

Expectation Maximization (EM)

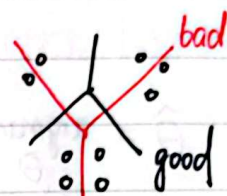
Unsupervised Learning: Kmeans & Gaussian Mixture Models (GMM)

考虑下述任务:

Clustering (informal goals): Goal: data 为 unlabeled, 将它们分成数组, 每组中点相似。

这个任务极为常用!

问题在于: 如何建模出 partitions 的好坏?

Input: unlabeled data: $D = \{x^{(i)}\}_{i=1}^N$, $x^{(i)} \in \mathbb{R}^M$ \Rightarrow cluster centers: $C = [c_1, \dots, c_k]$ $c_j \in \mathbb{R}^M$ cluster assignments: $z = [z^{(1)}, z^{(2)}, \dots, z^{(N)}]$, $z^{(i)} \in \{1, \dots, k\}$

$$\text{Objective: } \hat{C} = \underset{C}{\operatorname{argmin}} \sum_{i=1}^N \min_j \|x^{(i)} - c_j\|_2^2$$

$$= \underset{C}{\operatorname{argmin}} \sum_{i=1}^N \min_{z^{(i)}} \|x^{(i)} - c_{z^{(i)}}\|_2^2$$

$$\Leftrightarrow \hat{C}, \hat{z} = \underset{C, z}{\operatorname{argmin}} \sum_{i=1}^N \|x^{(i)} - c_{z^{(i)}}\|_2^2 \quad \mathcal{J}(C, z)$$

Algorithm: 1) Initialize $C = \{c_1, \dots, c_k\}$. * 随机k个样本作初始均值向量

2) Repeat until Convergence:

a) for i in $\{1, \dots, N\}$:

$$z^{(i)} \leftarrow \underset{j}{\operatorname{argmin}} (\|x^{(i)} - c_j\|_2)^2 \leftarrow \text{每个 } x_i \text{ 选最近 } c_j$$

b) for j in $\{1, \dots, k\}$:

$$c_j \leftarrow \underset{c_j}{\operatorname{argmin}} \sum_{i: z^{(i)}=j} (\|x^{(i)} - c_j\|_2)^2$$

$$\text{此优化 i.e.: } \mu_i = \frac{1}{|C_i|} \sum_{x \in C_i} x, \quad C_i \leftarrow \mu_i$$

* 不止随机, 也可 Furthest Point Heuristic: 先随机选一个, 之后每一个选取尽可能离先前选的点们远! (但 outlier 处理是问题)

也可 K-means++: 先随机选 c_1 , 则 $j = 2, \dots, K$, pick c_j from distribution:

$$P(c_j = x^{(i)}) \propto \min_{j' < j} \|x^{(i)} - c_{j'}\|^2, \text{ where } x^{(i)} \text{ are not picked before}$$

