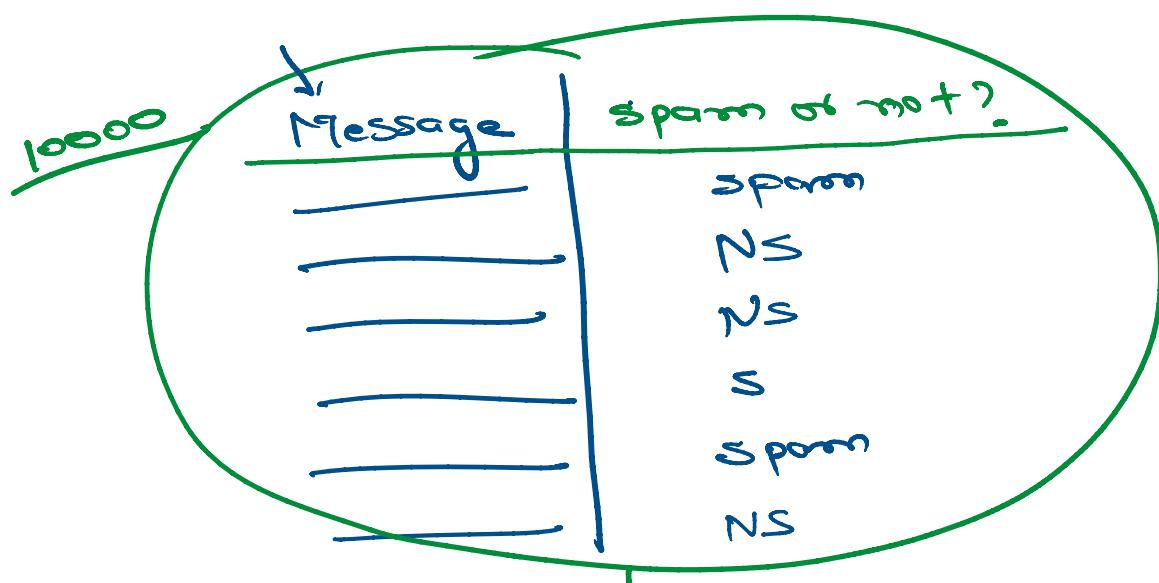


Based on Bayes Theorem (statistical concept)

Naive Bayes

Used specially for spam detection/classification



Toucanelli

Naive Bayes

Learn the patterns in spam and Non-spam msgs.

Future msgs and classify them on the basis of patterns we have learned.

Male

-

0

↑

Male	0
F	1
M	0
F	1
F	1
M	0

Example →

Can we <sup>have</sup> a quick meeting tomorrow at 5 PM?

↓  
Tokenization

(separate each word from the msg and make a list of words)

[can, we, have, <sup>a</sup>, quick, meeting, tomorrow,  
at, 5, PM]

↓  
Remove stop words (a, too, of, the, at etc)

↓  
Vectorization

(converting text to numbers)

Naïve Bayes →

$$P(C|x) = \frac{P(x|C) \times P(C)}{P(x)}$$

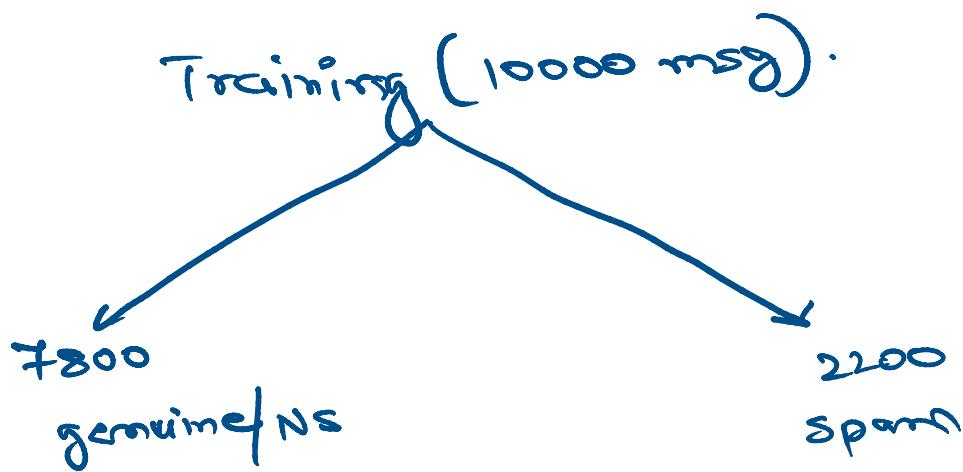
$P(x) \rightarrow$

$P(\text{spam}) | \text{con, we, have, quick, meeting, tomorrow, 5, PM}$

$$= P(\text{con}|\text{spam}) \times P(\text{we}|\text{spam}) \times P(\text{have}|\text{spam}) \times \dots \times P(\text{PM}|\text{spam}) \times P(\text{spam})$$

$$= 0.25 \times \dots \times \frac{2200}{10000}$$

$$= \underline{\underline{0.000156}}$$



$$P(NS) = \frac{7800}{10000}$$

$$= 0.78$$

$$= 78\%$$

$$P(\text{spam}) = \frac{2200}{10000}$$

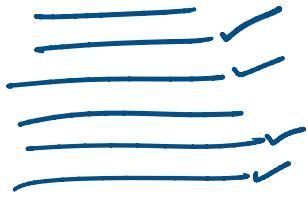
$$= 0.22$$

$$\approx 22\%$$

Let's take spam msgs in the training data.

2200

Training data  
---  $\rightarrow$  10,000


 $\frac{2200}{\text{Training Data}}$ 
 $P(\text{com}|\text{spam}) = \frac{576}{2200}$ 
 $= \underline{0.26}$

$$P(\text{we}|\text{spam}) = \frac{180}{2200}$$
 $=$

$$P(\text{Nonspam}|\text{com, we...}) = 1 - P(\text{spam}) - \dots$$
 $= 1 - 0.18$ 
 $= 0.82$ 
 $= \underline{82\%}$

Example →

Congratulation, you have a pre-approved loan of 50000 USD. No documentation required.

Fill the form to claim the loan.

$$\underline{P(\text{spam}|\text{congratulation, you, have, pre-approved, ... loan})} =$$

$$= P(\text{congratulation}|\text{spam})$$

our email  
Regular

✓ Regular