

Paper Code	Examiner	Department	Office
INT309	Fangyu Wu	Department of Intelligence Science	SD555

Based on the data in Table 1 below:

	docID	contents	in $c = \text{China?}$
Training Set	1	Taipei Taiwan	yes
	2	Macao Taiwan Shanghai	yes
	3	Japan Sapporo	no
	4	Sapporo Osaka Taiwan	no
Test Set	5	Taiwan Taiwan Sapporo	?

Table Q1

1. Estimate a **multinomial Naive Bayes classifier**, and apply the classifier to the test document.
2. Estimate a **Bernoulli Naive Bayes classifier** and apply the classifier to the test document.
3. Compute the **tf-idf vector (normalized)** representations of the documents in the table below. By using **Rocchio Classification Algorithm** to determine the class of the document.

———— *End of paper* ———

Multinomial:

Bernoulli:

$$P(c|d) \propto P(c)\prod_{1 \leq k \leq N_d} P(t_k|c)$$

$$P(t_k|c) = \frac{T_{ct} + 1}{\sum_{t' \in V(T_{ct}+1)}$$

$$P(c|d) \propto P(c)\prod_{t_k \in Q} P(t_k|c)\prod_{t_k \notin Q} [1 - P(t_k|c)] \quad (5)$$

$$P(t_k|c) = \frac{df_{ct} + 1}{N_c + \text{Number of classes}} \quad (6)$$

	docID	contents	in $c = \text{China?}$
Training Set	1	Taipei Taiwan	yes
	2	Macao Taiwan Shanghai	yes
	3	Japan Sapporo	no
	4	Sapporo Osaka Taiwan	no
Test Set	5	Taiwan Taiwan Sapporo	?

Table Q1

$$1. \hat{P}(c) = \frac{2}{4} = \frac{1}{2}$$

$$\hat{P}(\text{Taiwan}|c) = \frac{2+1}{5+7} = \frac{1}{4}$$

$$\hat{P}(\text{Sapporo}|c) = \frac{0+1}{5+7} = \frac{1}{12}$$

$$\therefore \hat{P}(c|d_5) = \hat{P}(c) \cdot (\hat{P}(\text{Taiwan}|c))^2 \cdot \hat{P}(\text{Sapporo}|c)$$

$$= \frac{1}{2} \cdot \frac{1}{16} \cdot \frac{1}{12}$$

$$\approx 0.0026.$$

$$\hat{P}(\bar{c}) = \frac{2}{4} = \frac{1}{2}$$

$$\hat{P}(\text{Taiwan}|\bar{c}) = \frac{1+1}{5+7} = \frac{1}{6}$$

$$\hat{P}(\text{Sapporo}|\bar{c}) = \frac{2+1}{5+7} = \frac{1}{4}$$

$$\therefore \hat{P}(\bar{c}|d_5) = \hat{P}(\bar{c}) \cdot (\hat{P}(\text{Taiwan}|\bar{c}))^2 \cdot \hat{P}(\text{Sapporo}|\bar{c})$$

$$= \frac{1}{2} \cdot \frac{1}{36} \cdot \frac{1}{4}$$

$$\approx 0.0035$$

$\because \hat{P}(c|d_5) < \hat{P}(\bar{c}|d_5) \therefore \text{Test document doesn't belong to } c.$

Taipei, Taiwan, Macao, Shanghai, Japan, Sapporo, Osaka.

Multinomial:

$$P(c|d) \propto P(c)\prod_{1 \leq k \leq N_d} P(t_k|c)$$

$$P(t_k|c) = \frac{T_{ct} + 1}{\sum_{t' \in V(T_{ct}+1)}}$$

Bernoulli:

$$P(c|d) \propto P(c)\prod_{t_k \in Q} P(t_k|c)\prod_{t_k \notin Q} [1 - P(t_k|c)] \quad (5)$$

$$P(t_k|c) = \frac{df_{ct} + 1}{N_c + \text{Number of classes}} \quad (6)$$

	docID	contents	in $c = \text{China?}$
Training Set	1	Taipei Taiwan	yes
	2	Macao Taiwan Shanghai	yes
	3	Japan Sapporo	no
	4	Sapporo Osaka Taiwan	no
Test Set	5	Taiwan Taiwan Sapporo	?

