

# NAIVE BAYES CLASSIFIER

1) Queremos saber si un comentario pertenece a la clase China o no.

	docID	words in document	in clase China?
Training Set	1	Chinese Beijing Chinese	yes
	2	Chinese Chinese Shanghai	yes
	3	Chinese Museo	yes
	4	Tokyo Japan Chinese	no
Test set	5	Chinese Chinese Chinese Tokyo Japan	?

2) Queremos saber si el cambio test (5) es + o -

$$\hat{C} = \max_c \Pr(c) \prod_{i=1}^n \Pr(f_i, c)$$

3) Calculamos la probabilidad a priori.

$$\Pr(c) = \frac{n^{\text{documentos de la clase China}}}{n^{\text{total de documentos in Training set}}} = \frac{N_c}{N_{\text{docs}}}$$

$$\Pr(\text{yes}) = 3/4, \quad \Pr(\text{no}) = 1/4$$

4) Calcular las probabilidades condicionales

$$\Pr(t/c) = \frac{T_{ct} + 1}{T'_{ct} + B}$$

$T_{ct}$  = n° ocurrencias del término  $t$  en la clase o grupo  $c$ .

$T'_{ct}$  = n° total de términos en cada clase (training set)

$B$  = n° total de términos diferentes en todo el grupo entrenado.  
Es n° palabras diferentes.

$$\Pr(\text{Chinese} / \text{yes}) = \frac{5+1}{8+6} = 3/7$$

$$T_{\text{yes}-t} = 8$$

$$T'_{\text{no}-t} = 3$$

$$B = 6$$

$$\Pr(\text{Tokyo} / \text{yes}) = \Pr(\text{Japan} / \text{yes}) = \frac{0+1}{8+6} = 1/14$$

$$\Pr(\text{Chinese} / \text{no}) = \frac{1+1}{3+6} = 2/9$$

$$\Pr(\text{Tokyo} / \text{no}) = \Pr(\text{Japan} / \text{no}) = \frac{1+1}{3+6} = 2/9$$

5) Multiplicamos la probabilidad a priori de los documentos de cada clase por las probabilidades condicionales de las palabras incluidas en la consulta de test.

$$\Pr(\text{yes} / \text{test doc}) = \text{Prob. Priori} (\text{Chinese}^3 \times \text{Japan} \times \text{Tokyo}) = 3/4 \cdot (3/7)^3 \cdot 1/14 \cdot 1/14 = [0,0003]$$

$$\Pr(\text{no} / \text{test doc}) = 1/4 \cdot (2/9)^3 \cdot 2/9 \cdot 2/9 = [0,0006]$$

[Sol: Se asigna el documento 5 a: