# Probability

— study of events

— frequentist - concern with events as frequencies

— Bayesian - prior beliefs/info.

A - Event humidity will run today

B - Event that my cat will meow today

$$P(A) = \frac{800}{1000} = \frac{8}{10} = \boxed{\frac{4}{5}}$$

$$P(B) = \frac{900}{1000} = \boxed{\frac{9}{10}}$$

$$P(A \text{ and } B) = P(A \cap B) = \underline{P(A) P(B)}$$

$$\underline{A \perp B} = \left(\frac{4}{5}\right)\left(\frac{9}{10}\right) = \frac{36}{50}$$

$$= 72\%$$

$$P(A \text{ or } B) = P(A \cup B) = P(A) + P(B)$$

$$- P(A \cap B)$$

$$= \frac{8}{10} + \frac{9}{10} - \frac{36}{50} = \boxed{98\%}$$

$$P(A^c) = 1 - P(A) = \frac{1}{5}$$

$$P(B^c) = 1 - P(B) = \frac{1}{10}$$

## Bayes Theorem

$$P(A \mid B) = \frac{\overbrace{P(B \mid A) \, P(A)}^{likelihood}}{\underbrace{P(B)}_{}}\Big]\;prior$$

posterior probability

marginal likelihood

Likelihood: $P(B|A) = \dfrac{100}{200} = \boxed{0.5}$

Prior: $P(A) = \dfrac{800}{1000} = \boxed{\dfrac{4}{5} \text{ (empirical prior)}}$

$\dfrac{9}{10}$ (subjective prior) ✓

$\dfrac{1}{2}$ (flat prior)

Marginal Likelihood: $\boxed{P(B) = \dfrac{9}{10}}$

ShowMe.com

Marginal Likelihood =

$$P(B) = \underbrace{P(B|A)\ P(A)} + \underbrace{P(B|A^c)\ P(A^c)}$$

$$P(A|B) = \frac{\left(\frac{1}{2}\right)\left(\frac{4}{5}\right)}{9/10} = \frac{4/10}{9/10}$$

$$= \frac{40}{90} = \boxed{12\%} \quad \text{empirical prior}$$

$$P(A|B) = \frac{\left(\frac{1}{2}\right)\left(\frac{1}{2}\right)}{9)_W} = \frac{1/4}{9)_{10}}$$

$$= \frac{10}{36} \approx \boxed{28\%}$$

<span style="color:red">] flat prior</span>

<span style="color:red">subjective</span>

$$\Rightarrow \frac{\left(\frac{1}{2}\right)\left(\frac{9}{10}\right)}{9/10} = 50\%$$ <span style="color:red">]</span>

# Text-as-Data

$N = 100$ emails $\rightarrow$ documents

$\underbrace{\phantom{N = 100 \text{ emails}}}_{\text{Corpus}}$

Spam $\begin{cases} 1 & \text{email is spam} \\ 0 & \text{email } '' \text{ not spam} \end{cases}$
(target)

Text $\rightarrow$ Text of emails
(features)

2 steps to turn text into data

① **Pre-processing**

NLP {

  a) Tokenization

  b) **Cleaning**

    ↳ Stemming

    ↳ lowercase, remove #'s, punctuation

    ↳ Stop word removal

② Document-term matrix
   (DTM)

# Tokenization

## Spam  Text

| Spam | Text |
|---|---|
| 1 | "You win million dollars give bank info to claim prize" |

n-grams → Unigrams (words)

n-grams → 2-grams (phrases)

$n > 1 \Rightarrow$ phrase

$$w_1 = \{ You, win, million, dollars\_\_\_ , \\ You win, win million, millio.. milli\$$$

$$W_1 = \{ you, win, minien, dollars, .. \}$$

→ <u>stemming</u> - removing prefixes & suffixes and replacing words with their root

dollars ⟶ dollar

groups ⟶ group
graphs
grouped

→ <u>stop word removal</u>
the ⟶ for
you
a

$w_1 = \{ wm, million, dollar, \ldots \}$

Document term matrix

$$X \in \mathbb{R}^{N \times W}$$

$N \equiv \#$ of docs in a corpus

$W \equiv \#$ of terms in a corpus

$\boxed{\text{TF ?}} \rightarrow$ # of terms in a doc

$w_1 = \{\underline{\text{win}}, \underline{\text{miMion}}, \underline{\text{dol}}\text{lar} \ldots \}$

win million dollar today yury ..

$1 \quad \left( \dfrac{1}{60} \quad \dfrac{1}{2}\checkmark \quad \dfrac{1}{4} \quad 0 \quad 0 \right.$

$\text{TF-IPF} \quad \text{IDF}(wm) = \dfrac{1}{50}$

$\text{IDF}(\text{Dollar}) \quad \text{IDF}(\text{million})$

$= \dfrac{1}{4} \qquad = \dfrac{1}{2}$

## Naive Bayes

$N = 100$ $\longrightarrow$ 20 test

$\longrightarrow$ 80 Train

$$P(S=1 \mid w_i) = \frac{\overbrace{P(w_i \mid S=1)}^{\text{Likelihood}} P(S=1)}{P(w_i)}$$

$$= \frac{P(\text{win}, \text{million}, \text{dollar} \mid S=1)}{P(\text{win} \mid S=1) \, P(\text{million} \mid S=1) \, P(\text{dollar})}$$

$80 \begin{cases} 30 \text{ spam} \\ 50 \text{ no spam} \end{cases}$

1000 words "Win" appears 300

$$P(\text{win} \mid \text{spam}) = \frac{300}{1000}$$

$$P(S=1) = \frac{30}{80} = \frac{3}{8}$$