

COL774

Machine Learning

Oct 26, 2021

Last class:-

K-Means

GMM

$P(x^u; \theta)$

$P(x^u, z^u; \theta)$

hidden  
class  
labels

$$P(z^u) \equiv P(z^u | x^u; \theta)$$

$\theta \equiv \argmax_{\theta} P(x^u, z^u; \theta)$   
M-step given  $z^u$  fitted in missing values

$\{x^u\}_{i=1}^m$   
 $\{c^u\}_{i=1}^m$  given  $m$  -  $m \times K$   
 $\{m_1, \dots, m_K\}$   
 $\Downarrow$  cluster  
centroids  
giving  $\{c^u\}_{i=1}^m$

do

E-step: fill in missing values (given parameters)  
M-step: estimate the parameters / given  
filled in values

} while (!converged)

$\Rightarrow$  A more general framework for dealing with  
the problem of missing data

$\{x^u\}_{i=1}^m$

$P(x^u, z^u; \theta)$  :- underlying  
distribution  
 $\hookrightarrow$  hidden

$\Rightarrow$  an optimal set of  $\theta$   
parameters which describe  
the observed data

in some  
parametric  
form

$$\theta^* = \arg \max_{\theta} L(\theta) \\ \hookrightarrow \log \left[ \prod_{i=1}^n p(x^{(i)}; \theta) \right] \\ = \sum_{i=1}^n \log p(x^{(i)}; \theta)$$

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

How do we optimize over  $\theta$ :-

$\Rightarrow$  Directly gradient descent

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

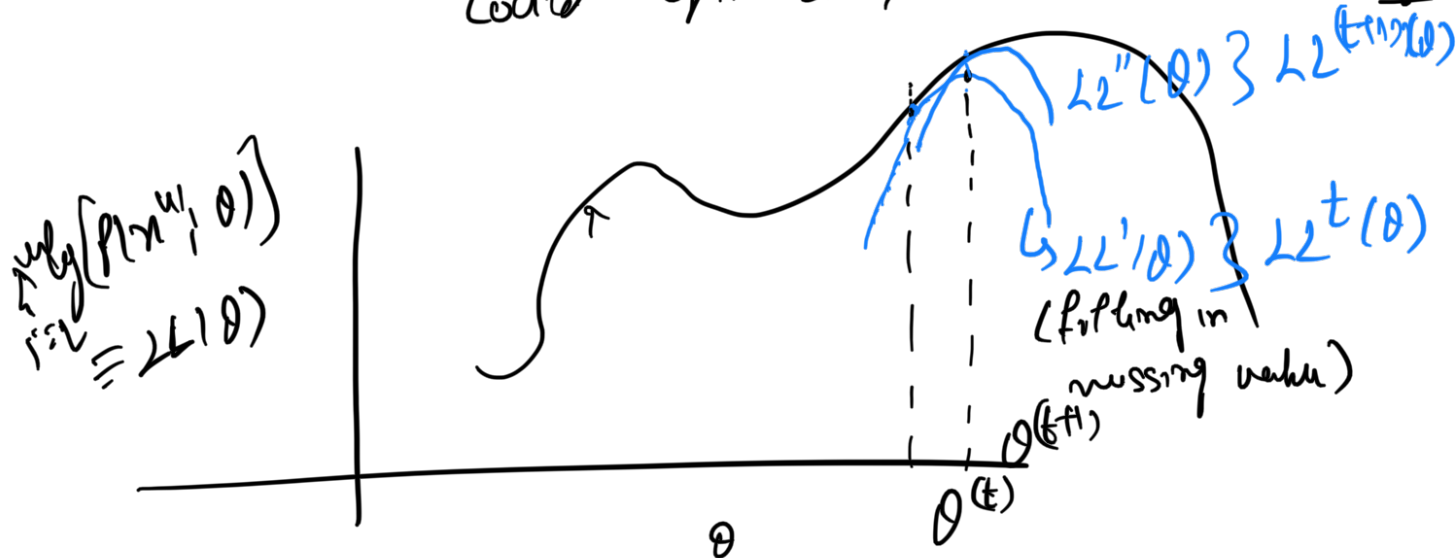
To  $L(\theta)$  :- perform gradient descent

Computing this way intractable

Assumption:- it is easy to optimize

$$\log \left[ \prod_{i=1}^n p(x^{(i)}, z^{(i)}; \theta) \right]$$

$\Rightarrow$  if we knew  $z^{(i)}$  then we could optimize the likelihood efficiently



$$① \quad L(\theta) \geq L'(\theta) \quad \forall \theta$$

$$② \quad L(\theta) \big|_{\theta^{(t)}} = L'(\theta^{(t)})$$

⇒ (A) E-step: - estimate  $L(\theta)$  (filling in of missing values)

