

Esercizio: costruire un Naive Bayes classifier per classificare documenti in argomenti. Come features si considerano la presenza o meno di parole di un vocabolario nel documento da classificare.

Vocabolario: (baseball calcio Moratti riforma sport studio)

La funzione di classificazione $f(x)$ assume due valori: sport, scuola.

Le x_i (i documenti) sono rappresentati mediante vettori a 6 valori binari. Il valore "1" della feature x_{ij} indica che il termine j-esimo ($j=1, 2, \dots, 6$) del vocabolario è presente nel documento i-esimo.

Insieme di addestramento D:

$x_1 = ((0,0,1,0,1,0), \text{sport})$

$x_2 = ((0,0,1,0,0,1), \text{scuola})$

$x_3 = ((1,1,1,0,0,0), \text{sport})$

$x_4 = ((1,0,0,0,1,0), \text{sport})$

$x_5 = ((0,0,0,1,0,1), \text{scuola})$

$x_6 = ((0,0,0,1,0,1), \text{scuola})$

Si vuole assegnare una classe all'istanza $x = (1,0,1,0,0,0)$ sulla base di un Naive Bayes classifier:

$$f(x) = \arg \max_{c_i \in C} P(c_i) \prod P(val_j | c_i)$$

$f(x)$: assume i valori in $C = \{\text{sport}, \text{scuola}\}$

a_j (valore del j-esimo attributo) = 0,1 (tutti gli attributi sono binari)

Nell'esempio x da classificare ho i seguenti valori per i 6 attributi:

baseball=1 calcio=0 Moratti=1 riforma=0 sport=0 studio=0

Le (stime delle) probabilità a priori delle due classi sono uguali, cioè $P(\text{sport})=P(\text{scuola})=3/6=1/2$. Devo stimare sull'insieme di addestramento D le seguenti probabilità condizionate:

$P(\text{baseball}=1|\text{sport})$ $P(\text{calcio}=0|\text{sport})$ $P(\text{Moratti}=1|\text{sport})$

$P(\text{riforma}=0|\text{sport})$ $P(\text{sport}=0|\text{sport})$ $P(\text{studio}=0|\text{sport})$

$P(\text{baseball}=1|\text{scuola})$ $P(\text{calcio}=0|\text{scuola})$ $P(\text{Moratti}=1|\text{scuola})$ $P(\text{riforma}=0|\text{scuola})$

$P(\text{sport}=0|\text{scuola})$ $P(\text{studio}=0|\text{scuola})$

avrò le seguenti stime per le probabilità condizionate (senza m-stime):

$P(\text{baseball}=1|\text{sport})=2/3$ $P(\text{calcio}=0|\text{sport})=2/3$ $P(\text{Moratti}=1|\text{sport})=2/3$

$P(\text{riforma}=0|\text{sport})=3/3$ $P(\text{sport}=0|\text{sport})=1/3$ $P(\text{studio}=0|\text{sport})=3/3$

$P(\text{baseball}=1|\text{scuola})=1/3$ $P(\text{calcio}=0|\text{scuola})=1/3$ $P(\text{Moratti}=1|\text{scuola})=1/3$

$P(\text{riforma}=0|\text{scuola})=0$ $P(\text{sport}=0|\text{scuola})=2/3$ $P(\text{studio}=0|\text{scuola})=0$

Dunque,

$P(\text{sport}) \prod P(a_j|\text{sport}) = 1/2 * 2/3 * 2/3 * 2/3 * 1 * 1/3 * 1 = 4/81 = 0,05$

$P(\text{scuola}) \prod P(a_j|\text{scuola}) = 1/2 * 1/3 * 1/3 * 1/3 * 0 * 2/3 * 0 = 0$

quindi l'argmax fornisce $c_{NB}=\text{sport}$