

## WEEK 2

### SENTIMENT ANALYSIS WITH NAÏVE BAYES.

#### • Probability and Bayes' Rule

$A \rightarrow$  Positive tweet

$$P(A) = P(\text{Positive}) = N_{\text{pos}} / N = 0.65$$

$$P(\text{Negative}) = 1 - P(\text{Positive}) = 0.35$$

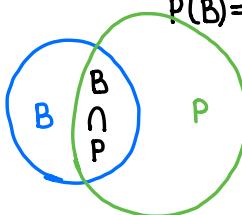
Tweets containing the word "happy"

$B \rightarrow$  tweet contains "happy"

$$P(B) = P(\text{happy}) = N_{\text{happy}} / N$$

$$P(B) = 4/20 = 0.2$$

$$P(B \cap P) = P(B, P) = \frac{3}{20} = 0.15$$



$P(P|B) = P(\text{Positive} | \text{"happy"}) = 3/4 = 0.75 \rightarrow 75\% \text{ of likelihood of being positive if tweet contains the word "happy".}$

$P(B|P) = P(\text{"happy"} | \text{Positive}) = 3/13 = 0.231 \rightarrow 23.1\% \text{ of likelihood of containing the word "happy" if the tweet is positive.}$

$$P(\text{Positive} | \text{"happy"}) = \frac{P(\text{Positive} \cap \text{"happy"})}{P(\text{"happy"})}$$

$$P(\text{"happy"} | \text{Positive}) = \frac{P(\text{"happy"} \cap \text{Positive})}{P(\text{Positive})}$$

**BAYES'  
RULE**

$$P(\text{Positive} | \text{"happy"}) = P(\text{"happy"} | \text{Positive}) \times \frac{P(\text{Positive})}{P(\text{"happy"})}$$

$$P(X|Y) = P(Y|X) \frac{P(X)}{P(Y)}$$

## • Naïve Bayes

this method makes the assumption that the features we are using for classification are all independent, which in reality is rarely the case.

$$P(w_i | \text{class})$$

word	Pos	Neg		word	Pos	Neg	
I	3	3		I	0.24	0.24	$\rightarrow 3/12$
am	3	3		am	0.24	0.24	
happy	2	1		happy	0.15	0.08	
because	1	0		because	0.08	0	
learning	1	1		learning	0.08	0.08	
NLP	1	1		NLP	0.08	0.08	
sad	1	2		sad	0.08	0.15	
not	1	2		not	0.08	0.15	
<u>Nclass</u>		13	13				

Naïve Bayes inference condition rule for binary classification

"I am happy today; I am learning"

smooth the values.

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

$\frac{0.2}{0.2} \times \frac{0.2}{0.2} \times \frac{0.14}{0.10} \times \frac{0.2}{0.2} \times \frac{0.2}{0.2} \times \frac{0.1}{0.1} \leftarrow = 1.4 > 1 \text{ (positive class).}$

## • Laplacian Smoothing

(to avoid probabilities being zero)

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

(smoothing)

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

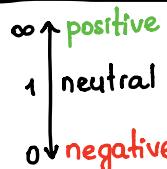
V: # of unique words in a vocabulary.

word	Pos	Neg
I	3	3
:	:	:
Nclass	13	13

$$P(I|Pos) = \frac{3+1}{13+8} = 0.19$$

## • Log Likelihood.

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



Naïve Bayes' inference:

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

prior ratio      likelihood

- products bring risk of underflow
- $\log(a \cdot b) = \log(a) + \log(b)$
- $\log \left[ \frac{P(\text{pos})}{P(\text{neg})} \prod_{i=1}^m \frac{P(w_i | \text{pos})}{P(w_i | \text{neg})} \right]$
- $\log \frac{P(\text{pos})}{P(\text{neg})} + \sum_{i=1}^m \log \frac{P(w_i | \text{pos})}{P(w_i | \text{neg})}$

## - calculating Lambda

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

word	Pos	Neg
I	0.05	0.05
am	0.04	0.04
happy	0.09	0.01

$$\lambda \rightarrow \log \frac{0.05}{0.05} = \log 1 = 0$$

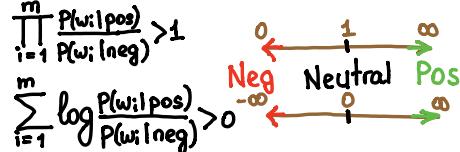
$$0 \rightarrow \log \frac{0.05}{0.01} = \log 5 = 2.2$$

word sentiment

$$\left\{ \begin{array}{l} \text{ratio}(w) = \frac{P(w | \text{pos})}{P(w | \text{neg})} \\ \lambda(w) = \log \frac{P(w | \text{pos})}{P(w | \text{neg})} \end{array} \right.$$

tweet: "I am happy because I am learning."

$$\log \text{likelihood} = 0 + 0 + 2.2 + 0 + 0 + 0 + 1.1 = 3.3 > 0$$



## • Training Naïve Bayes

0. Collect and annotate corpus.

1. Preprocess (lowercase, remove punctuation, urls, names, stop words, stemming, tokenize sentences)

2. Word count

3.  $P(w|class)$  by using Laplacian smoothing  $\frac{freq(w,class)+1}{N_{class}+V_{class}}$
4. Get lambda  $\lambda(w) = \log \frac{P(w|pos)}{P(w|neg)}$
5. Get the log prior:  $\text{logprior} = \log \frac{D_{pos}}{D_{neg}} \rightarrow \# \text{ of positive tweets}$

## • Testing Naïve Bayes.

**Predict**

- log-likelihood dictionary  $\lambda(w) = \log \frac{P(w|pos)}{P(w|neg)}$   
- logprior =  $\log \frac{D_{pos}}{D_{neg}}$   
- "I passed the NLP interview"  
 ↓  
 [I, pass, the, NLP, interview]  
 ↓  
 $\text{score} = -0.01 + 0.5 - 0.01 + 0 + \text{logprior}^0 = 0.48 > 0 \rightarrow \text{positive.}$

•  $X_{val}, Y_{val}, \lambda, \text{logprior}$

$\text{score} = \underline{\text{predict}}(X_{val}, \lambda, \text{logprior})$

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

$\text{accuracy} = \frac{1}{m} \sum_{i=1}^m (\text{pred}_i == Y_{val,i})$

(note: the words that do not appear in  $\lambda(w)$  are treated as neutral words)

## • Applications of Naïve Bayes

$P(\text{pos}|\text{tweet}) \approx P(\text{pos}) P(\text{tweet}|\text{pos}) \rightarrow \text{Author identification}$

$P(\text{neg}|\text{tweet}) \approx P(\text{neg}) P(\text{tweet}|\text{neg}) \rightarrow \text{Spam filtering}$

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

priors      likelihood

$\rightarrow \text{Information retrieval}$   
 $\rightarrow \text{Word disambiguation}$

## • Naïve Bayes Assumptions.

- Independence btw the predictors or features.
- Relative frequency in corpus

- Independence assumption could lead us to under or over estimate the conditional probabilities of individual words.
- Relative frequency assumption relies on the distribution of the training data sets: proportion of positive and negative classes. → affect the model.

## • Error Analysis

- removing punctuation and stop words → semantic may lost
- word order → affects the meaning of a sentence
- adversarial attacks → confuse naive Bayes models.

"My beloved grandmother :("

"I am happy because I did **not** go." ↔ "I am **not** happy because I did go."