

K-Means

Ha Phuong

January 21, 2022

1 Formula

Suppose we have a data set $\{x_1, \dots, x_N\}$ consisting of N observations of a random D -dimensional Euclidean variable x . We introducing a set of D -dimensional vectors μ_k , where $k = 1, \dots, K$, in which μ_k is a prototype associated with the k th cluster.

An objective function

$$J = \sum_{n=1}^N \sum_{k=1}^K r_{nk} \|x_n - \mu_k\|^2 \quad (1)$$

$$\text{with } \begin{cases} r_{n,k} \in \{0, 1\}, \forall n, k \\ \sum_{j=1}^K r_{n,j} = 1 \forall n \end{cases}$$

Our goal is to find values for the r_{nk} and the μ_k so as to minimize J . Now consider the optimization of the μ_k with the r_{nk} held fixed. The objective function J is a quadratic function of μ_k , and it can be minimized by setting its derivative with respect to μ_k to zero giving

$$2 \sum_{n=1}^N r_{nk} (x_n - \mu_k) = 0 \quad (2)$$

$$\mu_k = \frac{\sum_{n=1}^N r_{nk} x_n}{\sum_{n=1}^N r_{nk}} \quad (3)$$

So μ_k equal to the mean of all of the data points x_n assigned to cluster k