# Lecture 16 : Clustering

# Clustering

Given points $\underline{x}_i \in \mathbb{R}^p$, $i = 1, 2, \ldots, n$, group them into clusters $C_1, \ldots, C_K$ so that

- any two points in same cluster are "close"

- any two points in different clusters are 'distant'

- "hard" clustering   — each $\underline{x}_i$ is in one and only one cluster
- a "soft" clustering   — each $\underline{x}_i$ may be in multiple clusters

## Different cluster analysis results on "mouse" data set:



Original Data          k-Means Clustering          EM Clustering

For instance, let $\underline{\mu}_k \in \mathbb{R}^p$, $k = 1, \ldots, K$, be a prototypcal point for the $k^{th}$ cluster.

Then we want to assign the $\underline{x}_i$'s to clusters so that the sum of distances (or squared distances) from each data point to its assigned cluster to be minimized

# K-means clustering

$$\{\hat{C}_1, ..., \hat{C}_k\} = \underset{\{C_1, ..., C_k\}}{\text{argmin}} \sum_{k=1}^{K} \sum_{\underline{x}_i \in C_k} \| \underline{x}_i - \mu_k \|^2$$

$$= \underset{\{C_1, ..., C_k\}}{\text{argmin}} \underbrace{\sum_{k=1}^{K} \frac{1}{|C_k|} \sum_{\underline{x}_i, \underline{x}_j \in C_k} \| \underline{x}_i - \underline{x}_j \|^2}_{=: \text{obj} \ (\text{objective})}$$

where $\mu_k = \frac{1}{|C_k|} \sum_{\underline{x}_i \in C_k} \underline{x}_i$ = cluster center

# Lloyd's algorithm

Start with initial set of $K$ means $\underline{\mu}_1^{(1)}, \underline{\mu}_2^{(1)}, ..., \underline{\mu}_K^{(1)}$ (centers)

for $t = 1, 2, 3, ...$

    for $i = 1, 2, ..., n$

        # find nearest mean (center) to $\underline{X}_i$

$$\hat{k}_i = \underset{k}{\arg\min} \| \underline{X}_i - \underline{\mu}_k^{(t)} \|^2$$

    end

    for $k = 1, ..., K$

$$\hat{C}_k^{(t+1)} = \{ \underline{X}_i : \hat{k}_i = k \} \quad \text{(set cluster estimates)}$$

$$\underline{\mu}_k^{(t+1)} = \frac{1}{|\hat{C}_k^{(t+1)}|} \sum_{i \in \hat{C}_k^{(t+1)}} \underline{X}_i \quad \text{(find each cluster mean)}$$

    end

end

Generally this algorithm converges, but not to the optimal clustering — results depend on initial clustering

Handwritten annotations: (top) $E$ step, M — step; (middle) $E$ — step, M — step

In panel (a): $\hat{\mu}_1^{(1)}$, $\hat{\mu}_2^{(1)}$

**Figure 9.1** Illustration of the $K$-means algorithm using the re-scaled Old Faithful data set. (a) Green points denote the data set in a two-dimensional Euclidean space. The initial choices for centres $\mu_1$ and $\mu_2$ are shown by the red and blue crosses, respectively. (b) In the initial E step, each data point is assigned either to the red cluster or to the blue cluster, according to which cluster centre is nearer. This is equivalent to classifying the points according to which side of the perpendicular bisector of the two cluster centres, shown by the magenta line, they lie on. (c) In the subsequent M step, each cluster centre is re-computed to be the mean of the points assigned to the corresponding cluster. (d)–(i) show successive E and M steps through to final convergence of the algorithm.

N=200, K=3
K-means with random initialization

N=200, K=3
K-means with random initialization

$\Longleftarrow$ Importance of good initialization

**Forgy**: choose $k$ random $\underline{x}_i$'s as initial points

**K means++**:

- Choose one $\underline{x}_i$ uniformly at random as $\mu_1^{(1)}$

- for $k = 2, 3, \ldots, K$

  find remaining data point furthest from any $\mu_j^{(1)}$, $j < k$, selected so far, make this $\mu_k^{(1)}$

Provably good! Solution found by Kmeans++ may not be optimal, but

$$obj\,(kmeans++) \leq O(\log K)\, obj\,(optimal)$$

**original data (with ground truth)**

**original data (with kmeans clustering)**

⟸ k-means not always sensible.

