

Lec 24 EM algo (soft decision)

Recall: GMM & EM algo:

① Initialization:

Arbitrary bump membership $\{B_j\}_{j=1}^k$

② Given $\{B_j\}$ estimate $\Omega = \{w_j, \mu_j, \Sigma_j\}_{j=1}^k$

$$w_j = \frac{N_j}{N} = \frac{|B_j|}{N}$$

$$\mu_j = \frac{1}{N_j} \sum_{\mathbf{x}_n \in B_j} \mathbf{x}_n$$

$$\Sigma_j = \frac{1}{N_j} \sum_{\mathbf{x}_n \in B_j} (\mathbf{x}_n - \mu_j)(\mathbf{x}_n - \mu_j)^T$$

③ Given Ω , estimate $\{B_j\}$

$$j^* = \underset{j}{\operatorname{argmax}} N(\mathbf{x}_n; \mu_j, \Sigma_j) w_j$$

Repeat step 2 and 3 until converge.

1. EM algo (soft decision)

(Assume a datapoint is divisible)

→ At each step, let γ_{nj} denote the fraction (prob.) of \mathbf{x}_n that belongs to B_j .

$$\therefore \sum_{j=1}^k \gamma_{nj} = 1, \quad \gamma_{nj} \geq 0 \quad (\gamma_{nj} : \text{No. of fraction go to bound } j)$$

↳ subproblem 1 becomes:

Given $\{\gamma_{nj}\}_{1 \leq n \leq N, 1 \leq j \leq k}$, update Ω

$$N_j = \sum_{n=1}^N \gamma_{nj}$$

↑
Effective (expected) No. of data points in B_j .

$$w_j = \frac{N_j}{N}$$

$$\underline{\mu}_j = \frac{1}{N_j} \sum_{n=1}^N \gamma_{nj} \underline{x}_n$$

$$\underline{\Sigma}_j = \frac{1}{N_j} \sum_{n=1}^N \gamma_{nj} (\underline{x}_n - \underline{\mu}_j)(\underline{x}_n - \underline{\mu}_j)^T$$

↳ Subproblem 2:

Given Ω , update $\{\gamma_{nj}\}$

$$\gamma_{nj} = \Pr(j | \underline{x}_n)$$

$$= \frac{P(\underline{x}_n | j) \cdot \Pr(j)}{\sum_{i=1}^K P(\underline{x}_n | i) \cdot \Pr(i)}$$

$$= \frac{\mathcal{N}(\underline{x}_n; \underline{\mu}_j, \underline{\Sigma}_j) w_j}{\sum_{i=1}^K \mathcal{N}(\underline{x}_n; \underline{\mu}_i, \underline{\Sigma}_i) w_i}$$

↳ Summary of soft-EM:

① Initialization:

arbitrary $\{\gamma_{nj}\}_{1 \leq n \leq N}$
 $1 \leq j \leq K$

(e.g. k-means \Rightarrow binary γ_{nj})

② Given $\{\gamma_{nj}\}$, update Ω

③ Given Ω , update $\{\gamma_{nj}\}$

Repeat step ② until convergence

{ convergence guaranteed

hidden variables $\{\gamma_{nj}\}$ can be thrown out at the end

2. Relation to $E_{in}(\Omega)$:

(Hard decision case)

$$\hookrightarrow \text{Recall: } E_{in}(\Omega) = -\log \hat{p}_{\Omega}(\mathcal{D}) = -\log \prod_{n=1}^N \hat{p}(x_n)$$

$$= -\log \prod_{n=1}^N \underbrace{\sum_{j=1}^K w_j \mathcal{N}(x_n; \mu_j, \Sigma_j)}_{\text{GMM}}$$

In subproblem 1, $\{B_j\}$ is given

Let j_n be the bump that x_n belongs to.

$$\hat{p}(x_n | x_n \in B_{j_n}) = \mathcal{N}(x_n; \mu_{j_n}, \Sigma_{j_n})$$

$$\begin{aligned} \hat{p}(x_n, x_n \in B_{j_n}) &= \hat{p}(x_n | x_n \in B_{j_n}) \Pr\{x_n \in B_{j_n}\} \\ &= w_{j_n} \mathcal{N}(x_n; \mu_{j_n}, \Sigma_{j_n}) \end{aligned}$$

$$E_{in}^{\{B_j\}}(\Omega) = -\log \prod_{n=1}^N \hat{p}(x_n, x_n \in B_{j_n})$$

(likelihood of \mathcal{D} and $\{B_j\}$)

$$\begin{aligned} &= -\log \prod_{n=1}^N w_{j_n} \mathcal{N}(x_n; \mu_{j_n}, \Sigma_{j_n}) \\ &= \sum_{n=1}^N -\log(w_{j_n}) + \sum_{n=1}^N -\log[\mathcal{N}(x_n; \mu_{j_n}, \Sigma_{j_n})] \\ &= \sum_{j=1}^K -N_j \log w_j + \sum_{j=1}^K \sum_{x_n \in B_j} -\log[\mathcal{N}(x_n; \mu_j, \Sigma_{j_n})] \end{aligned}$$

Seperately optimize those two parts:

$$\textcircled{1} \text{ Lagrange } \Rightarrow w_j^* = \frac{N_j}{N}$$

$$\textcircled{2} \text{ Gradient } = 0 \Rightarrow \mu_j = \frac{1}{N_j} \sum_{x_n \in B_j} x_n \Rightarrow \Sigma_j = \dots$$

(same as what we had)

↳ In subproblem 2

$$\text{Given } \Omega, \hat{j}^* = \underset{\{B_j\}}{\operatorname{argmax}} \Pr\{j | \underline{x}_n\}$$

$$\equiv \text{Minimize } E_{\Omega}(\Omega)$$