

Lecture 6, 7: Maximum Likelihood, Model Selection, Bayesian Learning

General EM algorithm:

1. Initialize θ^{old}

2. E step: $P(z|x, \theta) \Rightarrow$ guess the value of z

3. M step: $\theta \leftarrow \arg \max_{\theta} Q(\theta, \theta^{\text{old}})$

$$Q(\theta, \theta^{\text{old}}) \leftarrow \sum_z P(z|x, \theta^{\text{old}}) \cdot \ln P(x, z|\theta) \Rightarrow \text{discrete}$$

$$\text{or } Q(\theta, \theta^{\text{old}}) \leftarrow \int_z P(z|x, \theta^{\text{old}}) \cdot (\ln P(x, z|\theta)) \cdot dz \Rightarrow \text{continuous}$$

\Rightarrow MLE to optimize parameters.

4. Repeat 2, 3 until converged.

MLE and MAP:

Maximum Likelihood Estimator:

$$\theta \leftarrow \arg \max_{\theta} P(X|\theta)$$

means model
params.

Maximum A Posterior:

$$\theta \leftarrow \arg \max_{\theta} P(\theta|X) = \arg \max_{\theta} \frac{P(X|\theta) \cdot P(\theta)}{P(X)}$$

iff $\propto \arg \max_{\theta} P(X|\theta) \cdot P(\theta)$

Model Selection:

Step 1: For each model k , compute the MLE:

$$\hat{\theta}_{ML}(k) = \underset{\theta}{\operatorname{argmax}} \log P(x | \theta, k), k = 1, 2, \dots, m$$

Step 2:

$$k^* = \underset{k}{\operatorname{argmax}} J(\hat{\theta}_{ML}(k))$$

Information criterion

$$J(\hat{\theta}_{ML}(k)) = \begin{cases} AIC: \ln \underbrace{P(x | \hat{\theta}(k))}_{\text{Maximum Likelihood}} + \underbrace{d_k}_{\# \text{params}} \\ BIC: \ln \underbrace{P(x | \hat{\theta}(k))}_{\# \text{samples}} + \frac{1}{2} \cdot d_k - \frac{\ln N}{N} \end{cases}$$

Bayesian Model Selection: VBEM

$$\text{MAP form} \Leftarrow \boxed{\underset{\text{model class}}{\overbrace{P(m|y)}}} = \boxed{\frac{P(y|m) \cdot P(m)}{P(y)}} \Rightarrow \text{Bayesian form}$$

data set

$$\underbrace{P(y|m)}_{\text{Bayesian Evidence}} = \int_{\theta_m} \underbrace{P(y | \theta_m, m)}_{\text{Evidence}} \cdot \underbrace{P(\theta_m | m)}_{\text{Prior}} \cdot \underbrace{d\theta_m}_{\text{Posterior}}$$

Variational form

E-step: compute $P(z|x, \theta_m) \Rightarrow$ as general EM

$$P(z|x, \theta_m) \propto e^{-E[\log p(x, z|\theta)]}$$

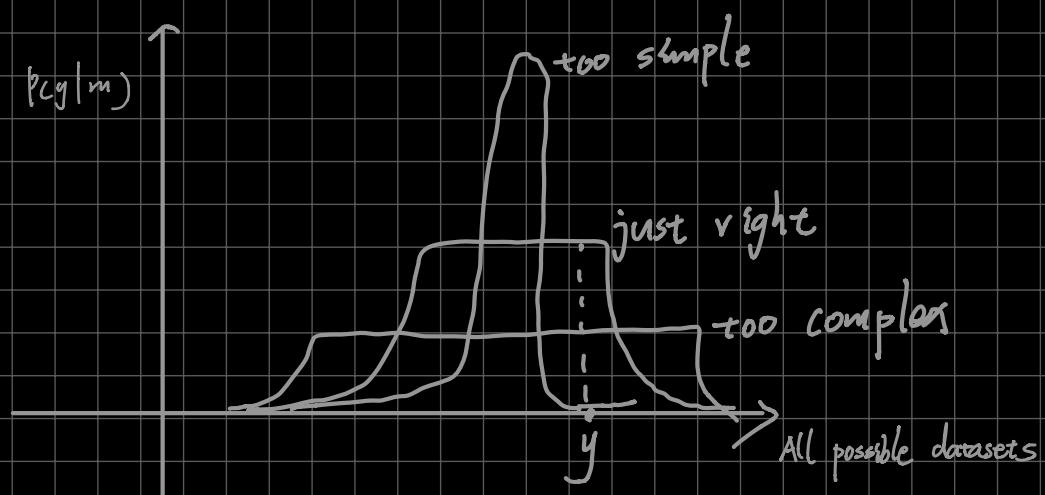
M-step: $\theta_m \propto e^{\left[\int_z P(z|x, \theta) \cdot \ln p(x, z|\theta) dz \right]}$

