# Relevance feedback + Text classification (Chapter 9+13)



## Definition 1 (Rocchio relevance feedback)
*Rocchio relevance feedback has the form*

$$q_m = \alpha q_0 + \beta \frac{1}{|D_r|} \sum_{\vec{d_r} \in D_r} \vec{d_r} - \gamma \frac{1}{|D_{nr}|} \sum_{\vec{d_{nr}} \in D_{nr}} \vec{d_{nr}}$$

*where $q_0$ is the original query vector, $D_r$ is the set of relevant documents, $D_{nr}$ is the set of non-relevant documents and the values $\alpha$, $\beta$, $\gamma$ depend on the system setting.*

## Exercise 9/1

What is the main purpose of Rocchio relevance feedback?

*PSEUDO - RELEVANCE FEEDBACK*

## Exercise 9/2

A user's primary query is *cheap CDs cheap DVDs extremely cheap CDs*. The user has a look on two documents: doc1 a doc2, marking doc1 *CDs cheap software cheap CDs* as relevant and doc2 *cheap thrills DVDs* as non-relevant. Assume that we use a simple *tf* scheme without vector length normalization. What would be the restructured query vector after considering the Rocchio relevance feedback with values $\alpha = 1$, $\beta = 0.75$, and $\gamma = 0.25$?

We rewrite the exercise to the table for an easier processing.

| terms | relevant doc1 | non-relevant doc2 | query |
|---|---|---|---|
| CDs | 2 | 0 | 2 |
| cheap | 2 | 1 | 3 |
| software | 1 | 0 | 0 |
| thrills | 0 | 1 | 0 |
| DVDs | 0 | 1 | 1 |
| extremely | 0 | 0 | 1 |

Table 1:

## Text classification and Naive Bayes (Chapter 13)

### Definition 2 (Naive Bayes Classifier)
*Naive Bayes (NB) Classifier assumes that the effect of the value of a predictor $x$ on a given class $c$ is class conditional independent. Bayes theorem provides a way of calculating the posterior probability $P(c|x)$ from class prior probability $P(c)$, predictor prior probability $P(x)$ and probability of the predictor given the class $P(x|c)$*

$$P(c|x) = \frac{P(x|c)P(c)}{P(x)}$$

*and for a vector of predictors $X = (x_1, \dots, x_n)$*

$$P(c|X) = \frac{P(x_1|c) \dots P(x_n|c)P(c)}{P(x_1) \dots P(x_n)}.$$

*The class with the highest posterior probability is the outcome of prediction.*

## Exercise 13/2

Considering the table of observations, use the Naive Bayes classifier to recommend whether to *Play Golf* given a day with *Outlook = Rainy, Temperature = Mild, Humidity = Normal* and *Windy = True*. Do not deal with the zero-frequency problem.

| Outlook | Temperature | Humidity | Windy | Play Golf |
|---|---|---|---|---|
| Rainy | Hot | High | False | No |
| Rainy | Hot | High | True | No |
| Overcast | Hot | High | False | Yes |
| Sunny | Mild | High | False | Yes |
| Sunny | Cool | Normal | False | Yes |
| Sunny | Cool | Normal | True | No |
| Overcast | Cool | Normal | True | Yes |
| Rainy | Mild | High | False | No |
| Rainy | Cool | Normal | False | Yes |
| Sunny | Mild | Normal | False | Yes |
| Rainy | Mild | Normal | True | Yes |
| Overcast | Mild | High | True | Yes |
| Overcast | Hot | Normal | False | Yes |
| Sunny | Mild | High | True | No |

Table 1: Exercise.

### Handwritten notes

$\alpha = 1$  $+0.7 \cdot \frac{1}{n} \sum \dots$

$-0.15 \cdot \frac{1}{n} \sum$

$q_0 \rightarrow \begin{bmatrix} d_1 \\ d_2 \end{bmatrix} D_r \rightarrow q_n = \alpha q_0 + \beta \dots + \dots$

$\begin{bmatrix} d_{n+1} \\ d_n \end{bmatrix} D_{nr}$

$\vec{x} = [x_1 \ x_2 \ x_3 \ x_4 \ \dots x_n]$

$P(x_1, x_2) = P(x_1) \cdot P(x_2)$

$P(x_1, x_2, x_3) = P(x_1, x_2) \cdot P(x_3)$

$q = [2 \ 3 \ 0 \ 0 \ 1 \ 1]$

$doc_1 = [2 \ 2 \ 1 \ 0 \ 0 \ 0]$

$doc_2 = [0 \ 1 \ 0 \ 0 \ 1 \ 1 \ 0]$

$p(x_i \mid x_{i+1}, \dots, x_n, C_k) = p(x_i \mid C_k)$

$P(x_i \mid x_{i+1}) = \frac{P(x_i, x_{i+1})}{P(x_{i+1})} = \frac{P(x_i) \cdot P(x_{i+1})}{P(x_{i+1})} = P(x_i)$

$q_m = \alpha q_0 + \beta \, doc_1 - \gamma \, doc_2$

$= 1 \cdot q_0 + 0.75 \, doc_1 - 0.15 \, doc_2$

$= [3.5 \ 4.25 \ 0.75 \ -0.25 \ 0.75 \ 1]$

$\cos \sin(q_m, \dots) \in [-1, 1]$

$P(Yes \mid Rainy, Mild, Normal, True)$

$\propto P(Yes) \cdot P(Rainy \mid Yes) \cdot P(Mild \mid Yes) \cdot P(Normal \mid Yes) \cdot P(Yes) \cdot P(True \mid Yes)$

$\propto \frac{9}{14} \cdot \frac{2}{9} \cdot \frac{4}{9} \cdot \frac{6}{9}$

$\cdot \frac{3}{9} = 0.014$

$P(W_0 \mid \dots) \propto \frac{5}{14} \cdot \frac{3}{5} \cdot \frac{2}{5}$

$\cdot \frac{1}{5} \cdot \frac{3}{5} = 0.010$

$P(Yes \mid \dots) = \frac{0.014}{0.014 + 0.010}$

$= 57.84\%$

$P(W_0 \mid \dots) = 42.16\%$