# Document Classification Problem

We have 2 classes of documents.

(Example: class 1 : spam, class 2 : important, class 3: personal)
class 4 : work.

$$k = 4$$

Dataset: class 1 has 5000 documents.
has class 2 " 6000 "
$25K$ " 3 " 4000 "
documents " 4 " 10000 "

Prior probabilities: $\pi_1 = \frac{5000}{20000} = \frac{5}{25}$, ... $\pi_4 = \frac{10}{25}$

$$\pi_2 = \frac{6}{25}, \quad \pi_3 = \frac{4}{25}.$$

$V = \{word_1, word_2, \ldots, word_{|V|}\}$ has $|V|$ elements.

$\bar{p}^{(1)} = (P_{11}, \ldots, P_{1|V|})$ is the model for class 1.

$\bar{p}^{(2)} = (P_{021}, \ldots, P_{2|V|})$ " " " " " 2.

$$\hat{P}_{21} = \frac{\# \text{ of word 1 in class 2}}{(\# \text{ of word 1} + \cdots + \# \text{ word}|V|) \text{ in class 2}}$$

$\bar{p}^{(3)} = (P_{31}, \ldots, P_{3|V|})$, $\qquad \sum\limits_{i=1}^{|V|} P_{3i} = 1$

$\bar{p}^{(4)} = (P_{41}, \ldots, P_{4|V|})$, $\qquad \sum\limits_{i=1}^{|V|} P_{4i} = 1$

$$\hat{P}_{43} = \frac{\# \text{ of word 3 in class 4}}{(\# \text{ of word 1} + \cdots + \# \text{ word}|V|) \text{ in class 4}}$$

We have 4 models:

Model 1 for class 1

$$\text{Likelihood}_1 = P_{11}^{x_1} \cdots \cdots P_{1|V|}^{x_{|V|}}$$

Model 2 for class 2

$$\text{Likelihood}_2 (x_1, x_2, \cdots, x_{|V|}) = P_{21}^{x_1} P_{22}^{x_2} \cdots P_{2|V|}^{x_{|V|}}$$

Model 3    "    3 (    "    ) $= P_{31}^{x_1} \cdots \cdots P_{3|V|}^{x_{|V|}}$

Model 4    "    4 (    "    ) $= P_{41}^{x_1} \cdots \cdots P_{4|V|}^{x_{|V|}}$

Training process is over

Testing (Classification or Recognition) phase

Given a document

$$\begin{array}{|c|} \hline \text{Word1, word2} \\ \text{word1, word2} \\ \text{word3, word4} \\ \text{word5, word6} \\ \hline \end{array} \xrightarrow[\text{extraction}]{\text{word1 word2 feature}} \quad x = \begin{bmatrix} 3 \\ 2 \\ 1 \\ 1 \\ 1 \\ 1 \end{bmatrix}$$

$|V| = 6$

Feature vector is obtained using the "bag of words" model.

for $i = 1:4$
$$\alpha_i = \pi_i \, P_{i1}^3 \, P_{i2}^2 \, P_{i3}^1 \, P_{i4}^1 \, P_{i5}^1 \, P_{i6}^1$$
end.

Answer = pick the $i$ that gives the max probability (likelihood).

$$\boxed{\text{Class Answer} = \arg\max_i (\alpha_i)}$$