| 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) |

10/7 ① Total unique words in spam: 14

Total unique words in not spam: 14

Vocabulary size = 28

using smoothing:

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

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

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

$$P\left(\frac{tomorrow}{meeting} \mid spam\right) = \frac{0+1}{14+28} = \frac{1}{42}$$

$$P(free \mid Notspam) = \frac{0+1}{14+28} = \frac{1}{42}$$

$$P(meeting \mid Notspam) = \frac{1+1}{14+28} = \frac{2}{42}$$

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

# Prior probabilities

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

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

$\Big\{$ (3 not spam and 3 spam emails, total 6) $\Big\}$

Now,

P(spam| free, meeting, tomorrow) $\propto$ P(spam) $\times$ P(free|spam) $\times$
P(meeting|spam) $\times$
P(tomorrow|spam)

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

$$\approx 0.00001349$$

P(notspam | free, meeting, tomorrow) $\propto$ P(notspam) $\times$ P(free|notspam) $\times$
P(meeting|notspam) $\times$ P(tomorrow|-notspam)

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

$$\approx 0.00002699$$

Since P(notspam) > P(spam), the email is not spam.

Normalization $= \dfrac{0.00002699}{0.06002699 + 0.00001349} \times 100\%$

$$= \frac{0.00002699}{0.00004048} \times 100\% = 66.67\% \text{ not spam}$$

② for email "claim your free prize"

using smoothing

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

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

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

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

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

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

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

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

Now,

$$P(\text{spam} \mid \text{claim, your, free, prize}) \propto P(\text{spam}) \times P(\text{claim} \mid \text{spam}) \times P(\text{your} \mid \text{spam}) \times P(\text{free} \mid \text{spam}) \times P(\text{prize} \mid \text{spam})$$

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

$$\propto 0.00000026$$

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

P(notspam|claim,your,free,prize) $\alpha$ P(notspam), P(claim|
notspam) $\times$
P(your|notspam)
$\times$ P(free|notspam)
$\times$ P(prize|notspam)

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

$$= \frac{1}{3111696}$$

$$\approx 0.0000003213$$

Since P(spam) > P(notspam), the email "claim
your free prize" is spam.

Normalize $= \dfrac{0.0000012855}{0.0000012855 + 0.0000003213}$

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

$$= 0.800 \, \text{AllD}$$

$$= 80.0\%. \quad \underline{\underline{spam}}$$