

Naive Bayes

- ① Intro
- ② Spam data
- ③ Mathematical Intuition
- ④ Naive Bayes Assumption
- ⑤ Train / Test time complexity
- ⑥ Space complexity

Google

Mail 1 : " I, Nigerian prince need your help.
Send money "

Mail 2 : " Meeting scheduled at 8pm. Kindly revert "

lottery , Million dollars , Jackpot

Sentiment

"terrible" , "bad" , "pathetic" → (-ve)

"great" , "good" → (+ve)

"This is shivank"



(1) Tokenisation ["This", "is", "shivank"]

good

Good

2 Lowercase

⇒ I, Shivank Agrawal wants to teach!



removing all punctuation ⇒ regex

(4) text → remove stopwords

"I", "and", "the"

(5)

grametical errors ←

Lema

stunning

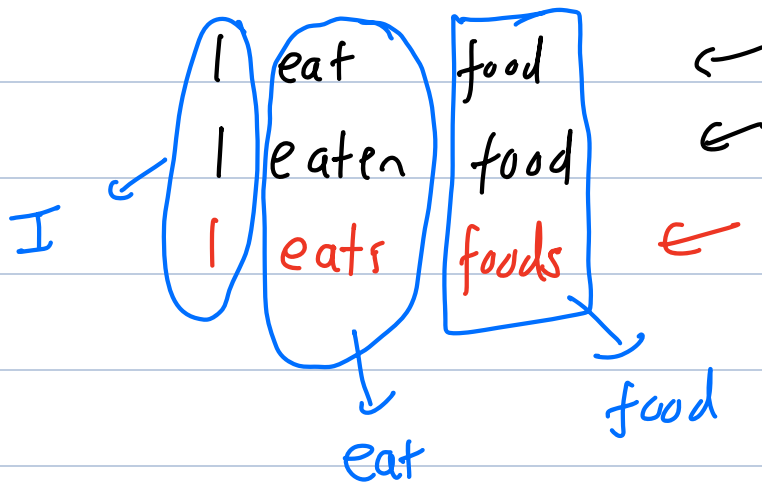
eating
eat

money

play
plays
playing

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

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



Mathematical Notation

email, \rightarrow spam or not spam

email, \rightarrow 0/1

$$P(y=1 | \text{email}_i)$$

$$P(y=0 | \text{email}_i)$$

$$P(y=0 \mid w_1^i, w_2^i, w_3^i, w_4^i, \dots, w_d^i)$$

$$P(y=1 \mid w_1^i, w_2^i, w_3^i, w_4^i, \dots, w_d^i)$$

Bayes Theorem

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

$$P(y=1 \mid w_1^i, w_2^i, w_3^i, w_4^i, \dots, w_d^i)$$

$$P_1(y=1 \mid x) = \frac{P(x \mid y=1) * P(y=1)}{P(x)}$$

$$P_2(y=0 \mid x) = \frac{P(x \mid y=0) * P(y=0)}{P(x)}$$

$$\frac{20}{35}$$

$$\frac{40}{31}$$

$$y=1$$

$$P(y=1|w_1, \dots, w_d) = \frac{P(w_1, \dots, w_d|y=1) \cdot P(y=1)}{\cancel{K}}$$

$$P(y=0|w_1, \dots, w_d) = \frac{P(w_1, \dots, w_d|y=0) \cdot P(y=0)}{\cancel{K}}$$

$$P(DF) = 0.01$$

$$P(smoke) = 0.1$$

$$P(smoke | DF) = 0.9$$

$$P(DF | smoke) \Rightarrow ? \Rightarrow \frac{P(DF) \times P(smoke|DF)}{P(smoke)}$$

$$\Rightarrow \frac{0.01 \times 0.9}{0.1}$$

$$\Rightarrow 0.09$$

$$P(y=1|w_1, \dots, w_d) = \frac{P(w_1, \dots, w_d | y=1) \cdot P(y=1)}{\cancel{K}}$$

$$P(y=0|w_1, \dots, w_d) = \frac{P(w_1, \dots, w_d | y=0) \cdot P(y=0)}{\cancel{K}}$$

$$P(y=1) \rightarrow \frac{\# \text{ train points with } y_i=1}{\text{total } \# \text{ of train pt.}}$$

$$P(y=0) \rightarrow \frac{\# \text{ of train point with } y=0}{\text{total } \# \text{ of train pt.}}$$

prince | spam

viagra | spam

$$P(w_1, w_2 | \text{spam}) \approx P(w_1 | \text{spam}) \cdot P(w_2 | \text{spam})$$

