

The EM Algorithm - in General

The EM Algorithm is a general algorithm for finding maximum likelihood solutions for probabilistic models having latent variables.

- Proof that heuristic algorithm does maximize likelihood function
- Basis for variational inference

Consider prob model with

- \bar{X} observed variables
- Z hidden variables
- Θ parameters

Joint distribution $P(X, Z | \Theta)$

Goal: Maximize likelihood function

$$P(X | \Theta) = \sum_z P(X, z | \Theta)$$

Assumption - 1) Z discrete, all args hold for continuous vars

- 2) Direct optimization of $P(X | \Theta)$ w.r.t. Θ hard
- 3) Optimization of complete likelihood function - easier

Z can be missing data, marginalized params, etc.

usually true because of sum

$$P(X, z | \Theta)$$

no summation, log passes through

$$\begin{aligned}
&= \sum_z q(z) \ln \frac{p(z|x, \theta) p(x|\theta)}{q(z)} - \sum_z q(z) \ln \frac{p(z|x, \theta)}{q(z)} \\
&= \sum_z q(z) \ln \frac{p(z|x, \theta)}{q(z)} + \sum_z q(z) \ln \frac{p(x|\theta)}{q(z)} - \sum_z q(z) \ln q(z) \\
&\quad - \left[\sum_z q(z) \ln p(z|x, \theta) - \sum_z q(z) \ln q(z) \right] \\
&= \sum_z q(z) \ln p(x|\theta) \\
&= \ln p(x|\theta) \cdot 1 \quad \square
\end{aligned}$$

$$\text{KL}(q \parallel p) = 0 \quad \text{iff} \quad \frac{q(x)}{p(x)} = 1 \quad \forall x$$

$\ln(1) = 0$

General EM

Introduce $q(z)$, function over latent vars,

Note the following decomposition holds for all q

$$\ln p(x|\theta) = \mathcal{L}(q, \theta) + \text{KL}(q||p)$$

where

$$9.71 \quad \mathcal{L}(q, \theta) = \sum_z q(z) \ln \left\{ \frac{p(x, z|\theta)}{q(z)} \right\}$$

and

$$9.72 \quad \text{KL}(q||p) = - \sum_z q(z) \ln \left\{ \frac{p(z|x, \theta)}{q(z)} \right\}$$

Note

- a) signs not equal
- b) $\mathcal{L}(q, \theta)$ contains complete-data likelihood
- c) $\text{KL}(q||p)$ is the KL divergence between $q(z)$ and $p(z|x, \theta)$ (and contains $p(z|x, \theta)$).

- Recall $\text{KL}(q||p) \geq 0$

$\text{KL}(q||p) \neq \text{KL}(p||q)$ in general
and $\text{KL}(q||p) = 0$ if $q = p$

Because $\text{KL}(q||p) \geq 0$

$$\mathcal{L}(q, \theta) \leq \ln p(x|\theta)$$

i.e. $\mathcal{L}(q, \theta)$ is a lower bound on $\ln p(x|\theta)$.

EM is two stage procedure
Suppose current param vector is θ^{old} .

E step: maximize $\mathcal{L}(q, \theta^{old})$ wrt. $q(z)$
- i.e. find $q(z)$ that maximizes $\mathcal{L}(q, \theta^{old})$,
this will be by making KL as small as possible
- ideal, set $q(z) = p(z|x, \theta^{old})$,
KL vanishes

M step: $q(z)$ held fixed and
 $\mathcal{L}(q, \theta)$ is maximized wrt. to θ
yielding θ^{new}

- this causes the lower bound to increase
and thereby $\ln P(x|\theta)$ to increase as well.

- since q is learned using θ^{old} , in
general $p(z|x, \theta^{new})$ will be different
from q and KL div will be non-zero

Substituting $q(z) = p(z|x, \theta^{old})$
into

$$\mathcal{L}(q, \theta) = \sum_z q(z) \ln \left\{ \frac{p(x, z|\theta)}{q(z)} \right\}$$

we have

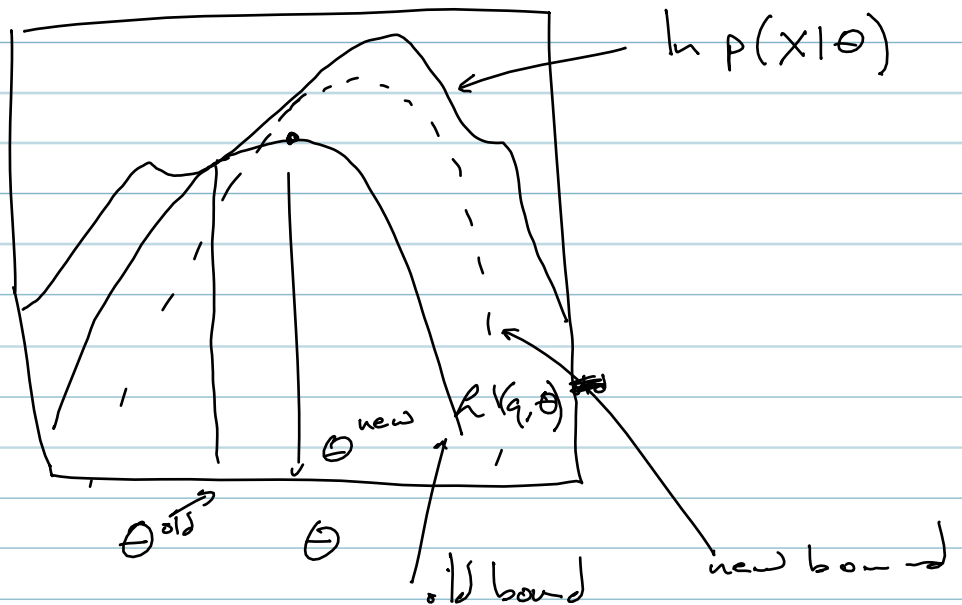
$$\mathcal{L}(q, \theta) = \underbrace{\sum_z p(z|x, \theta^{old}) \ln P(x, z|\theta)}_{Q(\theta, \theta^{old})} + \text{const.}$$

where

$$Q(\theta, \theta^{old})$$

involves the log of the complete data
likelihood, a quantity assumed to be easy to work
with.

Graphically



Note that in the case of n iid. data, i.e.,

$$p(X, z) = \prod_n p(x_n, z_n)$$

the posterior over z has a nice form

$$\begin{aligned} p(z|X, \theta) &= \frac{p(X, z|\theta)}{\sum_z p(X, z|\theta)} \\ &= \frac{\prod_{n=1}^n p(x_n, z_n|\theta)}{\sum_z \prod_{n=1}^n p(x_n, z_n|\theta)} \\ &= \prod_{n=1}^n p(z_n|x_n, \theta) \end{aligned}$$

So in E step, each data point's posterior "responsibility" can be computed independently from the others.

Note:

Any moves that increase $\mathcal{L}(q, \theta)$ will increase $\ln p(x|\theta)$.

Possibilities include:

- sampling z 's
- numerical gradient ascent of θ 's
- one data point at a time
- etc.