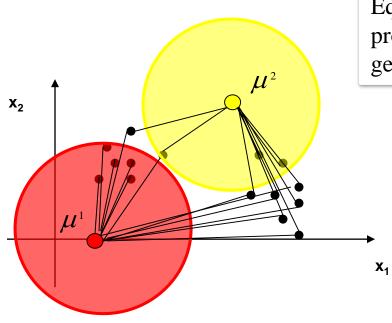


APPLIED MACHINE LEARNING

From soft K-means Clustering to Density Modeling of Data Clusters with Mixture of Gaussian PdF



Soft K-means Clustering (probabilistic interpretation)



Equivalent to computing the relative probability that the data point has been generated by the k-th cluster (d:norm-2).

$$r_i^k = \frac{e^{\left(-\beta \cdot d\left(\mu^k, x^i\right)\right)}}{\sum_{k'} e^{\left(-\beta \cdot d\left(\mu^{k'}, x^i\right)\right)}} \in [0, 1], \quad \beta \in \mathbb{R}_+$$

The likelihood of the k-th model is:
$$L(\mu^{k}; X) \sim \prod_{i=1}^{M} e^{\left(-\beta \cdot d(\mu^{k}, x^{i})\right)}$$

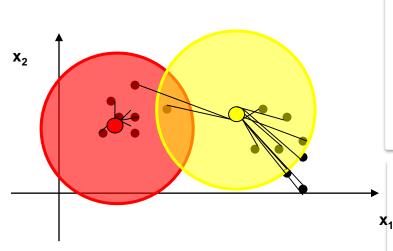
Assignment Step (~E-step)

Assume density under cluster k follows an (un-normalized) Gaussian pdf with $\sigma=1/\beta$ variance, centered on the cluster's centroid. Un-normalized Gaussian pdf are called Radial Basis Function – RBF.

We can derive the likelihood that each cluster has generated the dataset.



K-means Clustering (probabilistic interpretation)



The new centroid is closer to the datapoints after the update step,

$$\mu^{k} = \frac{\sum_{i} r_{i}^{k} x^{i}}{\sum_{i} r_{i}^{k}} \implies d(x^{i}, \mu^{k})$$

→ The likelihood of the k-th model increases

$$L(\mu^{k};X) = \prod_{i=1}^{M} e^{\left(-\beta \cdot d\left(\mu^{k},x^{i}\right)\right)}$$

Update Step (~M-step):

Update the position of the centroids.

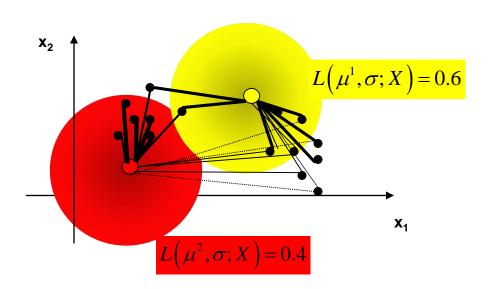
Modify the parameters to maximize likelihood of the pdf.



Soft-K-means (probabilistic interpretation)

Soft K-means is similar to fitting the data distribution with a *mixture of isotropic* (spherical) unnormalized Gaussian pdf-s and same variance (the stiffness).

Assignment-update ~ E-M on the parameters of each Gaussian to optimize the likelihood that the Gaussians represent the distribution of the datapoints.



Should we weight equivalently the likelihoods of each cluster?



From Soft-K-means to GMM

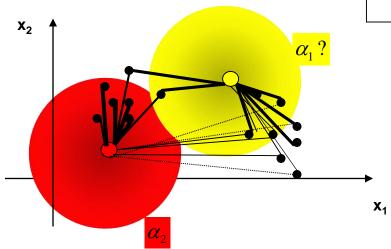
The responsibility factor gives a measure of the likelihood that cluster k generated the dataset.

 r_i^k : responsibility of cluster k for point x^i

$$r_{i}^{k} = \frac{\alpha_{k} p(x^{i}; \mu^{k}, \sigma)}{\sum_{k'} \alpha_{k} p(x^{i}; \mu^{k'}, \sigma)}, \qquad \alpha_{k} \in [0, 1]$$

 $p(x^i; \mu^k, \sigma) \in [0,1]$: Gauss pdf evaluated at x^i

Normalized over clusters: $\sum_{k} r_i^k = 1$



Relative importance of each of the K clusters.

