

Derivation of K-means update step

To make matters simple, we define a loss for a cluster rather than for all clusters:

$$\begin{aligned}\text{Loss}(z, \mu) &= \sum_{i=1}^n \|\phi(x_i) - \mu\|^2 \\ &= \sum_{i=1}^n (\phi(x_i) - \mu)^T (\phi(x_i) - \mu)\end{aligned}$$

where μ is the feature vector of the cluster's centroid, x_i is a data example and n is the number of data examples. Then, we can compute the gradient with respect to μ_j , which is an element of μ :

$$\begin{aligned}\frac{\partial}{\partial \mu_j} \text{Loss}(z, \mu) &= \frac{\partial}{\partial \mu_j} \sum_{i=1}^n \|\phi(x_i) - \mu\|^2 \\ &= \sum_{i=1}^n \frac{\partial}{\partial \mu_j} \|\phi(x_i) - \mu\|^2 \\ &= \sum_{i=1}^n \frac{\partial}{\partial \mu_j} \{(\phi(x_i) - \mu)^T (\phi(x_i) - \mu)\} \\ &= \sum_{i=1}^n \frac{\partial}{\partial \mu_j} (\phi(x_i)^T \phi(x_i) - 2\phi(x_i)^T \mu + \mu^T \mu) \\ &= \sum_{i=1}^n \left\{ \frac{\partial}{\partial \mu_j} (\phi(x_i)^T \phi(x_i)) + \frac{\partial}{\partial \mu_j} (-2\phi(x_i)^T \mu) + \frac{\partial}{\partial \mu_j} (\mu^T \mu) \right\} \\ &= \sum_{i=1}^n \left\{ \frac{\partial}{\partial \mu_j} \left(\sum_k^d \phi(x_i)_k^2 \right) + \frac{\partial}{\partial \mu_j} \left(-2 \sum_k^d \phi(x_i)_k \mu_k \right) + \frac{\partial}{\partial \mu_j} \left(\sum_k^d \mu_k^2 \right) \right\} \\ &\quad ; d \text{ is the dimension of } \phi(x_i) \text{ and } \mu \\ &= \sum_{i=1}^n \left\{ 0 + \frac{\partial}{\partial \mu_j} (-2\phi(x_i)_j \mu_j) + \frac{\partial}{\partial \mu_j} (\mu_j^2) \right\} \\ &= \sum_{i=1}^n \{-2\phi(x_i)_j + 2\mu_j\} \\ &= 2 \sum_{i=1}^n \{\mu_j - \phi(x_i)_j\}\end{aligned}$$

Now, we can get the optimal μ_j where $\frac{\partial}{\partial \mu_j} \text{Loss}(z, \mu) = 0$ (0 is a scalar):

$$\begin{aligned}\frac{\partial}{\partial \mu_j} \text{Loss}(z, \mu) &= 2 \sum_{i=1}^n (\mu_j - \phi(x_i)_j) \\ &= 0\end{aligned}$$

$$\implies \sum_{i=1}^n \phi(x_i)_j = \sum_{i=1}^n \mu_j$$

$$\implies \sum_{i=1}^n \phi(x_i)_j = n\mu_j$$

$$\therefore \mu_j = \frac{1}{n} \sum_{i=1}^n \phi(x_i)_j$$