

COL774

Machine Learning

Oct 27, 2021

Latent Class: EM expectation maximization

$$\{x^{(i)}\}_{i=1}^m$$

$$p(x^{(i)}, z^{(i)}; \theta)$$

↳ hidden

$$LL(\theta) = \sum_{i=1}^m \log p(x^{(i)}; \theta)$$

↳ intractable (using gradient descent)

$$\sum_{i=1}^m \log \sum_{z^{(i)}} p(x^{(i)}, z^{(i)}; \theta)$$

Jensen's inequality: for f concave

$$f(E[X]) \geq E[f(X)]$$

$$\Rightarrow LL(\theta) = \sum_{i=1}^m \log \sum_{z^{(i)}} \left[\frac{p(x^{(i)}, z^{(i)}; \theta)}{Q_i(z^{(i)})} \right] Q_i(z^{(i)})$$

$$= \sum_{i=1}^m \log E_{Q_i} \left[\frac{p(x^{(i)}, z^{(i)}; \theta)}{Q_i(z^{(i)})} \right]$$

↳ concave

E-step

$$\geq \sum_{i=1}^m E_{Q_i} \left[\log \frac{p(x^{(i)}, z^{(i)}; \theta)}{Q_i(z^{(i)})} \right]$$

$$= \sum_{i=1}^m \sum_{z^{(i)}} Q_i(z^{(i)}) \log \frac{p(x^{(i)}, z^{(i)}; \theta)}{Q_i(z^{(i)})} = LL'(\theta)$$

$$LL(\theta) \geq LL'(\theta) \rightarrow \text{Lower bound}$$

$$LL'(\theta) = \sum_{i=1}^n \sum_{z^{(i)}} Q_i(z^{(i)}) \left[\log \frac{P(x^{(i)}, z^{(i)}; \theta)}{Q_i(z^{(i)})} \right]$$

↳ not very different to optimize.

$$\arg \max_{\theta} LL'(\theta) = \arg \max_{\theta} \left[\sum_{i=1}^n \sum_{z^{(i)}} Q_i(z^{(i)}) \log P(x^{(i)}, z^{(i)}; \theta) \right]$$

+ ~~$\sum_{i=1}^n \sum_{z^{(i)}} Q_i(z^{(i)}) \log Q_i(z^{(i)})$~~ does not depend on θ
Entropy

$$= \arg \max_{\theta} \sum_{i=1}^n \sum_{z^{(i)}} Q_i(z^{(i)}) \log P(x^{(i)}, z^{(i)}; \theta)$$

(Contrast with EM)

Repeat
EM Algorithm:-

$$t \leftarrow 0;$$

$$Q^t \leftarrow \text{init}()$$

$$\{x^{(i)}\}_{i=1}^n \quad x^{(i)} \in \mathbb{R}^n$$

$$P(x^{(i)}, z^{(i)}; \theta)$$

~~to be deleted~~

do {

E-step $\forall i:- Q_i^t(z^{(i)}) \leftarrow P(z^{(i)} | x^{(i)}; \theta^{(t)})$

M-step:-

$$\theta^{(t+1)} \leftarrow \arg \max_{\theta} LL^t(\theta);$$

$$\equiv \arg \max_{\theta} \left[\sum_{i=1}^n \sum_{z^{(i)}} Q_i^t(z^{(i)}) \log P(x^{(i)}, z^{(i)}; \theta) \right]$$

$$t \leftarrow t+1;$$

~~$Q_i^t(z^{(i)})$~~

} while ! converged)

Convergence:- $\| \theta^{(t+1)} - \theta^{(t)} \| \leq \epsilon$

$$\hookrightarrow \nabla_{\theta} L(\theta) |_{\theta^{(t)}} \approx 0 \quad \Rightarrow \quad \theta^{(t+1)} \approx \theta^{(t)}$$

Note:- ① Converges to local optima of

$$L(\theta)$$

$$\textcircled{2} \arg \max_{\theta} L(\theta)$$

$$\equiv \arg \max_{\theta, \theta_1} \hat{L}(\theta, \theta_1)$$

EM can be seen as optimizing $\hat{L}(\theta, \theta_1)$ w.r.t block coordinate descent

$$\hat{L}(\theta, \theta_1) =$$

$$\sum_{i=1}^n \{ Q_i(z^{(i)}) \} \log \left[\frac{P(x^{(i)}, z^{(i)}; \theta)}{Q_i(z^{(i)})} \right]$$

