

First, we have a probabilistic model  $P(X, Z|\theta)$ , where  $X$  is observed data,  $Z$  is unobserved data,  $\theta$  is the parameter we would like to learn from data.

We want to get the maximum likelihood estimate for  $\theta$ , this should be done by maximizing the likelihood  $P(X, Z|\theta)$ , unfortunately we don't know values for  $Z$ , so this is impossible.

Then we find that we can calculate  $P(X|\theta)$  by summing out  $Z$ ,  $P(X|\theta) = \sum_Z P(X, Z|\theta)$ . So the problem becomes how to find a  $\theta$  that maximizes  $P(X|\theta)$ .

We are going to maximize  $P(X|\theta)$  iteratively. At every iteration, we begin with a  $\theta_{old}$ . We hope we can find a  $\theta_{new}$  so that  $P(X|\theta_{new}) > P(X|\theta_{old})$ .

To do this, we need to prove an important inequality, that for any  $\theta$ , we have

$$\begin{aligned} \log P(X|\theta) &= \log \sum_Z P(X, Z|\theta) \\ &= \log \sum_Z P(Z|X, \theta_{old}) \frac{P(X, Z|\theta)}{P(Z|X, \theta_{old})} \\ &\geq \sum_Z P(Z|X, \theta_{old}) \log \frac{P(X, Z|\theta)}{P(Z|X, \theta_{old})} \\ &= \mathbb{E}_{Z|X, \theta_{old}} \log \frac{P(X, Z|\theta)}{P(Z|X, \theta_{old})} \\ &= \mathbb{E}_{Z|X, \theta_{old}} \log P(X, Z|\theta) - \mathbb{E}_{Z|X, \theta_{old}} \log P(Z|X, \theta_{old}) = Q(\theta) . \end{aligned}$$

So we have proved that for any  $\theta$ ,  $\log P(X|\theta) \geq Q(\theta)$ .

We can also prove that  $Q(\theta_{old}) = \log P(X|\theta_{old})$ .

$$\begin{aligned} Q(\theta_{old}) &= \mathbb{E}_{Z|X, \theta_{old}} \log P(X, Z|\theta_{old}) - \mathbb{E}_{Z|X, \theta_{old}} \log P(Z|X, \theta_{old}) \\ &= \mathbb{E}_{Z|X, \theta_{old}} \log \frac{P(X, Z|\theta_{old})}{P(Z|X, \theta_{old})} \\ &= \mathbb{E}_{Z|X, \theta_{old}} \log \frac{P(X, Z, \theta_{old})/P(\theta_{old})}{P(X, Z, \theta_{old})/P(X, \theta_{old})} \\ &= \mathbb{E}_{Z|X, \theta_{old}} \log \frac{P(X, \theta_{old})}{P(\theta_{old})} = \mathbb{E}_{Z|X, \theta_{old}} \log P(X|\theta_{old}) = \log P(X|\theta_{old}) . \end{aligned}$$

Then we define  $\theta_{new}$  to be

$$\theta_{new} = \arg \max_{\theta} Q(\theta),$$

so we know that for any  $\theta$ ,  $Q(\theta_{new}) \geq Q(\theta)$ .

As a summary, we have proved that

- for any  $\theta$ ,  $\log P(X|\theta) \geq Q(\theta)$
- for any  $\theta$ ,  $Q(\theta_{new}) \geq Q(\theta)$ , where  $\theta_{new} = \arg \max_{\theta} Q(\theta)$
- $Q(\theta_{old}) = \log P(X|\theta_{old})$

Therefore,  $\log P(X|\theta_{new}) \geq Q(\theta_{new}) \geq Q(\theta_{old}) = \log P(X|\theta_{old})$ . Note that when  $P(X|\theta_{new}) = P(X|\theta_{old})$  the algorithm stops, so in other cases,  $P(X|\theta_{new}) > P(X|\theta_{old})$ .

As a final remark, since  $\mathbb{E}_{Z|X, \theta_{old}} \log P(Z|X, \theta_{old})$  in  $Q(\theta)$  is a constant, in practice we define  $\theta_{new}$  as

$$\theta_{new} = \arg \max_{\theta} \mathbb{E}_{Z|X, \theta_{old}} \log P(X, Z|\theta) .$$