(B) M-step:  $\arg \max_{\theta} L'(\theta)$

⇒ Converges  $\| \theta^{(t+1)} - \theta^{(t)} \| \leq \epsilon \quad \checkmark$

At convergence - local optimal of  $L(\theta)$

$t = 0$

$\theta^{(t)} \leftarrow \text{init}();$

do {

E-M Expectation  
Maximization  
to deal with missing data

E-step

$L^t(\theta) \leftarrow \text{construct lower bound } (L(\theta));$

M-step

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

$t \leftarrow t+1;$

} while ! converged)

⇒ In contrast, with directly optimizing  $L(\theta)$  using gradient ascent (may not be tractable)

Next:-

① How to construct the "lower" bound

② Prove convergence

③ EM algorithm can be seen as:-  
an instance of Block coordinate descent over a suitably defined

objective function

$$\mathcal{L}(\theta) \equiv \sum_{i=1}^m \log p(x_i; \theta)$$

$\Rightarrow \hat{\mathcal{L}}(\theta, \theta_1)$  performing Block coordinate descent

$$\arg\max_{\theta} \arg\max_{\theta_1} [\hat{\mathcal{L}}(\theta, \theta_1)]$$

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

$\hookrightarrow$  optimize our  $\theta_1$  given  $\theta$  (E-step)

$\hookrightarrow$  optimize our  $\theta$  given  $\theta_1$  (M-step)

} (while ! converged)

Jensen's Inequality:-

$f: \mathbb{R}^n \rightarrow \mathbb{R}$  be convex



$$f(\alpha x_1 + (1-\alpha)x_2) \leq \alpha f(x_1) + (1-\alpha)f(x_2) \quad \checkmark$$

- (1)

$$0 \leq \alpha \leq 1$$

$$f(E[x]) \leq E[f(x)]$$

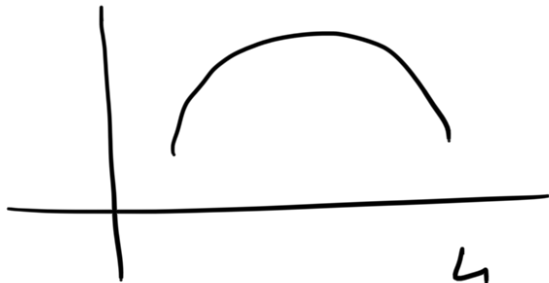
$P(x)$ :- our  $\underline{x}$

computed using

Think about: how to prove Jensen's inequality

Hint:- Generalize Eq ① to a convex combination of  $K$  points ( $K \geq 2$ )

For concave functions:  $f: \mathbb{R}^n \rightarrow \mathbb{R}$  concave



$$f(\alpha x_1 + (1-\alpha)x_2) \geq \alpha f(x_1) + (1-\alpha)f(x_2)$$

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

$\Rightarrow$  if  $f: \mathbb{R}^n \rightarrow \mathbb{R}$  is strictly convex (concave)  
(strictly) convex!

$$f(E[X]) < E[f(X)]$$

unless  $X$  is a constant (given  $p(X)$ )

$$f(E[C]) = E[f(C)]$$

(strictly)  
Concave:-

$$f(E[X]) > E[f(X)]$$

unless  $X$  is a constant (given  $p(X)$ )

$$f(E[C]) = E[f(C)]$$

Coming back:-

$$\mathcal{L}(\theta) = \log \prod_{i=1}^n [p(x^i; \theta)] = \sum_{i=1}^n \log p(x^i; \theta)$$

under iid

$$\log [P(x^u; \theta)] = \log \left( \sum_{z^u} P(x^u, z^u; \theta) \right) \quad \text{some expression}$$

$$\mathcal{L}(\theta) = \sum_{i=1}^u \log \left( \sum_{z^u} P(x^u, z^u; \theta) \right)$$

strictly concave



Let  $Q_i(z^u)$  be some distribution over  $z^u$  (non-zero)

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

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

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

concave

$$\geq \sum_{i=1}^u E_{Q_i(z^u)} \log \frac{P(x^u, z^u; \theta)}{Q_i(z^u)}$$

↳ Jensen's inequality

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

$$= \mathcal{L}'(\theta) \quad \text{exactly the same}$$

$$\mathcal{L}(\theta) \geq \mathcal{L}'(\theta)$$