Can we do better with kernels?

# Regular k-means

Start with initial set of K means

$$\underline{\mu}_1^{(1)}, \underline{\mu}_2^{(1)}, \ldots, \underline{\mu}_K^{(1)}$$

for $t = 1, 2, 3, \ldots$

  for $i = 1, 2, \ldots, n$

    # find nearest mean to $\underline{x}_i$

$$\hat{k}_i = \underset{k}{\arg\min} \; \| x_i - \underline{\mu}_k^{(t)} \|^2$$

  end

for $k = 1, \ldots, K$

$$\hat{C}_k^{(t+1)} = \{ \underline{x}_i : \hat{k}_i = k \}$$

$$\mu_k^{(t+1)} = \frac{1}{|\hat{C}_k^{(t)}|} \sum_{i \in \hat{C}_k^{(t)}} x_i$$

# Towards Kernel k-means

$$\hat{k}_i = \underset{k}{\arg\min} \; \| \phi(\underline{x}_i) - \phi(\underline{\mu}_k^{(t)}) \|_2^2$$

$$= \underset{k}{\arg\min} \; \underbrace{\phi(\underline{x}_i)^\top \phi(\underline{x}_i)}_{k(\underline{x}_i, \underline{x}_i)} - 2\phi(\underline{x}_i)^\top \phi(\underline{\mu}_k^{(t)}) + \phi(\underline{\mu}_k^{(t)})^\top \phi(\underline{\mu}_k^{(t)})$$

$$\phi(\underline{\mu}_k^{(t)}) = \frac{1}{|C_k^{(t)}|} \sum_{\underline{x}_j \in C_k^{(t)}} \phi(\underline{x}_j)$$

$$\Rightarrow \phi(\underline{x}_i)^\top \phi(\underline{\mu}_k^{(t)}) = \frac{1}{|C_k^{(t)}|} \sum_{\underline{x}_j \in C_k^{(t)}} \phi(\underline{x}_i)^\top \phi(\underline{x}_j)$$

$$= \frac{1}{|C_k^{(t)}|} \sum_{\underline{x}_j \in C_k^{(t)}} k(\underline{x}_i, \underline{x}_j)$$

$$\phi(\underline{\mu}_k^{(t)})^\top \phi(\underline{\mu}_k^{(t)}) = \frac{1}{|C_k^{(t)}|^2} \sum_{\substack{x_j, x_j' \\ \in C_k^{(t)}}} \phi(\underline{x}_j)^\top \phi(\underline{x}_{j'})$$

$$= \frac{1}{|C_k^{(t)}|^2} \sum_{\substack{x_j, x_j' \\ \in C_k^{(t)}}} k(\underline{x}_j, \underline{x}_{j'})$$

$x_i \in \mathbb{R}$

$0$

Regular k-means

Kernel k-means, $\phi(x_i) = \begin{bmatrix} x_i \\ x_i^2 \end{bmatrix}$

$x_i^2$

$x_i$

# Kernel k-means

Start with initial set of $K$ cluster
  assignments $\hat{C}_1^{(1)}, \hat{C}_2^{(1)}, \ldots, \hat{C}_K^{(1)}$

for $t = 1, 2, 3, \ldots$

  for $i = 1, 2, \ldots, n$

    \# find nearest mean to $\underline{x}_i$

$$\hat{k}_i = \underset{k}{\arg\min} \; \frac{1}{|C_k^{(t)}|^2} \sum_{\underline{x}_j, \underline{x}_{j'} \in C_k^{(t)}} k(\underline{x}_j, \underline{x}_{j'}) - \frac{2}{|C_k^{(t)}|} \sum_{\underline{x}_j \in C_k^{(t)}} k(\underline{x}_i, \underline{x}_j) + k(\underline{x}_i, \underline{x}_i)$$

  end

  for $k = 1, \ldots, K$

$$\hat{C}_k^{(t+1)} = \{\underline{x}_i : \hat{k}_i = k\}$$

  end

end



## k-means Vs. Kernel k-means