Bag of Words

– Anurag Dhaipule.

(1) Document Similarity

let the representation be

$$[woof \quad meow \quad squeak]$$

Then

$$D1 = [2 \quad 1 \quad 0]$$

$$D2 = [2 \quad 0 \quad 1]$$

Using $\quad Sim(A, B) = \dfrac{\sum_i a_i b_i}{\sqrt{\sum_i a_i^2} \sqrt{\sum_i b_i^2}}$ — ①

a) $\quad Sim(D1, D2) = \dfrac{2 \cdot 2 + 1 \cdot 0 + 0 \cdot 1}{\sqrt{2^2 + 1^2 + 0^2} \cdot \sqrt{2^2 + 0^2 + 1}}$

$$= \dfrac{4}{5}$$

b) Including idf weights

$$idf_i = \log\left(\dfrac{N}{n_i}\right)$$

$$idf_{woof} = \log \dfrac{2}{2} = 0$$

$$idf_{meow} = \log \dfrac{2}{1} = \log 2$$

$$idf_{speak} = \log \dfrac{2}{1} = \log 2.$$

Now, the new representations for the documents will
be

$$D_1' = [2 \times idf_{woof} \quad , \quad 1 \times idf_{meow} \quad \quad 0 \times idf_{squeak}]$$
$$D_2' = [2 \times idf_{woof} \quad , \quad 0 \times idf_{meow} \quad \quad 1 \times idf_{squeak}]$$

$$D_1' = [\; 0 \quad \log 2 \quad 0 \;]$$
$$D_2' = [\; 0 \quad \quad 0 \quad \log 2 \;]$$

$\rightarrow$ Using eq①

$$Sim(D_1', D_2') = \frac{0 \cdot 0 + 0 \cdot \log 2 + \log 2 \cdot 0}{\sqrt{0 + (\log 2)^2 + 0} \;\; \sqrt{0 + 0 + (\log 2)^2}}$$

$$= 0.$$

b) $\quad D_3 = [\; 0 \quad 1 \quad 1 \;]$

Now $\quad N = 3.$

New idfs, $\quad idf_{woof} = \log \dfrac{3}{2}$

$\quad\quad\quad\quad idf_{meow} = \log \dfrac{3}{2}$

$\quad\quad\quad\quad idf_{squeak} = \log \dfrac{3}{2}$

New representations

$$D_1' = [2 \times idf_{woof} \;, \; 1 \times idf_{meow} \;, \; 0 \times idf_{squeak}]$$

$$= \left[\; 2 \log \dfrac{3}{2} \;, \; \log \dfrac{3}{2} \;, \; 0 \;\right]$$

similarly,

$$D_2' = \left[\; 2 \log \dfrac{3}{2} \;, \; 0 \;, \; \log \dfrac{3}{2} \;\right]$$

Using ①,

$$\text{Sim}(D_1', D_2') = \frac{2\log\frac{3}{2} \times 2\log\frac{3}{2} + 0 + 0}{\sqrt{\left(2\log\frac{3}{2}\right)^2 + \left(\log\frac{3}{2}\right)^2 + 0}\sqrt{\left(2\log\frac{3}{2}\right)^2 + 0 + \left(\log\frac{3}{2}\right)^2}}$$

$$= \frac{4\left(\log\frac{3}{2}\right)^2}{5\left(\log\frac{3}{2}\right)^2} = \frac{4}{5}$$

**(2) Naive Bayes and Smoothing.**

(a)
$$P(t) = \frac{\text{Count of docs labelled } t}{\text{Total docs}}$$

$$P(+) = \frac{5}{10} = \frac{1}{2} \qquad\qquad P(-) = \frac{5}{10} = \frac{1}{2}$$

$$P(w_i|t) = \frac{\text{Count of docs labelled } t \text{ containing } w_i}{\text{Count of docs labelled } t}$$

$$P(\text{great}|+) = \frac{5}{5} = 1 \qquad\qquad P(\text{great}|-) = \frac{0}{5} = 0$$

$$P(\text{food}|+) = \frac{5}{5} = 1 \qquad\qquad P(\text{food}|-) = \frac{5}{5} = 1$$

$$P(\text{Served}|+) = \frac{0}{5} = 0 \qquad\qquad P(\text{Served}|-) = \frac{1}{5} = \frac{1}{5}$$

$$P(\text{terrible}|+) = \frac{0}{5} = 0 \qquad\qquad P(\text{terrible}|-) = \frac{5}{5} = 1$$

**(a)** Using Bernoulli's Naive Bayes

$P(+|\text{"great food Served"}) = P(\text{great}|+) \times P(\text{food}|+) \times$
$$P(\text{Served}|+) \times (1 - P(\text{terrible}|+) \times P(+)$$
$$= 1 \times 1 \times 0 \times (1-0) \times \frac{1}{2} =$$
$$= 0.$$

$P(-|\text{"great food Served"}) = P(\text{great}|-) \times P(\text{food}|-) \times$
$$P(\text{Served}|-) \times (1 - P(\text{terrible}|-) \times P(-)$$
$$= 0 \times 1 \times \frac{1}{5} \times (1-1) \times \frac{1}{2}$$
$$= 0.$$

**(b)** With laplace smoothing,

we add

$$P(t) = \frac{\text{Count of doc labelled } t + 1}{\text{total doc} + 2}$$

$$P(+) = \frac{6}{12} = \frac{1}{2} \qquad\qquad P(-) = \frac{1}{2}$$

similarly

$$P(\text{great}|+) = \frac{5+1}{5+2} = \frac{6}{7} \;;\quad P(\text{great}|-) = \frac{1}{7}$$

$$P(\text{food}|+) = \frac{5+1}{5+2} = \frac{6}{7} \;;\quad P(\text{food}|-) = \frac{6}{7}$$

$$P(\text{Served}|+) = \frac{0+1}{5+2} = \frac{1}{7} \;;\quad P(\text{Served}|-) = \frac{2}{7}$$

$$P(\text{terrible}|+) = \frac{0+1}{5+2} = \frac{1}{7} \qquad P(\text{Served}|-) = \frac{6}{7}$$

Using the equation ① & ②, we have.

$$P(+ \mid \text{"great food served"}) = \frac{6}{7} \times \frac{6}{7} \times \frac{1}{7} \times \left(1 - \frac{1}{7}\right) \times \frac{1}{2}$$

$$= 0.045.$$

$$P(- \mid \text{"great food serva}) = \frac{1}{7} \times \frac{6}{7} \times \frac{2}{7} \times \left(1 - \frac{6}{7}\right) \times \frac{1}{2}$$

$$= 0.002$$

$P(+ \mid \text{"great food served"})$ is greater.