

category	email text
Not spam	"Hi there, how are you?"
Not spam	"Meeting at 3PM tomorrow"
Not spam	"Please send the report"
spam	"Win a free prize now!"
spam	"Claim your discount today"
spam	"Limited time offer: click here"
?	"Free meeting tomorrow" (To classify)
?	"Claim your free prize" (To classify)

Q1 Total unique words in spam: 14

Total unique words in notspam: 14

$$\text{Vocabulary Size} = 28$$

using smoothing:

Now, To classify "free meeting tomorrow"
we need to ~~do~~ find;

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

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

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

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

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

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

Prior probabilities

$$P(\text{spam}) = \frac{3}{6} = 0.5$$

$$P(\text{not spam}) = \frac{3}{6} = 0.5$$

$\left. \begin{array}{l} \{3 \text{ not spam and } 3 \text{ spam}\} \\ \text{emails, total 6} \end{array} \right\}$

Now,

$$\begin{aligned} P(\text{spam} | \text{free, meeting, tomorrow}) &\propto P(\text{spam}) \times P(\text{free} | \text{spam}) \\ &\quad P(\text{meeting} | \text{spam}) \times \\ &\quad P(\text{tomorrow} | \text{spam}) \end{aligned}$$

$$\approx 0.5 \times \frac{2}{42} \times \frac{1}{42} \times \frac{1}{42}$$

$$\approx 0.00001349$$

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

$$\approx 0.5 \times \frac{2}{42} \times \frac{1}{42} \times \frac{2}{42}$$

$$\approx 0.0002699$$

Since $P(\text{not spam}) > P(\text{spam})$, the email is not spam.

$$\text{Normalization} = \frac{0.0002699}{0.0002699 + 0.0001349} \times 100\%$$

$$= \frac{0.0002699}{0.0004048} \times 100\% = 66.67\% \text{ not spam}$$

② for email "claim your free prize"

using smoothing

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

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

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

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

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

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

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

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

Now,

$$\begin{aligned} P(\text{spam} | \text{claim, your, free, prize}) &\propto P(\text{spam}) \times P(\text{claim} | \text{spam}) \times \\ &\quad P(\text{your} | \text{spam}) \times P(\text{free} | \text{spam}) \times \\ &\quad P(\text{prize} | \text{spam}) \\ &\approx 0.5 \times \frac{2}{42} \times \frac{1}{42} \times \frac{2}{42} \times \frac{2}{42} \end{aligned}$$

$$\approx 0.0000026$$

$$\approx \frac{4}{3111696} \approx 0.0000012855$$

$$\begin{aligned}
 & P(\text{not spam} | \text{claim, your, free, prize}) = P(\text{not spam}), P(\text{claim}) \\
 & \quad \times P(\text{your} | \text{not spam}) \\
 & \quad \times P(\text{free} | \text{not spam}) \\
 & \quad \times P(\text{prize} | \text{not spam}) \\
 & \approx 0.5 \times \frac{1}{42} \times \frac{1}{42} \times \frac{2}{42} \times \frac{1}{42} \\
 & = \frac{1}{3111696} \\
 & \approx 0.000003213
 \end{aligned}$$

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

$$\text{Normalize} = \frac{0.0000012855}{0.0000012855 + 0.0000003213}$$

$$= \frac{0.0000012855}{0.0000016068}$$

$$= 0.800 \times 100$$

→ 80.0% spam