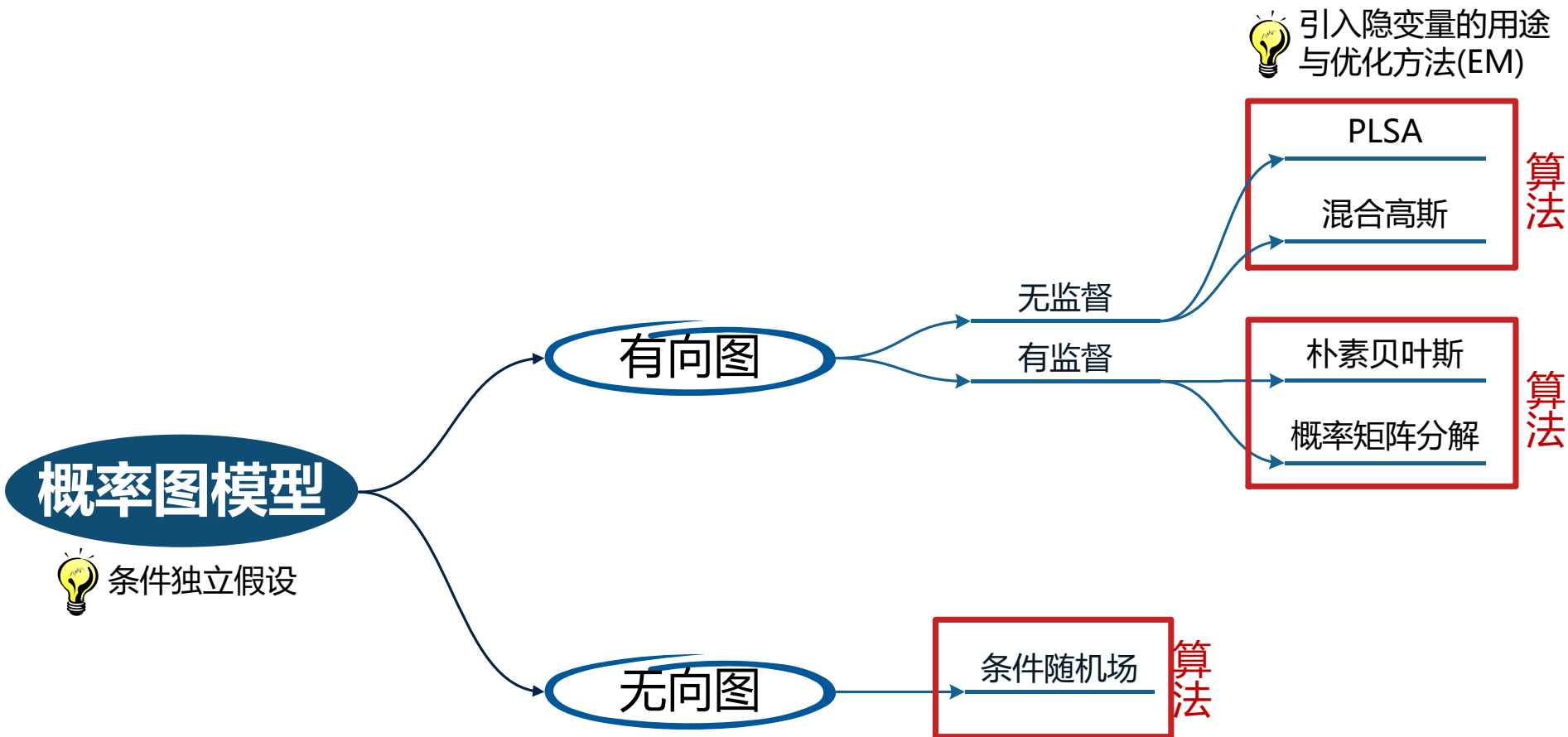


课程专题二：概率图模型方法



EM算法的一般形式

EM 算法推导

数据 $X = \{(x_i, c_i)\}_{i=1}^N$, x_i 出现了 c_i 次

$$L(\theta) = \ln P(X|\theta) = \sum_{i=1}^N c_i \ln P(x_i|\theta)$$

$$= \sum_{i=1}^N c_i \ln \sum_z P(x_i, z|\theta)$$

常数, θ' 是上一轮的参数

$$= \sum_{i=1}^N c_i \ln \sum_z \boxed{P(z|x_i, \theta')} \frac{P(x_i, z|\theta)}{P(z|x_i, \theta')}$$

$$\geq \sum_{i=1}^N c_i \sum_z P(z|x_i, \theta') \ln \frac{P(x_i, z|\theta)}{P(z|x_i, \theta')} = l(\theta|\theta')$$

等价于优化

$$\sum_{i=1}^N c_i \sum_z P(z|x_i, \theta') \ln P(x_i, z|\theta)$$

EM 算法

Iterate until convergence

- E step: Calculate $P(z|x_i, \theta')$ for each example x_i , Use this to construct

$$\sum_{i=1}^N c_i \sum_z P(z|x_i, \theta') \ln P(x_i, z|\theta)$$

- M step: Replace current θ' by

$$\theta' \leftarrow \operatorname{argmax}_{\theta} \sum_{i=1}^N c_i \sum_z P(z|x_i, \theta') \ln P(x_i, z|\theta)$$

