# Vector Space Classification (Chapter 14)

**Algorithm 1 (Rocchio classification)**
1: **function** TRAIN-ROCCHIO($\mathbb{C}, \mathbb{D}$)
2:   **for all** $c_j \in \mathbb{C}$ **do**
3:     $D_j \leftarrow \{d : \langle d, c_j \rangle \in \mathbb{D}\}$
4:     Centroid $\vec{\mu_j} \leftarrow \frac{1}{|D_j|} \sum_{d \in D_j} \vec{v}(d)$
5:   **end for**
6:   **return** $\{\vec{\mu_1}, \ldots, \vec{\mu_J}\}$
7: **end function**
8:
9: **function** APPLY-ROCCHIO($\{\vec{\mu_1}, \ldots, \vec{\mu_J}\}, d$)
10:   **return** $\arg\min_j |\vec{\mu_j} - \vec{v}(d)|$
11: **end function**

**Algorithm 2 ($k$ nearest neighbor classification)**
1: **function** TRAIN-KNN($\mathbb{C}, \mathbb{D}$)
2:   $\mathbb{D}' \leftarrow$ PREPROCESS($\mathbb{D}$)
3:   $k \leftarrow$ SELECT-$K$($\mathbb{C}, \mathbb{D}'$)
4:   **return** $\mathbb{D}', k$
5: **end function**
6:
7: **function** APPLY-KNN($\mathbb{C}, \mathbb{D}', k, d$)
8:   $S_k \leftarrow$ COMPUTENEARESTNEIGHBORS($\mathbb{D}', k, d$)
9:   **for all** $c_j \in \mathbb{C}$ **do**
10:     $p_j \leftarrow |S_k \cap c_j|/k$
11:   **end for**
12:   **return** $\arg\max_j p_j$
13: **end function**

k - odd

$\vec{\mu_1} = \begin{bmatrix} \frac{3}{2} & ; & \frac{1}{2} \end{bmatrix} = [1,5 ; 0,5]$

Rocchio:    q $\rightarrow$ class 2

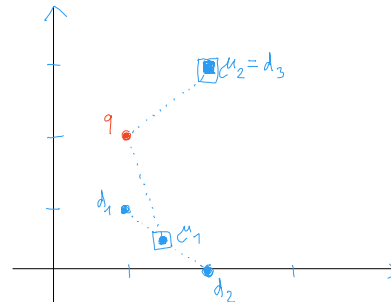$|\vec{\mu_1} - \vec{q}| = \sqrt{(1,5-1)^2 + (0,5-2)^2} = \sqrt{2,5}$

$|\vec{\mu_2} - \vec{q}| = \ldots = \sqrt{2}$

1NN:    k=1 ,    $\rightarrow$ class 1

$|\vec{d_1} - \vec{q}| = \ldots = \sqrt{1}$

$|\vec{d_2} - \vec{q}| = \ldots = \sqrt{5}$

$|\vec{d_3} - \vec{q}| = \ldots = \sqrt{2}$