## P61 ELBO + kL Divergence

$X$: observed data  $Z$: unobserved data (latent (隐藏) variable)

$(X, Z)$: complete data   $\theta$: parameter    $\theta_{MLE} = \arg\max p(x|\theta)$

EM: $\theta^{(t+1)} = \underset{\theta}{\arg\max} \int \underbrace{p(z|x,\theta^{(t)})}_{constant} \log \underbrace{p(x,z|\theta)}_{variable} dz$

E-step: $p(z|x,\theta^{(t)}) \Rightarrow E_{z|x,\theta^{(t)}}[\log p(x,z|\theta)] = ELBO$

M-step: $\theta^{(t+1)} = \underset{\theta}{\arg\max} E_{z|x,\theta^{(t)}}[\log p(x,z|\theta)]$

Proof

$\log p(x|\theta) = \log p(x,z|\theta) - \log p(z|x,\theta) = \log \frac{P(x,z|\theta)}{q(z)^{(t)}} - \log \frac{P(z|x,\theta)}{q(z)^{(t)}}$   $q(z)^{(t)}$ is the probability distribution of $z$, $\int q(z) dz = 1$

Direct calculation is too difficult $\Rightarrow$ Find the expectation  # $E_x[g(x)] = \int p(x)g(x) dx$ or $\sum_x p(x)g(x)$

Left side: $\int q(z) \log p(x|\theta) dz = \log p(x|\theta) \int q(z) dz = \log p(x|\theta)$

Right side: $\underbrace{\int q(z)^{(t)} \log \frac{P(x,z|\theta)}{q(z)^{(t)}} dz}_{ELBO: \text{evidence lower bound}} - \underbrace{\int q(z)^{(t)} \log \frac{P(z|x,\theta)}{q(z)^{(t)}} dz}_{kL(q(z)^{(t)} || P(z|x,\theta))} \Rightarrow$ relative entropy $kL(P||Q) = \int p(x) \log \frac{P(x)}{Q(x)} dx$ or $\sum p(x) \log \frac{P(x)}{Q(x)}$

$\therefore \log p(x|\theta) = ELBO + kL(q||P)$    $kL(q||P) \geq 0$   $\log(x|\theta) \geq ELBO$   Equality holds if and only if $q = P$

When the equality holds $q(z)^{(t)} = P(z|x,\theta^{(t)}) \to$ posterior    $\log(x|\theta) = ELBO$

$\therefore \hat{\theta} = \underset{\theta}{\arg\max} ELBO = \underset{\theta}{\arg\max} \int p(z|x,\theta^{(t)}) \log \frac{P(x,z|\theta)}{P(z|x,\theta^{(t)})} dz = \underset{\theta}{\arg\max} \int p(z|x,\theta^{(t)}) [\log p(x,z|\theta) - \log p(z|x,\theta^{(t)})] dz$

$= \underset{\theta}{\arg\max} \int p(z|x,\theta^{(t)}) \log p(x,z|\theta) dz$

## P62 ELBO + Jensen's inequlity

# The secant line (割线) of a convex function ($f''>0$) lies above the graph of the function.

Jensen's inequality: $f(tx_1 + (1-t)x_2) \leq t f(x_1) + (1-t) f(x_2)$

$\varphi$ is convex function, $\varphi(E[X]) \leq E[\varphi(X)]$    $\varphi(\sum_{i=1}^{N} \lambda_i x_i) \leq \sum_{i=1}^{N} \lambda_i \varphi(x_i)$ s.t. $\sum_{i=1}^{N} \lambda_i = 1$    Jensen $\alpha$: $E[\varphi(X)] - \varphi(E[X])$

Equlity holds if and and only if $x_1 = x_2 = \cdots = x_N$ or $\varphi$ is linear on domain (域) containing $x_1, x_2, \cdots, x_N$

Proof (finite form) (mathematical induction)