解此最优化问题：有可能找不到解析（Closed Form）的最优解；
也可以朝增大的方向优化，这就是Generalized EM算法

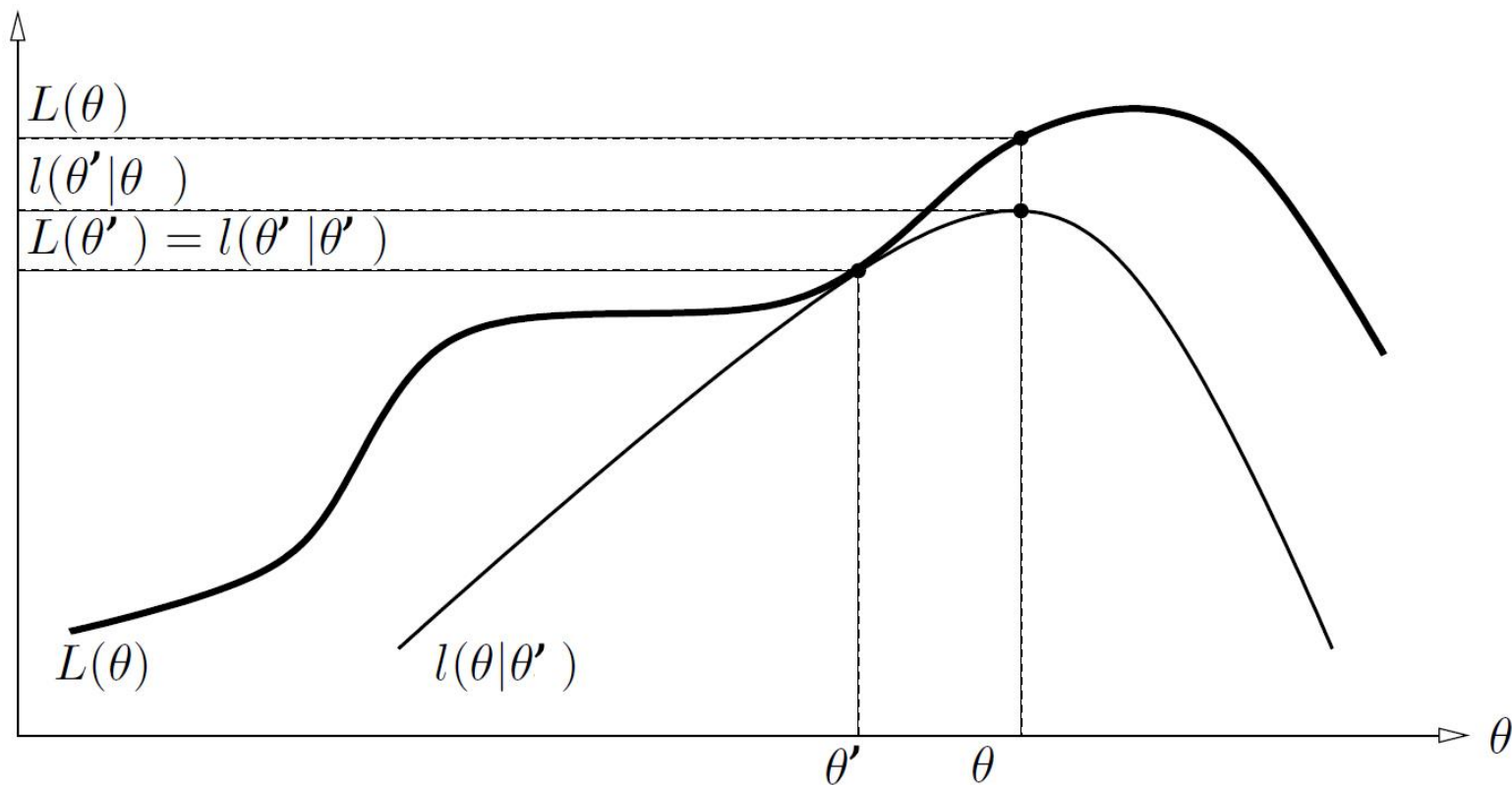
EM 算法推导：为什么使用 $P(z|x_i, \theta')$

最大似然的下界为

$$\begin{aligned}l(\theta|\theta') &= \sum_{i=1}^N c_i \sum_z P(z|x_i, \theta') \ln \frac{P(x_i, z|\theta)}{P(z|x_i, \theta')} \\l(\theta'|\theta') &= \sum_{i=1}^N c_i \sum_z P(z|x_i, \theta') \ln \frac{P(x_i, z|\theta')}{P(z|x_i, \theta')} \\&= \sum_{i=1}^N c_i \sum_z P(z|x_i, \theta') \ln P(x_i|\theta') \\&= \sum_{i=1}^N c_i \ln P(x_i|\theta') \sum_z P(z|x_i, \theta') \\&= \sum_{i=1}^N c_i \ln P(x_i|\theta') = L(\theta')\end{aligned}$$

下界在 θ' 时，等于 $L(\theta')$

EM 算法推导：为什么使用 $P(z|x_i, \theta')$



构造出一个“紧”的下界

EM 应用于PLSA

- 数据 $X = \{(x_i, c_i)\}_{i=1}^N$, x_i 出现了 c_i 次

$$\sum_{i=1}^N c_i \sum_z P(z|x_i, \theta') \ln P(x_i, z|\theta)$$

- 应用于PLSA, 数据 $X = \{((d_i, w_j), c_{ij})\}$

$$\sum_{i,j} c_{ij} \sum_z P(z|d_i, w_j, \theta') \ln P(d_i, w_j, z|\theta)$$

$$= \sum_{i,j} c_{ij} \sum_z P(z|d_i, w_j, \theta') \ln P(d_i|\theta) P(z|d_i, \theta) P(w_j|z, \theta)$$

EM 应用于PLSA

- 模型的输出结果: $P(z|d_i)$ $P(w_j|z)$
- 引入隐变量的目的
 - **获得可观察变量与隐变量的关系**