# K means within-point scatter deviation
——Elements of Statistical Learning Chapter 14 Reading Note

## Shuai

## December 29, 2013

The overall within-point dissimilarity can be replaced with the summation of dissimilarity between each point in the cluster and the cluster mean, meaning:

$$W(C) \quad = \quad \frac{1}{2}\sum_{k=1}^{K}\sum_{C(i)=k}\sum_{C(i')=k}||x_i - x_{i'}||^2 \tag{1}$$

$$= \quad \sum_{k=1}^{K}N_k\sum_{C(i)=k}||x_i - \bar{x_k}||^2 \tag{2}$$

where, there are K clusters; $C(i) = k$ means the label given to sample i is k; $N_k$ means the overall number of points in cluster k; $x_i$ means sample i.

In this doc, detailed proof will be presented. During the proof, norm will be just written as square.

For expression 1:

For clarity, we replace $\sum_{C(i)=k}\sum_{C(i')=k}$ with $\sum_{i}\sum_{j}$.

$$\sum_{C(i)=k}\sum_{C(i')=k}||x_i - x_{i'}||^2 \quad = \quad \sum_{i}\sum_{j}||x_i - x_{i'}||^2 \tag{3}$$

$$= \quad \sum_{i}\sum_{j}(x_i^2 + x_j^2) - \sum_{i}\sum_{j}x_i x_j \tag{4}$$

$$= \quad N_k\sum_{i}x_i^2 + N_k\sum_{j}x_j^2 - 2\sum_{i}\sum_{j}x_i x_j \tag{5}$$

$$= \quad 2(N_k\sum_{i}x_i^2 + \sum_{i}\sum_{j}x_i x_j) \tag{6}$$

For expression 2:

For clarity, we replace $\sum_{C(i)=k}$ with $\sum_i$.

$$\sum_{C(i)=k} ||x_i - \bar{x_k}||^2 = \sum_i ||x_i - \bar{x_k}||^2 \tag{7}$$

$$= \sum_i (x_i^2 + \bar{x_k}^2) - \sum_i x_i \bar{x_k} \tag{8}$$

$$= \sum_i x_i^2 + N_k \bar{x_k}^2 - \sum_i x_i \bar{x_k} \tag{9}$$

$$= \sum_i x_i^2 + N_k \bar{x_k}^2 - 2N_k \bar{x_k}^2 \tag{10}$$

$$= \sum_i x_i^2 - N_k \bar{x_k}^2 \tag{11}$$

$$= \sum_i x_i^2 - \sum_i x_i \sum_j x_j / N_k \tag{12}$$

$$\tag{13}$$

Since $\sum_i x_i \sum_j x_j = \sum_i \sum_j x_i x_j$, proof is done.