EM Algorithm

M-step corresponds to: $\max_{\theta} \hat{L}(\theta, \theta_1)$ (for a given θ_1)
 E-step corresponds to: $\max_{\theta_1} \hat{L}(\theta, \theta_1)$ (for a given θ)

Result:- $L(\theta) \geq \hat{L}(\theta, \theta_1) \quad \forall \theta_1$
 (By construction / Jensen's inequality)

\Rightarrow EM is nothing but block coordinate

descent $\nabla L(\theta, \alpha)$
EM converges to local optima of $L(\theta)$

$$L(\theta^{(1)}) \geq L(\theta^{(2)}) \geq L(\theta^{(3)}) \geq L(\theta^{(4)}) \dots \text{--- (i)}$$

$$\forall t \quad L(\theta^{(t)}) = L(\theta^{(t+1)}) \quad \text{--- (ii)}$$

Suppose EM has converged.

$$\theta^{(t+1)} \approx \theta^{(t)}$$

$$\Rightarrow \nabla_{\theta} L(\theta) |_{\theta^{(t)}} = 0 \quad \text{--- (iii)}$$

$$\therefore \nabla_{\theta} L(\theta) |_{\theta^{(t+1)}} = 0 \quad \text{By limit}$$

$$\text{Concl.:- } \nabla_{\theta} L(\theta) |_{\theta^{(t)}} = 0 \Rightarrow \text{Local optima}$$

Assume o.w.:- $\nabla_{\theta} L(\theta) |_{\theta^{(t)}} > 0$

$$\Delta \theta = \eta \cdot \nabla_{\theta} L(\theta) |_{\theta^{(t)}}$$

$$\frac{L(\theta^{(t)} + \Delta \theta) - L(\theta^{(t)})}{\Delta \theta} > 0$$

$$\frac{L(\theta^{(t)} - \Delta \theta) - L(\theta^{(t)})}{\Delta \theta} < 0 \quad \text{--- (4)}$$

$$\frac{L(\theta^{(t)} - \Delta \theta) - L(\theta^{(t)})}{\Delta \theta} = 0 \quad \text{--- (5)}$$

→

④-⑤

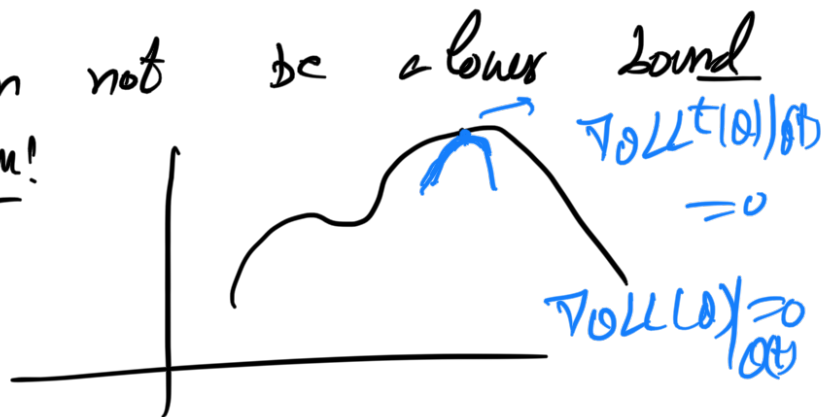
$$\frac{\mathcal{L}(\theta^{(t)} - \Delta\theta) - \mathcal{L}(\theta^{(t)}) - [\mathcal{L}^t(\theta^{(t)} - \Delta\theta) - \mathcal{L}^t(\theta^{(t)})]}{\Delta\theta} < 0$$

$$\Rightarrow \frac{\mathcal{L}(\theta^{(t)} - \Delta\theta) - \mathcal{L}^t(\theta^{(t)} - \Delta\theta)}{\Delta\theta} < 0$$

$$\Rightarrow \mathcal{L}(\theta^{(t)} - \Delta\theta) - \mathcal{L}^t(\theta^{(t)} - \Delta\theta) < 0$$

⇒ \mathcal{L}^t can not be a lower bound

⇒ contradiction!



⇒ EM algorithm:

PCA: Principal Component Analysis
(next class)