general EM: $\theta_m \leftarrow \underset{\theta}{\operatorname{argmax}} Q(\theta, \theta)$

change $\underset{\theta}{\operatorname{argmax}}$ to $e^{\underset{\theta}{\operatorname{argmax}}}$

△ Use EM algorithm to optimize $P(y|\theta)$.

Occam's Razor in Bayesian Model Selection:



\Rightarrow [m is too simple: can't generate y
 m is too complex: can generate y , but $P(y|m)$ is too low.

△ Bayesian inference automatically supplements

Occam's Razor principle.

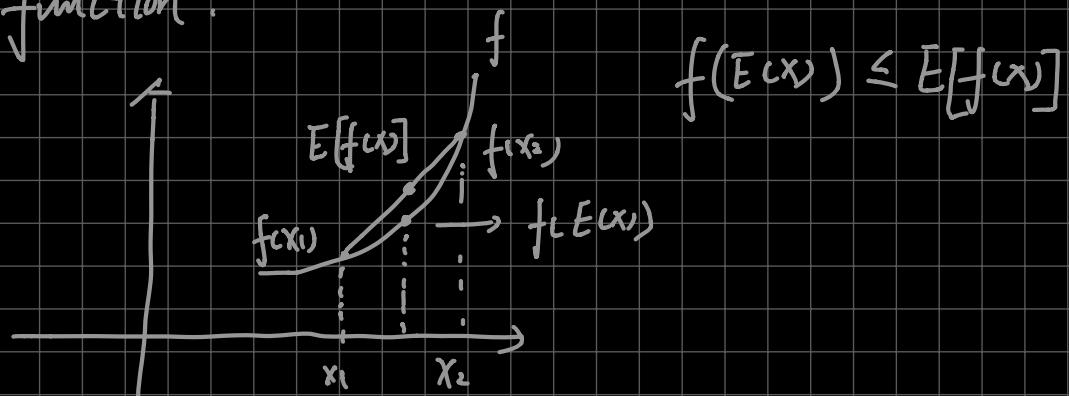
EM algorithm derivation

sometimes called bootstrap.

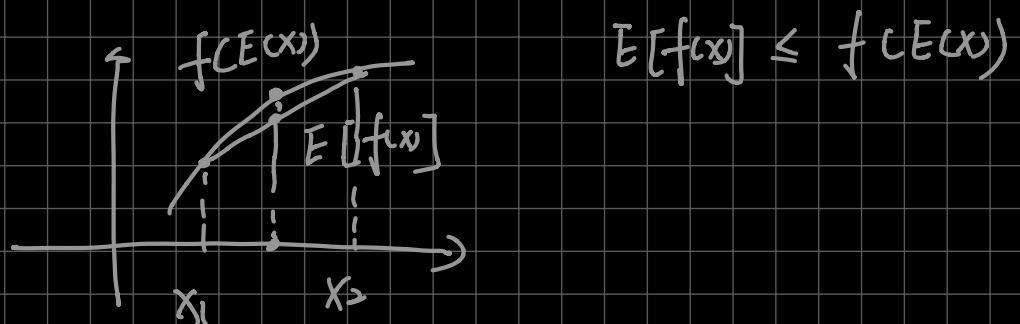
— Stanford Ng

Jensen's inequality:

Convex function:



Concave function:



EM process:

Have model for $P(X, z | \theta)$, only observe $X = \{x^{(1)}, x^{(2)}, \dots, x^{(m)}\}$.