$N = 1$ or $2$    $\varphi(\lambda_1 x_1 + \lambda_2 x_2) \leq \lambda_1 \varphi(x_1) + \lambda_2 \varphi(x_2)$ is clearly true

Assume the induction hypothesis that for a single case $N = n$ is true: $\varphi(\sum_{i=1}^{n} \lambda_i x_i) \leq \sum_{i=1}^{n} \lambda_i \varphi(x_i)$

It follows that: $\varphi((1-\lambda_{n+1}) \sum_{i=1}^{n} \frac{\lambda_i}{1-\lambda_{n+1}} x_i + \lambda_{n+1} x_{n+1})$   Let $\frac{\lambda_i}{1-\lambda_{n+1}} = \eta_i$   $\because \sum_{i=1}^{n+1} \lambda_i = 1$   $\therefore \sum_{i=1}^{n} \frac{\lambda_i}{1-\lambda_{n+1}} = 1$   $\therefore \sum_{i=1}^{n} \eta_i = 1$

$\therefore = \varphi((1-\lambda_{n+1}) \sum_{i=1}^{n} \eta_i x_i + \lambda_{n+1} x_{n+1}) \leq (1-\lambda_{n+1}) \varphi(\sum_{i=1}^{n} \eta_i x_i) + \lambda_{n+1} \varphi(x_{n+1})$

$\because \varphi(\sum_{i=1}^{n} \eta_i x_i) \leq \sum_{i=1}^{n} \eta_i \varphi(x_i)$   $\therefore \varphi(\sum_{i=1}^{n+1} \lambda_i x_i) \leq (1-\lambda_{n+1}) \sum_{i=1}^{n} \eta_i \varphi(x_i) + \lambda_{n+1} \varphi(x_{n+1}) = \sum_{i=1}^{n} \lambda_i \varphi(x_i) + \lambda_{n+1} \varphi(x_{n+1})$

$\therefore \varphi(\sum_{i=1}^{n+1} \lambda_i x_i) \leq \sum_{i=1}^{n+1} \lambda_i \varphi(x_i)$

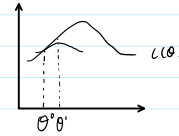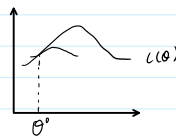$\varphi(\int x d\mu_n(x)) \leq \int \varphi(x) d\mu_n(x)$     $\mu_n(x) = \sum_{i=1}^{n} \lambda_i \delta_{x_i}$ ?

EM: $\log p(x|\theta) = \log \int p(x,z|\theta) dz = \log \int \frac{P(x,z|\theta)}{q(z)^{(t)}} q(z)^{(t)} dz \geq \int q(z)^{(t)} \log \frac{P(x,z|\theta)}{q(z)^{(t)}} dz \to ELBO$

$f(E[\frac{P(x,z|\theta)}{q(z)^{(t)}}]) \geq E[f(\frac{P(x,z|\theta)}{q(z)^{(t)}})]$

When the equlity holds: $\frac{P(x,z|\theta^{(t)})}{q(z)^{(t)}} = C$

$\therefore q(z)^{(t)} = \frac{1}{C} p(x,z|\theta^{(t)})$   $\therefore 1 = \int q(z)^{(t)} dz = \frac{1}{C} \int p(x,z|\theta^{(t)}) dz$   $\therefore p(x|\theta^{(t)}) = C$

$\therefore q(z)^{(t)} = \frac{P(x,z|\theta^{(t)})}{P(x|\theta^{(t)})} = P(z|x,\theta^{(t)})$



① initialize $\theta^0$   ② find the expectation $\log \frac{P(x,z|\theta)}{P(z|x,\theta^0)}$ (ELBO)   ③ find $\theta'$ that maximizes the expectation   ④ repeat ②③ until $\theta$ converges

## P60 EM Convergence

$\log(x|\theta) = \log(x,z|\theta) - \log(z|x,\theta)$   $\therefore q(z)^{(t)} \log(x|\theta) = q(z)^{(t)} \log(x,z|\theta) - q(z)^{(t)} \log(z|x,\theta)$

