

Assignment - 1B

category	Email text
Not spam	"Hi there how are you?"
Not spam	"Meeting at 3 PM tomorrow"
Not spam	"Please send the report"
spam	"Win a free prize now!"
spam	"Get claim your discount today"
spam	"Limited time offer click here"
?	"Free meeting tomorrow"
?	"Claim your free prize"

Prior probabilities

$$P(\text{spam}) = 3/6 = 1/2$$

$$P(\text{not spam}) = 3/6 = 1/2$$

Total unique word in spam = 14

Total unique word in not spam = 14

Total vocabulary size = 28

using laplace smoothing ("Free meeting tomorrow")

$$P(\text{free}|\text{spam}) = \frac{1+1}{14+28} = \frac{2}{42}$$

from formula :-

$$P(w|s) = \frac{\text{Count}(w \text{ in } s) + 1}{\text{Total words in } s + \text{Vocabulary size}}$$

$$P(\text{meeting}|\text{spam}) = \frac{0+1}{14+28} = \frac{1}{42}$$

$$P(\text{tomorrow}|\text{spam}) = \frac{0+1}{14+28} = \frac{1}{42}$$

$P(\text{spam} | \text{free, meeting, tomorrow})$

$$\begin{aligned}
 &= P(\text{spam}) \times P(\text{free}|\text{not spam}) \times P(\text{meeting}|\text{not spam}) \times P(\text{tomorrow}|\text{not spam}) \\
 &= \frac{1}{2} \times \frac{2}{42} \times \frac{1}{42} \times \frac{1}{42} \\
 &\approx 0.0000185
 \end{aligned}$$

$$P(\text{free}|\text{not spam}) = \frac{0+1}{14+28} = \frac{1}{42}$$

$$P(\text{meeting}|\text{not spam}) = \frac{0+1}{14+28} = \frac{1}{42}$$

$$P(\text{tomorrow}|\text{not spam}) = \frac{1+1}{14+28} = \frac{2}{42}$$

$P(\text{not spam} | \text{free, meeting, tomorrow})$

$$\begin{aligned}
 &= P(\text{not spam}) \times P(\text{free}|\text{not spam}) \times P(\text{meeting}|\text{not spam}) \times \\
 &\quad P(\text{tomorrow}|\text{not spam})
 \end{aligned}$$

$$= \frac{1}{2} \times \frac{1}{42} \times \frac{2}{42} \times \frac{2}{42}$$

$$\approx 0.000027$$

Since $P(\text{spam} | \text{free, meeting, tomorrow}) < P(\text{not spam} | \text{free, meeting, tomorrow})$ this email is 'not spam'

$$\begin{aligned}
 \text{Normalization} &= \frac{0.000027}{0.000027 + 0.000013} \times 100\% = 67.5\% \text{ not spam}
 \end{aligned}$$

$$\text{Normalization (spam)} = \frac{0.000013}{0.000013 + 0.000027} \times 100\% \\ \approx 32.5\% \text{ spam}$$

2 For email "claim your free prize"

Using Laplace smoothing,

$$P(\text{claim} | \text{spam}) = \frac{1+1}{14+28} = \frac{2}{42}$$

$$P(\text{your} | \text{spam}) = \frac{1+1}{14+28} = \frac{2}{42}$$

$$P(\text{free} | \text{spam}) = \frac{1+1}{14+28} = \frac{2}{42}$$

$$P(\text{prize} | \text{spam}) = \frac{1+1}{14+28} = \frac{2}{42}$$

$$P(\text{claim} | \text{not spam}) = 0+1/14+28 = 1/42$$

$$P(\text{your} | \text{not spam}) = 0+1/14+28 = 1/42$$

$$P(\text{free} | \text{not spam}) = 0+1/14+28 = 1/42$$

$$P(\text{prize} | \text{not spam}) = 0+1/14+28 = 1/42$$

Now,

$$P(\text{spam} | \text{claim, your, free, prize}) \\ = P(\text{spam}) \times P(\text{claim} | \text{spam}) \times P(\text{your} | \text{spam}) \times P(\text{free} | \text{spam}) \\ \times P(\text{prize} | \text{spam})$$

$$= \frac{1}{2} \times \frac{2}{42} \times \frac{2}{42} \times \frac{2}{42} \times \frac{2}{42}$$

$$\approx 0.0000025$$

$P(\text{not spam} | \text{claim, your, free, prize})$

$$= P(\text{not spam}) \times P(\text{not claim} | \text{not spam}) \times P(\text{your} | \text{not spam}) \\ \times P(\text{free} | \text{not spam}) \times P(\text{prize} | \text{not spam})$$

$$= \frac{1}{2} \times \frac{1}{2} \times \frac{1}{2} \times \frac{1}{2} \times \frac{1}{2}$$

$$= 0.00000016$$

Since $P(\text{email claim, your, free, prize} | \text{spam}) > P(\text{not spam} | \text{claim, your, free, prize})$ so this email is "spam".

$$\text{Normalization (spam)} = \frac{0.00000025}{0.0000025 + 0.0000016}$$

$\approx 90.5\% \text{ spam}$

$$\text{Normalization (not spam)} = \frac{0.0000016}{0.0000016 + 0.0000025}$$

$\approx 5.7\% \text{ not spam}$

$\therefore \text{"free meeting tomorrow"} \rightarrow \text{not spam}$

$\therefore \text{"claim your free prize"} \rightarrow \text{spam}$