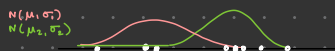


Oct 4th

Hard Clustering - 1 pt \rightarrow 1 cluster

Soft Clustering - helps w/ overlapping clusters, to give each pt a probability of being in a given cluster.



For each pt, we can compute probability of it being generated from any cluster.

MIXTURE MODEL

$$P(X=x) = \sum_{j=1}^K P(c_j) P(X=x | c_j)$$

MIXTURE PROPORTION

K clusters
so probability of $X=x$ is weighted by the prob. of each cluster and $X=x$ given the cluster

ASSUMING discrete

PROBABILITY OF SEEING x WHEN SAMPLING FROM c_j

$$\left\{ \frac{1}{\sigma_j \sqrt{2\pi}} e^{-\frac{1}{2} \left(\frac{x - \mu_j}{\sigma_j} \right)^2} \right\}$$

NORMAL DISTRIBUTION μ_j, σ_j

GMM

$$P(X=x | c_j) \sim N(\mu, \sigma)$$

GMM Clustering

GOAL: find the GMM that maximizes the probability of seeing the data we have.

Probability of seeing data we saw (assume iid) = the product of the probabilities of observing each data point.

PARAMETERS WE HAVE: $P(c_j), \mu_j, \sigma_j \quad \forall K$ clusters. $\rightarrow \theta = \{ \mu_1, \dots, \mu_K, \sigma_1, \dots, \sigma_K, P(c_1), \dots, P(c_K) \}$

$$\theta^* = \arg \max_{\theta} \prod_{i=1}^n \sum_{j=1}^K P(c_j) P(x_i | c_j)$$

(argument that maximizes θ)

NOTE: taking log = transform, doesn't change optimal pt's

$$\ell(\theta^*) = \sum_{i=1}^n \log \sum_{j=1}^K P(c_j) P(x_i | c_j)$$

Solve for

$$\frac{d}{d\mu} \ell(\theta) = 0$$

$$\frac{d}{d\sigma} \ell(\theta) = 0$$

$$\hat{\mu}_j = \frac{\sum_{i=1}^n P(c_j | x_i) x_i}{\sum_{i=1}^n P(c_j | x_i)}$$

$$\hat{\sigma}_j = \frac{\sum_{i=1}^n P(c_j | x_i) (x_i - \hat{\mu}_j)^T (x_i - \hat{\mu}_j)}{\sum_{i=1}^n P(c_j | x_i)}$$

$$\hat{P}(c_j) = \frac{1}{n} \sum_{i=1}^n P(c_j | x_i)$$

mle estimates \Rightarrow

$$P(c_j | x_i) = \frac{P(x_i | c_j) P(c_j)}{\sum_{i=1}^K P(c_i) P(x_i | c_i)}$$

BAYES' RULE

LOCAL PROBLEM \rightarrow NEED ESTIMATORS to estimate.

EXPECTATION-MAXIMIZATION ALGO

1. Start w/ random θ
2. Compute $P(c_j | x_i)$ for all x_i using θ
3. Compute / update θ from $P(c_j | x_i)$
4. Repeat 2-3 until convergence.

ESTIMATE WILL NOT GET WORSE.

MUST CONVERGE

NOT NECESSARILY THE OPTIMAL SOLN UPON CONVERGENCE

Clustering - a group of clusters output by a clustering algo

Cluster - a group of points.

Clustering Aggregation

GOALS: Compare clusterings, Combine info from mult clusterings to create a new clustering.

To compare clusterings, for each pair of points are they clustered to each other in each k clusters.

→ Disagreement Distance

$$D(P, C) = \sum_{x, y} \mathbb{I}_{P, C}(x, y) \quad \text{where} \quad \mathbb{I}_{P, C}(x, y) = \begin{cases} 1 & \text{if } P \neq C \text{ disagreed on which clusters } x, y \text{ belong to.} \\ 0 & \end{cases}$$

distance function on 2 clusterings P & C for all pairs of points. [IT'S A METRIC!]

$\sum_{i=1}^m D(C^i, C)$ generate C^* from set of clusterings and look to minimize

NOTE: PROBLEM IS EQUIVALENT TO CLUSTERING.

Benefits

- can identify best num of clusters
- can handle/detect outliers
- improve robustness of clustering algo
- privacy preserving clustering

BUT IT IS

NP-HARD

MAJORITY RULE → ISSUE IS DOES NOT ALWAYS PRODUCE A CLUSTERING.