$\log(x|\theta) = \underbrace{\int p(z|x,\theta^{(t)}) \log p(x,z|\theta) dz}_{Q(\theta, \theta^{(t)})} - \underbrace{\int p(z|x,\theta^{(t)}) \log p(z|x,\theta) dz}_{H(\theta, \theta^{(t)})}$

$\because \theta^{(t+1)} = \underset{\theta}{\arg\max} \int p(z|x,\theta^{(t)}) \log p(x,z|\theta) dz = \underset{\theta}{\arg\max} Q(\theta, \theta^{(t)})$

$\therefore Q(\theta^{(t+1)}, \theta^{(t)}) \geq Q(\theta^{(t)}, \theta^{(t)})$

$\because H(\theta^{(t+1)}, \theta^{(t)}) - H(\theta^{(t)}, \theta^{(t)}) = \int P(z|x,\theta^{(t)}) \log \frac{P(z|x,\theta^{(t+1)})}{P(z|x,\theta^{(t)})} dz \overset{Jensen}{\leq} \log \int \underbrace{P(z|x,\theta^{(t+1)}) dz}_{1} = \log 1 = 0$

$\therefore H(\theta^{(t+1)}, \theta^{(t)}) \leq H(\theta^{(t)}, \theta^{(t)})$

$\therefore \log(x|\theta^{(t+1)}) \geq \log(x|\theta^{(t)})$   $\therefore$ EM converges

## P63 EM Review

EM mainly solves problems arising from probabilities generative models:

MLE: $\hat{\theta} = \arg\max P(x|\theta)$

But in some cases, $P(x)$ is difficult to find. Because the observable variable $X$ does satify a certain distribution. So we assume that $X$ is generated by the latent variable $Z$ that satisfies a certain distribution.

$P(X) = \int_z P(x,z) dz = \frac{P(x,z)}{P(z|x)}$

## P64 Generalized EM

Narrow EM is a special case of generalized EM

The EM algorithm is designed to solve the parameter estimation problems

$\log P(x|\theta) = ELBO + kL(q||P) \leq ELBO$ $E_q[\log \frac{P(x,z|\theta)}{q}] - f(q,\theta)$

Narrow EM is a special case of generalized EM

The EM algorithm is designed to solve the parameter estimation problems

$$\log P(x|\theta) = ELBO + KL(q||P) \quad \begin{cases} ELBO = E_{q(z)}\left[\log \frac{P(x,z|\theta)}{q(z)}\right] = \mathcal{L}(q,\theta) \\ KL(q||P) = \int q(z) \log \frac{q(z)}{P(z|x,\theta)} dz \end{cases}$$

$KL \geq 0$   Equality holds if and only if $q = P \Rightarrow$ Narrow EM: $q(z) = P(z|x,\theta)$

But in some cases, $P(x|z,\theta)$ is hard to find $\Rightarrow$ Generalized EM: Fix $\theta$, then $P(x|\theta)$ is fixed $\Rightarrow \hat{q} = \underset{q}{\arg\min} \, KL(q||P) = \underset{q}{\arg\max} \, \mathcal{L}(q,\theta)$

$\therefore$ Generalized EM:

$$\begin{cases} E\text{-step}: \quad q^{(t+1)} = \underset{q}{\arg\max} \, \mathcal{L}(q, \theta^{(t)}) \\ M\text{-step}: \quad \theta^{(t+1)} = \underset{\theta}{\arg\max} \, \mathcal{L}(q^{(t+1)}, \theta) \end{cases}$$

$$\mathcal{L}(q,\theta) = E_q\left[\log P(x,z|\theta) - \log q(z)\right]$$

$$= E_q\left[\log P(x,z|\theta)\right] - \underbrace{E_q\left[\log q(z)\right]}_{H(q) \text{ entropy}}$$

Narrow EM: $H(q) = 0$

Generalized EM: $H(q) \neq 0$