\downarrow \downarrow
 prince viagra

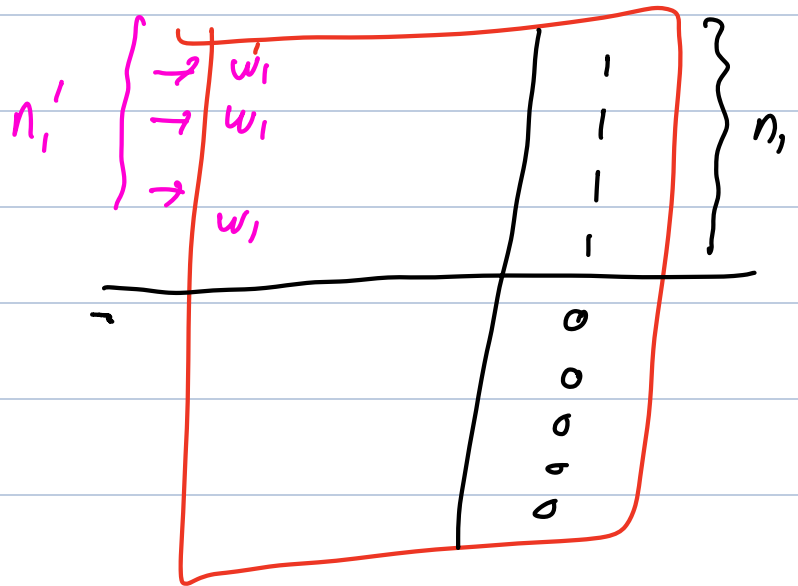
$$P(w_1, w_2 | y=1) = P(w_1 | y=1) * P(w_2 | y=1)$$

Naive Assumption

$$P(w_1, w_2, w_3, \dots, w_d | y=1)$$

$$= P(w_1 | y=1) \cdot P(w_2 | y=1) \cdot \dots \cdot P(w_d | y=1)$$

$$P(w_1 | y=1) = \frac{n_1'}{n_1}$$



$$P(y=1 | \text{text}) \approx P(w_1, w_2, \dots, w_d | y=1) \cdot \frac{P(y=1)}{K}$$

$$= P(w_1 | y=1) * P(w_2 | y=1) \cdot \dots \cdot P(w_d | y=1) \cdot \frac{P(y=1)}{K}$$

$$\approx \prod_{i=1}^d P(w_i | y=1) \cdot \frac{P(y=1)}{K}$$

$$\begin{aligned}
 P(y=1 | \text{text}) &\approx \prod_{i=1}^n \underbrace{p(w_i | y=1)}_{\text{likelihood}} \cdot \underbrace{P(y=1)}_{\text{class prior}} / K \\
 P(y=0 | \text{text}) &\approx \prod_{i=1}^n p(w_i | y=0) \cdot P(y=0) / K
 \end{aligned}$$

compare

This is Shivank

Shivank teaching python

This	Shivank	teaching	Python
1	1	0	0
0	1	1	1

Train time Complexity

Likelihood

$$\rightarrow P(y=1)$$

$$\rightarrow P(y=0)$$

$$O(nd+2)$$

$$\rightarrow p(w_i | y_i=1) \rightarrow d$$

$$\rightarrow O(nd)$$

$$\rightarrow p(w_i | y_i=0) \rightarrow d$$

$$O(nd)$$

$$\Downarrow$$

$$O(2nd)$$

$$\downarrow$$

$$n''$$

$$\cancel{n^2}$$

$$p(\tilde{w}_1, \tilde{w}_2 | y=1)$$

$$\approx p(w_1 | y=1) \cdot p(w_2 | y=1)$$

Test time complexity

$$w_1 \ w_2 \ w_3 \ w_4 \ - \ - \ - \ w_k$$

$$O(k)$$

Space Complexity

$$P(y=0) \\ P(y=1)$$

likelihood : $2d$

$$2d+2$$

$$O(2d) \\ \Downarrow \\ O(d)$$

Booth: $q: 16 \times 17$

\Rightarrow Laplace Smoothing

This is Shivank

Shivank teaching python

This	Shivank	teach	Python
1	1	0	0
0	1	1	1

Ahant teaches Python

$$x_q = w_1 w_2 w_6 w_7 w'$$

$$w' \in \{w_1, w_2, \dots, w_d\}$$

$$P(y=1 | w_1 w_2 w_6 w_7 w') = P(y=1) \times P(w_1 | y=1) \dots P(w' | y=1)$$

$$\frac{\# \text{ word } w' \text{ is part of } y=1}{\# \text{ values in } y=1} \approx 0$$

$$P(y=0 | w_1 w_2 \dots w') = 0$$

Laplace smoothing or Additive smoothing

$$P(w_j | y=1) = \frac{n_j + \alpha}{n_1 + \alpha \cdot C}$$

↑
distinct values
 w_j can take

100 text msg

$$y=1$$

$$d=1$$

$$c=2$$

$$P(w' | y=1) = \frac{0 + 1}{100 + (1 \times 2)} = \frac{1}{102} \approx 0.0098$$

$$\Rightarrow \frac{\cancel{0} + 1}{\cancel{100} + 1 \times 2} = \frac{1}{2}$$

HyperParam Tuning

Train Time

$$n_{j1} = 10$$

$$n_j = 100$$

Case 1

$$P(w_j | y=1) = \frac{n_{j1} + \alpha}{n_{j1} + 2\alpha} = \frac{10 + 1}{100 + 2} \approx 0.1$$

