

احتمال =

$$P(A) = \frac{n_A}{n_{total}}$$

احتمال وقوع A به شرطی که B اتفاق افتاده باشد

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

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

$$P(A \cap B) = P(A|B) P(B) = P(B|A) P(A)$$

$$P(A|B) = \frac{P(B|A) P(A)}{P(B)} \rightarrow \text{Bayse Rule}$$

کلماتی که در کلاس + و - اعمال تقریباً یکسانی دارند معمولاً احساسی به هم وابسته اند. نکته:  $\rightarrow$

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

تعداد کلمات کلاس

نویسندگان الان:  $\text{class} \in \{+, -\}$

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

nominator

demonator

vocab

laplacian smoothing

اگر کلمه در tweet نباشد منفرجه. برای جلوگیری از این منفرج از تکنیک smoothing استفاده میکنیم. مرتبه فرمول رو بدو (چون منفرجه بی نهایت + و - رو میگیره)

$$\text{ratio}(w_i) = \frac{P(w_i | +)}{P(w_i | -)} \approx \frac{\text{freq}(w_i, 1) + 1}{\text{freq}(w_i, 0) + 1}$$

منفرجه

صفر

یک

نهایت

نهایت

غریب میکنیم تا کلمات متبعا  $\text{ratio} = 1$  به این برسیم naive bayesian

$$\frac{P(+)}{P(-)} \prod_{i=1}^m \frac{P(w_i | +)}{P(w_i | -)} > 1$$

زمانی که تعداد m زیاد باشد دچار underflow می شویم و یا لوگ گرفته شدن صفر جلوگیری از گرفت

$$\log \frac{P(+)}{P(-)} + \sum_{i=1}^m \log \frac{P(w_i | +)}{P(w_i | -)}$$

log prior

log likelihood

$= \chi(w)$

$\log(1) = 0$

منفرجه

صفر

یک

نهایت

To train your naïve Bayes classifier, you have to perform the following steps:

**1) Get or annotate a dataset with positive and negative tweets**

**2) Preprocess the tweets:  $\text{process\_tweet}(\text{tweet}) \rightarrow [w_1, w_2, w_3, \dots]$ :**

- Lowercase
- Remove punctuation, urls, names
- Remove stop words
- Stemming
- Tokenize sentences

**3) Compute  $\text{freq}(w, \text{class})$ :**

Positive tweets
[happi, because, learn, NLP]
[happi, not, sad]
Negative tweets
[sad, not, learn, NLP]
[sad, not, happi]

Step 2:  
Word  
count

word	Pos	Neg
happi	2	1
because	1	0
learn	1	1
NLP	1	1
sad	1	2
not	1	2
$N_{\text{class}}$	7	7

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

**4) Get  $P(w|pos), P(w|neg)$**

You can use the table above to compute the probabilities.

**5) Get  $\lambda(w)$**

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

مربوط به کلاس = نسبت در یک کلاس

**6) Compute  $\log\text{prior} = \log(P(pos)/P(neg))$**

$\log\text{prior} = \log \frac{D_{pos}}{D_{neg}}$ , where  $D_{pos}$  and  $D_{neg}$  correspond to the number of positive and negative documents respectively.

- log-likelihood dictionary  $\lambda(w) = \log \frac{P(w|pos)}{P(w|neg)}$

- $logprior = \log \frac{D_{pos}}{D_{neg}} = 0$

- Tweet: [I, pass, the NLP interview] 🍀

$$score = -0.01 + 0.5 - 0.01 + 0 + logprior = 0.48$$

$$pred = score > 0$$

word	$\lambda$
I	-0.01
the	-0.01
happi	0.63
because	0.01
pass	0.5
NLP	0
sad	-0.75
not	-0.75

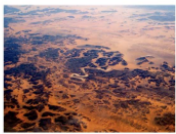
The example above shows how you can make a prediction given your  $\lambda$  dictionary. In this example the  $logprior$  is 0 because we have the same amount of positive and negative documents (i.e.  $\log 1 = 0$ ).

There are many applications of naive Bayes including:

- Author identification
- Spam filtering  $\rightarrow \frac{P(\text{spam} | \text{email})}{P(\text{non-spam} | \text{email})}$
- Information retrieval  $\rightarrow P(\text{document}_k | \text{query}) \propto \prod_{i=1}^{|\text{query}|} (query_i | \text{document}_k)$
- Word disambiguation  $\rightarrow \frac{P(\text{river} | \text{text})}{P(\text{memory} | \text{text})}$  (with Urdu text: "کے bank پر تیرا" vs "پانی کے کنارے پر")

This method is usually used as a simple baseline. It also really fast.

Naïve Bayes makes the independence assumption and is affected by the word frequencies in the corpus. For example, if you had the following



"It is sunny and hot in the Sahara desert."



"It's always cold and snowy in \_\_\_."

In the first image, you can see the word sunny and hot tend to depend on each other and are correlated to a certain extent with the word "desert". Naive Bayes assumes independence throughout. Furthermore, if you were to fill in the sentence on the right, this naive model will assign equal weight to the words "spring, summer, fall, winter".

Relative frequencies in corpus



On Twitter, there are usually more positive tweets than negative ones. However, some "clean" datasets you may find are artificially balanced to have to the same amount of positive and negative tweets. Just keep in mind, that in the real world, the data could be much noisier.

There are several mistakes that could cause you to misclassify an example or a tweet. For example,

- Removing punctuation

- Removing words

**Tweet:** This is not good, because your attitude is not even close to being nice.

**processed\_tweet:** [good, attitude, close, nice]

**Tweet:** My beloved grandmother :(

**processed\_tweet:** [belov, grandmoth]

- Word order

**Tweet:** I am happy because I did not go.



**Tweet:** I am not happy because I did go.



- Adversarial attacks

These include sarcasm, irony, euphemisms.