Table Q1

2.

$$\hat{P}(c) = \frac{2}{4} = \frac{1}{2}$$

$$\hat{P}(\bar{c}) = \frac{2}{4} = \frac{1}{2}$$

$$\hat{P}(\text{Taiwan}|c) = \frac{2+1}{2+2} = \frac{3}{4}$$

$$\hat{P}(\text{Sapporo}|c) = \frac{0+1}{2+2} = \frac{1}{4}.$$

$$\hat{P}(\text{Taipei}|c) = \hat{P}(\text{Macao}|c) = \hat{P}(\text{Shanghai}|c) = \frac{1+1}{2+2} = \frac{1}{2}$$

$$\hat{P}(\text{Japan}|c) = \hat{P}(\text{Osaka}|c) = \frac{0+1}{2+2} = \frac{1}{4}.$$

$$\begin{aligned} \therefore \hat{P}(c|d_5) &= \hat{P}(c) \cdot \prod_{t_k \in Q} \hat{P}(t_k|c) \cdot \prod_{t_k \notin Q} (1 - \hat{P}(t_k|c)) \\ &= \frac{1}{2} \times \frac{3}{4} \times \frac{1}{4} \times (1 - \frac{1}{2})^3 \times (1 - \frac{1}{4})^2 \\ &\stackrel{\text{Taiwan Sapporo}}{\approx} \stackrel{\text{Taipei, Macao, SH}}{=} \stackrel{\text{Japan, Osaka}}{\approx} 0.0066 \end{aligned}$$

$$\hat{P}(\text{Taiwan}|\bar{c}) = \frac{1+1}{2+2} = \frac{1}{2}$$

$$= \hat{P}(\text{Japan}|\bar{c}) = \hat{P}(\text{Osaka}|\bar{c})$$

$$\hat{P}(\text{Sapporo}|\bar{c}) = \frac{2+1}{2+2} = \frac{3}{4}$$

$$\hat{P}(\text{Taipei}|\bar{c}) = \hat{P}(\text{Macao}|\bar{c}) = \hat{P}(\text{Shanghai}|\bar{c}) = \frac{0+1}{2+2} = \frac{1}{4}.$$

Taipei, Taiwan, Macao,  
Shanghai, Japan, Sapporo  
Osaka.

	docID	contents	in $c = \text{China?}$
Training Set	1	Taipei Taiwan	yes
	2	Macao Taiwan Shanghai	yes
	3	Japan Sapporo	no
	4	Sapporo Osaka Taiwan	no
Test Set	5	Taiwan Taiwan Sapporo	?

Table Q1

$$\therefore \hat{P}(\bar{c} | d_5) = \hat{P}(\bar{c}) \cdot \prod_{t_k \in Q} \hat{P}(t_k | \bar{c}) \cdot \prod_{t_k \notin Q} (1 - \hat{P}(t_k | \bar{c})).$$

$$= \frac{1}{2} \times \frac{1}{2} \times \frac{3}{4} \times (1 - \frac{1}{4})^3 \times (1 - \frac{1}{2})^2$$

Taiwan Sapporo Taipei, Macao, SH Japan, Osaka.

$$\approx 0.0198$$

$$\because \hat{P}(c | d_5) < \hat{P}(\bar{c} | d_5)$$

$\therefore$  Test document doesn't belong to c.

3.

	docID	contents	in c = China?
Training Set	1	Taipei Taiwan	yes
	2	Macao Taiwan Shanghai	yes
	3	Japan Sapporo	no
	4	Sapporo Osaka Taiwan	no
Test Set	5	Taiwan Taiwan Sapporo	?

Table Q1

$$idf_i = \log_{10} \frac{N}{df_i}$$

$$N = 4$$

term	doc1	doc2	doc3	doc4	doc5	df	idf
Taipei	1	0	0	0	0	1	0.6
Taiwan	1	1	0	1	2	3	0.1
Macao	0	1	0	0	0	1	0.6
Shanghai	0	1	0	0	0	1	0.6
Japan	0	0	1	0	0	1	0.6
Sapporo	0	0	1	1	1	2	0.3
Osaka	0	0	0	1	0	1	0.6

tf-idf (unnormalized)

term	doc1	doc2	doc3	doc4	doc5
Taipei	0.6	0	0	0	0
Taiwan	0.1	0.1	0	0.1	0.2
Macao	0	0.6	0	0	0
Shanghai	0	0.6	0	0	0
Japan	0	0	0.6	0	0
Sapporo	0	0	0.3	0.3	0.3
Osaka	0	0	0	0.6	0

$$C \left\{ \begin{array}{l} \overrightarrow{doc1_N} = \frac{(0.6, 0.1, 0, 0, 0, 0, 0)}{\sqrt{0.6^2 + 0}} = (0.28, 0.16, 0, 0, 0, 0, 0) \\ \overrightarrow{doc2_N} = \frac{(0, 0.1, 0.6, 0.6, 0, 0, 0)}{\sqrt{0.1^2 + 0.6^2 + 0.6^2}} = (0, 0.12, 0.7, 0.7, 0, 0, 0) \\ \overrightarrow{doc3_N} = \frac{(0, 0, 0, 0, 0.6, 0.3, 0)}{\sqrt{0.6^2 + 0.3^2}} = (0, 0, 0, 0, 0.89, 0.45, 0) \\ \overrightarrow{doc4_N} = \frac{(0, 0.1, 0, 0, 0, 0.3, 0.6)}{\sqrt{0.1^2 + 0.3^2 + 0.6^2}} = (0, 0.15, 0, 0.17, 0.44, 0.88) \\ \overrightarrow{doc5_N} = \frac{(0, 0.2, 0, 0, 0, 0.3, 0)}{\sqrt{0.2^2 + 0.3^2}} = (0, 0.55, 0, 0, 0, 0.83, 0) \end{array} \right.$$

The centroid of the relevant documents is:  
(class C)

$$C = \frac{\overrightarrow{doc1N} + \overrightarrow{doc2N}}{2} = (0.50, 0.14, 0.35, 0.35, 0, 0, 0)$$

The centroid of the non-relevant documents is:  
(class  $\bar{C}$ )

$$\bar{C} = \frac{\overrightarrow{doc3N} + \overrightarrow{doc4N}}{2} = (0, 0.08, 0, 0, 0.45, 0.45, 0.44)$$

The distances from the query to the centroids are:

$$\begin{aligned} d_{q,C} &= \|(0, 0.55, 0, 0, 0, 0.83, 0) - (0.5, 0.14, 0.35, 0.35, 0, 0, 0)\| \\ &= \|(0.5, 0.41, -0.35, -0.35, 0, 0.83, 0)\| \\ &= \sqrt{0.5^2 + 0.41^2 + 0.35^2 \times 2 + 0.83^2} \approx 1.16 \end{aligned}$$

$$\begin{aligned} d_{q,\bar{C}} &= \|(0, 0.55, 0, 0, 0, 0.83, 0) - (0, 0.08, 0, 0, 0.45, 0.45, 0.44)\| \\ &= \|(0, 0.47, 0, 0, -0.45, 0.38, -0.44)\| \\ &= \sqrt{0.47^2 + 0.45^2 + 0.38^2 + 0.44^2} \approx 0.87 \end{aligned}$$

$$\because d_{q,C} > d_{q,\bar{C}}$$

$\therefore$  test document doesn't belong to C.