$$\alpha = 1$$

$$\approx \frac{n_{j1}}{n_j} = \frac{10}{100} \approx 0.1$$

$$\alpha = 10,000$$

$$P(w_j | y=1) = \frac{10 + 10000}{100 + 20000} \approx \frac{1}{2}$$

$\alpha \uparrow \rightarrow$ underfit

$$P(y=1) = 0.7$$

$$P(y=0) = 0.3$$

$$P(y=1 | w_1 \dots w_d) = P(y=1) \pi \dots$$

$$P(y=0 | w_1 \dots w_d) = P(y=0) \pi \dots$$

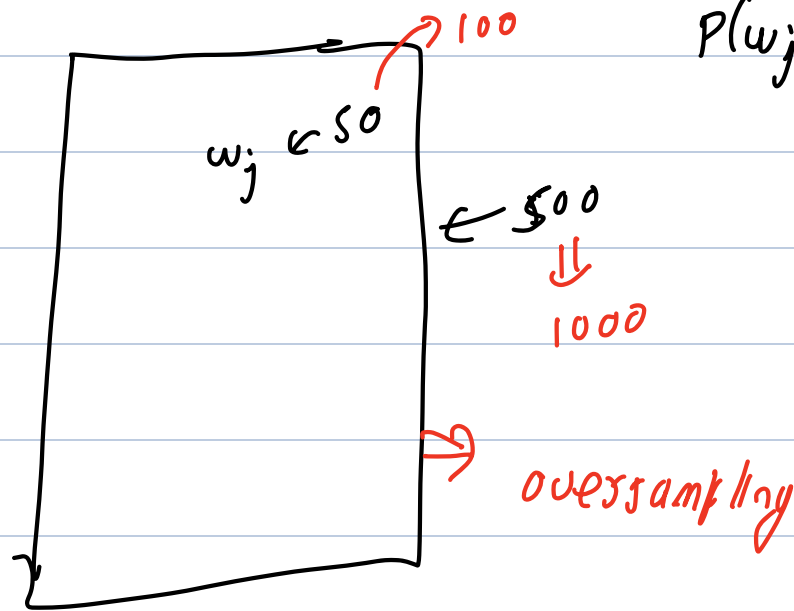
$\alpha \uparrow \rightarrow \text{underfit}$
 $\alpha \downarrow \rightarrow \text{overfit}$

$$\left. \begin{array}{l} P(y=1) = 0.7 \\ P(y=0) = 0.3 \end{array} \right\} \text{ imbalance}$$

\Downarrow rebalance

$$\parallel P(y=1 | \text{text}) \propto \frac{P(y=1)}{\cancel{0.7} \quad 0.5} \cdot \frac{\Pi \text{ likelihood}}{0.1} \approx 0.07$$

$$P(y=0 | \text{text}) \propto \frac{P(y=0)}{\cancel{0.3} \quad 0.6} \cdot \frac{\Pi \text{ likelihood}}{0.2} \approx 0.06$$



$$P(w_j | y=0) = \frac{50 + \alpha}{500 + 2\alpha}$$

$$\approx 10\%$$

$$P(w_j | y=1) = \frac{100 + \alpha}{1000 + 2\alpha} \approx 10\%$$

Break: 8:12 am

< 1

$P(y=0)$ ————
 $P(y=1)$ ————

$$0.2 \times 0.3 \times 0.4 \times 0.2 \times 0.2$$

$$\Rightarrow 0.000000002 \approx$$

↓
Floating Number

Underflow Problem

$\Rightarrow \prod$

$$\log(\prod w_1 \cdot w_2 \cdot w_3) \Rightarrow \log(w_1) + \log(w_2) + \log(w_3)$$

$$\log(p(y=1 | w_1 w_2 \dots w_d)) = \log(p(y=1)) + \sum_{j=1}^d \log(w_j | y=1)$$

$$\log(ab) = \log a + \log b$$

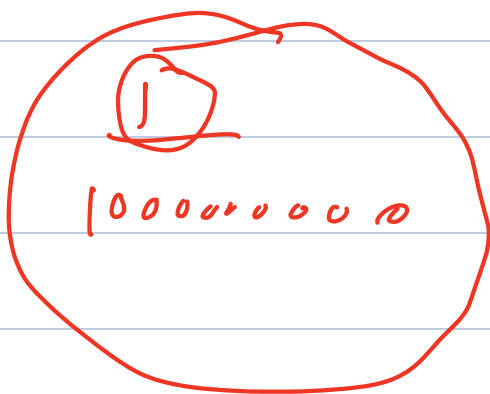
\Rightarrow Feature Importance

$$P(w_1 | y=1)$$

$$P(w_2 | y=1)$$

\Rightarrow

$$= w_1 w_2 \rightarrow w' \in D_{\text{Train}}$$



Bernouli NB

Multinomial NB

This is Shivank and shivank
Shivank teaching python

This	Shivank	teach	Python
1	2	0	0
0	1	1	1

1