

Exercise _ Text Classification_ Naïve Bayes_ big data

1.- Queremos entender si un comentario pertenece a la clase China o no (Análisis de sentimiento consiste en la asignación de un nuevo comentario a las clases positivas o negativas).

Anteriormente hemos entrenado nuestro sistema, y sabemos qué expresiones se consideran referidas a China o no. Nuestro léxico se muestra en la siguiente tabla

► **Table 13.1** Data for parameter estimation examples.

	docID	words in document	in $c = \text{China?}$
training set	1	Chinese Beijing Chinese	yes
	2	Chinese Chinese Shanghai	yes
	3	Chinese Macao	yes
	4	Tokyo Japan Chinese	no
test set	5	Chinese Chinese Chinese Tokyo Japan	?

$$P(\text{yes}) = \frac{n^{\circ} \text{yes}}{n^{\circ} \text{total}}$$

↳ 3/4

$$P(\text{No}) = \frac{n^{\circ} \text{no}}{n^{\circ} \text{total}}$$

↳ 1/4

Queremos saber si el conjunto test es positivo o negativo. Para eso, vamos a utilizar el clasificador denominado Naïve Bayes binomial

$$\hat{c} = \max_{c \in C} \Pr(c) \prod_{i=1}^n \Pr(f_i | c)$$

2.- Calcularlas probabilidades a priori

$$\begin{aligned} \Pr(c_i) &= \frac{\text{number of docs of class } c}{\text{total number of docs in training dataset}} \\ &= \frac{N_c}{N_{\text{docs}}} \end{aligned}$$

$$P(\text{yes}) = \frac{3}{4}$$

$$P(\text{no}) = \frac{1}{4}$$

3.- Calculamos las probabilidades condicionales

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

$\frac{\text{Nº veces término en el grupo de entrenamiento} + 1}{\text{Nº total de términos en cada clase} + \text{Nº total de términos en todo el grupo de entrenamiento}}$

- T_{ct} = número de ocurrencias del término t en el grupo de entrenamiento de los documentos clase c
- T'_{ct} = número total de términos en cada clase

$T'_{\text{yes-t}} = 8$ → mirar table.

$T'_{\text{no-t}} = 3$ → mirar table

- B = número total de términos diferentes en todo el grupo de entrenamiento (vocabulario)

• $B = 6$

Chinese, Beijing, Shanghai, Macao, Tokyo, Japan

¿Cuántas veces aparece la palabra chinese en la table? → 5

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

• $P(\text{Tokyo/yes}) = P(\text{Japan/yes}) = (\underbrace{10}_{\text{veces que estas palabras tienen yes}} + 1) / (8 + 6) = 1/14$

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

• $P(\text{Tokyo/no}) = P(\text{Japan/no}) = (\underbrace{1}_{\text{veces Tokyo no / Japan no}} + 1) / (3 + 6) = 2/9$

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

$P(\text{yes/ test document}) = \frac{3}{4} * (\frac{3}{7})^3 * \frac{1}{14} * \frac{1}{14} = 0.0003$

$P(\text{no/test document}) = \frac{1}{4} * (\frac{2}{9})^3 * \frac{2}{9} * \frac{2}{9} = 0.0001$

Entonces, ¿el clasificador asigna el documento 5 a ? Chinese yes