

Question 1

Based on the data in Table

- (i) Classify test document using Multinomial Naive Bayes Classifier
- (ii) Classify test document using Bernoulli NB classifier,
- (iii) Classify test document using Rocchio classifier
(Use Euclidean distance)

	docID	words in document	in $c = \textit{China}$?
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	?

Q 1 Solution

- Multinomial NB
- $$P(C | d5) = P(C) P(\text{Taiwan} | C) P(\text{Taiwan} | C) P(\text{Sapporo} | C)$$
$$= 0.5 * 0.25 * 0.25 * 1/12 = 0.0026$$
- $$P(\text{Not } C | d5) = P(\text{Not } C) P(\text{Taiwan} | \text{Not } C) P(\text{Taiwan} | \text{Not } C) P(\text{Sapporo} | \text{Not } C)$$
$$= 0.5 * 1/6 * 1/6 * 0.25 = 0.00347$$

d5 is classified in Class Not China

Q 1 Solution

- Bernoulli NB
- $$P(C | d5) = P(C) P(\text{Taiwan} | C) P(\text{Sapporo} | C) (1 - P(\text{Taipei} | C)) (1 - P(\text{Macao} | C)) (1 - P(\text{Taipei} | C)) (1 - P(\text{Shanghai} | C)) (1 - P(\text{Japan} | C)) (1 - P(\text{Osaka} | C))$$
$$= 0.00659$$

Q 1 Solution

- Bernoulli NB
- $$P(\text{Not } C \mid d5) = P(\text{Not } C) P(\text{Taiwan} \mid \text{Not } C) \\ P(\text{Sapporo} \mid \text{Not } C) (1 - P(\text{Taipei} \mid \text{Not } C)) \\ (1 - P(\text{Macao} \mid \text{Not } C)) (1 - P(\text{Taipei} \mid \text{Not } C)) \\ (1 - P(\text{Shanghai} \mid \text{Not } C)) (1 - P(\text{Japan} \mid \text{Not } C)) \\ (1 - P(\text{Osaka} \mid \text{Not } C)) \\ = 0.0197$$

Q 1 Solution

- Rocchio

	Taiwan	Sapporo	Japan	Taipei	Osaka	Shanghai	Macao
d_1	0.4	0	0	2	0	0	0
d_2	0.4	0	0	0	0	2	2
d_3	0	1	2	0	0	0	0
d_4	0.4	1	0	0	2	0	0
Centriod of C	0.4	0	0	1	0	1	1
Centriod of Not C	0.2	1	1	0	1	0	0
d_5	0.8	1	0	0	0	0	0

Distance of d_5 to centriod of C = 2.36

Distance of d_5 to centriod of Not C = 1.61

d_5 is classified in class Not C

Question 2

- (i) Compute the “export”/POULTRY contingency table for the “Kyoto”/JAPAN in the collection given below.
- (ii) Make up a contingency table for which MI is 0 – that is, term and class are independent of each other.

	$e_c = e_{poultry} = 1$	$e_c = e_{poultry} = 0$
$e_t = e_{EXPORT} = 1$	$N_{11} = 49$	$N_{10} = 27,652$
$e_t = e_{EXPORT} = 0$	$N_{01} = 141$	$N_{00} = 774,106$

Collection:

	docID	words in document	in $c = \text{Japan?}$
training set	1	Kyoto Osaka Taiwan	yes
	2	Japan Kyoto	yes
	3	Taipei Taiwan	no
	4	Macao Taiwan Shanghai	no
	5	London	no

Question 3

- Consider the following frequencies for the class *coffee for four terms in the total 100,000 documents*

term	N_{00}	N_{01}	N_{10}	N_{11}
brazil	98,012	102	1835	51
council	96,322	133	3525	20
producers	98,524	119	1118	34
roasted	99,824	143	23	10

- Calculate *mutual information* for the term brazil

Q 2 Solution

Part 1

$N_{11} = 2$	$N_{10} = 0$
$N_{01} = 0$	$N_{00} = 3$

Part 2: $N_{11} = N_{01}$ and $N_{10} = N_{00}$; term is present in half docs of each class

$N_{11} = 3$	$N_{10} = 5$
$N_{01} = 3$	$N_{00} = 5$

OR $N_{01} = 0$ and $N_{00} = 0$, term is present in all docs in all classes

$N_{11} = 2$	$N_{10} = 3$
$N_{01} = 0$	$N_{00} = 0$

Q 3 Solution

N11 = 51	N10 = 1835
N01 = 102	N00 = 98,012

$$N = 100,000$$

$$I(U; C) = \frac{N_{11}}{N} \log_2 \frac{NN_{11}}{N_{1.}N_{.1}} + \frac{N_{01}}{N} \log_2 \frac{NN_{01}}{N_{0.}N_{.1}} \\ + \frac{N_{10}}{N} \log_2 \frac{NN_{10}}{N_{1.}N_{.0}} + \frac{N_{00}}{N} \log_2 \frac{NN_{00}}{N_{0.}N_{.0}}$$

$$\frac{51}{100000} \log \frac{51 * 100000}{(51 + 1835)(51 + 102)} + \frac{102}{100000} \log \frac{102 * 100000}{(102 + 98012)(51 + 102)} +$$

$$\frac{1835}{100000} \log \frac{1835 * 100000}{(51 + 1835)(1835 + 98012)} + \frac{98012}{100000} \log \frac{98012 * 100000}{(98012 + 1835)(102 + 98012)}$$