_{\text{neg}})} \begin{cases} > 1: \text{Pos} \\ = 1: \text{neutral} \\ < 1: \text{Neg} \end{cases}$$

$$\text{Ex. } \frac{P(\text{'happi'}_{\text{pos}})}{P(\text{'happi'}_{\text{neg}})} = \frac{0.02}{0.001} = 20 > 1 \Rightarrow \text{Positive class.}$$

$$-\log \text{likelihood} \quad \log \left(\frac{P(\text{word}_{\text{pos}})}{P(\text{word}_{\text{neg}})} \right) = \log P(\text{word}_{\text{pos}}) - \log P(\text{word}_{\text{neg}}) \begin{cases} \geq 0: \text{Positive} \\ < 0: \text{Negative} \end{cases}$$

- Prior (of documents)

D: number of documents (in this case tweets)
 D_{Pos}: up up up clamed in database as Pos class (ex. Pos tweets)

D_{Neg}: up up up up Neg class (ex. Neg tweets)

$$\Rightarrow \text{Probability } P(D_{\text{pos}}) = D_{\text{pos}}/D ; \quad P(D_{\text{neg}}) = \frac{D_{\text{neg}}}{D}$$

$$\text{Prior probability} = \frac{P(D_{\text{pos}})}{P(D_{\text{neg}})} = \frac{D_{\text{pos}}}{D_{\text{neg}}} \Rightarrow \log \text{Prior} = \log D_{\text{pos}} - \log D_{\text{neg}}$$

- Probability of Pos and Neg words

$$P(\text{word}_{\text{pos}}) = \frac{\text{freq}_{\text{pos}} + 1}{N_{\text{pos}} + V}$$

$$P(\text{word}_{\text{neg}}) = \frac{\text{freq}_{\text{neg}} + 1}{N_{\text{neg}} + V}$$

where freq dictionary pos class

$$\text{freq} = \{("happi", 1): 99, ("happi", 0): 8, \dots\} \quad \text{neg. class.}$$

Python def train(freqs, train-x, train-y):

:
 return loglikelihood: dict (ex. {'word1': 1.7, 'word2': -0.3, ...})
 logprior: float (ex. 0.0)

def predict(tweet, logprior, loglikelihood):

words = process_tweet(tweet)

$$p = \text{logprior} + \sum_i^N \text{loglikelihood}[word]$$

p ≥ 0 : positive

p < 0 : negative