It should give a measure of the likelihood that the Gaussian k (or cluster k) generated the whole dataset.

$$\alpha_k = \frac{\sum_{i} r_i^k}{\sum_{k} \sum_{i} r_i^k}$$



One step towards Gaussian Mixture Model with Spherical Gaussians

 r_i^k : responsibility of cluster k for point x^i

$$r_i^k = \frac{\alpha_k p(x^i; \mu^k, \sigma^k)}{\sum_{k'} \alpha_{k'} p(x^i; \mu^{k'}, \sigma^k)}, \qquad \alpha_k \in [0, 1]$$

 $p(x^i; \mu^k, \sigma) \in [0,1]$: Gauss pdf evaluated at x^i

Normalized over clusters: $\sum_{k} r_i^k = 1$

$$\mu^k = \frac{\sum_{i} r_i^k x^i}{\sum_{i} r_i^k}$$

$$\left(\sigma^{k}\right)^{2} = \frac{\sum_{i} r_{i}^{k} \left\|x^{i} - \mu^{k}\right\|^{2}}{N \cdot \sum_{i} r_{i}^{k}}$$

This fits a mixture of *spherical* Gaussians. The variance of each Gauss pdf fits the spread of the data around its mean.

$$lpha_k = rac{\displaystyle\sum_i r_i^k}{\displaystyle\sum_k \displaystyle\sum_i r_i^k}$$



From spherical to diagonal Gaussian pdf-s.

Update Step (M-Step)

$$\left(\sigma_{j}^{k}\right)^{2} = \frac{\sum_{i} r_{i}^{k} \left(x_{j}^{i} - \mu_{j}^{k}\right)^{2}}{\sum_{i} r_{i}^{k}}$$

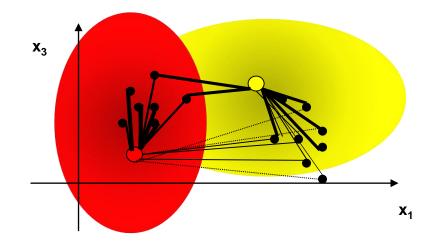
 r_i^k : responsibility of cluster k for point x^i

$$r_i^k = \frac{\alpha_k p(x^i; \mu^k, \sigma_j^k)}{\sum_{k'} \alpha_{k'} p(x^i; \mu^{k'}, \sigma_j^{k'})}$$

j = 1,...N: dimension of dataset

 $p(x^i; \mu^k, \sigma^k) \in [0,1]$: Gauss pdf evaluated at x^i

Normalized over clusters: $\sum_{k} r_i^k = 1$



One covariance element per dimension, aligned with the axes of the original frame of reference.



Clustering with Mixture of Gaussians

Likelihood of the mixture of Gaussians: $L\left(\Theta = \left\{\alpha_k, \mu^k, \Sigma^k\right\}_{k=1}^K; x\right) = \sum_{k=1}^K \alpha_k \cdot p\left(x; \mu^k, \Sigma^k\right)$

with $p(X; \mu^k, \Sigma^k) \sim \prod_{i=1}^M e^{-(x^i - \mu_k)^T (\Sigma^k)^{-1} (x^i - \mu_k)}$ (omit normalization factor)

 μ^k, Σ^k : mean and covariance matrix of Gaussian k

The mixing Coefficients are normalized.

$$\sum_{k=1}^{K} \alpha_k = 1$$

$$\mathbf{x}_3$$

$$\alpha_1 = 0.2$$

$$\alpha_2 = 0.8$$

$$\alpha_k \sim \frac{1}{M} \sum_{i=1}^{M} \frac{p(x^i; \mu^k, \Sigma^k)}{\sum_{k'} p(x^i; \mu^{k'}, \Sigma^{k'})}$$



E-M Steps for GMM

Initialization:

The priors $\alpha_1,...,\alpha_k$ can be uniform for starters.

The means $\mu^{\scriptscriptstyle 1},...,\mu^{\scriptscriptstyle K}$ can be initialized with K-means.

Calculate the initial value of the likelihood

$$p(X | \Theta^{(t)}) = \prod_{i=1}^{M} \sum_{k=1}^{K} \alpha_{k(t)} p(x^{i}; \mu^{k(t)}, \Sigma^{k(t)})$$

Expectation Step (E-step):

Evaluate responsibilities of each cluster k over each sample x^i using current parameters

$$r_i^k = \frac{\alpha_k p(x^i; \mu^k, \Sigma^k)}{\sum_{k'} \alpha_{k'} p(x^i; \mu^{k'}, \Sigma^{k'})}$$



Maximization (Update step) Step (M-step):

Recompute the means, covariances matrices and prior probabilities so as to maximize the log – likelihood of the current estimate : $\log \left(L\left(\Theta^{(t)} \mid X\right) \right)$

$$\mu_k^{(t+1)} = \frac{1}{M_k} \sum_{i=1}^M r_k^i x_i$$

$$\mu_k^{(t+1)} = \frac{1}{M_k} \sum_{i=1}^M r_k^i x_i$$

$$\sum_k^{(t+1)} = \frac{1}{M_k} \sum_{i=1}^M r_k^i (x_i - \mu_k^{(t+1)}) (x_i - \mu_k^{(t+1)})^T$$

$$lpha_{_{k}}^{^{(t+1)}}=rac{M_{_{k}}}{M}$$

$$\alpha_k^{(t+1)} = \frac{M_k}{M}$$
 where: $M_k = \sum_{i=1}^M r_k^i$

The E and M steps alternate until the log-likelihood reaches a plateau.