

### 3. K-means

~~3. Variational methods: from K-means (or EM algorithm) to~~

### 3. K-means and EM.

Let us consider a task of putting a set of data points

$\{\vec{x}^{(n)}\}_{n=1}^N$  into  $K$  clusters.

#### 1) Classical K-Means Algo:

1- Set  $K$  means  $\{\vec{m}^{(k)}\}$  to random values

2- For each  $x^{(n)}$ :  $\hat{k}^{(n)} = \underset{k}{\operatorname{argmin}} \{d(\vec{m}^{(k)}, \vec{x}^{(n)})\}$ ,  $d$ -distance  
predicted cluster  $k$

3- Set  $r_k^{(n)} = I(\hat{k}^{(n)} = k)$  and update  $\vec{m}^{(k)} \leftarrow \frac{\sum_{n=1}^N r_k^{(n)} \vec{x}^{(n)}}{\sum_{n=1}^N r_k^{(n)}} \quad \forall k$

4- Repeat 2 and 3 until convergence.

#### 2) Problems of K-Means

- too "hard" algorithm: points are assigned ~~to~~ to exactly one cluster, but ~~border~~ <sup>border</sup> points should be present better be presented in multiple

- the classic distance couldn't represent a narrow cluster with K-Means.

- doesn't have any ~~info~~ represent. of the weight and breadth of each cluster.

#### 3) Soft K-Means Clustering: [Improvement] $\beta$ -hyper-param.

1- set  $K$  means  $\{\vec{m}^{(k)}\}$  to random values

2- For each  $x^{(n)}$  and  $k$ :  $r_k^{(n)} = \frac{\exp(-\beta d(\vec{m}^{(k)}, \vec{x}^{(n)}))}{\sum_{k'} \exp(-\beta d(\vec{m}^{(k')}, \vec{x}^{(n)}))}$  - softmax over distances.

3- update  $\vec{m}^{(k)} \leftarrow \frac{\sum_{n=1}^N r_k^{(n)} \vec{x}^{(n)}}{\sum_{n=1}^N r_k^{(n)}}$

4- Repeat 2 and 3 until convergence.

4) Further soft k-means enhancements: ~~very~~

1- ~~set~~ Assign Initialize  $\{\vec{m}^{(k)}\}$  - cluster centers,  
 $\{\tau_k\}$  - cluster importances,  $\{\sigma_k^2\}$  - cluster dispersions.

2- For each  $x^{(n)}$  and  $k$ :  $r_k^{(n)} = \frac{\tau_k \left( \frac{1}{\sqrt{2\pi\sigma_k^2}} \right)^d \exp\left(-\frac{1}{2\sigma_k^2} d(\vec{m}^{(k)}, \vec{x}^{(n)})\right)}{\sum_{k'=1}^K \text{---}} \leftarrow \text{normalize. constant}$

3- Adjust  $m^{(k)}$ ,  $\tau_k$  and  $\sigma_k^2$  to match the data:

$$\vec{m}^{(k)} \leftarrow \frac{\sum_n r_k^{(n)} \vec{x}^{(n)}}{\sum_n r_k^{(n)}}, \quad \sigma_k^2 \leftarrow \frac{\sum_n r_k^{(n)} (\vec{x}^{(n)} - \vec{m}^{(k)})^2}{\sum_n r_k^{(n)}}, \quad \tau_k \leftarrow \frac{R^{(k)}}{\sum_k R^{(k)}}$$

$$\vec{m}^{(k)} \leftarrow \frac{\sum_n r_k^{(n)} \vec{x}^{(n)}}{R^{(k)}}, \quad \sigma_k^2 \leftarrow \frac{\sum_n r_k^{(n)} (\vec{x}^{(n)} - \vec{m}^{(k)})^2}{d \cdot R^{(k)}}, \quad \tau_k \leftarrow \frac{R^{(k)}}{\sum_k R^{(k)}},$$

where  $R^{(k)} = \sum_n r_k^{(n)}$ ,  $d$  - dimensionality of  $\vec{x}$ .

5) General Idea: EM.

Let  $X$  denote a set of observed data,  $Z$  - a set of latent data,

$\theta$  - a vector of unknown parameters.

The EM-algorithm is a way to find MLE of the  $\theta$ ,  
in other words to maximize  $p(X|\theta) = \int p(X, z|\theta) dz$

1- (Expectation step)  $Q(\theta|\theta^{(t-1)}) = \mathbb{E}_{z|X, \theta^{(t-1)}} [\log p(X, z|\theta)]$

2- (Maximization step)  $\theta^{(t)} = \arg \max_{\theta} Q(\theta|\theta^{(t-1)})$ .

In example 4)  $\theta^{(t)} = \{\vec{m}^{(k)}, \sigma_k^2, \tau_k\}_{k=1}^K$ ,  $z = \{r_k^{(n)}\}_{k,n=1}^{K,N}$

c) A fatal flaw of maximum likelihood

E.g., let's look at algorithm 4). If  $\vec{m}^{(k)} \equiv \vec{x}^{(n)}$  during iterations, then  $\sigma_k^2 \leftarrow 0$ . If  $\sigma_k^2$  is sufficiently small, then it becomes even smaller during iterations.