$$L(\theta) = \sum_{i=1}^m \log P(X^{(i)} | \theta)$$

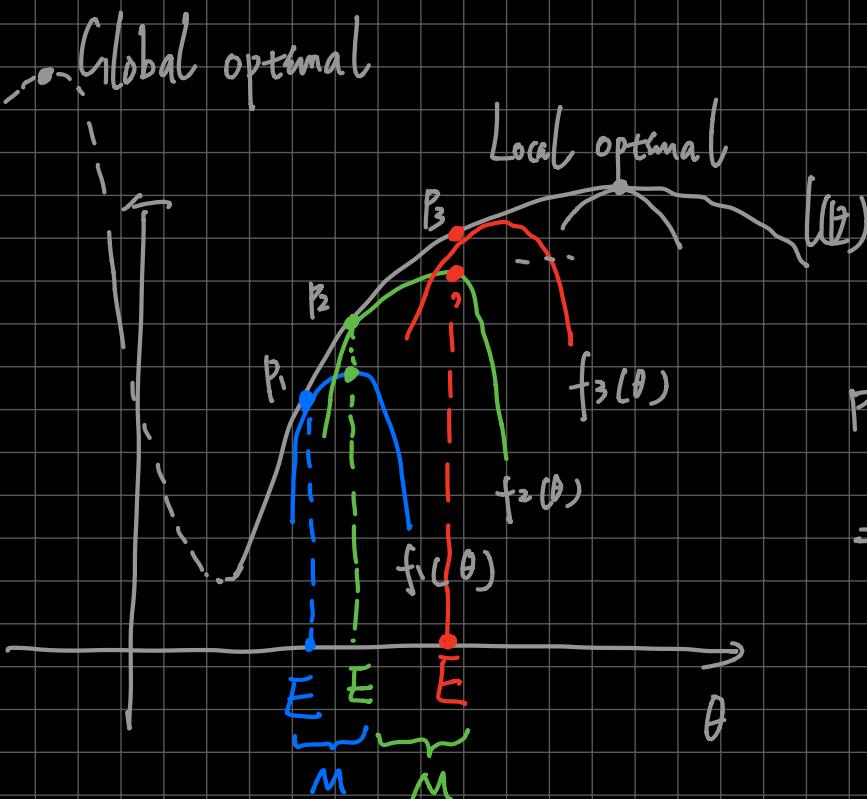
$$= \sum_{i=1}^m \cdot \log \sum_{k=1}^K P(X^{(i)}, z_k^{(i)} | \theta)$$

$$= \sum_{i=1}^m \cdot \log \sum_{k=1}^K Q(z_k^{(i)}) \left[\frac{P(X^{(i)}, z_k^{(i)} | \theta)}{Q(z_k^{(i)})} \right], \text{ where } \sum_{k=1}^K Q(z_k^{(i)}) = 1$$

$$\begin{aligned}
 \max_{\theta} \{ l(\theta) \} &\rightarrow \sum_{i=1}^m \cdot \log \sum_{k=1}^K Q(z_k^{(i)}) \left[\frac{P(x^{(i)}, z_k^{(i)} | \theta)}{Q(z_k^{(i)})} \right] \\
 &= \sum_{i=1}^m \log E \left[\frac{P(x^{(i)}, z_k^{(i)} | \theta)}{Q(z_k^{(i)})} \right] \\
 &\stackrel{\text{Jensen's inequality}}{\leq} \sum_{i=1}^m E \left[\log \left(\frac{P(x^{(i)}, z_k^{(i)} | \theta)}{Q(z_k^{(i)})} \right) \right] \\
 &= \sum_{i=1}^m \sum_{k=1}^K Q(z_k^{(i)}) \cdot \log \frac{P(x^{(i)}, z_k^{(i)} | \theta)}{Q(z_k^{(i)})} \\
 &= \sum_{i=1}^m \sum_{k=1}^K Q(z_k^{(i)}) \left(\log \frac{P(z_k^{(i)} | x^{(i)}, \theta) \cdot P(x^{(i)} | \theta)}{P(z_k^{(i)} | x^{(i)}, \theta)} \right) = \sum_{i=1}^m \sum_{k=1}^K Q(z_k^{(i)}) \log P(x^{(i)} | \theta)
 \end{aligned}$$

Graphical Explanation:

(lower bound), z, x is fixed
so it's a function of $\theta \Rightarrow f(\theta)$



E-step: Draw curve f_i that in point P_i , Jensen inequality satisfies.

M-step: Find the optimal point P_3 of curve f_i drew by E-step.

\Rightarrow Repeat EM until it converges
(Point found by M-step is the same as E-step)

△ E-step 即即 Jensen's inequality 相等，因此
这也是为什么 E-step 直观的解释为缩小 lower bound
function $f(\theta)$ 与 $l(\theta)$ 之间的 gap.

M-step 即 optimize θ 使得上述不等式左边变
大，这也是为什么 M-step 直观的解释为拔高
lower bound function $f(\theta)$ 的值。

